

1 **Dopamine waves as a mechanism for spatiotemporal credit assignment**

2 Arif A. Hamid¹, Michael J. Frank^{2,3*}, Christopher I. Moore^{1,3*}

3 ¹Department of Neuroscience, ²Department of Cognitive Linguistics & Psychological Sciences,
4 ³Carney Institute for Brain Science, Brown University, Providence, RI

5 * These authors contributed equally to this work; alphabetical order

6

7 **Abstract**

8 Significant evidence supports the view that dopamine shapes reward-learning by
9 encoding prediction errors. However, it is unknown whether dopamine decision-signals are
10 tailored to the functional specialization of target regions. Here, we report a novel set of wave-like
11 spatiotemporal activity-patterns in dopamine axons across the dorsal striatum. These waves
12 switch between different activational motifs and organize dopamine transients into localized
13 clusters within functionally related striatal subregions. These specific motifs are associated with
14 distinct task contexts: At reward delivery, dopamine signals rapidly resynchronize into
15 propagating waves with opponent directions depending on instrumental task contingencies.
16 Moreover, dopamine dynamics during reward pursuit signal the extent to which mice have
17 instrumental control, and interact with reward waves to predict future behavioral adjustments.
18 Our results are consistent with a computational architecture in which striatal dopamine signals
19 are sculpted by inference about instrumental controllability, and provide evidence for a
20 spatiotemporally “vectorized” role of dopamine in credit assignment.

21

22 -----

23 **Main text**

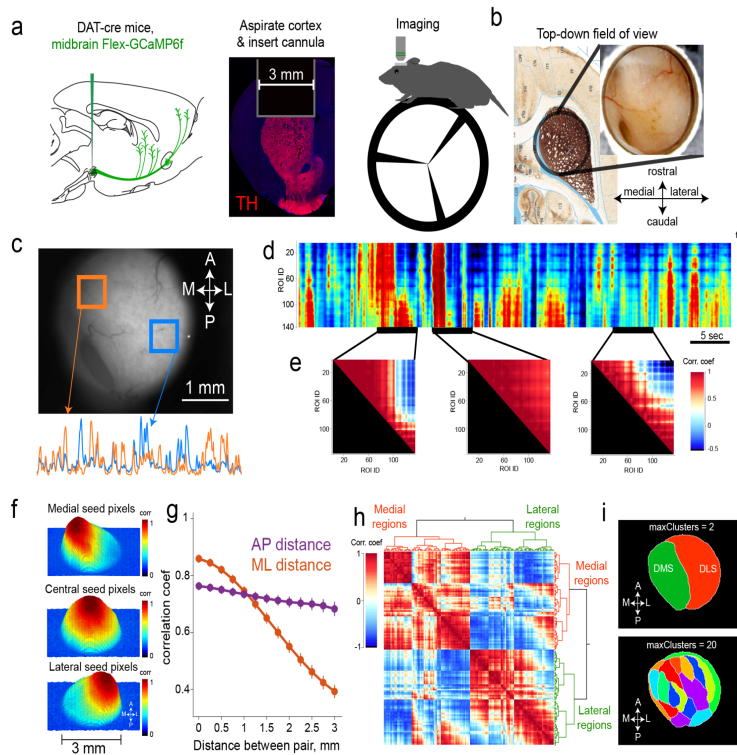
24 Dopamine supports reward learning and motivational activation, but details about what
25 decision variables are encoded, and how they are delivered to postsynaptic targets, continue to
26 be refined¹⁻³. The dopamine-reward prediction error (RPE) hypothesis emphasizes that
27 dopamine conveys deviations from reward expectation in reinforcement learning (RL) theory⁴.
28 This formulation generally treats dopamine as a “global” spatio-temporally uniform signal, a view
29 based on two key findings. First, extensively divergent dopamine axons^{5,6} provide an
30 architecture for broadcast-like communication. Second, dopamine cell spikes measured in the
31 midbrain are highly synchronized^{7,8}, putatively implementing a redundant population code⁹⁻¹¹ for
32 RPEs¹². These observations form the basis for an influential view^{13,14} of what dopamine
33 communicates and how it is delivered: scalar RPEs that are uniformly delivered to all recipient
34 subregions. The notion of uniform encoding also extends to alternative accounts for dopamine’s
35 role in motivation² by relaying scalar value signals¹⁵.

36 A key limitation of this global view is that scalar (or, spatio-temporally uniform) decision
37 variables are neither computationally advantageous, nor reflected in forebrain dopamine
38 dynamics. Postsynaptic striatal subregions are functionally specialized^{16,17}, receiving distinct
39 cortical and thalamic afferents^{18,19}, and express unique compliments of biomarkers²⁰.
40 Accordingly, rewards²¹, their motivated pursuit^{15,22} and predictive stimuli²³ produce vastly
41 different dopamine time courses across the dorsal-ventral and medial-lateral axes of the
42 striatum. While these observations indicate regional heterogeneity, the extent to which
43 dopamine inputs reflect the computational requirements of postsynaptic areas remains elusive.
44 For example, there is some theoretical motivation^{24,25} and empirical support^{26,27} for delivery of
45 vector-valued RPEs that depend on a target region’s computational specialty. Nonetheless, we
46 currently lack a clear understanding of organizing principles for striatal dopamine activity, and
47 what normative computational functions may be served by such heterogeneity.

48 ***Related striatal subregions get correlated dopamine input***

49 We set out to characterize the spatio-temporal organizational rules of dopamine activity
50 across the dorsal striatum. Standard methods for dopamine measurement typically survey small
51 territories (10s – 100s of micrometers) and are ill-suited to probing large-scale organization. To
52 overcome these limitations, we injected *cre*-dependent GCAMP6f into the midbrain of DAT-*cre*
53 mice, and imaged dopamine axons through a 7mm² chronic imaging window over the dorsal
54 striatum (DS)²⁸ (**Fig. 1a**). This approach provided optical access to 60-80% of the dorsal surface
55 of the mouse striatum, with a view of dorsomedial (DMS), dorsolateral (DLS) and partial access
56 to the posterior-tail (TS) region of the striatum (**Fig. 1b**). We imaged the activity of dopamine
57 axons at multiple levels of resolution with one or two photon microscopy.

Figure 1



58
59

60 **Figure 1: Dorsal striatal dopamine activity is spatiotemporally asynchronous and**
61 **clusters into contiguous territories.**

62 **a**, Schematic of methods for imaging dopamine axons over dorsal striatum. GCaMP6 was injected
63 into midbrain of DAT-cre mice. Cortex overlying dorsal striatum was aspirated, together with
64 insertion of 3 mm diameter imaging cannula, and activity was imaged in head-restrained mice. **b**,
65 Top-down field of view. **c**, Average delta f/F from two regions (top) that exhibit decorrelated activity
66 (bottom). **d**, Activity of several ROIs from the same session as **c**, time series are sorted such that
67 medial areas are top ROIs, and lateral regions are represented at the bottom. **e**, Correlation matrix
68 across ROIs for different 5sec epochs (highlighted in bottom of **d**), showing patterns of
69 correlations that evolve in time. For example, middle shows global correlation, whereas left and
70 right panels show instances of antagonistic activity patterns in top and bottom set of ROIs. **f**,
71 Results from spatial correlation from seed pixels, evaluating the Pearson's correlation of with all
72 other regions. *Top* panel shows medial seed pixels that are highly correlated with nearby regions,
73 and show graded decrease in correlation for distant regions. Same analysis was repeated for a
74 set of pixels in central striatum (*middle*) and lateral seed pixels (*bottom*). **g**, Quantification of
75 sessions-wide correlation between each pair of pixels as a function of distance, separated by
76 medio-lateral (orange) and antero-posterior distances. (n= 8 mice, p<0.001 wilcoxon signed-rank
77 test for difference of ML vs AP slopes)

78 **h**, Pairwise correlation matrix using hierarchical clustering summarizes similarity of dopamine
79 activity. **i**, *Top*, anatomical projection of pixels that share similarity at the highest cluster limit of
80 two outlining medial and lateral subregions of the dorsal striatum. Increasing the cluster threshold
81 to 20 (*bottom*) revealed smaller, but anatomically contiguous regions of the striatum.

82

83 Using a head-fixed preparation, we began by focusing on spontaneous activation of
84 dopamine axons in mice resting on a wheel in a dark chamber without external stimuli. To first
85 test if dopamine axons were globally activated, we compared fluorescence signals in DS
86 regions-of-interest (ROIs) (**Fig. 1c**). While ROIs were sometimes globally synchronized²⁸, we
87 observed decorrelated patterns across striatal subregions that evolved across time (**Fig.**
88 **1c,d,e**). These patterns of activation were observed across multiple anatomical scales (see
89 **Extended Data Fig. 1** for micron-scale organization), indicating that dopamine afferents can
90 become recruited asynchronously.

91 To examine how activity is spatially coordinated, we computed the Pearson's correlation
92 between pixels' fluorescence as a function of anatomical distance. Dopamine axons showed
93 strong local correlations that gradually decreased with distance (**Fig. 1f,g**), comparable to the
94 organization of striatal spiny-neuron activity²⁹. Strikingly however, this distance-dependent falloff
95 was selective to the medio-lateral axis, and was not present on the antero-posterior axis (**Fig.**
96 **1g**), suggesting an organization rule that promoted selective mediolateral decorrelation.

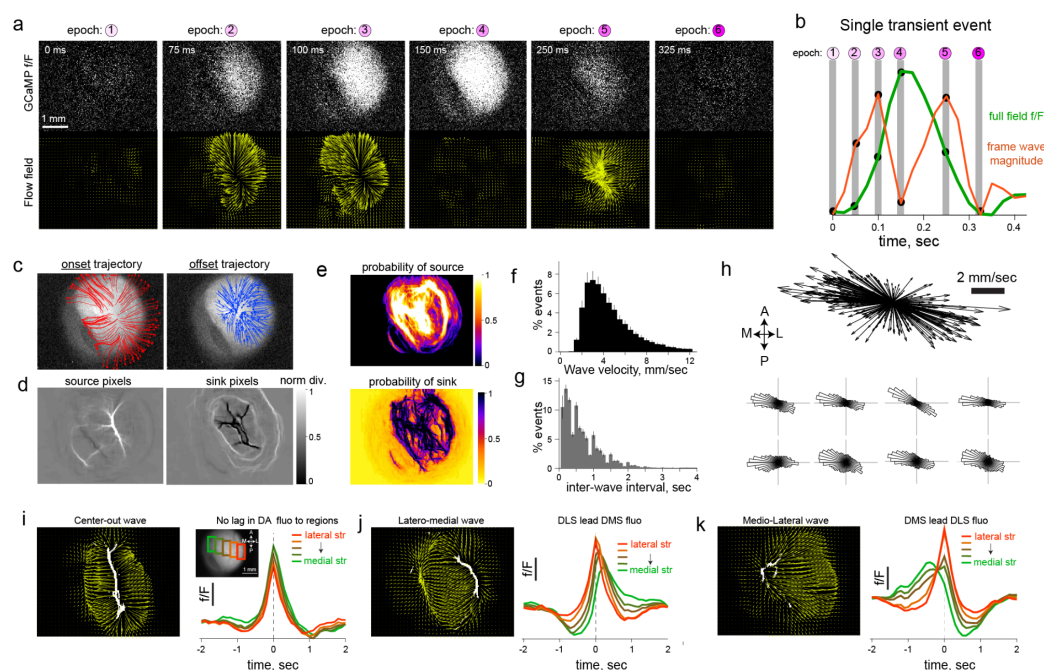
97 To further examine the topographical organization of dopamine signals, we leveraged
98 standard cluster analyses (**Fig. 1h**). In every dataset (n = 31 sessions, 8 mice), the highest
99 cluster threshold identified two contiguous subregions outlining well-established^{30,31} DS
100 subregions; medial (DMS) and lateral (DLS) striatum (**Fig. 1i top**). Further increasing cluster
101 limits progressively (**Extended Data Fig. 2**) revealed smaller subdomains of DS (**Fig. 1i**
102 **bottom**), resembling striatal sub-clusters previously identified based on glutamatergic input
103 patterns¹⁸. These areas had similar clustering patterns across days and animals (**Extended**
104 **Data Fig. 3**), with 25-30 optimal clusters identified in our field of view (**Extended Data**
105 **Fig. 4**), indicating a critical dependence on the underlying spatio-temporal activity pattern.
106 Together, these results provide evidence for regional coordination of dopamine transmission
107 and provided an initial basis for evaluating whether these signals are modulated by the
108 underlying subregion's computational speciality.
109

110 ***Wave-like patterns coordinate dopamine activity***

111 What spatiotemporal patterns produce systematically decorrelated dopamine signals?
112 We noticed that full-field fluorescence exhibited complex but spatially and temporally continuous
113 trajectories throughout the striatum, similar to travelling waves described in other cortical and
114 subcortical brain regions³²⁻³⁵ (**Fig. 2a,b, Supplementary Video 1**). To quantify these complex
115 trajectories, we used optic flow algorithms³⁶ to compute frame-by-frame flow fields (see
116 methods for details; **Supplementary Video 2**).

117

Figure 2



118
119
120 **Figure 2: Wave-like spatiotemporally continuous sequences of dopamine-axon activation**
121 **switch between motifs.**

122 **a**, *Top* row shows individual frames for different epochs of a transient as dopamine axon activity
123 emerges and extinguishes in DS. *Bottom* row displays the corresponding flow vector fields
124 computed for each pixel. Notice the divergence of vector fields during the rise phase of
125 fluorescence, and convergent vectors during fall phase. **b**, Average fluorescence (green) across
126 the entire field of view lasting ~300ms sampled at 40Hz, and corresponding, flow magnitude in
127 the fluorescence signal (red). **c**, Flow trajectory of fluorescence for 5 frames during the onset (*left*,
128 red lines), or offset (*left*, blue lines) phase of the wave from **b**. Each line shows the pattern of flow
129 from individually seeded pixels. **d**, Heatmap quantifying how divergent the vector fields are at
130 each pixel during onset or offset of activity (left and right respectively). A positive value indicates
131 that diverging pattern of flow at each pixel indicating that fluorescence is entering the striatum
132 from those locations. Conversely, negative values are sink regions with converging flow vectors.
133 **e**, Peak-normalized projection of the flow vector divergence for the onset (*top*) or offset (*bottom*)
134 of all transients in one session (n=1516 events). Note that a repeated configuration of pixels serve
135 as sinks and sources. **f**, Distribution of propagation velocity (n=8 mice, 1625 +/- 213 events per
136 mouse). Error bar denotes SEM. **g**, Distribution of interwave intervals for the same data **f**. **h**, *Top*,
137 quiver plot summarizing the direction and magnitude of waves in a single session, and distribution
138 of angle of wave propagation for each animal, *bottom* (n=8 mice, all p < 0.001, Omnibus test for
139 angular uniformity). **i**, *Left*, vector fields (yellow) superimposed onto source pixels (white) for
140 waves that are sourced at the midline and propagate bidirectionally outward. *Right*, corresponding
141 fluorescence time course in ROIs on a medio-lateral gradient of the striatum (*inset*). **j**, Same
142 format as **i**, for lateral source and medial flow or, **k**, medial source and laterally flowing wave.

143 -----

144 The onset of activity in GCaMP fluorescence originated from clustered “source”
145 locations, and rapidly migrated to other regions (**Fig. 2c,d, left**). By contrast, activity terminated
146 as a result of flow toward “sink” locations (**Fig. 2c,d, right**). A repeated configuration of pixels
147 had a high probability of serving as sinks and sources (**Fig. 2e, Extended Data Fig. 5**),
148 indicating that local rules may dictate the initiation and termination of dopamine activity.

149 Dopamine waves entered the dorsal striatum with exponentially decaying inter-wave-
150 intervals (**Fig. 2g**) and propagate with a range of velocities (median = 3.8 mm/s, interquartile
151 range = 2.5, **Fig. 2f**). The overall direction of flow is bimodally distributed, with a biased medial-
152 lateral propagation axis (**Fig. 2h**, all $p < 0.001$, Omnibus test for angular uniformity).

153

154 We next sought to determine if the collection of complex trajectories were made up of
155 simpler, repeated sequences that may influence the time course of dopamine arriving at
156 different parts of the striatum. Indeed, the combination of initiation loci and flow direction gave
157 rise to motif waves that were scaled by propagation velocity and extent of striatum covered. We
158 focused our attention on three motifs that produced most of the dopamine transients (**Extended**
159 **Data Fig. 5**).

160 First, source pixels clustered at the juncture of DMS and DLS (**Fig. 2i**) initiated
161 dopamine activity that rapidly spread bilaterally outward (Type-1, “Center-Out” or CO wave, **Fig.**
162 **2i, left**). These waves radiate across the stratum with the fastest velocities, arriving at all
163 subregions with almost zero lag (**Fig. 2i, right**). Second, source pixels in lateral DLS initiated a
164 wave that propagates medially (Type-2, “latero-medial” or LM wave). LM waves advanced
165 across the striatum relatively slowly and delivered dopamine transients to DMS that were
166 delayed relative to DLS in proportion to propagation speed (**Fig. 2j right**). Third, a medially
167 sourced wave propagated laterally (Type-3, “medio-lateral” or ML wave, **Fig. 2k, left**),
168 terminating in DLS. ML waves activate dopamine axons in the medial striatum first and recruit
169 lateral regions with substantial delay (**Fig. 2k, right**). Together, these results demonstrate that
170 wave-like patterns are a fundamental organizational principle of dopamine axonal activity,
171 prescribing how activity initiates, propagates and terminates across DS.

172 **Rewards evoke directional dopamine waves**

173 What is the functional role of dopamine waves in adaptive behavior? We set out to
174 determine the computational significance of wave-like trajectories in the context of the well-
175 studied role of striatal dopamine in instrumental behavior. The dorsal striatum exhibits graded
176 behavioral specialty, with the DMS orchestrating goal-directed behaviors involving action-
177 outcome contingencies, and the DLS implicated in stimulus-response behaviors^{30,31,37}.
178 Inactivation or manipulation of dopamine in DMS degrades goal-directed planning and action
179 due to an inability to learn whether rewards are under instrumental control^{38,39}.

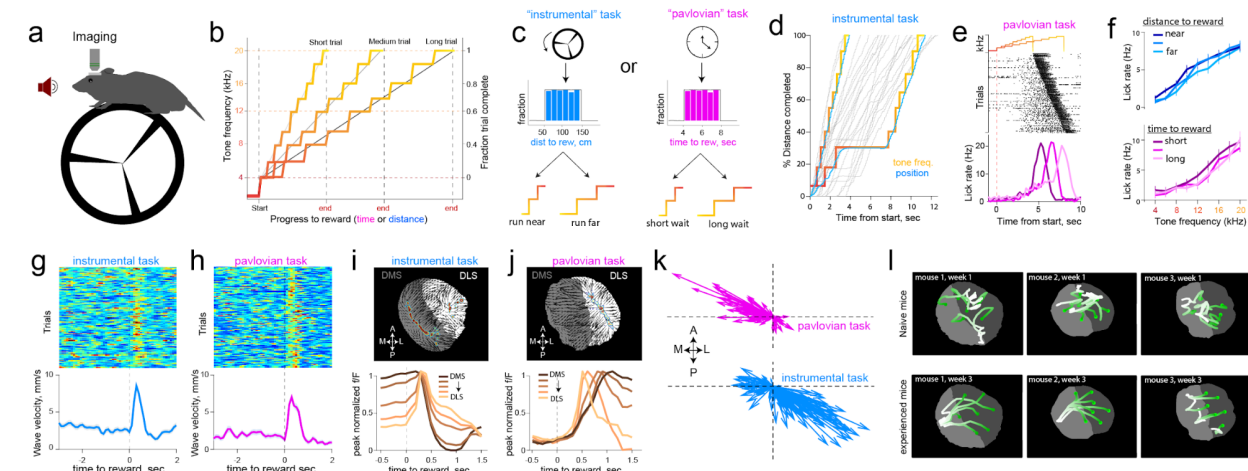
180 We thus designed two operant tasks intended to manipulate action-outcome
181 contingencies, and asked whether dopamine dynamics carry information about the degree of
182 instrumental controllability (**Fig. 3a,b,c**). First, in an ‘instrumental’ task, rewards were contingent

183 on mice running on a wheel to traverse linearized distance, with the progress to reward
184 indicated by an auditory tone that escalated in frequency (**Fig. 3b,d**). On each trial, the distance
185 that was needed to run for tone transitions (and ultimately, reward) was randomly selected from
186 a uniform distribution (50-150 cm, **Fig. 3c, left**). Thus, while the mouse was in control of tone
187 transitions, the specific contingencies varied across trials. A second ‘pavlovian’ task was
188 administered in separate sessions. The task structure was identical except tone-transitions
189 occurred after fixed durations within a trial (randomly drawn, 4-8 sec, **Fig. 3c right**), and
190 progress to reward was unrelated to running. Trained mice exhibited anticipatory lick trajectories
191 that increased with ascending tone frequency in both tasks (**Fig. 3e,f**), indicating that mice used
192 these tones to update their online judgment of progress to reward. Analysis of run bouts
193 (**Extended Data Fig. 6**) revealed that mice invested goal-directed effort to receive rewards
194 selectively in the instrumental task.

195 As in spontaneous conditions reported above, dopamine waves were ubiquitous during
196 task-performance. Notably, reward delivery immediately resynchronized irregular patterns into
197 smooth waves (**Fig. 3g,h**) that had opposite directions depending on task conditions.
198 Specifically, completion of a trial in the instrumental task triggered ML waves (**Fig. 3i,k bottom**,
199 **Supplementary video 3**), whereas rewards in the pavlovian task promoted LM waves (**Fig. 3j,k**
200 **top, Supplementary video 4**, $p < 0.001$ Watson-Wilson test for equality of mean directions in
201 two tasks, $n=6$ mice for instrumental task, $n=8$ mice for pavlovian task).

202 These patterns evolved dynamically with learning: Reward-waves were initially irregular
203 in naive animals but became progressively smooth and directional with experience in task (**Fig.**
204 **3l, Supplementary video 5**). The dynamic sculpting of the spatiotemporal patterns by training
205 and task demands ruled out explanations related to the intrinsic anatomy or biophysics of
206 dopamine axons that would constrain the array of observed activation patterns. Thus, we
207 conclude that dopamine waves carry behaviorally relevant decision signals and set out to
208 formalize their precise contribution. In particular, the continuous propagation of dopamine
209 across the striatum both in space and time motivated a revision of standard “temporal-
210 difference” models wherein a single reward-value influences learning about earlier events that
211 are predictive of rewards. We reasoned that these views could be expanded to include
212 “spatiotemporal differences” in which waves carry additional, graded information about structural
213 sub-circuits that are most likely to be responsible for rewards.

Figure 3



214
215

216 **Figure 3: Reward delivery promotes directional waves that depend on instrumental**
217 **requirement of task.**

218 **a**, Schematic of test chamber. **b**, Changes in tone frequency for short, medium and long trials
219 tiling fraction of trial complete. **c**, In the instrumental task (*left*), tone transitions are linked to
220 rotation of the wheel, and change in tone frequency. The total distance to travel on each trial is
221 drawn from a uniform distribution of 50-150 centimeters. In pavlovian task (*right*), the passage of
222 time escalated tones, and the duration to wait was also drawn from a uniform distribution of 4-8
223 seconds. **d**, Example trials in the instrumental task. When the mouse traverses linearized distance
224 rapidly, the tones also escalate quickly, but if the mouse pauses running, auditory tones signal
225 the completed fraction of distance. **e**, Example licking behavior in pavlovian, sorted by delay to
226 reward. Mice increase lickrate in anticipation of reward. **f**, These anticipatory licks were not
227 influenced by distance to run, or duration to wait, but increased in proportion to progress to reward
228 signaled by tones (two-way ANOVA effect of tones $F(8,683) = 3.32$, $p = 0.001$ and effect of 4
229 duration bins $F(3,683)=0.48$, $p = 0.7$ in pavlovian sessions. For instrumental sessions, effect of
230 tones $F(8,359) = 8.41$, $p<0.001$ and influence of four distance bins $F(3,359) = 0.13$, $p=0.9$). **g**,
231 Alignment of trial-by-trial wave velocity across the striatum. Rewards consistently resynchronized
232 dopamine axons into wave in the instrumental task ($n=123$ trials), but also in the pavlovian task **h**
233 (111 trials). **i**, *Top*, example flow vectors (black arrows) and source locations (contour plot
234 representing source regions) across pixels for a single trial. *Bottom*, peak-normalized
235 fluorescence time course across trials produced by mediolateral waves on the medio-lateral
236 gradient of the striatum. **j**, Same format as **i** for pavlovian session. **k**, Mean flow vectors for reward
237 epoch (0-1s post reward) for each trial in pavlovian and instrumental sessions shown in **h** and **i**.
238 **l**, Flow trajectory of fluorescence in response to reward as mice gained experience with the task.
239 *Top*, Naive mice had irregular responses in the first two days of reward exposure, and at *bottom*,
240 the same mice exhibit smooth waves after 3 weeks of learning reward contingency. See
241 supplementary video 5 for responses plotted.
242

243 **Dopamine waves implement spatio-temporal credit assignment**

244 Our functional interpretation of dopamine dynamics is that the opponent wave
245 trajectories at reward are relevant for spatiotemporal credit assignment. The key inference the
246 animal must make is whether it is in control of the reward-predictive tone transitions, and
247 moreover, which specific contingency applies in the current trial (i.e., distance to run to advance
248 tones). Thus, for mice to preferentially run in the instrumental task (and persist running for long-
249 distance trials), the extent of instrumental controllability should guide reward-evoked dopamine
250 to favor the DMS (i.e. strengthen action-outcome learning). Trial by-trial controllability is partly
251 ambiguous in the task because contingencies were stochastic (drawn from uniform
252 distributions), and mice natively run to varying levels. Nonetheless, we reasoned that task
253 contingencies could still be inferred within trials based on the extent that tone transitions are
254 congruent with locomotion, and dopamine signals can be informed by such congruency.

255 To formalize this notion, we constructed a multi-agent mixture of experts (MoE) model,
256 extending earlier hierarchically nested corticostriatal circuit models of learning and decision
257 making^{24,25} (**Fig. 4a, Extended Data Fig. 7**, see *Methods* for details). At the highest layer (*level*
258 *1*) is an expert, putatively corresponding to DMS, that computes the online evidence for action-
259 outcome contingencies and thus task controllability (**Fig. 4a**). Sub-experts within that area (*level*
260 *2*) represent specific contingencies (e.g., distance needed to run is short, medium or long)
261 based on previous exposure to the tone transition distributions, learned as a semi-markov
262 decision process via temporal difference learning⁴⁰. Sub-expert prediction errors (PEs; level 3)
263 occur at tone transitions and are used to compute evidence for (or against) the accuracy of each
264 sub-expert's prediction. This formulation allows an agent to flexibly adapt behavior based on
265 task contingencies (**Extended Data Fig. 7**)^{24,25} and expands the RL account of dopamine to
266 allow both RPE and value signals to be informed by their inferred causal contributions⁴¹⁻⁴⁴.

267 This architecture makes novel predictions at multiple levels which can potentially tie
268 together the separable roles of dopamine during reward pursuit (performance) and learning.
269 First, when reward waves initiate in the DMS (i.e. ML waves in the instrumental task), that
270 region will receive the most credit, and hence mice will be faster to initiate running on the next
271 trial and will persist in doing so until rewards are obtained. Second, the reward wave dynamics
272 should be informed by a trace of which circuit ("expert") was most responsible for the reward
273 (i.e., which circuit's predictions were most valuable). We posited that DA dynamics during the
274 tone transitions (anticipatory epoch) could provide such a responsibility signal; that is, the sub-
275 circuit that best predicts the action-outcome contingencies will exhibit increases in dopamine.
276 These levels of dopamine can facilitate mice's motoric output to be guided most-strongly by that
277 subexpert, while also signaling the degree to which it is responsible for future rewards. We thus
278 hypothesized that DA dynamics during anticipation would impact how reward-waves circulate
279 among striatal subregions and the behavioral expression of running in future trials. Finally, at
280 the most fine-grained level, our model predicts that RPEs should occur at tone transitions to
281 inform the extent of instrumental controllability. In the remaining sections we unpack and test
282 each of these predictions.

283 The first prediction is that dopamine waves experienced at reward outcome reflects a
284 measure of credit assignment across the striatal experts. ML waves deliver dopamine first to
285 medial subregions (**Fig. 2i,j, 3i,j**), and these DMS-biased signals would selectively strengthen
286 corticostriatal representations for action-outcome contingencies that compete for instrumental
287 control in future trials. As such, we predicted that stronger ML waves at reward would enhance
288 instrumental learning that will drive future running. Indeed, we found a significant correlation
289 between the trial-by-trial magnitude of reward wave and latency to start running on the next trial
290 ($n=6$ mice, mean $r = -0.32$, $p = 0.0019$ two-sided t-test on correlation coefficients). Furthermore,
291 these wave magnitudes predicted the velocity even late in the next trial, 10.2 ± 1.4 seconds after
292 the reward response (**Fig. 4c**). The influence of these waves in future-trial behavioral
293 adjustments indicated that they are used for learning functions. Further, these effects were
294 selective to instrumental sessions, indicating that DMS sourced ML waves promote learning
295 about instrumental contingency that is employed for future reward pursuit.

296 ***Anticipatory dopamine ramps provide eligibility for credit assignment***

297 If reward-waves reflect credit assignment, what determines which subregion should
298 receive the credit? Canonical accounts in RL invoke dopamine RPEs that have graded effects
299 on learning depending on “eligibility” signaled by recent MSN activity^{45,46}. As noted above, we
300 considered the possibility that local dopamine dynamics during the anticipation epoch
301 themselves signal a coarser measure of eligibility in terms of which subregion was responsible.
302 This trace would be in proportion to the value of the underlying subregion’s predictions,
303 providing a tag for a subregion’s credit.

304 We thus focused on the activity of dopamine axons during the anticipatory period as
305 mice drew closer to reward. In the instrumental task, we observed a buildup of activity in the
306 DMS (**Fig. 4d,f**), ramping in proportion to the progress to reward^{15,47}. Strikingly, the opposite
307 profile was observed in the pavlovian condition (**Fig. 4e,f**), with decreasing ramps even as the
308 mouse continued to increase licking in anticipation of rewards. These findings do not support
309 extant models of DA ramps in accumbens or midbrain, where they have been linked to value
310 functions or RPEs^{15,48-50}, none of which predict opposite profiles across the two tasks.

311 Instead, we posited that anticipatory dopamine dynamics in the DMS reflects the
312 evidence of agency or controllability, and that subregions within might differentially represent
313 distinct controllable transition functions (which vary from trial to trial). Escalating tones in our
314 tasks provide information about online action-outcome contingency. For example, if tone
315 transitions consistently follow locomotion (as in **Fig. 3d**), they signal evidence for control. The
316 opposite inference can be made in the pavlovian task when tone transitions diverge from
317 locomotion. Respectively increasing or decreasing ramps in the instrumental and pavlovian
318 tasks accumulate in MoE ‘distance’ expert-weights as controllability is confirmed (or
319 contradicted) with each tone transition (**Fig. 4a, right, Extended Data Fig. 7**).

320 Thus, according to our model, DA ramps do not reflect a monolithic value function, but
321 rather the value of the underlying sub-region’s agentic predictions for reward pursuit, as a
322 marker of that region’s responsibility. Consequently, we argue here that the computational
323 function of dopamine waves at reward-outcome is to assign spatio-temporal credit by delivering

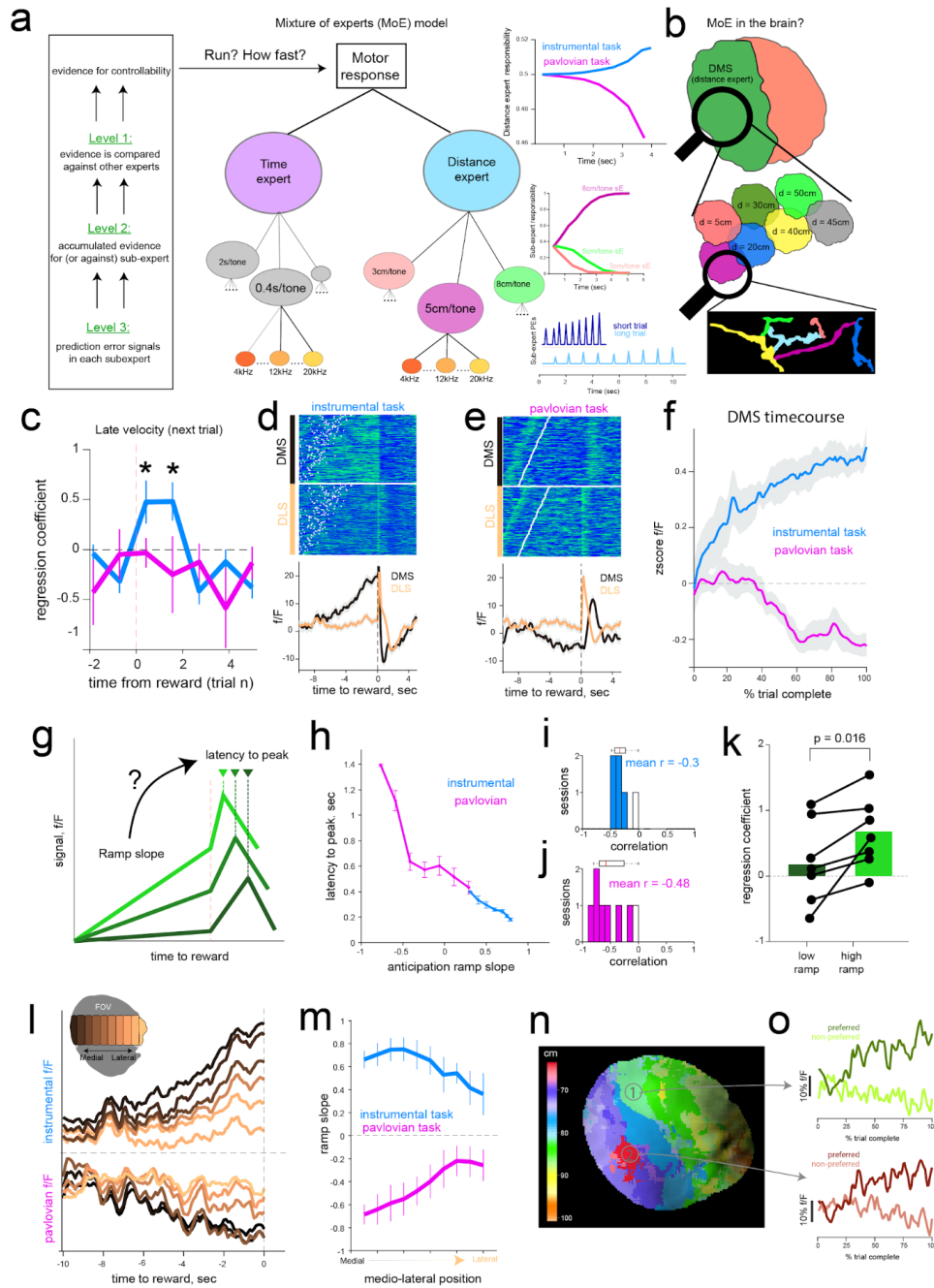
324 dopamine to striatal subregions with different latencies as a function of their graded
325 responsibility signals. This proposal is also motivated by theory and observations that
326 dopamine-mediated plasticity at striatal synapses is strongly attenuated with delayed dopamine
327 release^{45,46}.

328 This credit assignment interpretation makes additional testable predictions at both
329 physiological and behavioral levels. If dopamine ramps during reward anticipation hold
330 persistent information about a sub-region's prediction accuracy, they should modify the impact
331 of dopamine bursts at reward to focus preferentially on the sub-region with the highest
332 accuracy. As such, striatal areas that ramp with the steepest slopes during anticipation (highest
333 eligibility) should receive a reward response soonest (largest credit, **Fig. 4g**). Indeed,
334 anticipatory ramp slopes across pixels were significantly correlated with the fastest latency to
335 peak fluorescence following reward for both tasks (**Fig. 4h,i**).

336 Second, if DMS ramps signal responsibility for learning about instrumental control, then
337 trial-by-trial DMS ramp slopes should also modulate the impact of reward waves on next-trial
338 velocity. Indeed, we found that the impact of ML waves on future velocity in the instrumental
339 task (**Fig. 4c**) were dependent on the level of DMS ramps in the previous trial. When DMS
340 ramps were steep, reward waves strongly predicted speeded velocity in subsequent trials; this
341 effect was absent when ramps were weak (**Fig. 4k**, $p = 0.016$ Wilcoxon signed-rank test, $n=6$
342 mice). Together these results suggest that anticipatory dopamine ramps provide a tag for how
343 midbrain driven reward-credit circulate across the striatum to deliver a reinforcement signal for
344 future performance.

345

Figure 4



346
347

348 **Figure 4: Anticipatory epoch dopamine dynamics reflect inferred controllability, trial-**
349 **specific task statistics, and modulate reward responses in line with a mixture of striatal**
350 **experts.**

351 **a**, Schematic of hierarchical mixture of experts model as a framework for interpreting functional
352 relevance of dopamine dynamics applied to escalating tone tasks. At the lowest level, individual
353 states (representing auditory tones) induce reward prediction errors if they misalign with learned
354 contingencies. These prediction errors are accumulated within trials to provide evidence for or

355 against “sub-experts” specialized to represent local task contingencies (e.g. short, medium or long
356 durations/distances). At the highest level, instrumental or Pavlovian “experts” computes the
357 overall (weighted across sub-experts) evidence for instrumental task requirements, used to infer
358 controllability of the value function. These distance expert responsibility weights accumulate
359 within and across trials in the Instrumental task and decline in the Pavlovian task, and are used
360 to adjust model velocities. **b**, Proposed implementation of hierarchical task signals in striatal
361 dopamine activity. Widefield and 2-photon imaging at the micron level was used to test sub-
362 region-specific computations in dopamine terminals. **c**, Multiple regression predicting future
363 running speed of mice in the late phase of the next trial as a function of trial-by-trial wave
364 propagation-velocity (in 1 second bins) surrounding the reward from the previous trial. Reward-
365 induced wave velocity predicted future running speed in instrumental (blue) but not Pavlovian
366 (pink) sessions. Regression coefficients significantly different from zero (blue, asterix $p=0.005$,
367 two-tailed t-test). Error bars are S.E.M. **d**, Anticipatory and reward response in the medial and
368 lateral DS in a representative instrumental session. White points indicate start of trial. **e**, same
369 format as **d** for pavlovian session. **f**, Aggregate ramping profile during anticipatory epoch for the
370 DMS. Mean activity for each session was z-scored and averaged across mice. Activity in DMS
371 but not DLS showed task-dependent ramping profiles in line with inferred controllability. Shaded
372 regions represent S.E.M. **g**, Schematic for testing whether anticipatory epochs ramp slope is
373 related to the latency to peak dopamine in the outcome epoch. **h**, Results of the relationship from
374 two representative sessions, each from instrumental and pavlovian condition. For both tasks,
375 ramp slope was inversely related to subsequent latency to peak reward response. **i, j** summarize
376 the distribution of correlation coefficients across sessions. **k**, Multiple regression as in **c** for
377 instrumental sessions. In each session, trial-by-trial ramp slope was median split into low-ramp
378 or high-ramp trials. Across trials, steeper DMS ramps magnified the impact of reward waves on
379 subsequent trial running speed, in line with credit assignment. **l**, Anticipatory epoch ramps in a
380 sample instrumental (*top*) and pavlovian (*bottom*) sessions were expressed to varying extent on
381 the medio-lateral axis. ROIs were drawn with fixed distance from the edge of the field of view and
382 illustrated by inset. **m**, Quantification of ramp slope across sessions. Error bars represent S.E.M.
383 **n**, Local subregions within the dorsal striatum respond to distinct distance contingencies during
384 reward pursuit, reminiscent of sub-expert dynamics. Contingency specialization map in an
385 example session with color indicating the average distance that produces the steepest ramp
386 slopes for each pixel. **o**, Example time course of anticipatory ramps in two example subregions
387 for their respective preferred (high ramp trials), and non-preferred (low ramp trials). See Extended
388 Data Fig. 7c for similar pattern of activation in model simulations.

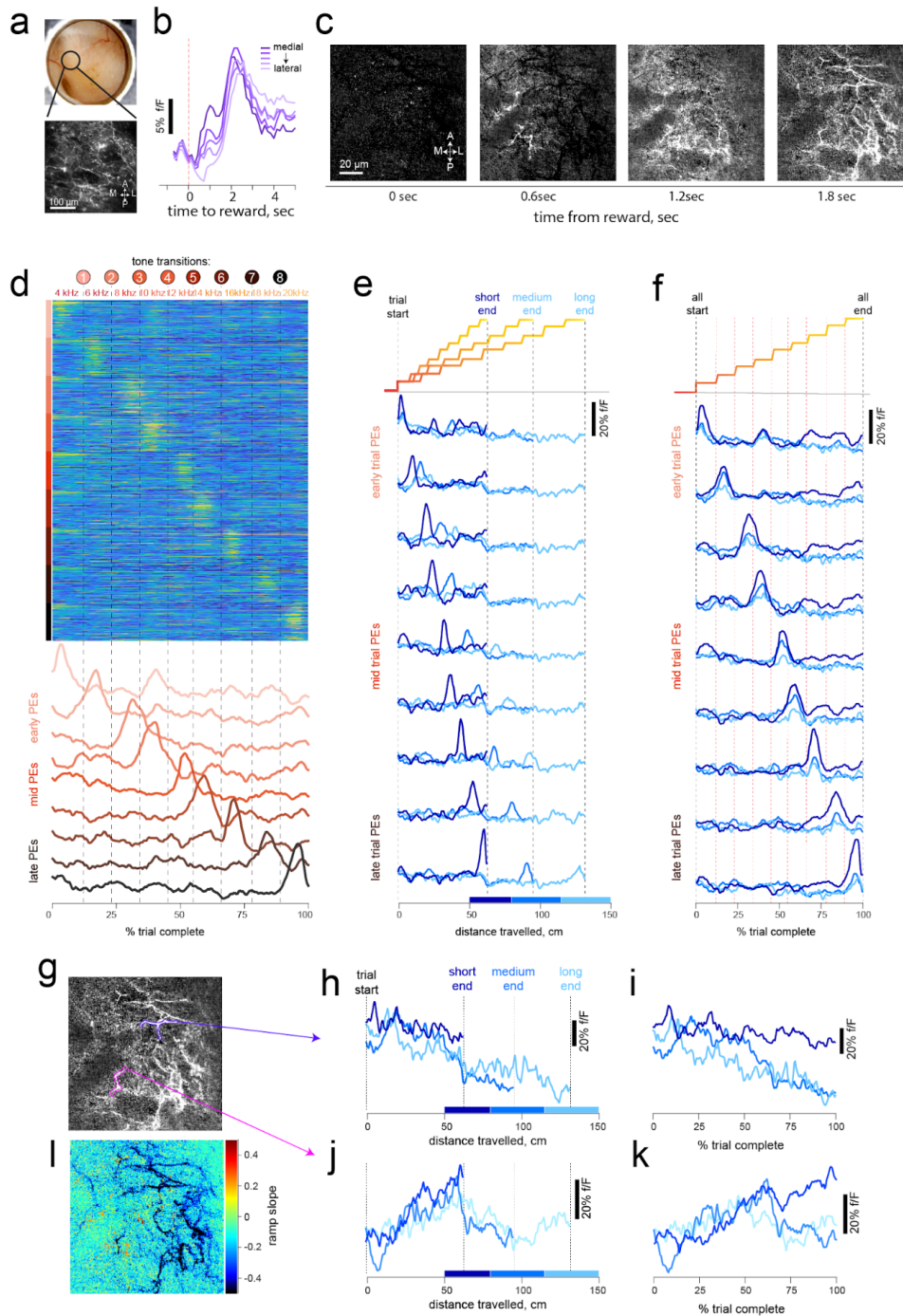
389 -----

390

391 Thus far, we have focused on the coarsest division of labor related to the highest level in
392 our model (controllability, level 1), but the agent's ability to infer control depends on underlying
393 sub-experts that learn distinct action-outcome contingencies (level 2, **Fig. 4a,b**). Such a
394 hierarchical scheme implies that striatal subregions should also differentially ramp for different
395 distance contingencies (**Fig. 4a**). Overall, we observed that DS dopamine ramps are expressed
396 in a gradient across both tasks, with the strongest ramps in the most medial portions (**Fig. 4l,m**).
397 These results are in line with previous work on progressive instrumental specialization of DS on
398 the mediolateral axis⁵¹. Moreover, contiguous territories within the DS exhibit varying ramp
399 profiles for different distance conditions (**Fig. 4n, Extended Data Fig. 8**), with each area
400 expressing the steepest dopamine ramps in preferred set of trials with related distance
401 requirements (**Fig. 4o**). On a trial-by-trial basis, we further observed a significant rank
402 correlation between each pixel's ramp slope and latency to peak response during reward (mean
403 $r = -0.13$, spearman's correlation $p < 0.001$ for all instrumental sessions). These results indicate
404 that the heterogeneously expressed anticipatory ramp gradients across the striatum modulate
405 the spread of reward waves, further strengthening the relationship between eligibility and credit
406 assignment. These findings further support our interpretation that is motivated by MoE account
407 by demonstrating that DMS consists of smaller sub-regions that learn, and express predictions
408 for a variety of potential instrumental contingencies.

409 These findings led us to ask whether waves organized the response of dopamine axons
410 on the micron scale, and functionally, how evidence for instrumental controllability accumulated
411 in single axon segments. The ramp-like responses we observed at the coarser scale using
412 widefield, one-photon imaging may emerge from trial-by-trial ramps within individual axons, or
413 from a weighted distribution of sharply tuned activation patterns. To directly address these
414 questions, we used 2-photon imaging in two mice to examine the behavior of individual axons in
415 the DMS (**Fig. 5a**).

Figure 5



416
417

418 **Figure 5: Single dopamine axons show wave-like reward dynamics, tone-specific**
419 **transients, and distance-dependent ramping during instrumental anticipation.**

420 **a**, Schematic of imaged region (*top*), and example field of view of dopamine axons in DMS. **b**,
421 Sequence of frames showing how individual dopamine axons respond to reward. Time relative to
422 solenoid click is shown below each frame. Note the activation of leftmost fibers first, then

423 progressive activation of more lateral axons. **c**, Average time course of reward response from
424 rectangular ROIs equally distributed along the ML axis. **d**, Activity of dopamine segments that
425 respond to tone transitions during anticipation. The activity in each trial is shown as percent trial
426 complete instead of alignment in time. *Top*, heatmap shows the trial-by-trial responses (106 total
427 trials) of groups of pixels in the 2-photon field of view that respond transiently at specific fraction
428 of trial completed. The responses of nine different types are concatenated. *Bottom*, average delta
429 f/F for each type. **E**, Transient response-peaks are not tuned to time or distance run within trial.
430 Each time series is aligned to trial-start and binned into 3 distances: short (50-80cm, dark blue),
431 medium (80-120cm), and long (120-150cm, light blue). Note that the peak location of response
432 arrives at different distances in each trial. By contrast, when aligned to % trial complete, in **f**, the
433 peak response arrives with fixed delay from tone transitions (illustrated at the top for both panels)
434 for all distance contingencies. The transient responses for each tone had larger amplitudes for
435 shorter trials (ie when rewards are predicted to occur sooner, in line with state-dependent reward
436 prediction errors within the lowest level of MoE). **g**, Individual axon segments highlighted to
437 demonstrate example ramp-like trajectories. **h**, Some axon segments ramped downward across
438 all distance trajectories when aligned to distance travelled or fraction of trial completed as in **i**. **j**,
439 An axon segment that progressively ramps upward only in short distance trials. **k**, same alignment
440 as **i**. **l**, pixelwise map of ramp slope during anticipation.
441 -----

442 Similar to our observations at the macro scale, reward delivery recruited dopamine
443 axons in a spatial sequence that was directional (**Fig. 5b,c**), demonstrating that wave-like
444 activation patterns also organize individual axon lattices on the micron scale.

445 The activity in individual dopamine axons were also modulated during the anticipation
446 epoch. Strikingly, segments of axons transiently responded to auditory tone transition, tiling the
447 full sequence of escalating tones (**Fig. 5d**). The timing of these responses was not affected by
448 distance travelled (**Fig. 5e**), but reliably responded to changes in tone frequency across a
449 variety of distance contingencies (**Fig. 5f**). The systematic tuning of these axons to tone-
450 transitions are consistent with PEs at the lowest level of our model (**Fig. 4a, Extended Data**
451 **Fig. 7**) that are used to update the online evidence of predictions within each sub-expert. Each
452 tone is represented as a unique state within a sub-expert's semi-markov process, and PEs arise
453 at tone transitions when they misalign with the predicted distance (or dwell time) until state
454 change. Furthermore, the model predicts larger PEs for state transitions indicative of rewards
455 that will arrive when the distance to run is shorter, due to temporal discounting (**Extended Data**
456 **Fig. 7**). Supporting this prediction, we observed that tone responses were largest in trials that
457 had shorter distance contingencies, and progressively decreased in amplitude for longer trials
458 (**Fig. 5e,f**, mean $r = -0.33$, and -0.14 $n=2$ mice; see **Extended Data Fig. 9**). These PE-like
459 responses were distributed throughout the 2-photon field of view, with equivalent fractions of
460 pixels selectively tuned to each tone transition (**Extended Data Fig. 9**).

461 We also noted that contiguous segments of axon lattices had single trial ramps that were
462 either upward (**Fig. 5h,i**) or downward (**Fig. 5j,k**) as mice get closer to reward. Instead of tuning
463 to tone-transitions reported above, these dopamine ramps were selectively expressed for
464 different contingencies, ramping to varying extents depending on the required distance in
465 separate trials. Together these results provide evidence for two, simultaneous classes of nearby
466 dopamine axon segments (**Extended Data Fig. 9**) used for sub-expert computations: transient
467 PE signals that respond to state transitions, and ramping segments that accumulate evidence
468 for controllability as predicted by a sub-expert.

469 Discussion

470 Our report of dopamine waves provides the earliest evidence for a foundational
471 organizational principle of dopamine axons that correlate activity within functionally related
472 striatal boundaries. In the cortex, travelling waves have been described to facilitate (or
473 constrain) computations that are topographically organized⁵²⁻⁵⁴. Similarly, we interpret the
474 computational significance of dopamine waves as orchestrating dopamine release to striatal
475 subregions that exhibit a graded functional specialization on the medio-lateral axis^{29,51,55}. Thus,
476 waves are a natural candidate for solving the spatiotemporal credit assignment problem when
477 multiple, topographically organized striatal actors/sub-experts compete to guide action selection
478 across multiple levels of abstraction^{24,25,56}. We used a very simple task to manipulate reward
479 and sensory statistics, requiring mice to resolve ambiguity about instrumental contingency by
480 comparing predicted and actual tone transitions. Consistent with the MoE account, wave
481 directions during reward were sensitive to controllability of task structure, and -- only in the

482 controllable task -- dopamine waves were related to future behavioral adjustment on a trial-by-
483 trial basis.

484 We also describe anticipatory epoch ramping dynamics that appear to signal the value of
485 a subregion's prediction about reward contingency. These dynamics may serve a dual purpose.
486 First, they could promote online behavioral flexibility (e.g., optimize reward-rate and minimize
487 energetic costs) according to the predictions of the most accurate subregions during reward
488 pursuit. Second, these activity patterns would also signal which subregion was most responsible
489 for behavioral output and hence provide a low dimensional tag for responsibility (akin to an
490 eligibility trace in RL⁵⁷), which would then allow for reward-driven RPEs to preferentially credit
491 the appropriate subregion and the eligible MSNs within it. While the two functions are not
492 mutually exclusive, our data provide strong support for the second interpretation: On a trial-by-
493 trial basis, the degree of ramping across regions was related to the latency to reward peak
494 elicited by the wave, and the combination of ramp slope and wave magnitude was predictive of
495 subsequent-trial behavioral adjustments. These findings accord with views that dopamine
496 signals can have different functions during reward pursuit and outcome, which can be gated by
497 local microcircuit elements that regulate plasticity windows^{3,58-61}. Moreover, we also interpret
498 transient and localized RPEs during reward pursuit as facilitating inference about the current
499 task state (i.e., determining credit), whereas RPEs during reward itself facilitates reinforcement
500 learning; a dual operation that can also be gated^{43,61-63}. Put together, the synthesis of our data
501 and computational simulations imply that dopamine signals are spatio-temporally vectorized
502 during both epochs, tailored to underlying region's computational specialty.

503 Although the dorsal striatal dopamine dynamics support the computations of the
504 Distance expert in the MoE, one limitation of our study is that we did not identify or assess the
505 dopamine dynamics with properties of the 'Time' expert. Many studies investigating RPEs
506 involve classical conditioning in which temporal representations are evident in the midbrain⁶⁴⁻⁶⁶.
507 Models and data have suggested that ramping signals related to timing may be present in other
508 regions upstream of the DA system⁶⁷⁻⁶⁹. Another limitation is that we did not deduce the origin
509 of dopamine waves, which may be inherited from sequential firing of midbrain dopamine cells
510 that have a topographical projection pattern⁷⁰. To date, such dynamics have not been reported
511 in the literature, potentially because limited studies investigated the activity of a large population
512 dopamine neurons simultaneously⁷¹. Another likely mechanism may involve local sculpting of
513 dopamine release within the striatum. Wave-like patterns have been reported in neocortex^{34,72}
514 and striatal cholinergic interneurons⁷³, both of which can potently regulate dopamine axon
515 activity⁷⁴⁻⁷⁶. Moreover, dopamine waves at reward outcome may also be a consequence of the
516 interaction between primed excitability of dopamine axons during the anticipatory epoch and
517 midbrain-sourced synchronous reward bursts. The combination of these two patterns may
518 produce sequential activation that propagates across the striatum in proportion to expressed
519 ramp during anticipation.

520 -----

521 **Methods**

522 **Animals and Surgery.** All procedures were conducted in accordance with the guidelines of the
523 NIH and approved by Brown University Institutional Animal Care and Use Committee. We used
524 17 DAT-cre mice (9 females, 8 males; *Jax Labs # 020080*) that were maintained on reversed
525 12hr cycle and all behavioral training and testing was performed during the dark phase. To
526 achieve selective expression of GCaMP6f in dopamine cells, we followed standard surgical
527 procedures for stereotaxic injection of cre-dependent virus. Briefly, mice were anesthetized with
528 isoflurane (2% induction and maintained at 0.75-1.25% in 1 liter/min oxygen). To attain
529 widespread infection of dopamine cells throughout the midbrain, we drilled two burr holes above
530 the midbrain (-3.2mm AP, 0.4mm and 1.0mm ML relative to bregma) and injected 0.1-0.2 μ L of
531 AAV-syn-Flex-GCaMP6f at two depths per burr hole (3.8 and 4.2 mm relative to brain surface).
532 We next secured a metal head post to the skull and implanted an imaging cannula over the
533 ipsilateral dorsal striatum. The cannula is a custom fabricated stainless-steel cylinder
534 (Microgroup; 3mm diameter and 2.5-3mm height) with a 3mm coverslip (CS-3R, Warner
535 Instruments) glued at the bottom with optical adhesive (Norad Optical #71). To insert the
536 cannula into the brain, a 3mm diameter craniotomy was first drilled over the striatum (at bregma,
537 centered on 2.0mm ML), and then gently removed the dura and slowly aspirated the overlying
538 cortex until white colossal fibers were clearly visible (~0.8-1.2mm from brain surface). These
539 fibers were also gently aspirated layer by layer until the underlying dorsal striatal tissue was
540 uniformly exposed. A sterile imaging cannula was progressively lowered until the coverslip
541 contacted striatal tissue uniformly. Dental cement was applied to secure implant to the skull and
542 mice received a single dose of slow-release buprenorphine and allowed to recover for 1-2
543 weeks with post-operative care.

544
545 **Behavioral Training.** After full recovery from surgery, mice underwent 2-3 days of habituation
546 in operant chambers outfitted with a 3D printed wheel (15 cm diameter), audio speakers and a
547 solenoid-gated liquid reward dispenser. Following acclimation, mice were water-restricted,
548 receiving 1ml/day in addition to water earned during task performance. We used custom
549 LabVIEW scripts to control operant boxes during training and testing in behavioral tasks. In the
550 first stage of training, mice received non-contingent rewards that were delivered randomly (3-15
551 second apart, uniform distribution) for 3 consecutive days. Next, training in a “pavlovian” task
552 began, wherein rewards were delivered after a variable delay from trial start. The start of each
553 trial is signaled by the onset of a 4.3kHz tone that continued to escalate in frequency in
554 proportion to fraction of trial completed. We used nine different frequencies that were selected
555 to minimize harmonic overlap; 4.3kHz, 6.2kHz, 8.3kHz, 10kHz, 12.4kHz, 14.1kHz, 16kHz
556 ,8.4kHz, 20kHz. Across trials, the duration to wait for reward is randomly drawn from a uniform
557 distribution (4-8 seconds). At the end of a trial, the auditory sound is turned off, and the solenoid
558 delivered 3 μ L of water reward to a spout in front of the mouse. Licking behavior is detected
559 using capacitive touch sensors (AT42QT1010, Sparkfun). In some catch trials, the initial 4.3kHz
560 tone turned off after 0.5s and the mouse did not have continuous information of progress to
561 reward. For clarity, we only focused on escalating-tone trials. The next trial started after a
562 variable inter-trial-interval of 3-8 seconds. After 2-3 weeks of the pavlovian task, activity of
563 dopamine axons in the striatum was imaged in a test chamber with a widefield and 2-photon

564 imaging system. The same animals were then further trained on a distance-variant of the same
565 task, where reward delivery is now contingent on mice running on the wheel. Mice were
566 exposed to the “instrumental” task in training chambers requiring them to run on the wheel to
567 traverse linearized distances, also randomly selected from a uniform distribution (50-150cm).
568 Progress to reward was indicated by the same tone frequencies, and the angular position of the
569 wheel recorded using a miniature rotary encoder (MA3A10250N, US Digital). All behavioral data
570 is digitized and stored to disc at 50Hz.

571
572 **Widefield and two-photon imaging.** All imaging was performed using a multi-photon
573 microscope with modular laser-scanning and light microscopy designed by Bruker/Prairie
574 Technologies. For widefield imaging, we used a full-spectrum LED illumination with FITC filter
575 cassette for illumination at 470nm and detection centered at 530nm. Images were acquired
576 using a CoolSnap ES2 CCD camera (global shutter, Photometrics) and synchronized with
577 behavioral events through TTL triggers. All widefield images during behavioral tasks were
578 acquired with a 4X objective (Olympus), 100ms exposure (10Hz) and 8X on-camera binning to
579 achieve a sample resolution of 40 μ m/pixel (unless indicated otherwise). Two-photon microscopy
580 was performed using a 20X air objective (Olympus) on the same imaging platform with a
581 femtosecond pulsed Ti:Sapphire laser source (MaiTai DeepSee, 980nm power measured at
582 objective was 20-50mW) that was scanned across the sample using a resonant (x-axis) and
583 non-resonant (y-axis) galvanometer scanning mirrors. Returning photons were collected through
584 an imaging path onto multi-alkali PMTs (R3896, Hamamatsu), and recorded frames were online-
585 averaged to achieve a sampling rate of 10-15Hz.

586
587 **Data Analysis and statistics.** All images were processed with custom routines in MATLAB.
588 Each session is preprocessed for image registration, and alignment to behavioral events based
589 on triggers. Movement artifacts and image drift in the XY plane were corrected using rigid-body
590 registration using a DFT-based method⁷⁷. To cluster the activity of dopamine axons, we used
591 the K-means algorithm in MATLAB. To compute robustness of clustering results, we used the
592 adjusted rand Index measure which computes the similarity of two clusters based on the
593 probability of member overlap (corrected for chance; 0=random clusters, 1=exact same
594 membership). To examine how robust the clustering results were, we re-clustered the same
595 dataset 100 times in K-means using random initialization and varied cluster limits. We compared
596 the extent that pixels were re-clustered into the same group using the adjusted rand index as an
597 indicator of robustness of underlying structure of the data that produced clusters (see Extended
598 Data Fig. 4. To additionally test the extent spatial relationship between the pixels, or the how
599 similarity in temporal activation influenced the identified clusters, we repeated the same analysis
600 but shuffled the spatial or temporal relationships between the pixels. To estimate the optimal
601 number of clusters within each dataset, we computed the Bayesian information criterion (BIC)
602 on the K-means algorithm.

603 We characterized flow patterns in dopamine waves by adapting standard optical flow
604 algorithms in machine vision that are adapted for imaging of fluorescence signals^{34,36,78}. Briefly,
605 flow trajectories were computed for any two successive frames as a displacement of intensity
606 across the pixels in time. This method allows us to evaluate a pixel-by-pixel velocity vector fields
607 that summarizes the direction and strength of flow at each pixel. While there are multiple

608 methods to achieve this calculation^{36,78}, we adapted a combined Global-Local (CGL)
609 algorithm^{79,80} that combines the Lucas-Kanade and Horn-Schunck methods. The frame-by-
610 frame vector fields calculated using the CGL method was further processed to extract sink and
611 source locations and also flow trajectories across multiple frames (**Fig. 2a, bottom**). The frame-
612 by-frame flow magnitude for each frame (or flow-velocity, with units of mm/second) is computed
613 by averaging the length of vectors at each pixel (e.g. **Fig. 2b, red**). The locations of sinks or
614 sources were estimated based on local vector orientations: i.e sinks are points of inward flow,
615 whereas sources are points of outward flow. We estimated the pixel-wise likelihood of sinks and
616 sources by simply computing the divergence of the vector field in each frame (“*divergence*”
617 function in MATLAB, **Supplementary Video 2**). The flow trajectory across frames were
618 calculated from vector fields using the “*stream3*” function in MATLAB from seeded pixels (e.g.
619 **Fig3I**).

620 For alignment of fluorescence time series, DMS and DLS masks were defined using one
621 of three methods: i) manual drawing, ii) boundaries using cluster boundaries (as in Fig. 1i, top)
622 and iii) uniformly spaced ROIs on the mediolateral axis (as in Fig. 4I, inset). Each animal
623 performed multiple behavioral sessions, and we used one session per animal (n=6 mice in
624 instrumental task, and n=8 mice in pavlovian task) that had the largest $\Delta f/F$ deviations to avoid
625 results from being dominated by a few animals.

626 To determine the influence of reward-wave on behavioral performance on the next trial,
627 we performed a multiple regression predicting the running velocity of mice late (i.e. 75-100% of
628 trial complete) in the next trial based on reward aligned wave magnitude (1-sec bins, **Fig. 4c**).
629 To determine whether DMS ramp slopes influenced how last-trial wave outcome, on the next
630 trial, we conditioned this analysis on the ramping profile in the DMS, median split into low, and
631 high ramp conditions (**Fig 4k**). We evaluated the correlation between the ramp slope and
632 latency to peak by first peak-normalizing the reward response in 2-sec window and finding the
633 time index (after reward) for which the fluorescence signal reached peak levels. To examine
634 whether anticipatory dopamine ramps had a preference for different distance conditions (**Fig.**
635 **4n**, also see **Extended Data Fig. 7**), we sorted the trials based on the expressed ramps in each
636 pixel and averaged the distance contingency in trials with top 90% ramping.

637 TIFF stacks of 2-photon images of dopamine axon segments were also pre-processed
638 for registration and alignment with behavioral data. To draw ROIs of these segments for
639 assessing organization of responses (**Extended Data Fig. 1**), we followed the Howe and
640 Dombeck²⁸. Otherwise, we generally used pixel-wise analyses.

641
642 **Computational model.** We modeled mouse behavior using a mixture of experts / multi-agent
643 RL architecture²⁵, extended here to accommodate the sequential tone structure with semi-
644 markov dynamics⁴⁰. We modeled the two task structures as separate “experts” that learned a
645 value function V as a function of either elapsed time as in classical temporal difference learning
646 applied to Pavlovian condition, or as a function of distance travelled. Because mice were trained
647 on both time and distance tasks, multiple sub-experts (representing clusters in mediolateral
648 coordinates of striatum) were pre-trained for 2000 trials to span a range of contingencies (e.g.,
649 400ms, 600ms, or 800ms per tone transition; or 5, 10 or 15cm). For simplicity, we modeled the
650 task with discrete sub-experts that specialized on (had been preferentially exposed to) particular
651 times/distances. However, one can easily generalize the framework to the continuous case

652 (e.g., using basis functions⁸¹) and the discrete space can be modeled with arbitrary resolution
653 by simply increasing the number of sub-experts. Moreover, various models have shown that
654 prediction errors can be used to segregate learning of different latent task states^{43,56}.

655 *Subexpert and expert learning.* The value function for each time sub-expert s estimates
656 the discounted future reward $V^s(X_{i,t}) = r(t) + \gamma V^s(X_{i,t+1})$ and was trained via temporal
657 differences⁸² based on reward prediction errors $\delta(X_{i,t}) = r(t) + \gamma V(X_{i,t+1}) - V(X_{i,t})$. Each auditory
658 tone was modeled as a distinct state $X_{i,t}$ or $X_{i,d}$ with semi-markov dynamics. That is, the onset of
659 each tone i would advance the state vector to the corresponding position even if the tone
660 occurred earlier or later in absolute time/distance. Thus the value function learned for each sub-
661 expert was tied to the current state (tone) and the (discretized) dwell time (t) or distance (d)
662 since it has been entered, and not to the absolute time or distance that passed from the onset of
663 the first state. This semi-markov process was based on the assumption that the tone stimuli
664 induce a neural state representation upon which TD is computed^{40,81} and evidence that rodents
665 are endowed with such a rich state representation⁸³. The value function was learned by
666 adjusting weights in response to the X state vector, with $V(X_{i,t}) = w_t X_{i,t}$ and $w_t \leftarrow w_t + \alpha \delta(t)$,
667 where α is a learning rate. The distance experts were trained analogously but with the X vector
668 advancing with discretized distance steps taken rather than passive time. Thus if the agent
669 stopped moving, the $X_{i,d}$ vector remained constant until it moved again, and if it moved faster
670 than usual, the $X_{i,d}$ vector would advance to later states accordingly. We fixed $\alpha=0.25$ and γ
671 $=.95$ for all experts but verified that the patterns were robust to other settings.

672 *Performance and inference.* After learning, the on-line evidence (responsibilities, fig 4B,
673 modeling the ramps) for each sub-expert was computed as an approximation to the likelihood of
674 the trial-wise tone transitions for that sub-expert given the dwelling time or distance since the
675 last tone. We adopted a hybrid Bayesian-RL formulation²⁵. Each expert learns a value function
676 associated with the tones, and inference about which expert is responsible is computed in
677 proportion to the log likelihood ratio of the observations (tone transitions at particular instances
678 following elapsed time or distance) given their predictions relative to the other experts.

679 From a Bayesian perspective, the attentional weights for each expert can be
680 evaluated by computing the posterior probability that each expert encompasses the best
681 account of the observed data x : $P(s|x) = P(x|s) P(s) / \sum P(x|s_i) P(s_i)$. Thus the evidence for each
682 expert is computed by considering its prior evidence and the likelihood that the observed tone
683 transitions or rewards (positive or negative) would have been observed under the expert's
684 model, relative to all other experts. For example, if there was a low probability for a tone
685 transition at a particular moment under a given expert, then the likelihood of that observation
686 given the expert's model is low. Once the posterior evidence for each expert is computed, one
687 can then apply Bayesian model averaging to allocate attentional weights to each expert in
688 proportion to their log evidence.

689 Rather than a fully Bayesian realization, we instead implemented an RL approximation
690 that we posited would more directly relate to corticostriatal DA mechanisms²⁵. Instead of
691 computing the likelihood directly, each expert was penalized as a function of its reward
692 prediction errors. In particular we updated the evidence for each sub-expert's predictions in
693 terms of a responsibility weight ω'_s which was decremented when the corresponding sub-expert
694 experienced a reward prediction error: $\omega'_s \leftarrow \omega'_s - \delta_s$, where δ_s is the positive reward
695 prediction error according to the corresponding sub-expert's value function at the state vector X .

696 (Similar results hold if using $|\delta_s|$ instead of only positive RPEs to decrement expert weights, but
697 positive RPEs dominate). Intuitively, experts with more prediction errors are less likely to have
698 been responsible for the outcome (tone transition or reward). These responsibility weights were
699 then normalized relative to all sub-experts as an approximation to the log evidence for a given
700 subexpert: $\omega_{si} = \exp(\beta \omega_{si}) / \sum_j \exp(\beta \omega_{sj})$, where β is an inverse temperature parameter.
701 Thus, in contrast to standard RL in which RPEs reinforce actions that yield rewards, during
702 inference, more frequent RPEs for a given subexpert are indicative that it is less responsible for
703 observations compared to those that predict these observations. Such a scheme is compatible
704 with extant models that use reward prediction errors for state creation and inference separate
705 from reinforcement per se^{25,43,56,63}. We posited that these RPEs correspond to the phasic
706 events observed at tone transitions in the 2p imaging data. The accumulation of these
707 responsibility weights were posited to relate to the 1p imaging data in discrete sub-regions of
708 DMS.

709 Finally, a second-level task selection process was implemented to arbitrate responsibility
710 between the overall distance expert and overall time expert (each of which constituted a
711 weighted combination of their subordinate experts). This inference process was identical to that
712 for the sub-experts, with responsibility updated based on their experienced prediction errors: $\omega_D \leftarrow \omega_D - \delta_D$,
713 where ω_D is the accumulated responsibility of the distance expert based on its
714 reward prediction errors, $\delta_D = r(t) + \gamma V_D(t+1) - V_D(t)$. The value function for the distance and
715 time experts V_D and V_T are in turn weighted averages according to the inferred responsibilities of
716 the subordinate experts within each structure: $V_D(t) = \sum \omega_{sD} V_{sD}(t)$ and $V_T = \sum \omega_{sT} V_{sT}(t)$. Similarly,
717 the value function of the agent as a whole is the weighted average value function across the two
718 experts $V(t) = \omega_D V_D(t) + \omega_T V_T(t)$. These responsibility weights for each task structure were
719 again normalized across tasks, $\omega_D = e^{\beta \omega_D} / (e^{\beta \omega_D} + e^{\beta \omega_T})$.

720 For each distance or time, 100 test trials were run with 10 tones each and an inter trial
721 interval was randomly drawn from 5-15s. The agent as a whole selects actions in terms of
722 speeds to run for a period of time at each tone transition or after it has completed it's previous
723 running. Speeds were selected in proportion to the inferred responsibility of the DMS expert,
724 together with some stochasticity: $\text{speed}(t) = 5 * (\omega_D(t) - 0.5) + \epsilon$, where ϵ was drawn from a
725 uniform distribution with a mean of 3. Stochasticity facilitates the agent ability to disambiguate
726 distance from time tasks within a trial (a constant speed would equate the prediction errors for
727 the two tasks given appropriate sub-experts). Increasing speed with inferred DMS expert
728 responsibility ω_D allows the model to capture the increased running with instrumental task
729 structure (extended Fig 7). More detailed investigation of how speeds may be optimized
730 according to reward/effort/delay tradeoffs will be examined in future work.

731 **References**

- 732 1 Schultz, W. Dopamine reward prediction-error signalling: a two-component response.
733 *Nat Rev Neurosci* **17**, 183-195, doi:10.1038/nrn.2015.26 (2016).
- 734 2 Berridge, K. C. The debate over dopamine's role in reward: the case for incentive
735 salience. *Psychopharmacology (Berl)* **191**, 391-431, doi:10.1007/s00213-006-0578-x
736 (2007).
- 737 3 Berke, J. D. What does dopamine mean? *Nat Neurosci* **21**, 787-793,
738 doi:10.1038/s41593-018-0152-y (2018).
- 739 4 Schultz, W., Dayan, P. & Montague, P. R. A neural substrate of prediction and reward.
740 *Science* **275**, 1593-1599, doi:10.1126/science.275.5306.1593 (1997).
- 741 5 Prensa, L. & Parent, A. The nigrostriatal pathway in the rat: a single-axon study of the
742 relationship between dorsal and ventral tier nigral neurons and the striosome/matrix
743 striatal compartments. *Journal of Neuroscience* **21**, 7247-7260 (2001).
- 744 6 Matsuda, W. *et al.* Single nigrostriatal dopaminergic neurons form widely spread and
745 highly dense axonal arborizations in the neostriatum. *J Neurosci* **29**, 444-453,
746 doi:10.1523/JNEUROSCI.4029-08.2009 (2009).
- 747 7 Hyland, B. I., Reynolds, J. N., Hay, J., Perk, C. G. & Miller, R. Firing modes of midbrain
748 dopamine cells in the freely moving rat. *Neuroscience* **114**, 475-492 (2002).
- 749 8 Li, W., Doyon, W. M. & Dani, J. A. Acute in vivo nicotine administration enhances
750 synchrony among dopamine neurons. *Biochemical pharmacology* **82**, 977-983 (2011).
- 751 9 Joshua, M. *et al.* Synchronization of midbrain dopaminergic neurons is enhanced by
752 rewarding events. *Neuron* **62**, 695-704, doi:10.1016/j.neuron.2009.04.026 (2009).
- 753 10 Kim, Y., Wood, J. & Moghaddam, B. Coordinated activity of ventral tegmental neurons
754 adapts to appetitive and aversive learning. *PLoS One* **7**, e29766,
755 doi:10.1371/journal.pone.0029766 (2012).
- 756 11 Mohebi, A. *et al.* Dissociable dopamine dynamics for learning and motivation. *Nature*
757 **570**, 65-70, doi:10.1038/s41586-019-1235-y (2019).
- 758 12 Eshel, N., Tian, J., Bukwich, M. & Uchida, N. Dopamine neurons share common
759 response function for reward prediction error. *Nat Neurosci* **19**, 479-486,
760 doi:10.1038/nn.4239 (2016).
- 761 13 Schultz, W. Predictive reward signal of dopamine neurons. *J Neurophysiol* **80**, 1-27,
762 doi:10.1152/jn.1998.80.1.1 (1998).
- 763 14 Glimcher, P. W. Understanding dopamine and reinforcement learning: the dopamine
764 reward prediction error hypothesis. *Proc Natl Acad Sci U S A* **108 Suppl 3**, 15647-
765 15654, doi:10.1073/pnas.1014269108 (2011).
- 766 15 Hamid, A. A. *et al.* Mesolimbic dopamine signals the value of work. *Nat Neurosci* **19**,
767 117-126, doi:10.1038/nn.4173 (2016).
- 768 16 Haber, S. N. The primate basal ganglia: parallel and integrative networks. *J Chem*
769 *Neuroanat* **26**, 317-330, doi:10.1016/j.jchemneu.2003.10.003 (2003).
- 770 17 Graybiel, A. M. Habits, rituals, and the evaluative brain. *Annu Rev Neurosci* **31**, 359-387,
771 doi:10.1146/annurev.neuro.29.051605.112851 (2008).
- 772 18 Hunnicutt, B. J. *et al.* A comprehensive excitatory input map of the striatum reveals novel
773 functional organization. *Elife* **5**, e19103 (2016).
- 774 19 Hintiryan, H. *et al.* The mouse cortico-striatal projectome. *Nature neuroscience* **19**, 1100
775 (2016).
- 776 20 Riedel, A. *et al.* Principles of rat subcortical forebrain organization: a study using
777 histological techniques and multiple fluorescence labeling. *Journal of chemical*
778 *neuroanatomy* **23**, 75-104 (2002).
- 779 21 Brown, H. D., McCutcheon, J. E., Cone, J. J., Ragozzino, M. E. & Roitman, M. F.
780 Primary food reward and reward-predictive stimuli evoke different patterns of phasic

- 781 dopamine signaling throughout the striatum. *Eur J Neurosci* **34**, 1997-2006,
782 doi:10.1111/j.1460-9568.2011.07914.x (2011).
- 783 22 Shnitko, T. A. & Robinson, D. L. Regional variation in phasic dopamine release during
784 alcohol and sucrose self-administration in rats. *ACS Chem Neurosci* **6**, 147-154,
785 doi:10.1021/cn500251j (2015).
- 786 23 Menegas, W., Babayan, B. M., Uchida, N. & Watabe-Uchida, M. Opposite initialization to
787 novel cues in dopamine signaling in ventral and posterior striatum in mice. *Elife* **6**,
788 e21886 (2017).
- 789 24 Doya, K., Samejima, K., Katagiri, K.-i. & Kawato, M. Multiple model-based reinforcement
790 learning. *Neural computation* **14**, 1347-1369 (2002).
- 791 25 Frank, M. J. & Badre, D. Mechanisms of hierarchical reinforcement learning in
792 corticostriatal circuits 1: computational analysis. *Cerebral cortex* **22**, 509-526 (2011).
- 793 26 Gershman, S. J., Pesaran, B. & Daw, N. D. Human reinforcement learning subdivides
794 structured action spaces by learning effector-specific values. *J Neurosci* **29**, 13524-
795 13531, doi:10.1523/JNEUROSCI.2469-09.2009 (2009).
- 796 27 Badre, D. & Frank, M. J. Mechanisms of hierarchical reinforcement learning in cortico-
797 striatal circuits 2: Evidence from fMRI. *Cerebral cortex* **22**, 527-536 (2011).
- 798 28 Howe, M. & Dombeck, D. A. Rapid signalling in distinct dopaminergic axons during
799 locomotion and reward. *Nature* **535**, 505 (2016).
- 800 29 Klaus, A. *et al.* The spatiotemporal organization of the striatum encodes action space.
801 *Neuron* **95**, 1171-1180. e1177 (2017).
- 802 30 Balleine, B. W., Delgado, M. R. & Hikosaka, O. The role of the dorsal striatum in reward
803 and decision-making. *Journal of Neuroscience* **27**, 8161-8165 (2007).
- 804 31 Yin, H. H. & Knowlton, B. J. The role of the basal ganglia in habit formation. *Nature*
805 *Reviews Neuroscience* **7**, 464 (2006).
- 806 32 Grinvald, A., Lieke, E. E., Frostig, R. D. & Hildesheim, R. Cortical point-spread function
807 and long-range lateral interactions revealed by real-time optical imaging of macaque
808 monkey primary visual cortex. *Journal of Neuroscience* **14**, 2545-2568 (1994).
- 809 33 Lubenov, E. V. & Siapas, A. G. Hippocampal theta oscillations are travelling waves.
810 *Nature* **459**, 534 (2009).
- 811 34 Mohajerani, M. H. *et al.* Spontaneous cortical activity alternates between motifs defined
812 by regional axonal projections. *Nature neuroscience* **16**, 1426 (2013).
- 813 35 Muller, L., Chavane, F., Reynolds, J. & Sejnowski, T. J. Cortical travelling waves:
814 mechanisms and computational principles. *Nature Reviews Neuroscience* **19**, 255
815 (2018).
- 816 36 Afrashteh, N., Inayat, S., Mohsenvand, M. & Mohajerani, M. H. Optical-flow analysis
817 toolbox for characterization of spatiotemporal dynamics in mesoscale optical imaging of
818 brain activity. *NeuroImage* **153**, 58-74 (2017).
- 819 37 Corbit, L. H. & Janak, P. H. Posterior dorsomedial striatum is critical for both selective
820 instrumental and Pavlovian reward learning. *Eur J Neurosci* **31**, 1312-1321,
821 doi:10.1111/j.1460-9568.2010.07153.x (2010).
- 822 38 Balleine, B. W. & O'doherty, J. P. Human and rodent homologues in action control:
823 corticostriatal determinants of goal-directed and habitual action.
824 *Neuropsychopharmacology* **35**, 48 (2010).
- 825 39 Wunderlich, K., Smittenaar, P. & Dolan, R. J. Dopamine enhances model-based over
826 model-free choice behavior. *Neuron* **75**, 418-424 (2012).
- 827 40 Daw, N. D., Courville, A. C. & Touretzky, D. S. Representation and timing in theories of
828 the dopamine system. *Neural computation* **18**, 1637-1677 (2006).
- 829 41 Chang, Y.-H., Ho, T. & Kaelbling, L. P. in *Advances in neural information processing*
830 *systems*. 807-814.

- 831 42 O'Reilly, R. C. & Frank, M. J. Making working memory work: a computational model of
832 learning in the prefrontal cortex and basal ganglia. *Neural computation* **18**, 283-328
833 (2006).
- 834 43 Gershman, S. J., Norman, K. A. & Niv, Y. Discovering latent causes in reinforcement
835 learning. *Current Opinion in Behavioral Sciences* **5**, 43-50 (2015).
- 836 44 Russell, S. J. & Zimdars, A. in *Proceedings of the 20th International Conference on*
837 *Machine Learning (ICML-03)*. 656-663.
- 838 45 Yagishita, S. *et al.* A critical time window for dopamine actions on the structural plasticity
839 of dendritic spines. *Science* **345**, 1616-1620, doi:10.1126/science.1255514 (2014).
- 840 46 Shindou, T., Shindou, M., Watanabe, S. & Wickens, J. A silent eligibility trace enables
841 dopamine-dependent synaptic plasticity for reinforcement learning in the mouse
842 striatum. *European Journal of Neuroscience* **49**, 726-736 (2019).
- 843 47 Howe, M. W., Tierney, P. L., Sandberg, S. G., Phillips, P. E. & Graybiel, A. M. Prolonged
844 dopamine signalling in striatum signals proximity and value of distant rewards. *Nature*
845 **500**, 575-579, doi:10.1038/nature12475 (2013).
- 846 48 Lloyd, K. & Dayan, P. Tamping Ramping: Algorithmic, Implementational, and
847 Computational Explanations of Phasic Dopamine Signals in the Accumbens. *PLoS*
848 *Comput Biol* **11**, e1004622, doi:10.1371/journal.pcbi.1004622 (2015).
- 849 49 Gershman, S. J. Dopamine ramps are a consequence of reward prediction errors.
850 *Neural Comput* **26**, 467-471, doi:10.1162/NECO_a_00559 (2014).
- 851 50 Morita, K. & Kato, A. Striatal dopamine ramping may indicate flexible reinforcement
852 learning with forgetting in the cortico-basal ganglia circuits. *Frontiers in neural circuits* **8**,
853 36 (2014).
- 854 51 Thorn, C. A., Atallah, H., Howe, M. & Graybiel, A. M. Differential dynamics of activity
855 changes in dorsolateral and dorsomedial striatal loops during learning. *Neuron* **66**, 781-
856 795 (2010).
- 857 52 Muller, L., Reynaud, A., Chavane, F. & Destexhe, A. The stimulus-evoked population
858 response in visual cortex of awake monkey is a propagating wave. *Nature*
859 *communications* **5**, 3675 (2014).
- 860 53 Nauhaus, I., Busse, L., Carandini, M. & Ringach, D. L. Stimulus contrast modulates
861 functional connectivity in visual cortex. *Nature neuroscience* **12**, 70 (2009).
- 862 54 Binguier, V., Chavane, F., Glaeser, L. & Frégnac, Y. Horizontal propagation of visual
863 activity in the synaptic integration field of area 17 neurons. *Science* **283**, 695-699 (1999).
- 864 55 Barbera, G. *et al.* Spatially compact neural clusters in the dorsal striatum encode
865 locomotion relevant information. *Neuron* **92**, 202-213 (2016).
- 866 56 Collins, A. G. & Frank, M. J. Cognitive control over learning: Creating, clustering, and
867 generalizing task-set structure. *Psychological review* **120**, 190 (2013).
- 868 57 Singh, S. P. & Sutton, R. S. Reinforcement learning with replacing eligibility traces.
869 *Machine learning* **22**, 123-158 (1996).
- 870 58 Morris, G., Arkadir, D., Nevet, A., Vaadia, E. & Bergman, H. Coincident but distinct
871 messages of midbrain dopamine and striatal tonically active neurons. *Neuron* **43**, 133-
872 143, doi:10.1016/j.neuron.2004.06.012 (2004).
- 873 59 Threlfell, S. & Cragg, S. J. Dopamine signaling in dorsal versus ventral striatum: the
874 dynamic role of cholinergic interneurons. *Front Syst Neurosci* **5**, 11,
875 doi:10.3389/fnsys.2011.00011 (2011).
- 876 60 Bradfield, L. A., Bertran-Gonzalez, J., Chieng, B. & Balleine, B. W. The thalamostriatal
877 pathway and cholinergic control of goal-directed action: interlacing new with existing
878 learning in the striatum. *Neuron* **79**, 153-166 (2013).
- 879 61 Franklin, N. T. & Frank, M. J. A cholinergic feedback circuit to regulate striatal population
880 uncertainty and optimize reinforcement learning. *Elife* **4**, e12029 (2015).

- 881 62 Schoenbaum, G., Stalnaker, T. A. & Niv, Y. How did the chicken cross the road? With
882 her striatal cholinergic interneurons, of course. *Neuron* **79**, 3-6 (2013).
- 883 63 Redish, A. D., Jensen, S., Johnson, A. & Kurth-Nelson, Z. Reconciling reinforcement
884 learning models with behavioral extinction and renewal: implications for addiction,
885 relapse, and problem gambling. *Psychological review* **114**, 784 (2007).
- 886 64 Pan, W.-X., Schmidt, R., Wickens, J. R. & Hyland, B. I. Dopamine cells respond to
887 predicted events during classical conditioning: evidence for eligibility traces in the
888 reward-learning network. *Journal of Neuroscience* **25**, 6235-6242 (2005).
- 889 65 Hollerman, J. R. & Schultz, W. Dopamine neurons report an error in the temporal
890 prediction of reward during learning. *Nature neuroscience* **1**, 304 (1998).
- 891 66 Soares, S., Atallah, B. V. & Paton, J. J. Midbrain dopamine neurons control judgment of
892 time. *Science* **354**, 1273-1277 (2016).
- 893 67 Brown, J., Bullock, D. & Grossberg, S. How the basal ganglia use parallel excitatory and
894 inhibitory learning pathways to selectively respond to unexpected rewarding cues.
895 *Journal of Neuroscience* **19**, 10502-10511 (1999).
- 896 68 Hazy, T. E., Frank, M. J. & O'Reilly, R. C. Neural mechanisms of acquired phasic
897 dopamine responses in learning. *Neuroscience & Biobehavioral Reviews* **34**, 701-720
898 (2010).
- 899 69 Mello, G. B., Soares, S. & Paton, J. J. A scalable population code for time in the
900 striatum. *Current Biology* **25**, 1113-1122 (2015).
- 901 70 Lerner, T. N. *et al.* Intact-Brain Analyses Reveal Distinct Information Carried by SNc
902 Dopamine Subcircuits. *Cell* **162**, 635-647, doi:10.1016/j.cell.2015.07.014 (2015).
- 903 71 Engelhard, B. *et al.* Specialized coding of sensory, motor and cognitive variables in VTA
904 dopamine neurons. *Nature*, 1 (2019).
- 905 72 Kasanetz, F., Riquelme, L. A., Della-Maggiore, V., O'Donnell, P. & Murer, M. G.
906 Functional integration across a gradient of corticostriatal channels controls UP state
907 transitions in the dorsal striatum. *Proceedings of the National Academy of Sciences* **105**,
908 8124-8129 (2008).
- 909 73 Rehani, R. *et al.* Activity patterns in the neuropil of striatal cholinergic interneurons in
910 freely moving mice represent their collective spiking dynamics. *Eneuro* **6** (2019).
- 911 74 Krebs, M. *et al.* Glutamatergic control of dopamine release in the rat striatum: evidence
912 for presynaptic N-methyl-D-aspartate receptors on dopaminergic nerve terminals.
913 *Journal of neurochemistry* **56**, 81-85 (1991).
- 914 75 Threlfell, S. *et al.* Striatal dopamine release is triggered by synchronized activity in
915 cholinergic interneurons. *Neuron* **75**, 58-64, doi:10.1016/j.neuron.2012.04.038 (2012).
- 916 76 Cachepe, R. *et al.* Selective activation of cholinergic interneurons enhances accumbal
917 phasic dopamine release: setting the tone for reward processing. *Cell Rep* **2**, 33-41,
918 doi:10.1016/j.celrep.2012.05.011 (2012).
- 919 77 Guizar-Sicairos, M., Thurman, S. T. & Fienup, J. R. Efficient subpixel image registration
920 algorithms. *Optics letters* **33**, 156-158 (2008).
- 921 78 Townsend, R. G. & Gong, P. Detection and analysis of spatiotemporal patterns in brain
922 activity. *PLoS computational biology* **14**, e1006643 (2018).
- 923 79 Bruhn, A., Weickert, J. & Schnörr, C. in *Joint Pattern Recognition Symposium*. 454-462
924 (Springer).
- 925 80 Liu, C. *Beyond pixels: exploring new representations and applications for motion*
926 *analysis*, Massachusetts Institute of Technology, (2009).
- 927 81 Ludvig, E. A., Sutton, R. S. & Kehoe, E. J. Stimulus representation and the timing of
928 reward-prediction errors in models of the dopamine system. *Neural computation* **20**,
929 3034-3054 (2008).
- 930 82 Sutton, R. S. & Barto, A. G. *Reinforcement learning : an introduction*. illustrated edn,
931 (MIT Press, 1998).

932 83 Gardner, M. P., Schoenbaum, G. & Gershman, S. J. Rethinking dopamine as
933 generalized prediction error. *Proceedings of the Royal Society B* **285**, 20181645 (2018).

934 **Acknowledgements.** We thank Matthew Nassar, Joshua Berke, Peter Dayan, Theresa
935 Desrochers, for valuable discussion of the project and feedback on an earlier version of the
936 manuscript, and members of the Frank and Moore laboratories for feedback at various stages of
937 the project. We also thank Ines Belghiti and Aneri Soni for help with adapting the MoE model
938 implementation to the tone task. GCaMP6 was developed and generously made available by
939 the HHMI Janelia Farm Research Campus GENIE Project. This work was supported by HHMI
940 Hanna Gray Fellowship to AAH, R01MH080066 and NSF grant 1460604 to MJF, and awards
941 from Carney Institute for Brain Sciences and Dean's office to CIM.

942 **Contributions.** AAH designed and performed the experiments, analyzed the data, and applied
943 computational model and wrote the paper. MJF designed and supervised the study, developed
944 computational model, and contributed to writing and revising of the paper. CIM designed and
945 supervised the study and revised the paper.

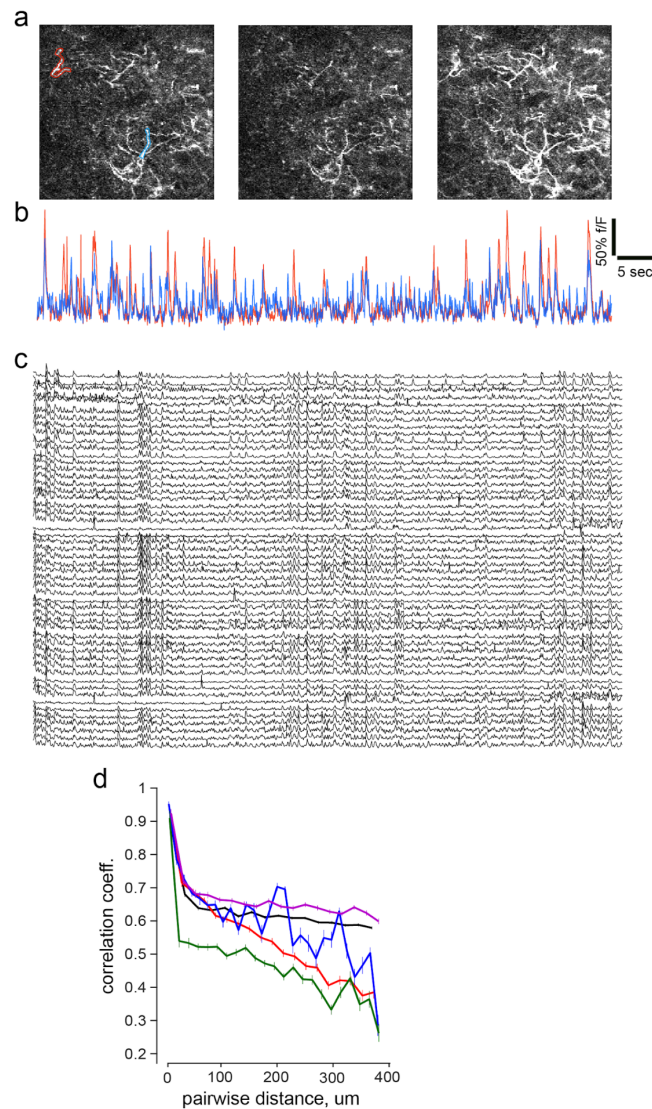
946 **Competing Interests:** The authors declare no competing interests.

947 **Data and code availability.** All data and code is available from corresponding author(s) upon
948 reasonable request.

949 **Materials and Correspondence:** Correspondence and material request is to be addressed to
950 Arif Hamid (Arif_Hamid@brown.edu), Christopher Moore (Christopher_Moore@brown.edu) or
951 Michael Frank (Michael_Frank@brown.edu).

952 -----

Extended Data Figure 1

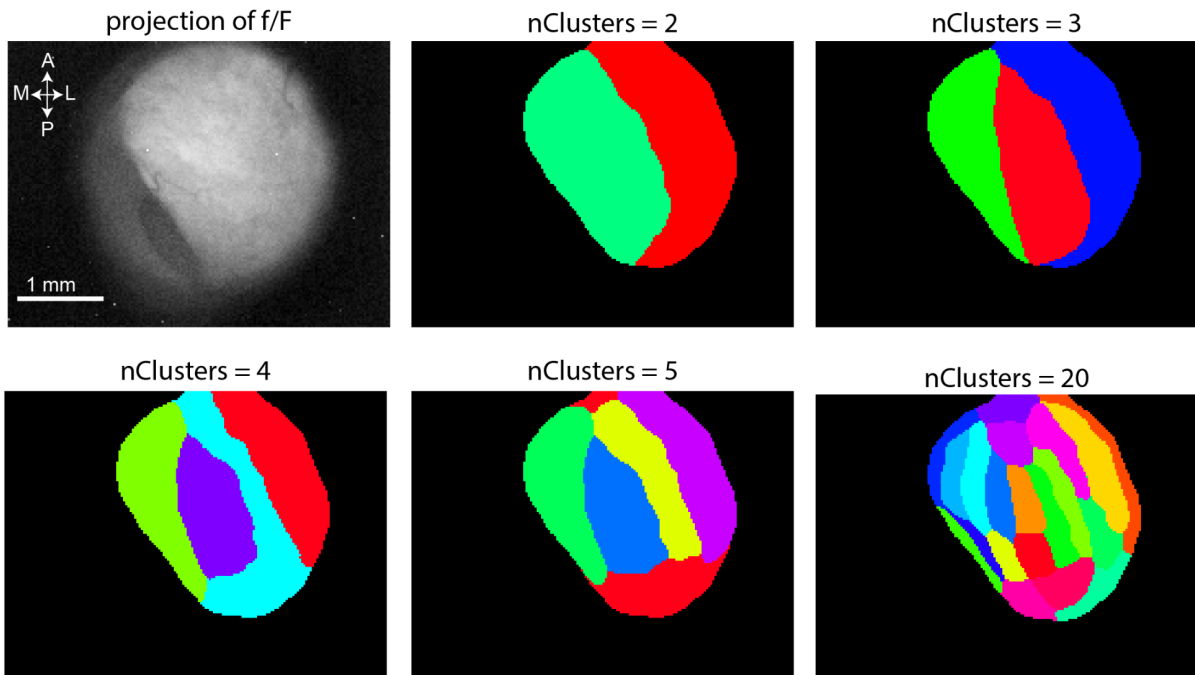


953

954 **Extended Data Fig. 1: Individual dopamine axons also exhibit decorrelated activity**
955 **patterns.**

956 **a**, Example frames illustrating that different portions of axon laticies are activated
957 asynchronously. **b**, Representative timeseries of fluorescence from two axon segments outlined
958 in blue and red at **a**. **c**, Additional examples of activity in dopamine axon segments. Data is
959 organized such that nearby axons are plotted closer. **d**, Quantification of correlation between the
960 sessionwide timeseries of axons based anatomical distances. Note that nearby axons are
961 highly correlated, but they exhibit a distance dependent falloff as reported in Fig. 1g, although a
962 different anatomical scale.

Extended Data Figure 2



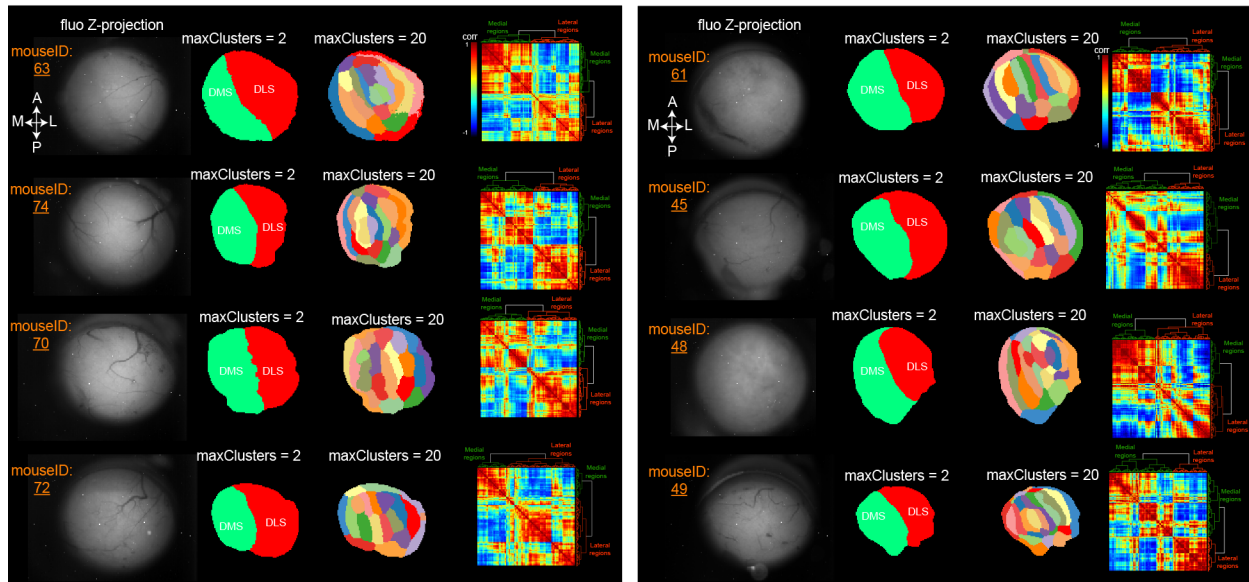
963

964 **Extended Data Fig. 2: Dopamine axon activity clusters into hierarchical domains.**

965 Data from one session, mean projection of fluorescence is displayed at the top left.

966 Progressively increasing cluster limits identifies contiguous striatal subregions that decompose
967 into sub-clusters.

Extended Data Figure 3



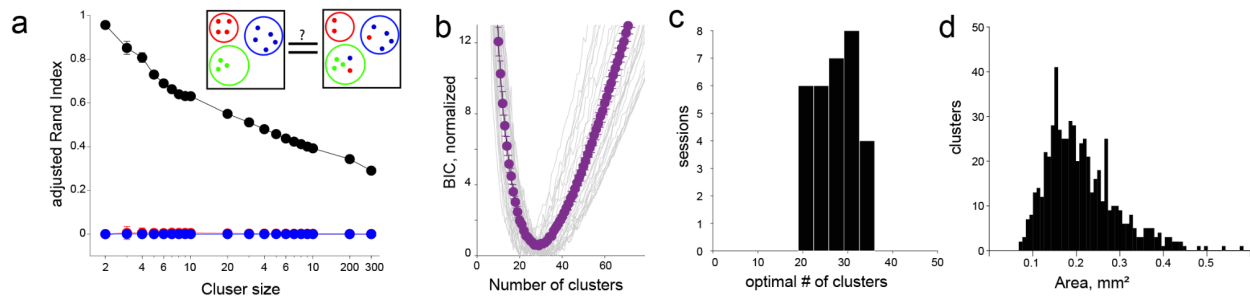
968

969 **Extended Data Fig. 3: Clustering patterns in all 8 mice.**

970 We provide the K-means cluster of each of the animals examined. Plotting format follows panels
971 in Fig 1.

972

Extended Data Figure 4

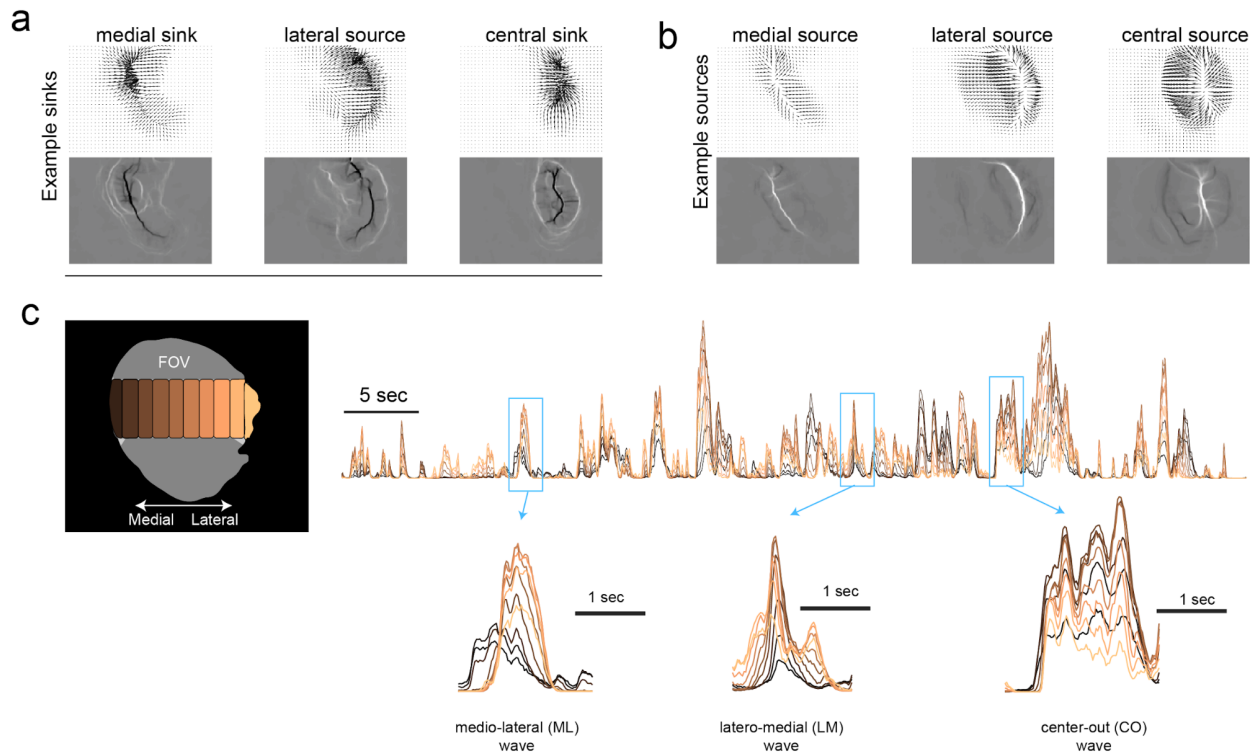


973

974 **Extended Data Fig. 4: Cluster patterns are robust.**

975 **a**, Adjusted rand-index score of cluster patterns determined using the K-means for re-clustering
976 (black) or shuffling the temporal (red) or spatial(blue) indices of pixels. Results are shown for
977 100 reclustering iteration with or without shuffling. Note that randomizing the temporal or
978 spatial relationships of fluorescence activity results in random clusters. **b**, BIC score for K-
979 means results all sessions examined (n=31 sessions, 8 mice; gray), and average (purple). **c**,
980 Distribution of optimal number of clusters identified using the BIC metric. **d**, Distribution of areas
981 of identified clusters across all mice.

Extended Data Figure 5



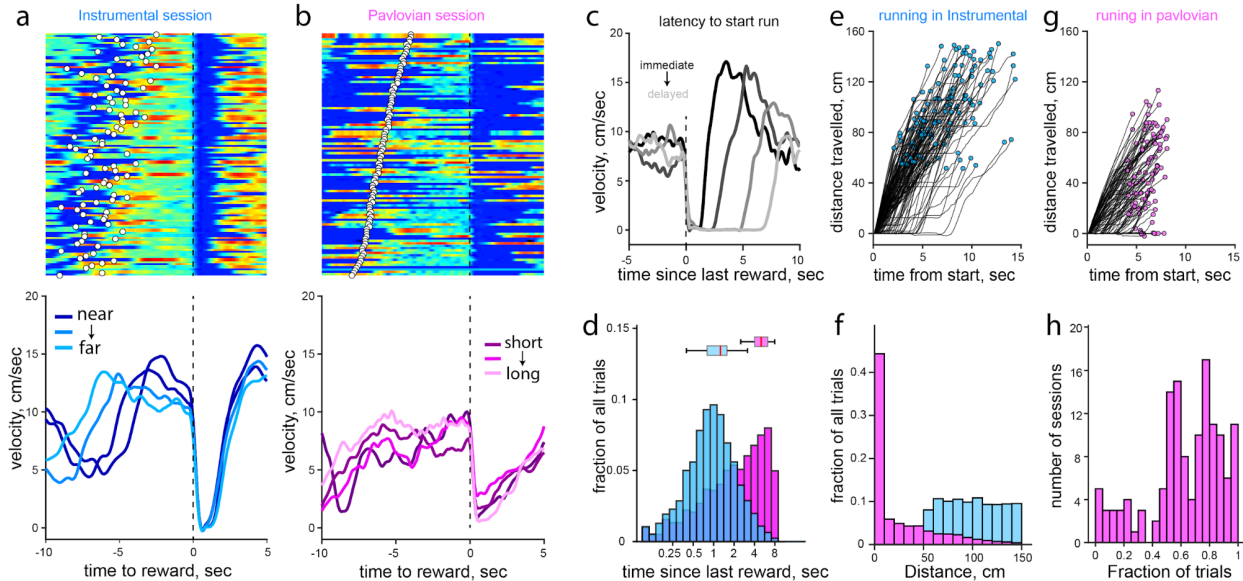
982

983 **Extended Data Fig. 5: Local sources and sinks initiate and terminate dopamine activity,**
984 **delivering temporally delayed dopamine to striatal subregions.**

985 **a**, Flow pattern (*top*) and divergence map (*bottom*) for sinks that are clustered in medial, lateral
986 or central striatal regions. **b**, same format as **a**, for source locations. **c**, Time Course of activity
987 across the mediolateral gradient for a one minute recording epoch. Blue boxes focus on
988 transient events that were produced by ML, LM, or CO waves that deliver dopamine to different
989 parts of the dorsal striatum with relative lags.

990

Extended Data Figure 6



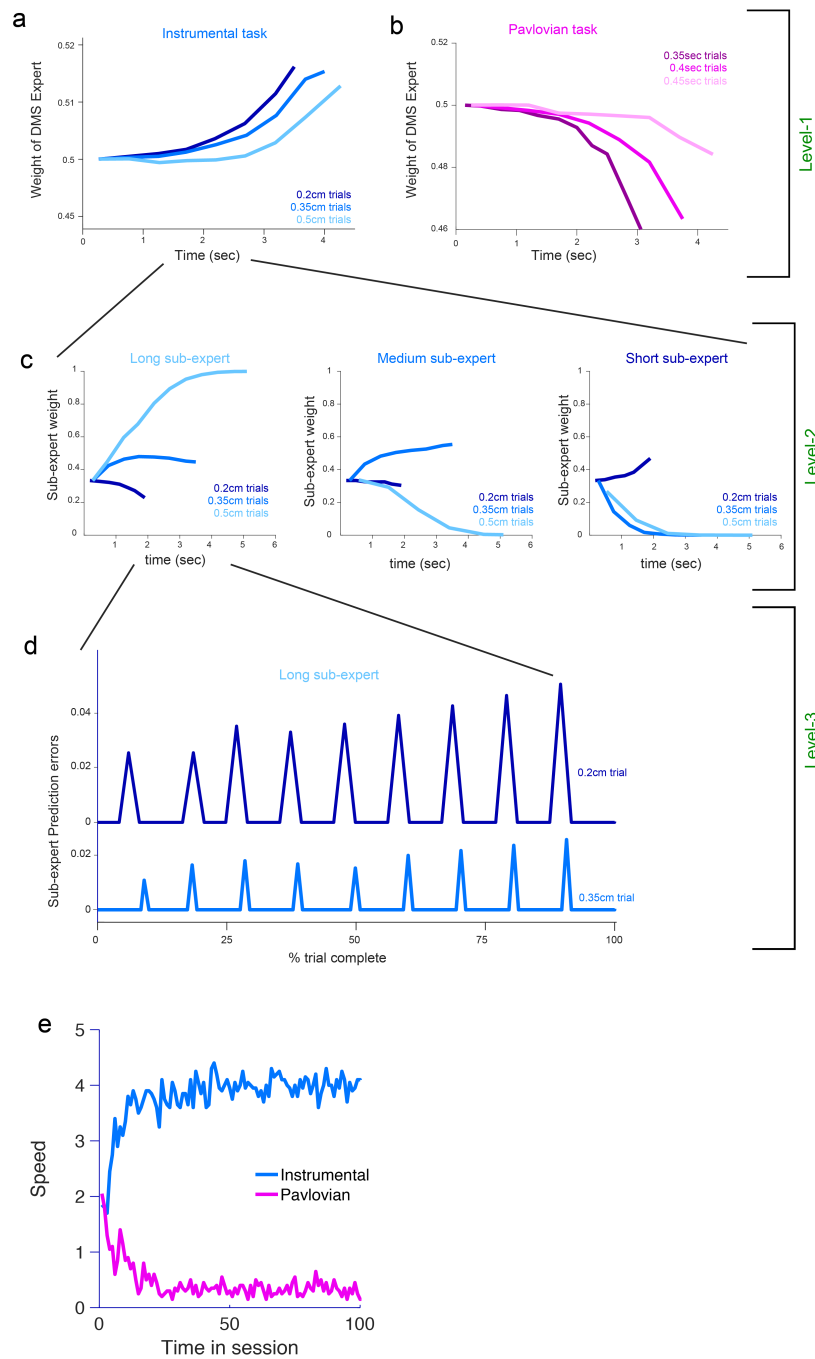
991

992

993 **Extended Data Fig. 6: The running behavior of mice is more structured and goal-directed**
 994 **in the instrumental task.**

995 **a**, Example velocity profile for an instrumental session. *Top* shows the trial by trial velocity
 996 aligned to the end of trial (reward receipt) White dots indicate the start of the trial. *Bottom*
 997 illustrates mean velocity trajectory for different distance contingencies. **b**, Same format as **a** but
 998 for pavlovian session. Note that the running behavior of the mouse is disorganized relative to
 999 task events (quantified in later sessions). **c**, Example session showing variability the latency to
 1000 start running on the next trial. **d**, We quantified this latency for training sessions across all mice
 1001 and observed a significantly shorter latency to initiate next trial running. X-axis is displayed in
 1002 log scale. **e**, Single trial trajectories of position from trial start during instrumental sessions, and
 1003 pavlovian sessions in **f**. Circles denote mouse position at the end of a trial. **g**, Overall, mice ran
 1004 less distance than the requirement in instrumental sessions (i.e. 50-150cm), and **h**, mice chose
 1005 not to run at all in a significant fraction of trials during the pavlovian task (note that running is
 1006 required for rewards in instrumental sessions).

Extended Data Figure 7



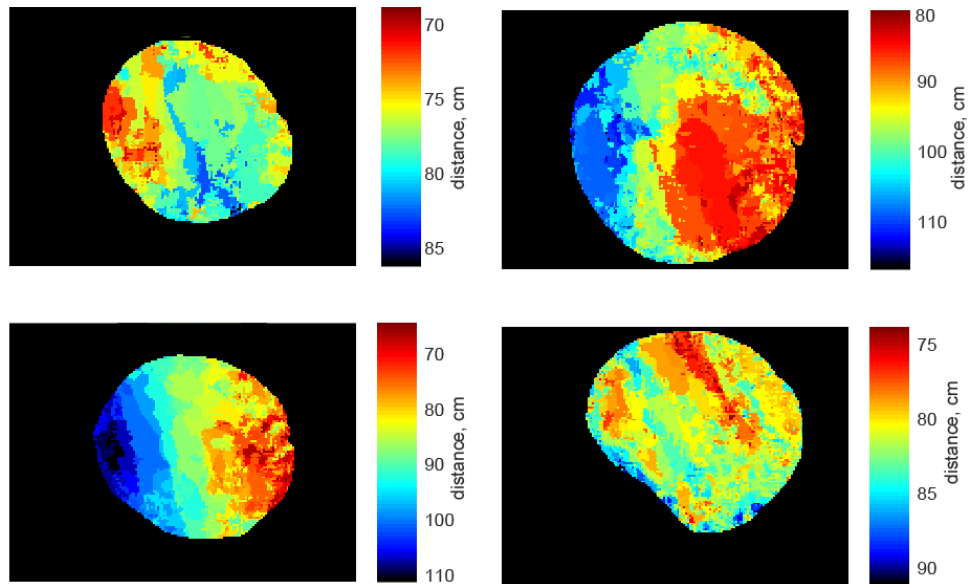
1007

1008 **Extended Data Fig. 7: Within-trial dynamics of model variables at all three levels under**
 1009 **different task conditions.**

1010 **a,b** Positive and negative accumulation of distance expert (level-1, equivalent to DMS) weights
 1011 under (a) instrumental and (b) Pavlovian task condition, for short, medium and long trial types.
 1012 Each trace is the average dynamics on the very first trial, averaged for 10 simulation sessions.
 1013 Similar dynamics accumulate across trials within a session when the task is repeated (not

1014 shown). Note that the ramp shape is convex in the first trial but concave for later trials. **c**, Within
1015 the distance expert, sub-experts (level-2) specialize on distinct contingencies and the weights
1016 ramp accordingly depending on task conditions. **d**, RPEs within a sub-expert in which tone
1017 transitions occur at unexpected times/distances (RPEs are zero for sub-experts that perfectly
1018 predict the current contingency; not shown). Note the larger magnitude RPEs for short
1019 compared to longer trials, as seen empirically (Fig. 5). Escalation of RPEs across the trial is due
1020 to temporal discounting. Similar to the empirical data, the impact of larger RPEs on short
1021 distances is more evident later in the trial. **e**) Model velocities (averaged across simulations)
1022 recapitulate increase in running in instrumental compared to pavlovian sessiona. The model
1023 selects speeds in proportion to inferred responsibility of the distance (“DMS”) expert, which
1024 accumulates across a session.

Extended Data Figure 8

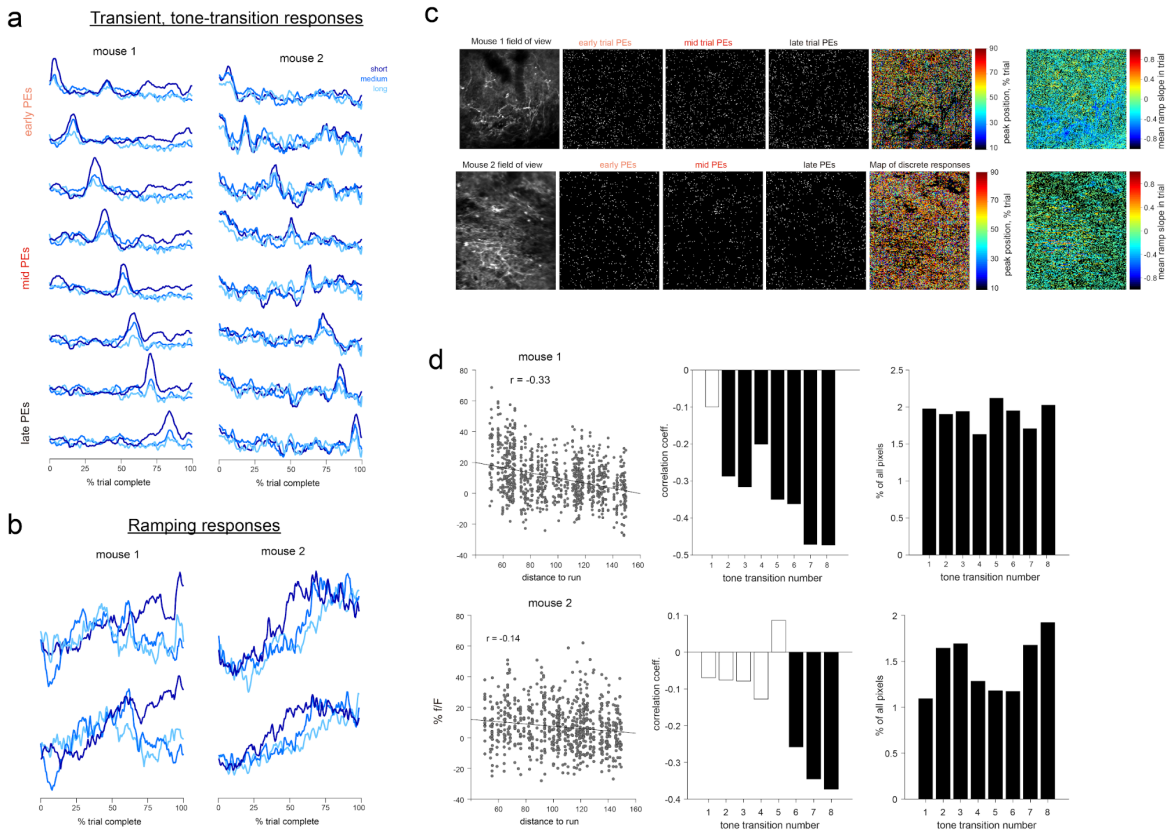


1025

1026 **Extended Data Fig. 8: Additional example session of striatal subregions that have**
1027 **preferred distance contingency.**

1028 The trial-by-trial ramp slope during anticipation epoch was distinctly modulated for different
1029 striatal subregions. Color maps are in same format as Fig. 4n.

Extended Data Figure 9



1030

1031 **Extended Data Fig. 9: Discrete and ramp-like responses in dopamine axon segments.**

1032 **a**, Discrete, tone responses in two mice. Format is as in Fig 5. **b**, Example ramping patterns in
 1033 the two animals. **c**, Examination of the anatomical distribution of pixels that exhibit tone-
 1034 transition tuning. *Top* row summarizes data from mouse-1, and *bottom* is for the second mouse.
 1035 Leftmost panels show the mean projection of field of view, and the next three panels show the
 1036 individual pixels that display discrete responses early (first transition), mid-trial (5th transition)
 1037 and late-trial (last-tone transition) responses. Moreover, the anatomical organization of these
 1038 pixels are intermingled. Fifth panel shows the anatomical position of all phasically responding
 1039 pixels color coded for which transition they respond to. Finally, Last panel on the right shows the
 1040 anatomical distribution and trial-averaged ramp slopes of pixels within the 2-photon field of view
 1041 that exhibit sustained upward or downward activity during the anticipatory epoch. Same format
 1042 for mouse-2 at *bottom* **d**, Quantification of how peak response at tone-transition is affected by
 1043 distance to needed to run on current trial. Our simulations predict (see Extended Data Fig. 7)
 1044 that shorter trials will elicit larger PEs. We found a significantly negative correlation overall in
 1045 both mice ($p < 0.001$, *left*), but the influence of distance was more prominent for later tones
 1046 (*middle*, filled bar are have $p < 0.05$) as in the model (Extended Data fig 7). Right panel shows
 1047 that similar fraction of pixels that were responsive to each tone transition.

- 1048 **Supplementary Video 1:** Example recording session demonstrating the activity pattern of
1049 dopamine axons in the dorsal striatum. Video playback is 1X.
- 1050 **Supplementary Video 2:** Video illustrating extraction of flow trajectories on a frame by frame
1051 basis. Video playback is slowed down 0.25X.
- 1052 **Supplementary Video 3:** Reward response in Instrumental task. Clock at top left displays time
1053 relative to reward. Video playback is 1X.
- 1054 **Supplementary Video 4:** Reward response in Pavlovian task. Clock at top left displays time
1055 relative to reward. Video playback is 1X.
- 1056 **Supplementary Video 5:** Progressively organized and continuous reward response to
1057 unpredicted reward deliver in naive mice (*top*), or animals that have received training in
1058 pavlovian sessions for 3 weeks. Clock at top left displays time relative to reward. Video
1059 playback is slowed down 0.5X.