

1 **Title**

2 Neural Basis of the Delayed Gratification

3 **Short Title**

4 Neural Basis of the Delayed Gratification

5 **Authors**

6 Zilong Gao^{1,2#}, Hanqing Wang^{3#}, Chen Lu⁴, Sean Froudish-Walsh³, Ming Chen⁴,

7 Xiao-jing Wang^{3*}, Ji Hu^{4,5*}, Wenzhi Sun^{1,6*}

8 **Affiliations**

9 ¹ Chinese Institute for Brain Research, Beijing 102206, China

10 ² Academy for Advanced Interdisciplinary Studies, Peking University, Beijing

11 100871, China

12 ³ Center for Neural Science, New York University, New York, NY 10003, USA.

13 ⁴ School of Life Science and Technology, ShanghaiTech University, Shanghai

14 201210, China

15 ⁵ Shanghai Key Laboratory of Psychotic Disorders, Shanghai Mental Health Center,

16 200030, China

17 ⁶ School of Basic Medical Sciences, Capital Medical University, Beijing 100069,

18 China

19 # These authors contributed equally to this work

20 * Corresponding authors. Email: xjwang@nyu.edu, huji@shanghaitech.edu.cn &

21 sunwenzhi@cibr.ac.cn

22

23 **Abstract**

24 Balancing instant gratification versus delayed, but better gratification is important for
25 optimizing survival and reproductive success. Although psychologists and neuroscientists
26 have long attempted to study delayed gratification through human psychological and brain
27 activity monitoring, and animal research, little is known about its neural basis. We successfully
28 trained mice to perform a waiting-and-water-reward delayed gratification task and used these
29 animals in physiological recording and optical manipulation of neuronal activity during
30 the task to explore its neural basis. Our results showed that the activity of DA neurons
31 in ventral tegmental area (VTA) increases steadily during the waiting period. Optical
32 activation vs. silencing of these neurons, respectively, extends or reduces the duration
33 of waiting. To interpret this data, we developed a reinforcement learning (RL) model
34 that reproduces our experimental observations. In this model, steady increases in
35 DAergic activity signal the value of waiting and support the hypothesis that delayed
36 gratification involves real-time deliberation.

37 **KEYWORDS**

38 Delayed Gratification, Dopaminergic, Ventral Tegmental Area, Ramping Activity,
39 Reinforcement Learning, Continuous Deliberation

40

41 **TEASER**

42 Sustained ramping dopaminergic activation helps individuals to resist impulsivity and
43 wait for laerger but later return.

44 **INTRODUCTION**

45 To optimize survival and reproductive success, animals need to balance instant
46 gratification versus delayed, but better gratification. Repeated exposure to instant
47 gratification may disrupt this balance, thereby increasing impulsive decisions. Such
48 decisions contribute to numerous human disorders, such as addiction and obesity(1, 2).
49 Delayed gratification is an important process that balances time delay with increased
50 reward (3). It is influenced by strengths in patience, will-power, and self-control(4).
51 Although psychologists and neuroscientists have long studied this important behavior through
52 human psychological and brain activity assessments as well as rodent-based studies, little is
53 known about its neurological basis.

54 During a well-controlled delayed gratification task, an individual must balance the
55 benefits vs. risks of delay in receiving an available reward. Sustaining choice requires
56 suppression of constant temptation by the expectation of enhanced reward in the future
57 (3, 5, 6). Midbrain dopaminergic neurons are well known to play central roles in
58 reward-related and goal-directed behaviors (7-12). Studies have revealed that DAergic
59 activity signals proximity to distant rewards, either spatially or operationally (7, 13, 14),
60 which has been postulated to sustain or motivate goal-directed behaviors while resisting

61 distractions. DAergic neurons play important roles in time judgment (15) and cost-
62 benefit calculations which are necessary for value-based decision making (13, 16-18).

63 We successfully trained mice to perform a waiting-and-water-reward delayed
64 gratification task. Recording and manipulation of neuronal activities during this task
65 allows us to explore the cellular regulation of delayed gratification. We found that the
66 activity of VTA DAergic neurons ramped up consistently while mice were waiting in
67 place for rewards. Transient activation of DAergic neurons extended and inhibition
68 reduced the duration of the waiting period. Then we adopted reinforcement learning
69 (RL) computational models to predict and explain our experimental observations.

70 **RESULTS**

71 **Mice can learn to wait for greater rewards by delayed gratification task training**

72 First, we trained mice to perform a one-arm foraging task (pre-training) in which
73 delay did not result in increased reward (19). The period the mouse in the waiting
74 zone was set as waiting duration and the time a mouse used in running from the
75 waiting zone to fetch the water reward was as running duration (Fig. 1A, left panel).
76 When a water-restricted mouse exited the waiting zone and licked the water port in
77 the reward zone, it could receive a 10 μ l of water drop regardless of the time spent in
78 the waiting zone (Fig. 1A, right panel, black line). In a week of training, the average
79 waiting and running durations both significantly decreased from Day 1 to Day 7 ($p <$
80 0.001 , $n = 7$ mice, Figs. 1C-E, Movie S1). All mice learned the strategy of reducing

81 durations of both waiting and running to maximize the reward rate (as ul of water per
82 second in a trial, fig. S1A).

83 Next, we trained the same mice using a delayed gratification paradigm, where the
84 size of the reward increased quadratically with time spent in the waiting zone (Fig.
85 1A, right panel, green line). Over the next three weeks, this resulted in shifting of the
86 distributions of waiting duration towards longer time durations. The averaged waiting
87 period significantly increased from Day 1 to Day 15 ($p < 0.001$, Figs. 1F, H, and Movie
88 S1), whereas the duration of running did not decrease beyond that observed initially
89 ($p = 0.97$, $n = 7$ mice, Fig.s 1G, H). The reward rate increased steadily, indicating that
90 the mice were learning to successfully delay gratification (fig. S1B).

91 **The activity of VTA DAergic neurons increases steadily during the waiting** 92 **period**

93 To monitor the activity of VTA DAergic neurons during the delayed gratification
94 task, we employed fiber photometry to record the calcium signals in VTA DAergic
95 neurons in freely moving mice for as long as one month (Fig. 2A-C, optical fiber
96 placement illustrated in fig S2). On the first day of pre-task training, the calcium signal
97 rose rapidly on reward and quickly reached a peak. A few days of training dramatically
98 reshaped the response pattern. Once the mice re-entered the waiting zone, the activities
99 of DAergic neurons started to rise and reached the highest level when the animal
100 received a reward (fig. 3SA).

101 We next analyzed the activity of these same neurons in the mice as they learned
102 the delayed gratification task. The recording traces showed that training gradually
103 reshaped the pattern and time course of activity (Fig. 2D). The activity started to ramp
104 up once the mice entered the waiting zone, and then reached its highest level when
105 animals exited. To investigate carefully the dynamical properties of the ramping
106 activity during waiting, we sorted the calcium signals from one training day of one
107 mouse by their length of waiting durations and plotted them with a heat map (Fig. 2E).
108 We divided trials according to the trial outcome (reward volume) and calculated the
109 calcium signals while the mouse exited waiting zone with different reward volumes.
110 Our results showed that the Z-scored calcium signals at 0.5 sec before exit were
111 significantly different while the reward volumes were different (Fig. 2F), but the mean
112 signal curves raised along a similar trajectory regardless of trial outcome (Fig. 2G).
113 Then, we calculated the slopes of signal curves with different outcomes over 4 time
114 windows (0~2, 2~4, 4~6, 6~8 s) by linear regression analysis. The slopes during the
115 same time window had no significant differences between reward groups (Fig. 2H). We
116 pooled and plotted the slopes of different waiting periods together and found the activity
117 curves kept rising steadily and almost saturated after 6 secs from the time the mice
118 entered the waiting zone. Besides, the ramping DAergic activity became less variable
119 along with delayed gratification task training in our experimental data (figs. S4A-D).
120 All these results indicated the VTA DAergic neurons consistently ramp up their activity
121 during waiting in as animals are trained in the delayed gratification task.

122 **Optogenetic manipulation of VTA DAergic activity altered the waiting durations**
123 **in delayed gratification task**

124 To determine whether VTA DAergic activity controls performance in the delayed
125 gratification task, we manipulated VTA DAergic neurons temporally within 20%
126 pseudo-randomly chosen trials utilizing optogenetic tools while the mice were waiting
127 during the delayed gratification task (Figs. 3A-C). Activating the VTA DAergic
128 neurons shifted the cumulative probability distribution to statistically significant longer
129 waiting duration (Fig. 3D, blue), while inhibiting these same neurons shifted this
130 distribution to significantly shorter periods(Fig. 3E, yellow). The impacts on the
131 cumulative probability duration distributions were only observed in the Laser-ON trials.
132 In contrast, the Laser-OFF trials, including the next trials after the Laser-ON as a single
133 group, were not significantly different from the trials from the previous day (Figs. 3D-
134 E). The optical manipulations didn't influence the running durations in ChR2 or eNpHR
135 3.0 expressing mice (figs. S5A-B), nor did it change the waiting duration distribution
136 of mice that expressed mCherry in DAergic neuron in delayed gratification task (figs.
137 S6A-B). To rule out the possibility of optogenetic manipulation-induced memory, we
138 performed a random place preference test with the same stimulation dosage. Nither
139 activating nor inhibiting VTA DAergic neurons significantly changed the transient
140 waiting duration and pattern in the location in which the laser was activated in all tested
141 mice (figs. S5C-F) as well as the mCherry expressing controls (figs. S6C-F).

142 **A reinforcement learning (RL) model suggests that ramping VTA DAergic**
143 **activity signals the value of waiting for delayed gratification**

144 How does a mouse manage to wait longer for a larger reward vs. smaller but more
145 immediate reward options? We propose two models to explain behavioral scenarios to
146 exemplify possible strategies a mouse may implement to achieve extended waiting
147 performances: setting a goal of expected waiting duration before initiation of waiting
148 or continuously deliberating during the waiting period. According to the first hypothesis,
149 we modeled a RL agent that keeps timing until the preset moment has passed(Fig. 4A,
150 *Decision Ahead*); In the second, we modeled a second RL agent that continuously
151 balances the values of waiting versus leaving to control the decision on waiting versus
152 leaving for the reward. Practically, we used a version of the state-action-reward-state-
153 action (SARSA) algorithm with a series of stares (20, 21)(Fig. 4A, *Continuous*
154 *Deliberation*, see methods). The behaviors of both models were able to replicate the
155 behavioral performance we observed in animal experiments (Figs. 4B-C). There was
156 no significant difference between the Kullback-Leibler (KL) divergence we chose to
157 quantitatively assess the divergence from the waiting distribution of simulated behavior
158 to that of the animals for both RL models. We couldn't determine which model is better
159 based on behavioral performance alone, given that both models well reproduced the
160 behavioral data (Fig. 4D).

161 What does the ramping DAergic activity mean in the delayed gratification task?
162 We tried to explain it with our RL model. In the *Decision Ahead* model, the agent keeps
163 timing until the preset moment has passed, which suggests the ramping DAergic
164 activity may relate to timing in delayed gratification task. Some studies have proposed
165 that the ramping activity is consistent with a role in the classical model of timing with

166 the movement initiated when the ramping activity reaches a nearly fixed threshold value,
167 following an adjusted slope of ramping activity(22-25). In contrast, our results showed
168 that the DAergic activity ramped up to different values with similar trajectories on a
169 nearly constant slope (Figs. 2F-H). This suggests that VTA DAergic neurons may not
170 implement a decision variable for the *Decision Ahead* scenario. In the *Continuous*
171 *Deliberation* RL model, we compared the curves of the value of waiting and leaving
172 with the ramping DAergic activity and found the behavioral performance of both
173 animals and model agents reached the asymptote. The values of waiting (Fig. 4E, light
174 purple) and the leaving (Fig. 4E, green) each correlated positively with the ramp of
175 DAergic activity during waiting (Fig. 4E, green, blue, Z-scored $\Delta F/F$, 0.5 sec before
176 exit from the last week of training). This detailed analysis suggested that the *Continuous*
177 *Deliberation* RL model agreed with previous studies (13, 26-28) and that ramping
178 DAergic activity signals the value of actions, either waiting or leaving, in the delayed
179 gratification task.

180 In the *Decision Ahead* RL model, if the agent keeps timing during waiting through
181 ramping DAergic activity to encode the elapse of time (29-31), extra VTA DAergic
182 activation should represent a longer time thus lead to an earlier stop of waiting. This
183 deduction is contrary to our optogenetics result, namely DAergic activation led to a
184 longer waiting (Fig. 3D). Instead, we reproduced the optogenetic manipulations in the
185 *Continuous Deliberation* RL model by either increasing or decreasing the value of
186 waiting (Q_{wait}) in pseudo-random 20% trials. The increase or decrease in waiting
187 durations only occurred in the Q_{wait} – manipulated trials, whereas the remaining trials,

188 including the next trials after value manipulation, had no significant difference with
189 control (Figs. 4F-G). Importantly, manipulating the value of leaving (Q_{leave}) in pseudo-
190 random 20% trials induced the opposite results (figs. S8A-B) compared with
191 experimental data of optogenetic manipulation. Our experimental data and *Continuous*
192 *Deliberation* RL model together indicated that the ramping VTA DAergic activity
193 profoundly influenced the waiting behavioral performance in the delayed gratification
194 task, which suggested ramping DAergic activity signal the value of waiting, rather than
195 the value of leaving. Our analysis conceptually revealed that the delayed gratification
196 involved real-time deliberation.

197 **VTA DAergic activity during waiting predicts the behavioral performance in the** 198 **delayed gratification task**

199 Our optogenetic manipulation experiments and *Continuous Deliberation* RL
200 model indicated that VTA DAergic activity during waiting influenced the waiting
201 durations while the mouse was performing the delayed gratification task (Figs. 3D-E
202 and Figs. 4F-G). Although the activity of VTA DAergic neurons ramped up
203 consistently during waiting (Figs. 2G-H), they still fluctuated to a certain extent
204 moment by moment. Therefore, we are next to determine whether this fluctuation
205 influences the waiting behavior in the delayed gratification task. A strong prediction
206 given by the *Continuous Deliberation* model is that, if DAergic activity signals the
207 value of waiting at each specific moment, the more likely the agent will keep waiting
208 in the next "time bin", but not in the later ones (fig. S8E-F, the value of waiting is only
209 positively correlated with the behavior of next time bin), which agrees with the Markov

210 property(32). We thus aimed to test the relationship between the amplitude of
211 momentary VTA DAergic signal and the behavior (i.e., waiting or leaving) within each
212 time bin to determine how the momentary DAergic activity (the calcium signal
213 amplitude in 0~1 sec, 1~2 sec, 2~3 sec, or 3~4 sec after waiting onset, shown as each
214 cluster of bars in Fig. 5B) affects the waiting performance in the subsequent periods
215 (behavior within 1~2 sec, 2~3 sec, 3~4 sec, and 4~5 sec for DA in 0~1 sec, behavior
216 within 2~3 sec, 3~4 sec and 4~5 sec for DA in 1~2 sec, and so on, Fig. 5A). To integrate
217 data from multiple sessions as well as multiple animals, we took the advantage of the
218 linear mixed model analysis (LMM, or linear mixed-effects, LME, see method) (33-
219 35). The regression coefficients between momentary DAergic activities and momentary
220 waiting (1 for waiting and 0 for not) were significantly positive between adjacent
221 DAergic and behavioral periods (Fig. 5B, 1sec for the bars on the most left of each
222 cluster/adjacent DAergic-behavior). The pair of momentary DAergic activity in 3~4
223 sec and waiting in 4~5 sec didn't show a significant correlation ($p = 0.61$, $n = 7$), which
224 may possibly result from insufficient data for those long trials. This result indicates that
225 the waiting decision of the current moment is only influenced by the most recent
226 DAergic signal but not by DAergic signal further in the past, which suggests that
227 deliberation for waiting in delayed gratification may be a Markov process as we
228 formalized in the *Continuous Deliberation* RL model(32).

229 In the *Continuous Deliberation* RL model, the probability of waiting (P_w)
230 positively correlates with the value of waiting (Q_{wait}). To explore the impact of DAergic
231 activity on the probability of waiting in our experimental data, we binned DAergic

232 activity of every trial and normalized data points (V_{DA}) in each momentary DAergic
233 period (1 sec started from 0 to 9 s). Then, we divided the trials into two groups by
234 setting a series of arbitrary thresholds (red, High DAergic activity, $V_{DA-Z} \geq Th$; green,
235 Low DAergic activity, $V_{DA-Z} \leq -Th$, Th was the threshold for the analysis of high/low
236 DAergic activity) from these trials (Fig. 5C, Th was set to 0.9). By analyzing the
237 probability of waiting of low and high DAergic activity at different thresholds for the
238 adjacent waiting period, we found that the probability of waiting increased rapidly as
239 the absolute value of threshold was set larger and larger. The probability of waiting was
240 significantly different between the high and low dopamine trials when the absolute
241 value of Th ($|Th|$) was greater than 2.0 (Fig. 5D).

242 Finally, we investigated the influence of fluctuations of intrinsic VTA DAergic
243 activity on the waiting performance of mice in the delayed gratification task. There
244 were trials whose DAergic activity in the whole waiting duration was significantly
245 higher (red, high-ramping) or lower (green, low-ramping) than the mean DAergic
246 activity (see Methods). Then we separated two groups of trials and found that the
247 cumulative distribution of waiting durations of high-ramping trials shifted to the right
248 with significantly higher normalized waiting durations compared with the normalized
249 waiting durations of low-ramping trials (Fig. 4E), but there was no difference between
250 the normalized waiting durations for the next-in series trials of high-ramping and low-
251 ramping trials (Fig. 4F). These results accorded with our optogenetic manipulation
252 experiment (Figs. 3D-E) that optogenetically manipulated VTA DAergic activity

253 transiently influences the behavioral performance of waiting in delayed gratification
254 task.

255 **DISCUSSION**

256 Here we reported a novel behavioral task to train the mice to learn a foraging task
257 with a delayed gratification paradigm. Mice learned to wait for bigger rewards with the
258 increase of waiting durations (Figs. 1F-H). Moreover, the calcium signal of VTA
259 DAergic neurons ramped up consistently when the mouse waited in place before taking
260 action to fetch an expected reward (Figs. 2G-H). Further data analysis showed that the
261 ramping VTA DA activity indeed influenced the behavioral performance of waiting
262 (Figs. 5B-E), which was confirmed with bi-directional optogenetic manipulations of
263 VTA DAergic activity (Figs. 3D-E). At last, a RL model well predicted our
264 experimental observations and consolidated the conclusion that the ramping VTA
265 dopaminergic activity signaled the value of waiting in the delayed gratification task,
266 which involves real-time deliberation (Figs. 4B-G)..

267 DA release in NAc was previously conjectured for sustaining or motivating the
268 goal-directed behavior as well as resisting distractions (13, 14). Here, we explicitly
269 implemented continuous 'distractions' or less-optimal options along the delayed
270 gratification process, in which, to achieve better performance, the mice need to sustain
271 waiting as well as prevent/control impulsivity (3, 6, 36, 37). We found remarkable and
272 sustaining DAergic activation when mice managed to wait longer, and further
273 demonstrated a causal link between DAergic activation and the increase in transient

274 waiting probability. Furthermore, we found DAergic activity ramps up in a consistent
275 manner during waiting, mimicking the value of waiting along with a series of states in
276 our *Continuous Deliberation* RL model, both of which are presumably resulted from
277 and contributed to resisting an increasing magnitude of distraction in our task.
278 Intriguingly, the momentary DAergic activity was found positively correlated to the
279 momentary waiting probability, which also suggested DAergic activity may be
280 involved in the continuous deliberation process. Therefore, we not only for the first
281 time to our knowledge demonstrated the behavioral significance of DAergic activity in
282 delayed gratification, but also depicted a "Continuous Deliberation" framework where
283 DAergic activity may participate and help achieve more flexible and sophisticated
284 performance.

285 Numerous works use Pavlovian conditioning in studying DA activity(10, 12, 38-
286 40). Some studies paired the reward with a cue (or cues), in which animals don't need
287 effortful work to obtain rewards. It is well known this kind of DA activity signals the
288 RPE via phasic firing. In the studies using operant conditioning or goal-directed
289 behavior, the animals have to perform actions and need effortful work to obtain
290 outcomes, and a ramping DA activity was reported to emerge while the animals were
291 approaching the reward (13, 14, 41, 42). The ramping activity is suggested to signal the
292 value of work (13) or distant rewards (14), but key evidence is lacking because the
293 change of sensory input flow remarkably alters the DA activity over time. Under such
294 mutual influence, it is impossible to identify RPEs or the value of work from external
295 cues. The RPE model of ramping activity assumes that the value increases

296 exponentially (or at least in a convex curve) as the reward is approached. Under this
297 model, sensory feedback is suggested to result in the RPE signal to ramp (41, 43, 44),
298 while a lack of sensory feedback is predicted to make a flat RPE signal. In contrast, the
299 ramping DAergic activity is well isolated from the external sensory inputs when
300 performing a delayed gratification task in our model. The mice continuously deliberate
301 the current state and future rewards without any external sensory inputs during waiting
302 in place. We still observed the calcium signal of VTA DAergic neurons ramped up in
303 a stable dynamic. This ramp may indicate an escalating value for the closer reward in
304 temporal and represent the 'willpower' of waiting.

305 Midbrain DAergic neurons play an important role in reinforcement learning(9, 11,
306 12, 45, 46), where activation of DAergic neurons usually produces a reinforcement
307 effect on associated action, stimulus, or place. But in our delayed gratification task,
308 optogenetic manipulation of DAergic activity substantially influenced the ongoing
309 behavior on the current trial without visible reinforcement effect on later trials. Notably,
310 this optogenetic manipulation was not sufficient to induce a reinforcement effect in the
311 random place performance test. These results revealed the distinct and potent
312 instantaneous effect of DAergic activity during delayed gratification. The observations
313 and analysis in our experiments integrate more reliable evidence for the value coding
314 in VTA DAergic neurons and significantly update the understanding of the coding
315 mechanisms and fundamental functions of the DAergic system in delayed gratification.
316 Our design of the delayed gratification task recapitulates the realistic situation where
317 **distractions and less-valuable choices lie in the way of pursuing a larger but later benefit.**

318 The deficit of resisting distractions (temptations), which disrupt the balance between
319 constant reward and delayed reward, is closely related to a variety of disorders like
320 obesity, gambling, or addiction(1, 47). The ramping VTA DAergic activity accords
321 with the model about NOW vs LATER decisions that tonic/stable DAergic signal have
322 a strong influence on dlPFC and favor LATER rewards(2). We proposed that the
323 sustained VTA DAergic activity during the delayed period could serve as a
324 conservative neural basis for the power to resist the ubiquitous distractions (temptations)
325 and improve reward rate or goal pursuit in the long run.

326

327 **MATERIALS AND METHODS**

328 **Mice.** Animal care and use strictly follow institutional guidelines and governmental
329 regulations. All experimental procedures were approved by the IACUC at the Chinese
330 Institute for Brain Research (Beijing) and ShanghaiTech University. Adult (8-10
331 weeks) DAT-IRES-Cre knock-in mice (Jax stock# 006660) were trained and
332 recorded. Mice were housed under a reversed 12/12 day/night cycle at 22–25°C with
333 free access to ad libitum rodent food.

334 **Stereotaxic viral injection and optical fiber implantation.** After deep anesthesia
335 with isoflurane in oxygen, mice were placed on the stereoscopic positioning
336 instrument. Anesthesia remains constant at 1~1.5% isoflurane supplied per anesthesia
337 nosepiece. The eyes were coated with aureomycin eye cream. The scalp was cut open,
338 and the fascia on the skull was removed with 3% hydrogen peroxide in saline. The
339 Bregma and Lambda points are used to level the mouse head. A small window of
340 300~500µm in diameter was drilled just above VTA (AP: -3.10 mm, ML: ±1.15mm,
341 and DV: -4.20 mm) for viral injection and fiber implantation. 300 nl of AAV2/9-
342 hSyn-DIO-GCamp6m (10^{12}) solution was slowly injected at 30nl /min unilateral for
343 fiber photometry recording. 300 nl either AAV2/9-EF1a-DIO-hChR2(H134R)-
344 mCherry (10^{12}) or AAV2/9-EF1a-DIO-eNpHR3.0-mCherry (10^{12}) was injected
345 bilaterally for optogenetic experiments. The injection glass pipette was tilted with an
346 angle of 8° laterally to avoid the central sinus. After injection, the glass pipette was
347 kept in place for 10 min and then slowly withdraw. An optical fiber (200 µm O.D.,
348 0.37 NA; Anilab) hold in a ceramic ferrule was slowly inserted into the brain tissue

349 with the tip slightly above the viral injection sites. The fiber was sealed to the skull
350 with dental cement. Mice were transferred on a warm blanket for recovery, then
351 housed individually in a new home until all experiments were done.

352 **Behavioral tasks.** One week after surgery, mice started a water restriction schedule to
353 maintain 85–90% of free-drinking bodyweight for 5 days. The experimenter petted
354 the mice 5 minutes per day for 3 days in a row and then started task training. All
355 behavioral tasks were conducted during the dark cycle of mice.

356 The foraging task shuttle box has two chambers (10×10×15 cm) connected by a
357 narrow corridor (45×5×15 cm, Fig 1a). A water port (1.2 mm O.D. steel tube, 3 cm
358 above the floor) is attached to the end of one chamber defined as the reward zone, the
359 other as the waiting zone. The position of the mouse in the shuttle box is tracked online
360 with a custom MATLAB (2016b, MathWorks) program through an overhead camera
361 (XiangHaoDa, XHD-890B). The experimental procedure control and behavioral event
362 acquisition are implemented with a custom MATLAB program and an IC board
363 (Arduino UNO R3).

364 **One-arm foraging task (pre-training):** A water-restricted mouse was put in the
365 shuttle box for free exploration for up to 1 hour. When the animal started from the
366 waiting zone through the corridor to the reward zone to lick the water port, 10 μ l water
367 was delivered by a step motor in 100 ms as a reward. A capacitor sensor monitors the
368 timing and duration of licking. Then the animal returned to the waiting zone to re-
369 initiate a new trial. Exiting from the waiting zone triggered an auditory cue (200 ms at

370 4 kHz sine wave with 90 dB) to signal the exit from the waiting zone. The time in the
371 waiting zone was defined as the waiting duration. The training was conducted every
372 day in a week. All mice learned to move quickly back and forth between two chambers
373 to maximize the reward rate within one week.

374 **Delayed gratification task.** From the second week, the volume of water reward is
375 changed to a function proportional to the waiting time: 0~2 s for 0 μ l; 2~4 s, 2 μ l; 4~6
376 s, 6 μ l; 6~8 s, 18 μ l; >8 s, 30 μ l as shown in Fig. 1A. There is no water delivered if the
377 animal waits less than 2 sec. The training was conducted five days a week, from
378 Monday to Friday.

379 **P_w Calculation.** We divided all trials into two groups: Waiting Trials and Leaving
380 Trials according to whether an animal to keep waiting or to leave in a given time
381 duration such as 1 sec after each behavioral period. And then, we calculated the
382 probability of waiting (P_w) in this given time duration by the number of 'Waiting Trials'
383 ($N_{w(n)}$) and the number of 'Leaving Trials' ($N_{L(n)}$) in the time window n:

384
$$P_{w(n)} = \frac{N_{w(n)}}{N_{w(n)} + N_{L(n)}}$$

385 Then we can get the P_w for a given time duration:

386
$$P_{w(n)} = \frac{\sum_0^9 N_{w(n)}}{\sum_0^9 N_{w(n)} + \sum_0^9 N_{L(n)}}$$

387 **Linear Mixed Model.** We implemented the Linear Mixed Model Analysis using the
388 open-source Python package "statsmodels" (https://www.statsmodels.org/stable/mixed_linear.html). The binary value of waiting or leaving during a specific behavioral
389 period t_{beh} was set as the dependent factor ($t_{beh}=[1, 2), [2, 3), [3, 4),$ or $[4, 5)$, unit:

391 second); the fluctuation of momentary DA signal from its mean during a preceding
392 period t_{DA} was set as a fixed effect ($t_{DA}=[0, 1), [1, 2), [2, 3), [3, 4)$, unit: second. Note
393 that t_{DA} is always smaller than t_{beh}); the animal identity and session numbers were set
394 as a random effect ($n=5$ for each animal from the third week). The parameters of the
395 model are estimated by restricted maximum likelihood estimation (REML).

396 **Optogenetic stimulation.** Lasers, 473 nm for activation and 589 nm for inhibition,
397 were coupled to the common end of a patchcord (200 μm O.D., 1-m long, 0.37 NA).
398 The patchcord split through an integrated rotatory joint into two ends connecting to
399 chronically implanted optical fibers (200 μm O.D., 0.37 NA) for bilateral light delivery.
400 First, the mice were trained for 3 weeks to learn the delayed gratification task. Optical
401 stimulation was delivered pseudo-randomly in $\sim 20\%$ of behavioral trials in the test
402 experiment. 20 ms square pulses at 10 Hz for activation or a continuous stimulation for
403 inhibition were delivered. The laser was set to ON when the animal entered the reward
404 zone and to OFF on the exit. The maximal laser stimulation was no longer than 16
405 seconds, even in the case a mouse stayed in the waiting zone longer than this time.
406 Continuous laser power at the tip of splitting patchcord was about 10 mW for 473 nm
407 laser and 8 mW for 589 nm laser, respectively.

408 **Random place performance test (RPPT).** After finishing optogenetic tests for
409 delayed gratification, all mice took an RPPT. RPPT is carried on in a rectangular
410 apparatus consisted of two chambers (30 \times 30 \times 30 cm) separated by an acrylic board.
411 With an 8 cm wide door open, the mice could move freely between the two chambers.
412 Before testing, each mouse was placed into the apparatus for 5-min free exploration.

413 RPPT consists of two rounds of 10-min tests. First, we randomly assigned one chamber
414 as a test chamber. Laser pulses were delivered in 20% possibility while the mouse
415 entered the test chamber. The delivery of light, no longer than 16 sec, stopped while the
416 mouse exited the test chamber. Next, we switched the chamber to deliver laser pulses.
417 The laser output power and pulse length were set the same as optogenetic manipulations
418 in the delayed gratification task.

419 **Fiber photometry recording.** During the behavioral task training and test, we
420 recorded the fluorescence signal of VTA dopaminergic (DAergic) neurons. The signal
421 was acquired with a fiber photometry system equipped with a 488 nm excitation laser
422 and a 505~544 nm emission filter. The GCaMP6m signal was focused on a
423 photomultiplier tube (R3896 & C7319, Hamamatsu) and then digitalized at 1 kHz and
424 recorded with a 1401 digitizer and Spike2 software (CED, Cambridge, UK). An
425 optical fiber (200 μ m O.D., 0.37 NA, 1.5-m long, Thorlabs) was used to transfer the
426 excitation and emission light between recording and brain tissue. The laser power
427 output at the fiber tip was adjusted to 5~10 μ W to minimize bleaching.

428 All data were analyzed with custom programs written in MATLAB (MathWorks).
429 First, we sorted the continuously recorded data by behavioral trials. For each trial, the
430 data spanned the range between 1 s before the waiting onset and 2 s after the reward.
431 Before hooking the fiber to the mouse, we recorded 20 s of data and averaged as F_b as
432 the ground baseline. For each trial, we averaged 1-sec data before the waiting onset as
433 baseline F_0 and then calculated its calcium transient as:

434
$$\Delta F/F (\%) = (F - F_0)/(F_0 - F_b) \times 100 (\%)$$

435 In the correlation analysis between VTA DAergic activity before waiting and
436 waiting duration of mice, we used averaged 1-sec data before the waiting onset as the
437 DAergic activity before waiting.

438 In the analysis of high-ramping and low-ramping DAergic activity, we compared
439 the whole calcium signal of every trial with the average curve (the same length as the
440 analyzed calcium signal) of all trials from one mouse in a single training day with paired
441 t-test.

442 To facilitate presenting the data, we divided each trial data into four segments,
443 including 1 s before waiting onset, waiting, running, and 2 s after rewarding. For
444 comparing the rising trends, we resampled the data segments at 100, 100, 50, and 100
445 data points, respectively. In the delayed gratification task, the trial data were aligned to
446 the waiting onset and presented by the mean plots with a shadow area indicating SEM
447 of fluctuations.

448 **Reinforcement learning model.** We investigate two potential scenarios. One was
449 that the mouse decided on a waiting duration before entering the waiting area, and
450 then waits according to the decided goal. The other scenario was that the mouse
451 entered the waiting zone, and determined whether to wait or leave as an ongoing
452 process throughout the whole waiting period. We called these two scenarios
453 "*Decision Ahead*" and "*Continuous Deliberation*", respectively, and formulated

454 corresponding reinforcement learning based models for simulation using Python
455 (Python Software Foundation, version 2.7. Available at <https://www.python.org/>).

456 **Decision Ahead.** Inspired by animal behavior, we simply set three optional "actions"
457 with different expected waiting durations that could empirically cover the main range
458 of animal's waiting duration across training ($T_{a1} = 1.65$ sec for action1, $T_{a2} = 2.72$ sec
459 for action2, $T_{a3} = 4.48$ sec for action3). These waiting durations were equally spaced on
460 the log-time axis, consistent with Weber's law (that is, $\ln(T_{a1}) = 0.5$, $\ln(T_{a2}) = 1$, $\ln(T_{a3})$
461 $= 1.5$). During the execution of action a_i , we imposed additional noise to the timing so
462 that the actual waiting time τ_{ai} for action a_i follows a Gaussian distribution on the log-
463 time axis centered at the T_{ai} , $\ln \sim \mathcal{N}(\ln(T_{ai}), 0.4^2)$, $i = 1, 2, 3$. These settings allowed
464 us to best capture the animal's waiting performance in the model. For each trial, the
465 agent chose action randomly based on three action values and a Boltzmann distribution
466 (Softmax):

467
$$P_{a_i} = \frac{e^{\beta Q_{a_i}}}{\sum_{j=1,2,3} e^{\beta Q_{a_j}}}$$

468 Where P_{a_i} was the probability of choosing action a_i and waiting for τ_{ai} . Q_{a_i} was the
469 value for a_i . β was the inverse temperature constant tuned to 5 according to our
470 experimental data. After waiting, the agent would get a reward according to the same
471 reward schedule used in our experiment. Each action value was updated separately
472 during the reward delivery:

473
$$\delta = r - Q_a$$

474
$$r = R/(\tau + 1)$$

475
$$Q_a \leftarrow Q_a + \alpha * \delta$$

476 Where the reward prediction error δ was calculated by the difference between the
477 hyperbolically discounted reward r (or "reward rate", given by the absolute reward R
478 dividing total time $\tau+1$ for obtaining the reward, where τ was the waiting duration and
479 the additional 1sec was the estimated delay for running between two zones) and the
480 chosen action value Q_a . The reward prediction error was then used to update the value
481 of the chosen action. We tuned the learning rate α to 0.002 to fit the animal behavioral
482 data.

483 ***Continuous Deliberation.*** In each trial, the agent would go through a series of hidden
484 states, each lasting for 0~2sec randomly according to a Gaussian distribution (mean at
485 1 sec). At each hidden state, the agent had two action options, either to keep waiting or
486 to leave. If it chose to keep waiting, the agent would transition to the next hidden state,
487 with the past time of the previous state cumulated to the whole waiting duration. If the
488 choice was to leave, the cumulation would cease and a virtual reward dependent on the
489 duration will be delivered, and then a new trial would begin from the initial state. The
490 reward schedule was identical to that used for the animals during the experiments.

491 The action choice for the future was determined randomly by a Boltzmann
492 distribution (SoftMax) and action values:

493
$$P_{a_w}^{(T+1)} = \frac{e^{\beta Q_{a_w}^{(T+1)}}}{e^{\beta Q_{a_w}^{(T+1)}} + e^{\beta Q_{a_L}^{(T+1)}}}$$

494 $P_{a_w}^{(T+1)}$ was the probability of choosing to wait for the next state $T + 1$. $Q_{a_w}^{(T+1)}$
495 and $Q_{a_L}^{(T+1)}$ were the value of waiting and leaving, respectively, for state $T + 1$. β
496 was the inverse temperature constant tuned to 5.

497 The action values for each hidden state T were updated by temporal difference
498 learning algorithm (SARSA):

$$499 \quad \delta = r + \gamma * Q_{a'}^{(T+1)} - Q_a^{(T)}$$

$$500 \quad r = R/(\tau + 1)$$

$$501 \quad Q_a^{(T)} \leftarrow Q_a^{(T)} + \alpha * \delta$$

502 Where the future action a' was determined by the Boltzmann distribution in the
503 previous step. The current action a and the future action a' could both be either
504 waiting or leaving. The prediction error δ was calculated by the sum of reward rate r
505 (r remained zero until the reward R was delivered. $\tau+1$ was the total time for obtaining
506 the reward, where τ was the waiting duration and the additional 1sec was the
507 estimated delay for running between two zones) and the future action value $\gamma *$
508 $Q_{a'}^{(T+1)}$ discounted by γ ($\gamma = 0.9$), minus the current action value $Q_a^{(T)}$. When a was
509 leaving, the future action value $Q_{a'}^{(T+1)}$ would always be zero. This error signal δ
510 was used to update $Q_a^{(T)}$ with learning rate $\alpha = 0.001$.

511 As a Markov process, each state would be identical to the agent no matter how
512 the state was reached or what the following actions are. So, we extracted the learned
513 value of waiting as a time series along all the hidden states to compare with the averaged
514 curve of VTA DAergic activity. For each trial, we also extracted the time series of the

515 transient waiting value for a trial-wise analysis. Apart from the value of waiting, we
516 could also extract the time series of RPE for each trial.

517 For optogenetics manipulation, we simulated it in the model after normal training
518 was accomplished as in the animal experiments.

519 **Value manipulation.** In 20% trials of the stimulation session, the future waiting value
520 throughout the whole waiting period was manipulated. The optogenetics activation was
521 simulated as an extra positive value added onto the future waiting value, and the
522 optogenetics inhibition corresponded to a proportional decrease of the future waiting
523 value as follows:

$$524 \quad Q_{a_w}^{(T+1)} \leftarrow \tilde{Q}_{a_w}^{(T+1)}, \quad \text{for the current trial}$$

$$525 \quad \text{where } \tilde{Q}_{a_w}^{(T+1)} = \begin{cases} Q_{a_w}^{(T+1)} + \Delta_{value-ext}, & \text{if "ChR2 - lighton"} \\ \kappa_{value-inh} Q_{a_w}^{(T+1)}, & \text{if "eNPHR - lighton"} \end{cases}$$

$$526 \quad \text{and, } \delta = r + \gamma * \tilde{Q}_{a_w}^{(T+1)} - Q_a^{(T)}, \quad \text{if } a' = a_w$$

527 Here we set $\Delta_{value-ext} = 0.15$, and $\kappa_{value-inh} = 0.9$, so that the change in
528 averaged waiting duration in the simulated "light-on" trials can capture the magnitude
529 of the instantaneous effect of optogenetic stimulations on the current trials. Using these
530 parameters "calibrated" by the current trial effect, we were able to compare the
531 stimulation effect on the light-off or the following trials in both real and simulated
532 situations. Also note that if the future action was chosen as waiting, the manipulated
533 value of waiting would be used in the RPE calculation and thus current action value
534 updating as well.

535 **RPE manipulation.** Under this situation, in 20% trials of the stimulation session,
536 instead of future waiting value, RPE (δ) was manipulated throughout the whole waiting
537 period as follows:

$$538 \quad \tilde{\delta} = \begin{cases} \delta + \Delta_{RPE-ext}, & \text{if "Chr2 - lighton"} \\ \delta - \Delta_{RPE-inh}, & \text{if "eNPHR - lighton"} \end{cases}$$

539 and, $Q_a^{(T)} \leftarrow Q_a^{(T)} + \alpha * \tilde{\delta}$

540 we set $\Delta_{RPE-ext} = 15$, and $\Delta_{RPE-inh} = 20$, which was calibrated by the current trial
541 effect of real light stimulation.

542 To simulate the fluctuation in real DAergic signal, we simply multiplied the future
543 waiting value during each state by a factor $\sigma \sim \mathcal{N}(1, 0.3^2)$ (determined by the averaged
544 signal-dependent noise magnitude / relative standard deviation for all momentary
545 DAergic amplitudes), additionally to the original model (this is only implemented for
546 figs. S8E~F).

547 **Electrophysiological recordings.** Adult (8-10 weeks) DAT-IRES-Cre knock-in male
548 mice 4 weeks after injection with AAV2/9-EF1a-DIO-ChR2(H134R)-mCherry or
549 AAV-DIO-eNpHR3.0-mCherry were anesthetized with an intraperitoneal injection of
550 pentobarbital (100 mg kg⁻¹) and then perfused transcardially with ice-cold oxygenated
551 (95% O₂/5% CO₂) NMDG ACSF solution (93 mM NMDG, 93 mM HCl, 2.5 mM
552 KCl, 1.25 mM NaH₂PO₄, 10 mM MgSO₄·7H₂O, 30 mM NaHCO₃, 25 mM glucose,
553 20 mM HEPES, 5 mM sodium ascorbate, 3 mM sodium pyruvate, and 2 mM
554 thiourea, pH 7.4, 295-305 mOsm). After perfusion, the brain was rapidly dissected out
555 and immediately transferred into an ice-cold oxygenated NMDG ACSF solution.

556 Then the brain tissue was sectioned into slices horizontally at 280 μ m in the same
557 buffer with a vibratome (VT-1200 S, Leica). The brain slices containing the VTA
558 were incubated in oxygenated NMDG ACSF at 32°C for 10~15 min, then transferred
559 to a normal oxygenated solution of ACSF (126 mM NaCl, 2.5 mM KCl, 1.25 mM
560 NaH_2PO_4 , 2 mM $\text{MgSO}_4 \cdot 7\text{H}_2\text{O}$, 10 mM Glucose, 26 mM NaHCO_3 , 2 mM CaCl_2) at
561 room temperature for 1h. A slice was then transferred to the recording chamber,
562 which was submerged and superfused with ACSF at a rate of 3 ml/min at 28°C. Cells
563 were visualized using infrared DIC and fluorescence microscopy (BX51, Olympus).
564 VTA DAergic neurons were identified by their fluorescence and other
565 electrophysiological characteristics. Whole-cell current-clamp recordings of VTA
566 DAergic neurons were made using a MultiClamp 700B amplifier and Digidata 1440A
567 interface (Molecular Devices). Patch electrodes (3-5 M Ω) were backfilled with
568 internal solution containing (in mM): 130 K-gluconate, 8 NaCl, 10 HEPES, 1 EGTA,
569 2 Mg \cdot ATP and 0.2 $\text{Na}_3\cdot$ GTP (pH:7.2, 280 mOsm). Series resistance was monitored
570 throughout the experiments. For optogenetic activation, blue light was delivered onto
571 the slice through a 200- μ m optical fiber attached to a 470 nm LED light source
572 (Thorlabs, USA). The functional potency of the ChR2-expressing virus was validated
573 by measuring the number of action potentials elicited in VTA DAergic neurons using
574 blue light stimulation (20 ms, 10 Hz, 2.7 mW) in VTA slices. For optogenetic
575 inhibition, yellow light (0.7 mW) was generated by a 590 nm LED light source
576 (Thorlabs, USA) and delivered to VTA DAergic neurons expressing eNpHR3.0
577 through a 200- μ m optical fiber. To assure eNpHR-induced neuronal inhibition,

578 whole-cell recordings were carried out in current-clamp mode and spikes were
579 induced by current injection (200 pA) with the presence of yellow light. Data were
580 filtered at 2 kHz, digitized at 10 kHz, and acquired using pClamp10 software
581 (Molecular Devices).

582 **Immunostaining.** Mice were deeply anesthetized with pentobarbital (100 mg/kg, i.p.),
583 following saline perfusion through the heart. After blood was drained out, 4%
584 paraformaldehyde (PFA) was used for fixation. Then the head was cut off and soaked
585 in 4% PFA at room temperature overnight. The brain was harvested the next day,
586 post-fixed overnight in 4% PFA at 4°C, and transferred to 30% sucrose in 0.1 M PBS,
587 pH 7.4 for 24~48 h. Coronal sections (20 μ m) containing the VTA were cut on a
588 cryostat (Leica CM3050 S). The slides were washed with 0.1 M PBS, pH 7.4,
589 incubated in blocking buffer (0.3% Triton X-100, 5% bovine serum albumin in 0.1 M
590 PBS, pH 7.4) for an hour, and then transferred into the primary antibody (rabbit anti-
591 tyrosine hydroxylase antibody, 1:1,000; Invitrogen) in blocking buffer overnight at
592 4°C. The sections were washed three times in 0.1 M PBS, then incubated with donkey
593 anti-rabbit IgG H&L secondary antibody (conjugated to fluor-488 or fluor-594,
594 1:1,000; Jackson ImmunoResearch) at room temperature for 2 h. The nucleus was
595 stained with DAPI (4',6-diamidino-2-phenylindole). Sections were mounted in
596 glycerine and covered with coverslips sealed in place. Fluorescent images were
597 collected using a Zeiss confocal microscope (LSM 880).

598 **Quantification and statistics.** All statistics were performed by MATLAB (R2016b,
599 MathWorks) and Python (V2.7, Python Software Foundation) routines. Data were

600 judged to be statistically significant, while the P-value less than 0.05. Asterisks denote
601 statistical significance * $p < 0.05$; ** $p < 0.01$; *** $p < 0.001$. Unless stated otherwise,
602 values were presented as Mean \pm s.e.m..

603 REFERENCES AND NOTES

- 604 1. D. Tomasi, N. D. Volkow, Striatocortical pathway dysfunction in addiction and
605 obesity: differences and similarities. *Crit Rev Biochem Mol Biol* **48**, 1-19 (2013).
- 606 2. N. D. Volkow, R. D. Baler, NOW vs LATER brain circuits: implications for obesity
607 and addiction. *Trends Neurosci* **38**, 345-352 (2015).
- 608 3. W. Mischel, Y. Shoda, M. I. Rodriguez, Delay of gratification in children. *Science*
609 **244**, 933-938 (1989).
- 610 4. J. E. Maddux, J. P. Tangney, *Social psychological foundations of clinical*
611 *psychology*. (Guilford Press, New York, 2010), pp. xv, 555 p.
- 612 5. J. Grosch, A. Neuringer, Self-control in pigeons under the Mischel paradigm. *J*
613 *Exp Anal Behav* **35**, 3-21 (1981).
- 614 6. B. Reynolds, H. de Wit, J. B. Richards, Delay of gratification and delay
615 discounting in rats. *Behav Process* **59**, 157-168 (2002).
- 616 7. B. Engelhard *et al.*, Specialized coding of sensory, motor and cognitive variables
617 in VTA dopamine neurons. *Nature*, (2019).
- 618 8. C. K. Starkweather, B. M. Babayan, N. Uchida, S. J. Gershman, Dopamine reward
619 prediction errors reflect hidden-state inference across time. *Nat Neurosci* **20**,
620 581-589 (2017).
- 621 9. P. W. Glimcher, Understanding dopamine and reinforcement learning: the
622 dopamine reward prediction error hypothesis. *Proc Natl Acad Sci U S A* **108**
623 **Suppl 3**, 15647-15654 (2011).
- 624 10. G. Morris, A. Nevet, D. Arkadir, E. Vaadia, H. Bergman, Midbrain dopamine
625 neurons encode decisions for future action. *Nat Neurosci* **9**, 1057-1063 (2006).
- 626 11. J. R. Hollerman, W. Schultz, Dopamine neurons report an error in the temporal
627 prediction of reward during learning. *Nat Neurosci* **1**, 304-309 (1998).
- 628 12. W. Schultz, P. Dayan, P. R. Montague, A neural substrate of prediction and
629 reward. *Science* **275**, 1593-1599 (1997).
- 630 13. A. A. Hamid *et al.*, Mesolimbic dopamine signals the value of work. *Nat Neurosci*
631 **19**, 117-126 (2016).
- 632 14. M. W. Howe, P. L. Tierney, S. G. Sandberg, P. E. M. Phillips, A. M. Graybiel,
633 Prolonged dopamine signalling in striatum signals proximity and value of distant
634 rewards. *Nature* **500**, 575-+ (2013).
- 635 15. S. Soares, B. V. Atallah, J. J. Paton, Midbrain dopamine neurons control
636 judgment of time. *Science* **354**, 1273-1277 (2016).

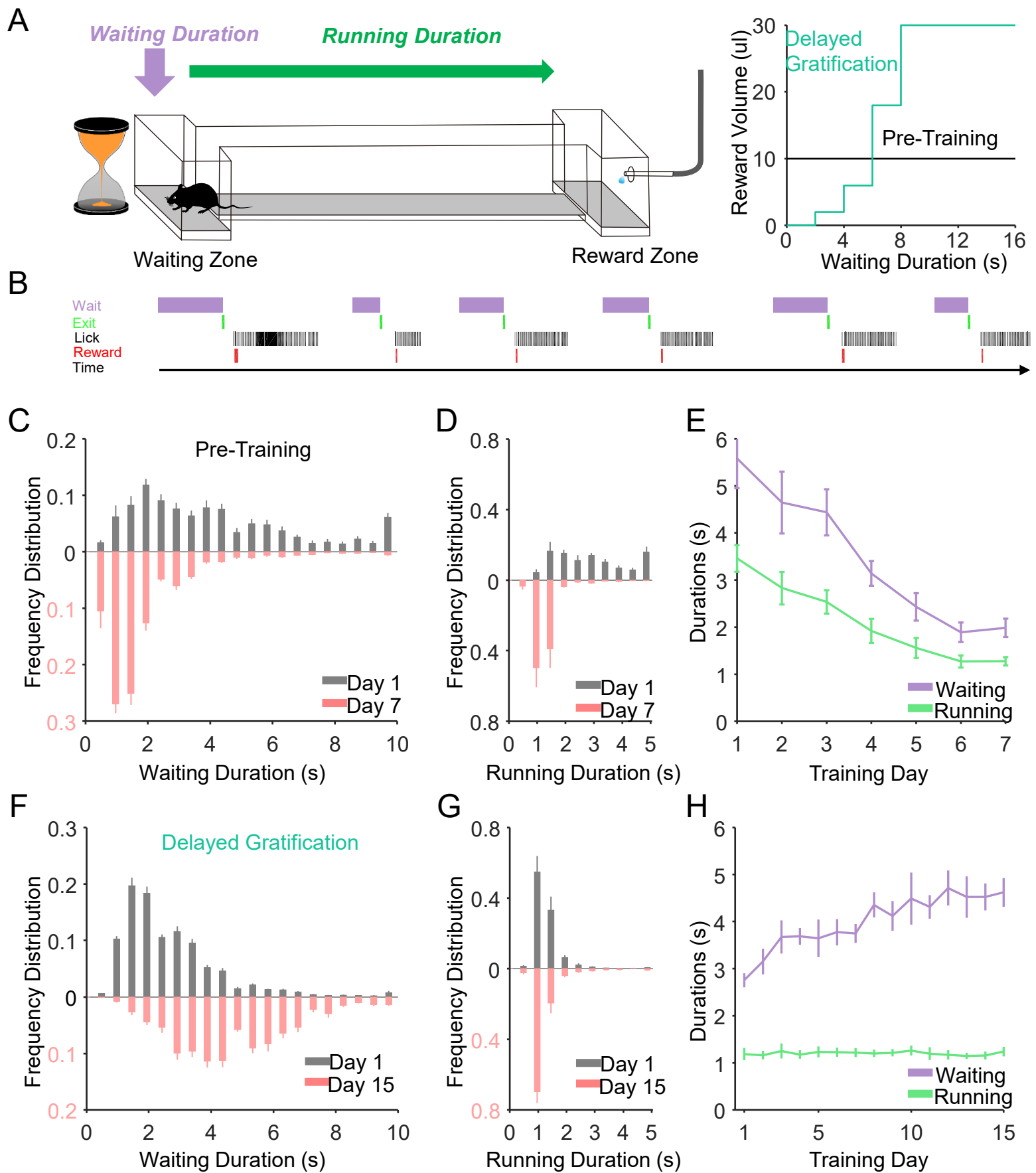
- 637 16. M. Guitart-Masip, U. R. Beierholm, R. Dolan, E. Duzel, P. Dayan, Vigor in the face
638 of fluctuating rates of reward: an experimental examination. *J Cogn Neurosci* **23**,
639 3933-3938 (2011).
- 640 17. Y. Niv, Cost, benefit, tonic, phasic: what do response rates tell us about
641 dopamine and motivation? *Ann N Y Acad Sci* **1104**, 357-376 (2007).
- 642 18. Y. Niv, N. D. Daw, P. Dayan, Choice values. *Nat Neurosci* **9**, 987-988 (2006).
- 643 19. Y. Li *et al.*, Serotonin neurons in the dorsal raphe nucleus encode reward signals.
644 *Nat Commun* **7**, 10503 (2016).
- 645 20. G. A. Rummery, M. Niranjan, On-Line Q-Learning Using Connectionist
646 Systems. *Technical Report CUED/F-Infeng/TR 166*, (1994).
- 647 21. R. S. Sutton, A. G. Barto, *Reinforcement learning : an introduction*. Adaptive
648 computation and machine learning (MIT Press, Cambridge, Mass., 1998), pp.
649 xviii, 322 p.
- 650 22. M. Treisman, Temporal discrimination and the indifference interval. Implications
651 for a model of the "internal clock". *Psychol Monogr* **77**, 1-31 (1963).
- 652 23. P. R. Killeen, J. G. Fetterman, A behavioral theory of timing. *Psychol Rev* **95**, 274-
653 295 (1988).
- 654 24. W. H. Meck, Neuropharmacology of timing and time perception. *Brain Res Cogn*
655 *Brain Res* **3**, 227-242 (1996).
- 656 25. M. Jazayeri, M. N. Shadlen, A Neural Mechanism for Sensing and Reproducing a
657 Time Interval. *Curr Biol* **25**, 2599-2609 (2015).
- 658 26. Y. Niv, N. D. Daw, D. Joel, P. Dayan, Tonic dopamine: opportunity costs and the
659 control of response vigor. *Psychopharmacology (Berl)* **191**, 507-520 (2007).
- 660 27. R. S. Sutton, A. G. Barto, *Reinforcement learning : an introduction*. Adaptive
661 computation and machine learning series (The MIT Press, Cambridge,
662 Massachusetts, ed. Second edition., 2018), pp. xxii, 526 pages.
- 663 28. B. A. Bari *et al.*, Stable Representations of Decision Variables for Flexible
664 Behavior. *Neuron*, (2019).
- 665 29. P. Simen, F. Balci, L. de Souza, J. D. Cohen, P. Holmes, A model of interval timing
666 by neural integration. *J Neurosci* **31**, 9238-9253 (2011).
- 667 30. D. Durstewitz, Self-Organizing Neural Integrator Predicts Interval Times through
668 Climbing Activity. **23**, 5342-5353 (2003).
- 669 31. F. Balci, P. Simen, A decision model of timing. *Current Opinion in Behavioral*
670 *Sciences* **8**, 94-101 (2016).
- 671 32. S. I. Gass, C. M. Harris, in *Encyclopedia of Operations Research and*
672 *Management Science*, S. I. Gass, C. M. Harris, Eds. (Springer US, New York, NY,
673 2001), pp. 490-490.
- 674 33. T. K. Koerner, Y. Zhang, Application of Linear Mixed-Effects Models in Human
675 Neuroscience Research: A Comparison with Pearson Correlation in Two Auditory
676 Electrophysiology Studies. *Brain Sci* **7**, (2017).
- 677 34. M. P. Boisgontier, B. Cheval, The anova to mixed model transition. *Neurosci*
678 *Biobehav Rev* **68**, 1004-1005 (2016).
- 679 35. S. N. Chettih, C. D. Harvey, Single-neuron perturbations reveal feature-specific
680 competition in V1. *Nature* **567**, 334-340 (2019).

- 681 36. K. Jimura, M. S. Chushak, T. S. Braver, Impulsivity and self-control during
682 intertemporal decision making linked to the neural dynamics of reward value
683 representation. *J Neurosci* **33**, 344-357 (2013).
- 684 37. B. Schmidt, C. B. Holroyd, S. Debener, J. Hewig, Why Is It So Hard to Wait? Brain
685 Responses to Delayed Gratification Predict Impulsivity and Self-Control.
686 *Psychophysiology* **52**, S42-S42 (2015).
- 687 38. C. D. Fiorillo, W. T. Newsome, W. Schultz, The temporal precision of reward
688 prediction in dopamine neurons. *Nat Neurosci* **11**, 966-973 (2008).
- 689 39. B. M. Babayan, N. Uchida, S. J. Gershman, Belief state representation in the
690 dopamine system. *Nat Commun* **9**, 1891 (2018).
- 691 40. J. Y. Cohen, S. Haesler, L. Vong, B. B. Lowell, N. Uchida, Neuron-type-specific
692 signals for reward and punishment in the ventral tegmental area. *Nature* **482**,
693 85-88 (2012).
- 694 41. J. G. Mikhael, H. R. Kim, N. Uchida, S. J. Gershman, Ramping and State
695 Uncertainty in the Dopamine Signal. 805366 (2019).
- 696 42. A. Guru *et al.*, Ramping activity in midbrain dopamine neurons signifies the use
697 of a cognitive map. 2020.2005.2021.108886 (2020).
- 698 43. S. J. Gershman, Dopamine ramps are a consequence of reward prediction errors.
699 *Neural Comput* **26**, 467-471 (2014).
- 700 44. H. R. Kim *et al.*, A Unified Framework for Dopamine Signals across Timescales.
701 *Cell* **183**, 1600-1616 e1625 (2020).
- 702 45. W. X. Pan, R. Schmidt, J. R. Wickens, B. I. Hyland, Dopamine cells respond to
703 predicted events during classical conditioning: evidence for eligibility traces in
704 the reward-learning network. *J Neurosci* **25**, 6235-6242 (2005).
- 705 46. H. C. Tsai *et al.*, Phasic firing in dopaminergic neurons is sufficient for behavioral
706 conditioning. *Science* **324**, 1080-1084 (2009).
- 707 47. A. E. Goudriaan, M. Yucel, R. J. van Holst, Getting a grip on problem gambling:
708 what can neuroscience tell us? *Front Behav Neurosci* **8**, 141 (2014).

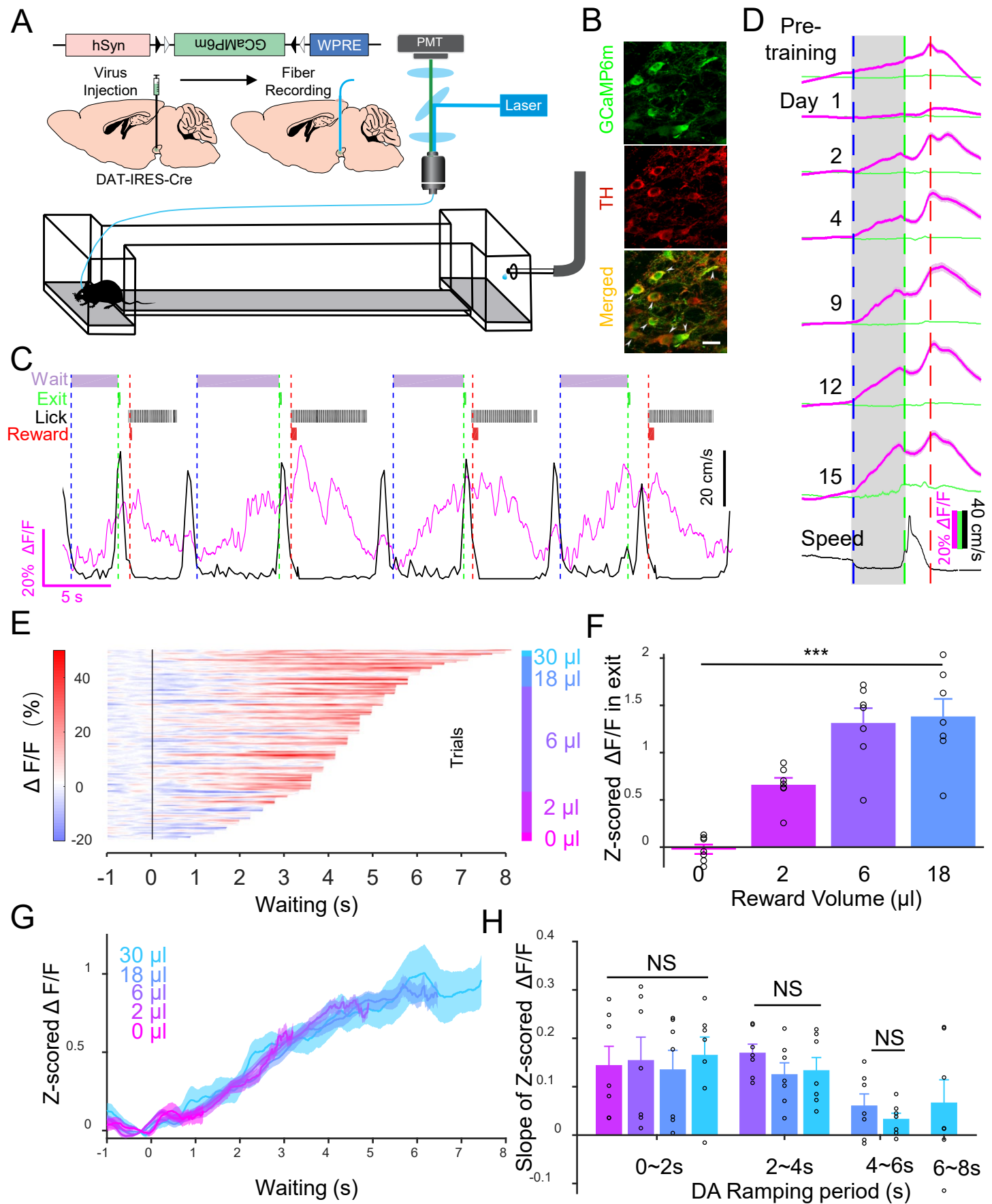
709 **Acknowledgments:** We thank Drs. M. Lou, W. Ge, Y. Rao, W Zhou, and B Min for
710 comments on the manuscript. This work was supported by the National Natural Science
711 Foundation of China (grant nos. 31922029, 31671086, 61890951, and 61890950 to J.H.), a
712 Shanghai Pujiang Talent Award (grant no. 2018X0302-101-01 to W.S.). **Author**
713 **contributions:** W.S. and J. H. oversaw the whole project. W.S., J. H. & Z. G. designed the
714 experiments. Z. G. & C. L. performed all animal experiments. Z. G. & H. W. analyzed the
715 data. H. W. & S. F. performed the computational modeling under the supervision of X. J. W..
716 M. C. performed the electrophysiological recordings. W.S., J. H., H. W. & Z. G. wrote the

717 paper with the participation of all other authors. **Competing interests:** The authors declare no

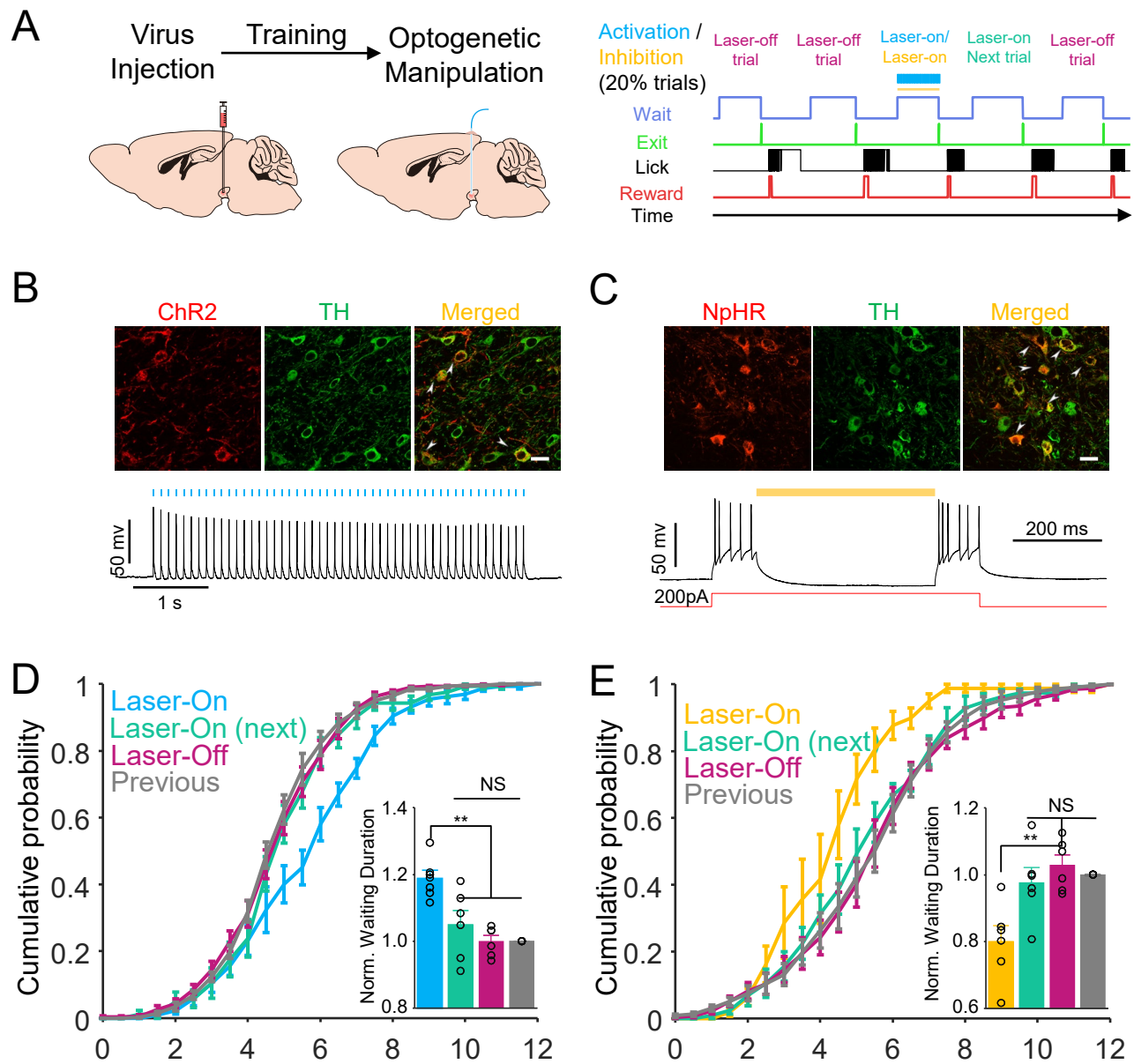
718 competing interests.



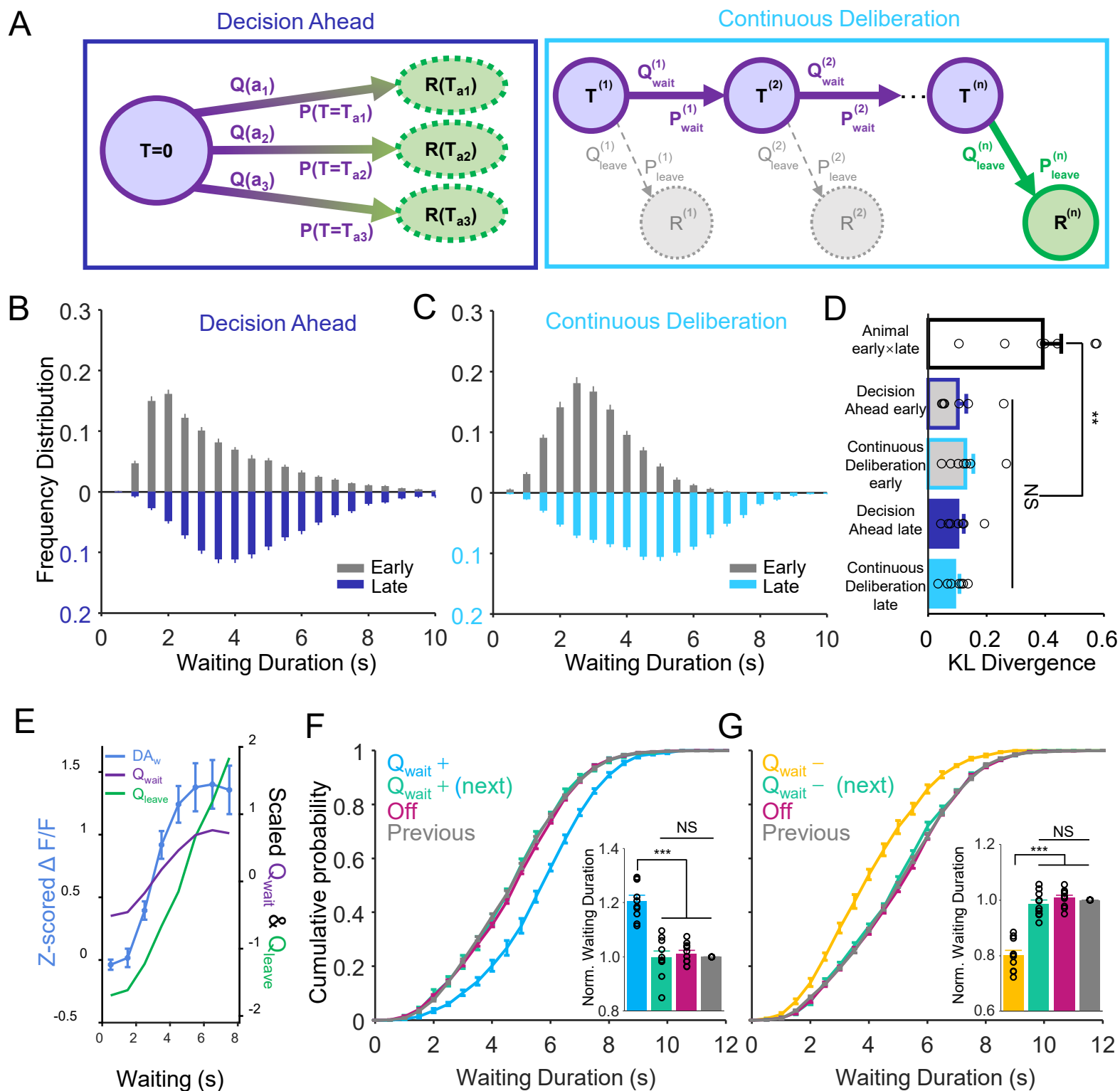
1 **Fig. 1. The behavioral performance of mice during a delayed gratification task**
2 **learning.** (A) Left panel, schematic of the delayed gratification task. Right panel, the
3 relationship between reward volumes and waiting durations in the two behavioral tasks.
4 (B) This plot presents the Transistor-Transistor Logic (TTL) signals for the
5 chronological sequence of behavioral events in the tasks. (C-E). The waiting duration
6 and running duration both decreased with the training process in the pre-training phase
7 (Day1, Waiting: 5.58 ± 0.63 sec; Running: 3.46 ± 0.28 sec; $p < 0.001$; Day 7, Waiting:
8 1.99 ± 0.19 sec; Running: 1.28 ± 0.09 sec, $p < 0.001$, $n = 7$ mice, Friedman test). (F) The
9 distribution of waiting durations from the behavioral session on the last analyzed day
10 (Day 15, light red), revealing significantly longer waiting durations compared with that
11 from day 1 (D 1, gray, $n = 7$ mice). (G) The distribution of running durations from D1
12 and D15 did not differ with training. (H) The plots show that the continuous training
13 steadily increased the averaged waiting duration from 2.76 ± 0.15 s on Day 1 to $4.62 \pm$
14 0.30 s on Day 15 ($p < 0.001$, $n = 7$ mice, Friedman test), whereas the training did not
15 change the average running duration from 1.19 ± 0.13 s on Day 1 to 1.24 ± 0.10 s on Day
16 15 ($p = 0.97$, $n = 7$ mice, Friedman test). All error bars represent the s.e.m..



17 **Fig. 2. VTA DAergic activity ramps up consistently while the mice are waiting for**
18 **the reward.** (A) Schematic of stereotaxic virus injection procedures. (B) Confocal
19 images illustrating GCaMP6m (green) expression in VTA TH⁺ neurons (red). Scale bar:
20 20 μm . (C) An example of a live-recording trace (magenta line) of Ca²⁺ signal in VTA
21 ^{DA} neurons and running speed (black line) when a Dat-Cre: GCaMP6m mouse was
22 performing the delayed gratification task. Delayed gratification task events over time
23 (top): the dashed vertical lines indicated waiting onset (blue), waiting termination
24 (green), and reward onset (red). (D) The scaled Ca²⁺ signals curves (magenta) and GFP
25 signals (green) curves of VTA^{DA} neurons from the last day in pre-training and day 1 to
26 day 15 in the delayed gratification task training (black line, speed). (E) Sorted Ramping
27 Ca²⁺ signal data from one mouse on the last day (D15) of the delayed gratification task
28 training (150 trials). The signal traces were aligned to waiting onset, sorted in waiting
29 duration length, and separated into five groups of the reward outcomes (0, 2, 6, 18, and
30 30 μl). f. Z-scored $\Delta F/F$ values at 0.5s before exit were significantly different while the
31 reward volumes were different ($F=24.67$, $p<0.01$, $n=7$, one-way ANOVA). (G)
32 Averaged Ca²⁺ signal curves with different outcomes from Fig. E. Slopes of Ca²⁺
33 signals for every outcome, showing that there were no differences in all DAergic
34 ramping periods throughout the last week of training (0~2s, $F=0.10$, $p=0.96$; 2~4s,
35 $F=1.03$, $p=0.38$; 4~6s, $F=1.00$, $p=0.34$, $n=7$, one-way ANOVA). All error bars represent
36 the s.e.m.. For (D) and (G), the shaded region represents s.e.m..

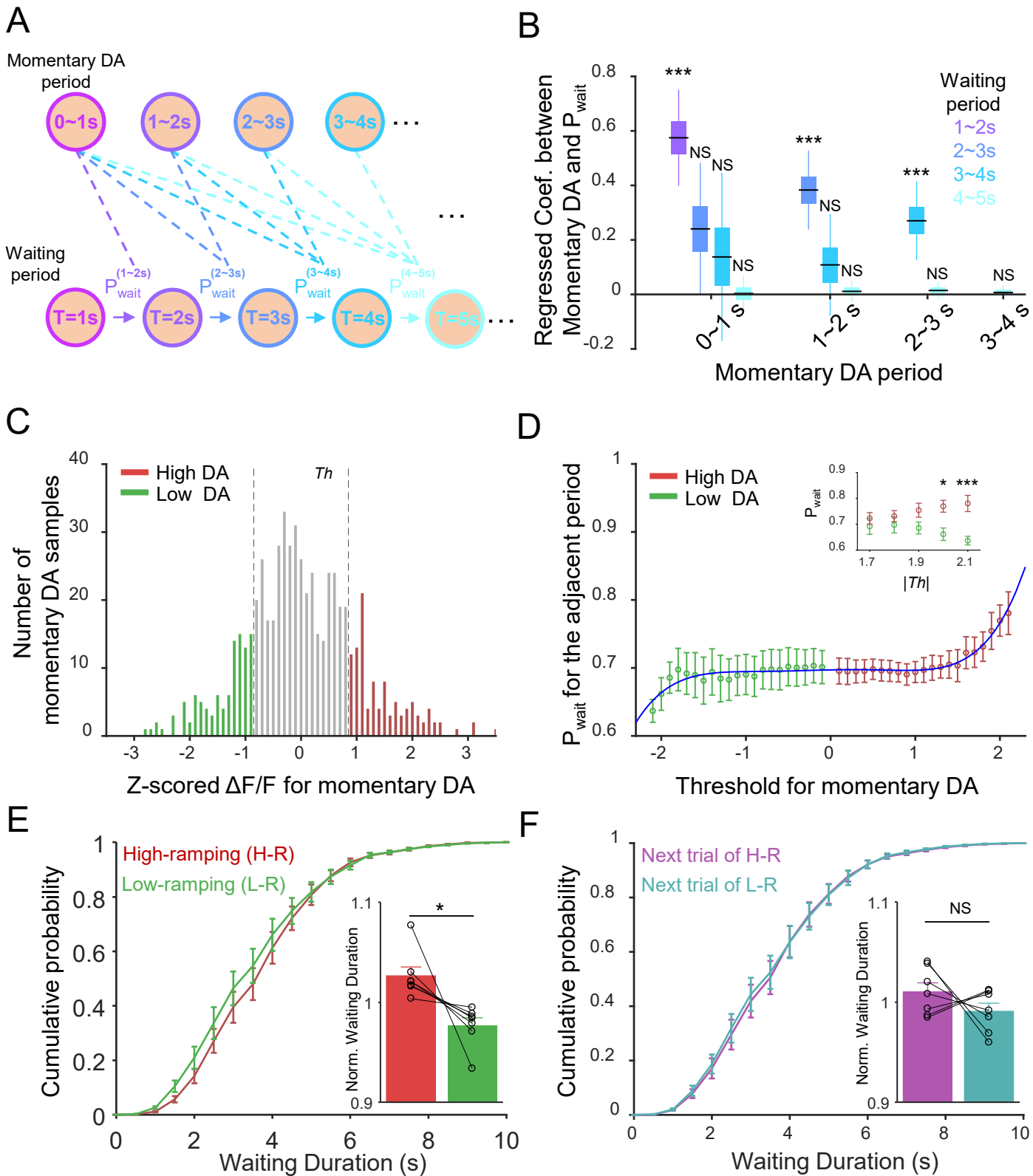


37 **Fig. 3. Optogenetic manipulation of VTA DAergic activity altered the waiting**
38 **durations.** (A) Left panels: schematic of stereotaxic virus injection and surgical
39 procedure. Right panels: the behavioral events and optogenetic manipulation protocol.
40 (B) Top panels: confocal image showing ChR2-mCherry (red) expression in VTA TH+
41 neurons (green). Bottom panels: whole-cell recording of VTA TH+ neurons in brain
42 slice showing action potentials evoked by 10-Hz 473 nm laser flash sequences (50
43 flashes, 20ms interval). (C) Top panels: confocal image showing eNpHR3.0-mCherry
44 (red) expression in VTA TH+ neurons (green). Bottom panels: whole-cell recording of
45 VTA TH+ in brain slice showing that action potentials evoked by 200pA current
46 injection were inhibited by continued 589nm laser. Scale bar: 20 μ m. (D) Cumulative
47 probabilities of waiting durations. The waiting durations of optogenetically activated
48 trials were significantly increased (blue, $F=12.93$, $p=0.002$, $n=6$ mice, one-way
49 ANOVA) than that of the previous day's trials (gray); note that the waiting duration of
50 unstimulated trials (red) did not differ from that of the previous day's trials (magenta,
51 $p=0.96$), or the next trials following photoactivation (green, $p=0.63$). Insert: a bar graph
52 of the normalized waiting durations from lasing stimulation (blue), the previous day's
53 trials (gray), photoactivated trials (blue, 1.19 ± 0.03), unstimulated trials (red,
54 1.00 ± 0.02), and the next trials following the photoactivation (green, 1.05 ± 0.04). (E)
55 The same experimental configuration as in (D), but VAT TH⁺ neurons were
56 optogenetically inhibited by a yellow laser. Optogenetic inhibition decreased the
57 waiting duration (yellow, 0.80 ± 0.05 , $F = 7.76$, $p=0.008$, $n=6$ mice, one-way ANOVA),
58 whereas there was no difference between uninhibited trials (red, 1.03 ± 0.03 , $p=0.80$),
59 the trials following the photoinhibition (green, 0.98 ± 0.04 , $p=0.80$), and the previous
60 day's trials (gray). All error bars represent the s.e.m..



61 **Fig. 4. Behavioral performances and ramping VTA DAergic activity are explained**
62 **by RL model. (A)** Two reinforcement learning computational models, *Decision Ahead*
63 in the dark blue box and *Continuous Deliberation* in the light blue box, simulating the
64 decision processes and variables under the delayed gratification task. In *Decision Ahead*:
65 $Q(a_n)$ is the value for action $a(n)$ and $P(T_{an})$ is the probability of action $a(n)$; $Q(a_n)$ was
66 used to compute the probability of action $a(n)$. In *Continuous Deliberation*: probability
67 of waiting, $P_{wait}^{(n)}$; waiting action value, $Q_{wait}^{(n)}$; probability of leaving, $P_{leave}^{(n)}$; leaving
68 action value, $Q_{leave}^{(n)}$; $R^{(n)}$, received reward; $Q_{wait}^{(n)}$ and $Q_{leave}^{(n)}$ were used to compute
69 the probability of waiting or leaving. **(B-C)** The distributions of waiting durations from
70 the early session and late session simulated in *Decision Ahead* model **(B)** and
71 *Continuous Deliberation* model **(C)** both displayed a similar distribution with our
72 experiment data **(Fig. 1F)**. **(D)** The distributions of behavioral performances between
73 early training days and late training days from experiment data were very different, in
74 which the Kullback-Leibler (KL) divergence was big enough (0.39 ± 0.06). The KL
75 divergences between the distributions of simulated behavioral performances from both
76 models in early or late session and experiment data of every mouse cross whole training
77 sessions were significant small ($p = 0.005$, $n = 7$ mice, Friedman test), in which there
78 was no difference ($p > 0.99$) between *Decision Ahead* RL model and *Continuous*
79 *Deliberation* RL model in late and early session. Data are represented as mean \pm SEM.
80 **(E)** Plots of Z-scored $\Delta F/F$ values (DA_w , light blue) at 0.5s before the waiting ended,
81 the scaled values of waiting (Q_{wait} , light purple), and the value of leaving (Q_{leave} , green)
82 from *Continuous Deliberation* model in last training session. The Q_{wait} and Q_{leave} both
83 predicted the experimental observation well (Q_{wait} : $r = 0.99$, $p < 0.001$; Q_{leave} : $r = 0.91$, p
84 $= 0.002$, Pearson correlation). **(F-G)** Computational reinforcement learning model
85 (continuous deliberation)-simulated data, dependent on manipulating the value of
86 waiting (Q_{wait}) in delayed gratification task. As with the experimental data in **Fig. 3D-**
87 **E**, the model-simulated data also shows that increased Q_{wait} only increases the waiting
88 durations of Q_{wait} increased trials **((F)** , $p < 0.001$, Friedman test, $n=10$) whereas
89 decreased Q_{wait} can decrease the waiting durations of Q_{wait} decreased trials **((G)**,

90 $p < 0.001$, Friedman test, $n=10$). The unstimulated trials including the next trials after
91 Q_{wait} manipulation had no difference with the last round regular running ((**F-G**),
92 $p > 0.999$, Friedman test, $n=10$). All error bars represent the s.e.m..



93 **Fig. 5. VTA DAergic activity during waiting predicts the behavioral performance**
94 **in the delayed gratification task.** (A) Schematic of waiting probability (P_{wait}) in
95 waiting periods after momentary DAergic periods from the experimental data. For
96 momentary DAergic activity in each period, the mouse has P_{wait} , which is calculated by
97 the waiting durations and trial number, in any waiting periods after the momentary
98 DAergic period. (B) Relationship between momentary VTA DAergic activity (Ca^{2+}
99 signals) and its waiting probability. For each momentary DAergic period, its DAergic
100 activity is only highly correlated ($p < 0.001$, $n = 7$ mouse, black lines, regressed
101 coefficient median; boxes, 50% confidence interval; whisker, 95% confidence interval)
102 with P_{wait} in the adjacent waiting period (the left bar of each cluster). (C) The
103 distribution of Z-scored mean $\Delta F/F$ of momentary DAergic periods. Three colors
104 illustrate high dopamine activity (High DA: red, greater than the threshold value, gray
105 dash line, while the threshold value is positive) trial numbers, low dopamine activity
106 (Low DA: green, less than the threshold value, while the threshold value is negative)
107 trial numbers, and all other (gray) dopamine activity trial numbers. (D) The waiting
108 probability of High DA (red) and Low DA (green) activity trials for the adjacent period
109 after the momentary DA periods. The P_w of High DA and Low DA activity trials fit
110 well with a fifth-degree polynomial function ($R^2 = 0.93$, $-2.1 \leq \text{threshold} \leq 2.1$). While the
111 absolute values of the threshold are big enough ($|\text{Th}| \geq 1.7$), the P_w of the High DA
112 activity trails is significantly ($p = 0.04$, $F(1,12) = 5.483$, Two-way ANOVA) higher than
113 the P_{wait} of the Low DA-ramping activity trials in adjacent waiting periods
114 ($|\text{threshold}| = 2.0$, $p = 0.02$; $|\text{threshold}| = 2.1$, $p < 0.001$, Sidak's multiple comparisons test,
115 $n = 7$). (E-F) Cumulative probabilities of waiting durations for the high DA-ramping
116 trials (e, H-R, red), the lower-DA-ramping trials (E, L-R, green) and their "next-in-
117 series" trials (F). Bar graph showing that the normalized waiting durations (1.03 ± 0.01)
118 of the higher-DA-ramping trials are significantly longer than that of the lower-DA-
119 ramping trials ((E), 0.98 ± 0.01 , $p = 0.024$, $n = 7$, paired Student's t-test), but have no
120 difference between their "next-in-series" trials ((F), next trial of H-R, 1.01 ± 0.01 ; next

121 trial of L-R, 0.99 ± 0.01 ; $p=0.290$, $n=7$, paired Student's t-test). All error bars represent
122 the s.e.m..

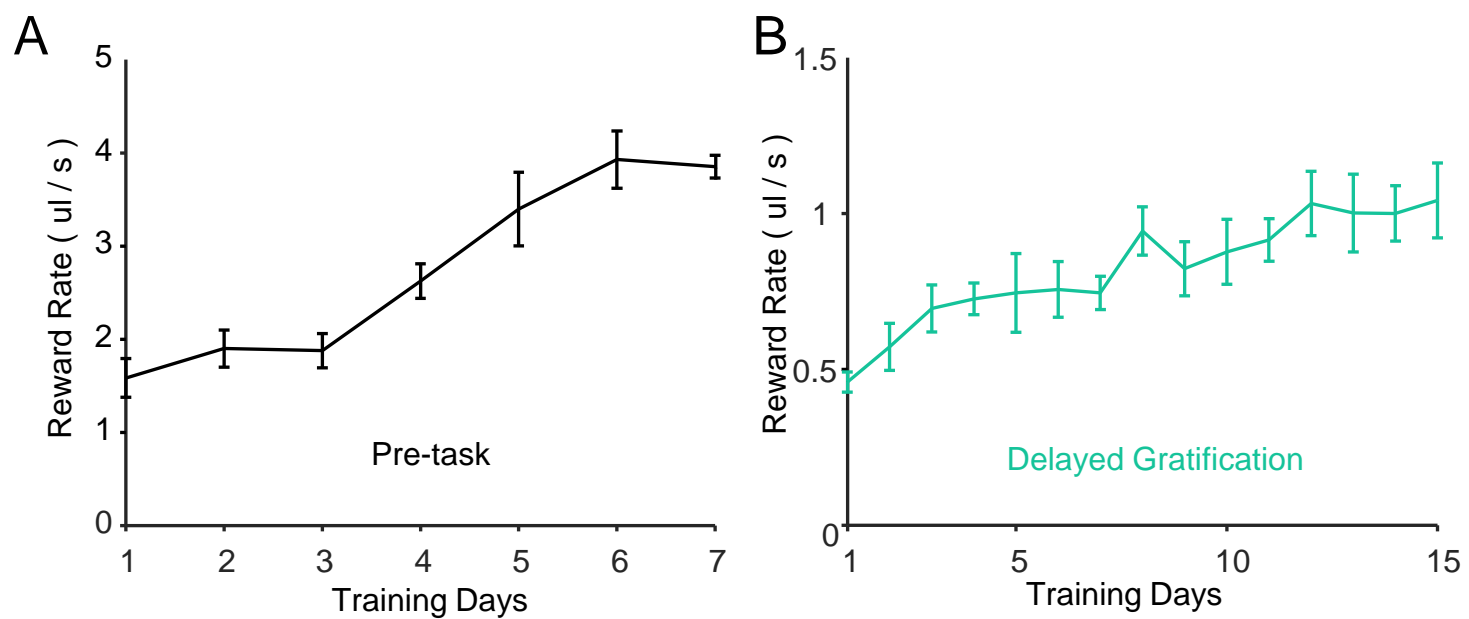


Fig. S1. The reward rate during behavioral training. (A-B) The reward rate both increased in pre-task training (A) and delayed gratification training (B). All error bars represent the s.e.m..

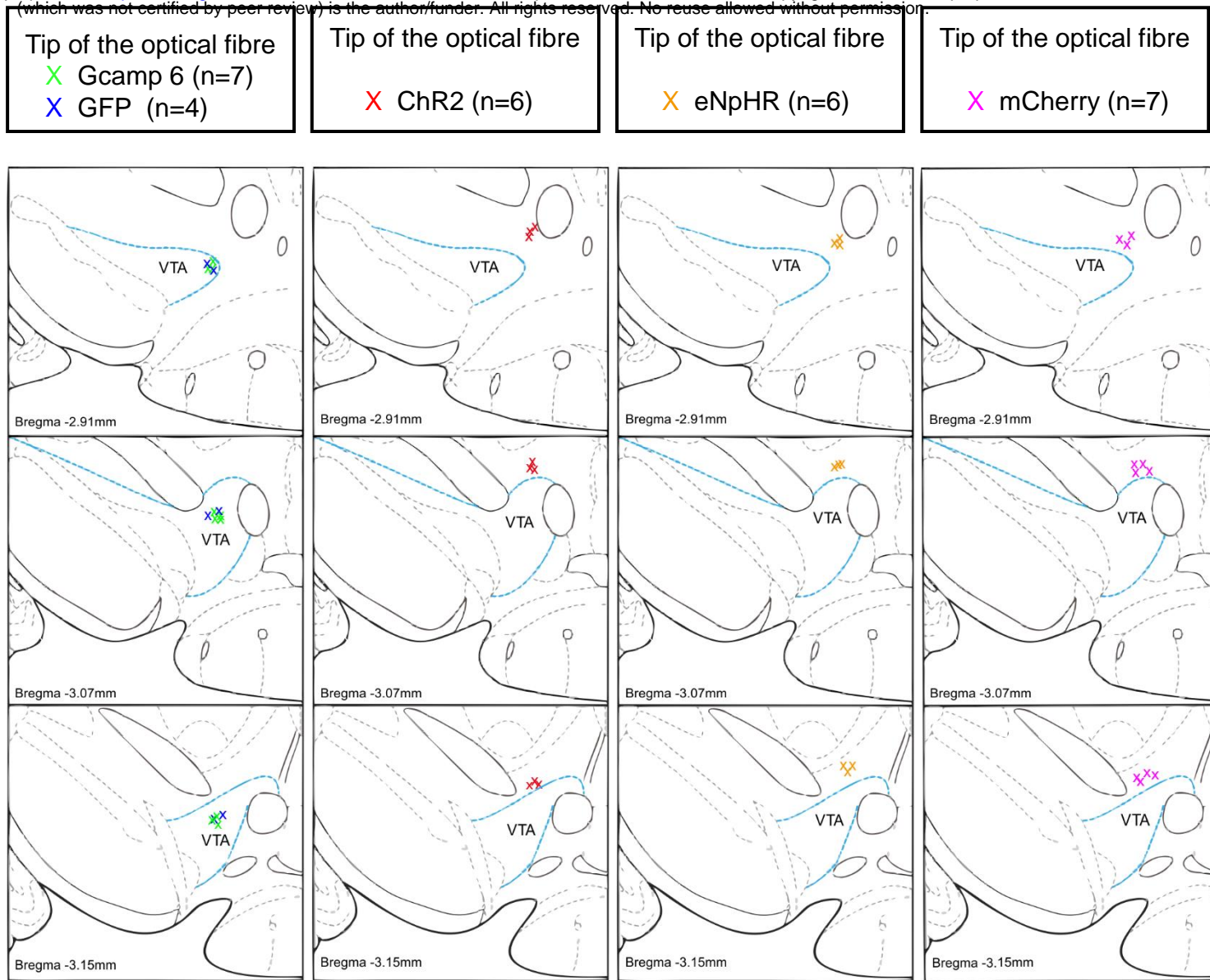


Fig. S2. Tip Positions of optical fibre in GCaMP6m (green, n = 7), GFP (blue, n = 4), ChR2 (red, n = 6), ENpHR (yellow, n = 6) and mCherry (magenta, n = 7) shown as coordinates in the mouse brain atlas.

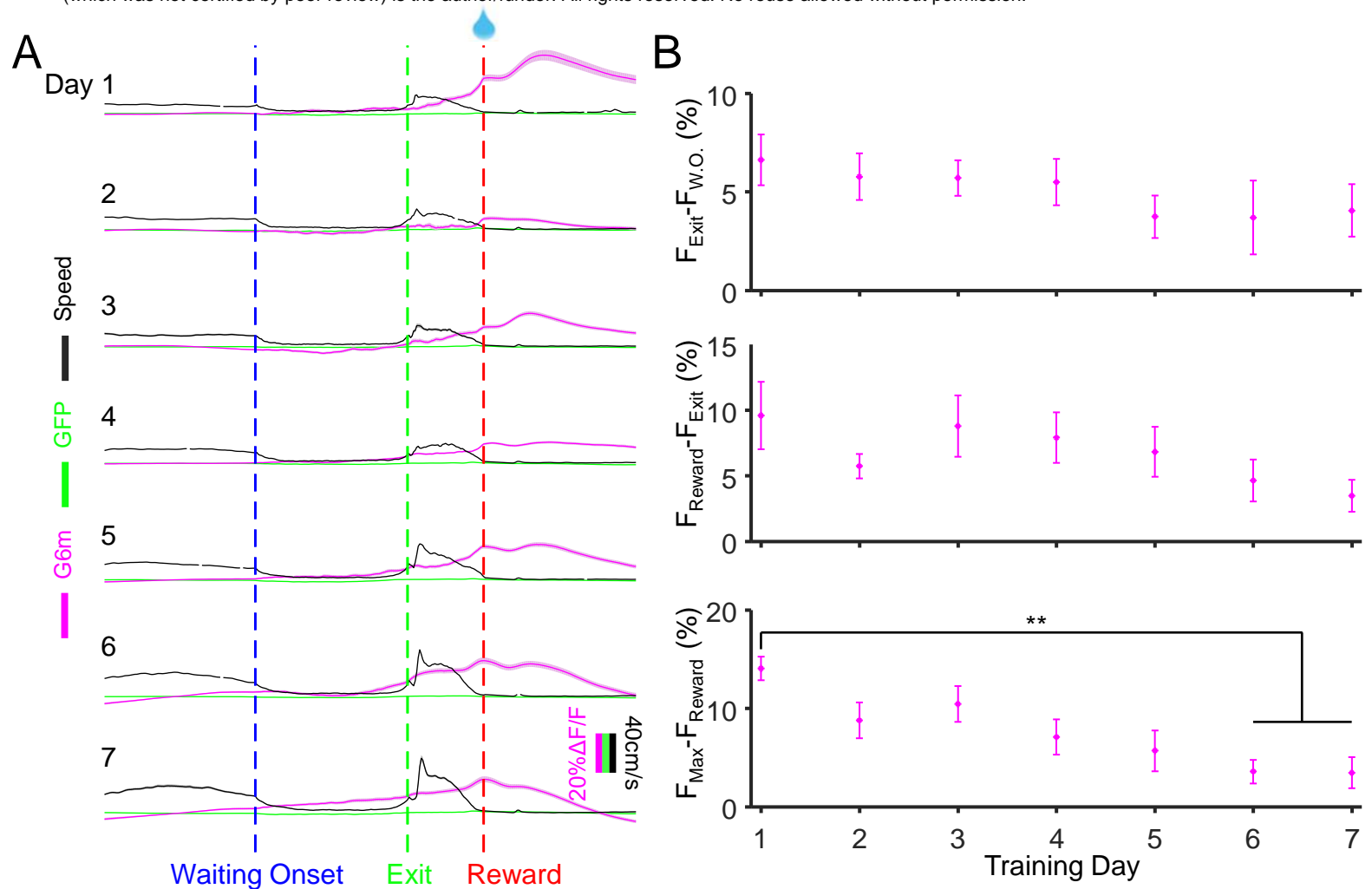


Fig. S3. The calcium and GFP signals of VTA DA neurons in behavioral tasks. (A) The GFP (green), calcium signals (magenta) of VTA DA neurons, and mouse running speed (black) in 7 days of pre-training. (B) The calcium signals of VTA DA neurons were changing with pre-training progress. The max calcium signals after the mouse received water rewards significantly decreased in day 6&7 ($p < 0.01$, Friedman test). All error bars represent the s.e.m..

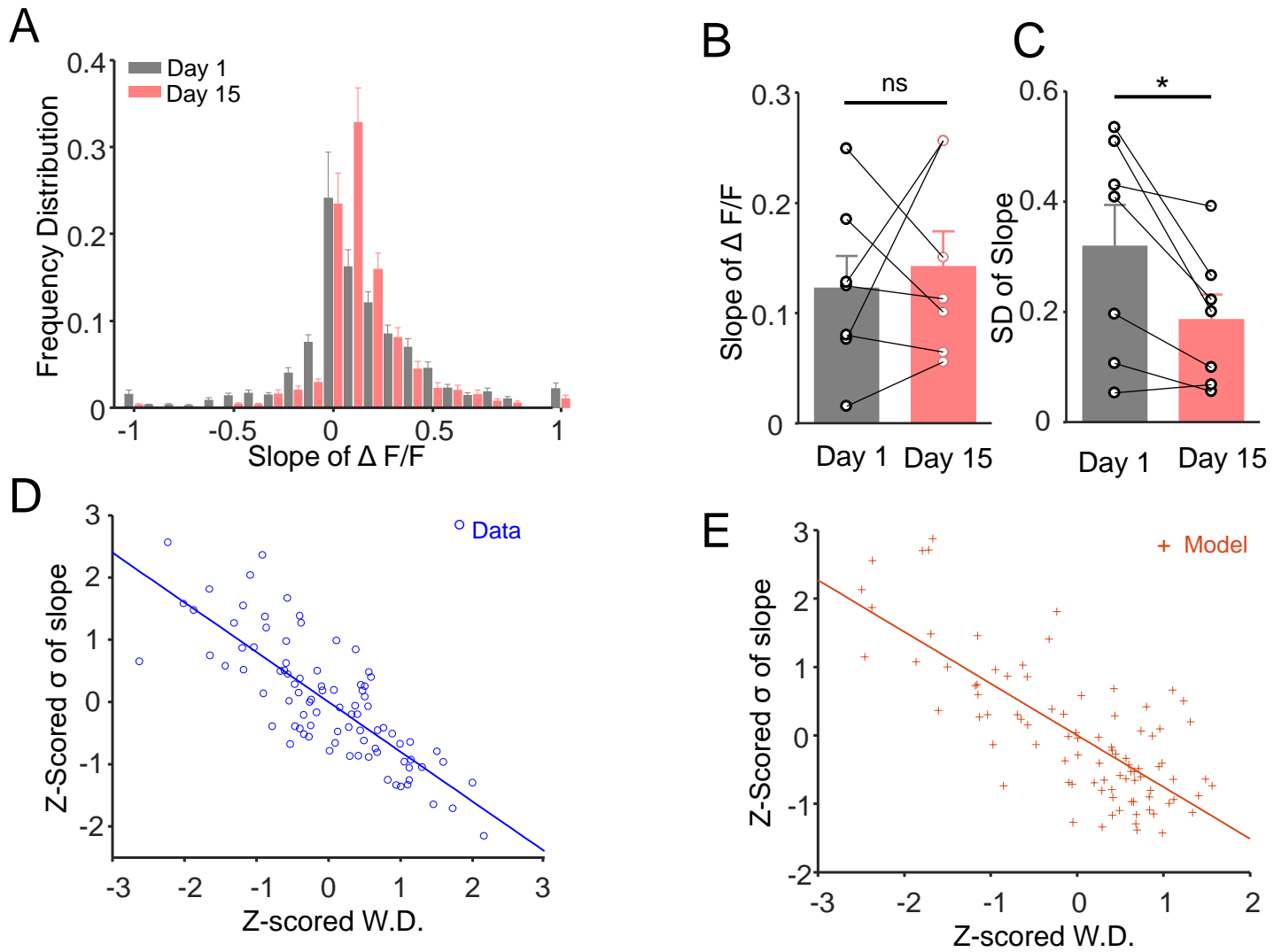


Fig. S4. The ramping dopamine activity became more stable with delayed gratification training. (A) The Frequency distribution of the slope of Z-scored $\Delta F/F$ during waiting (gray: day 1, red: day 15, $n = 7$). (B) The slope of $\Delta F/F$ had no difference between day 1 and day 15 ($p=0.63$, paired Student's t-test). (C) The standard deviation(σ) of the slope of $\Delta F/F$ in day 15 (0.32 ± 0.07) was significantly decreased than that in day 1 (0.17 ± 0.05 , $p=0.03$, paired Student's t-test). (D-E) The z-scored σ of $\Delta F/F$ and z-scored waiting durations were negatively correlated both in the experimental data (d, blue, $r = -0.80$, $p < 0.001$) and RL model (e, red, $r = -0.76$, $p < 0.001$). W.D. is for waiting duration. All error bars represent the s.e.m..

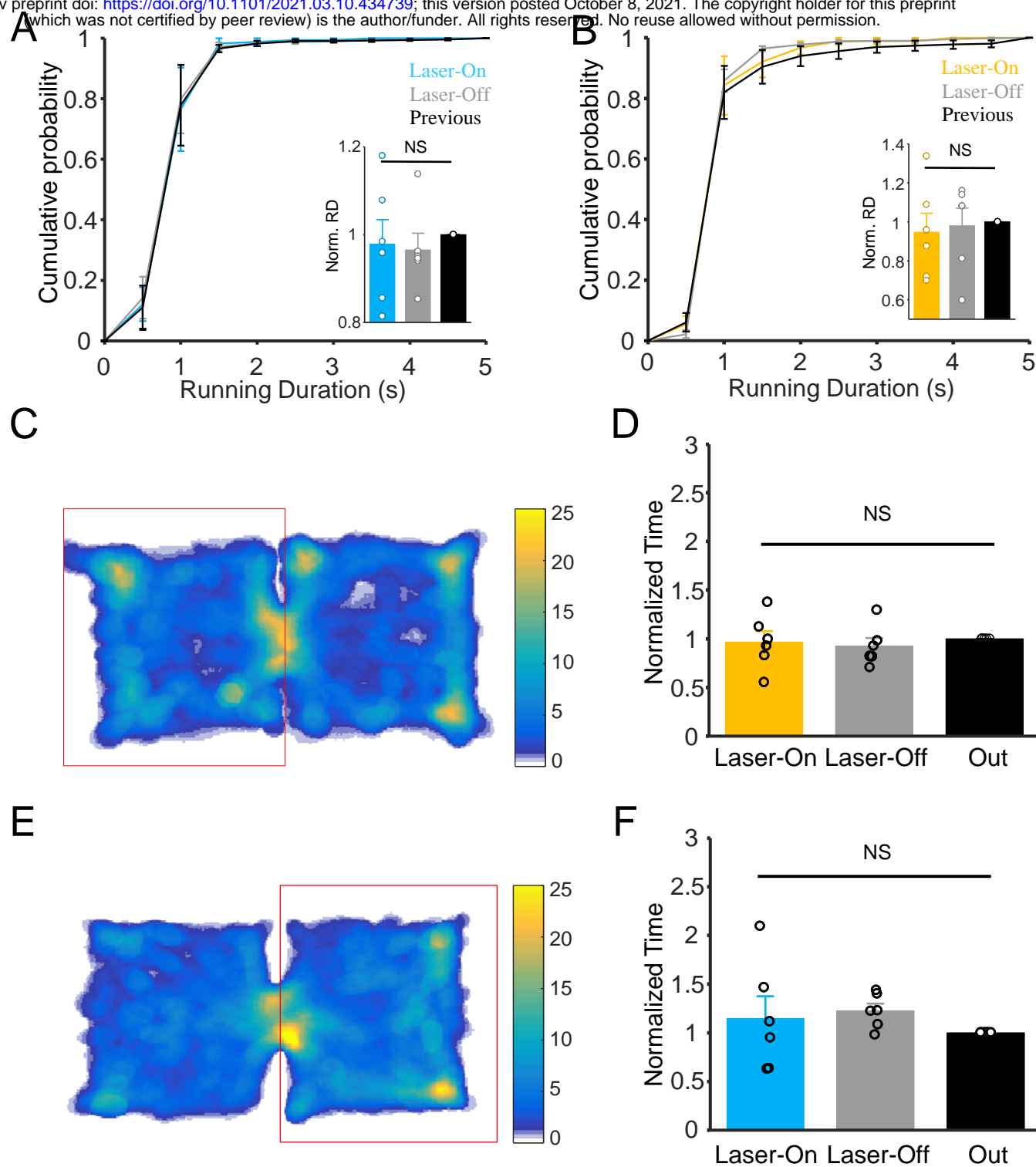


Fig. S5. Optical activation or inhibition of VTA DA activity didn't change the motivation of mouse in the delayed gratification tasks and the duration the mouse stayed in the given box in RPPT. (A) The running duration didn't change while the VTA DA neurons were optically activated ($F=0.20$, $p=0.82$, one-way ANOVA). (B) The result was the same as Figure A while optical inhibiting the VTA DA neurons ($F=0.12$, $p=0.88$, one-way ANOVA). (C) The heatmap of mouse traces in RPPT in which the VTA DA neurons were optically activated pseudo-randomly in 20% probability while the mouse entered into the given box (red rectangle). (D) The Z-scored duration that the mouse stayed in the given box while the VTA DA neurons were activated (Laser-On In) had no significant difference ($F=0.75$, $p=0.44$, one-way ANOVA, $n=6$) with the uninhibited durations (Laser-Off In) and durations in another box (Out). (E) The heatmap of mouse traces as shown in Figure C while inhibiting the VTA DA neurons in the given box (red rectangle). (F) Optical inhibiting the VTA DA neurons also didn't change the duration the mouse stayed in the given box ($F = 0.17$, $p = 0.73$, one-way ANOVA, $n=6$). All error bars represent the s.e.m..

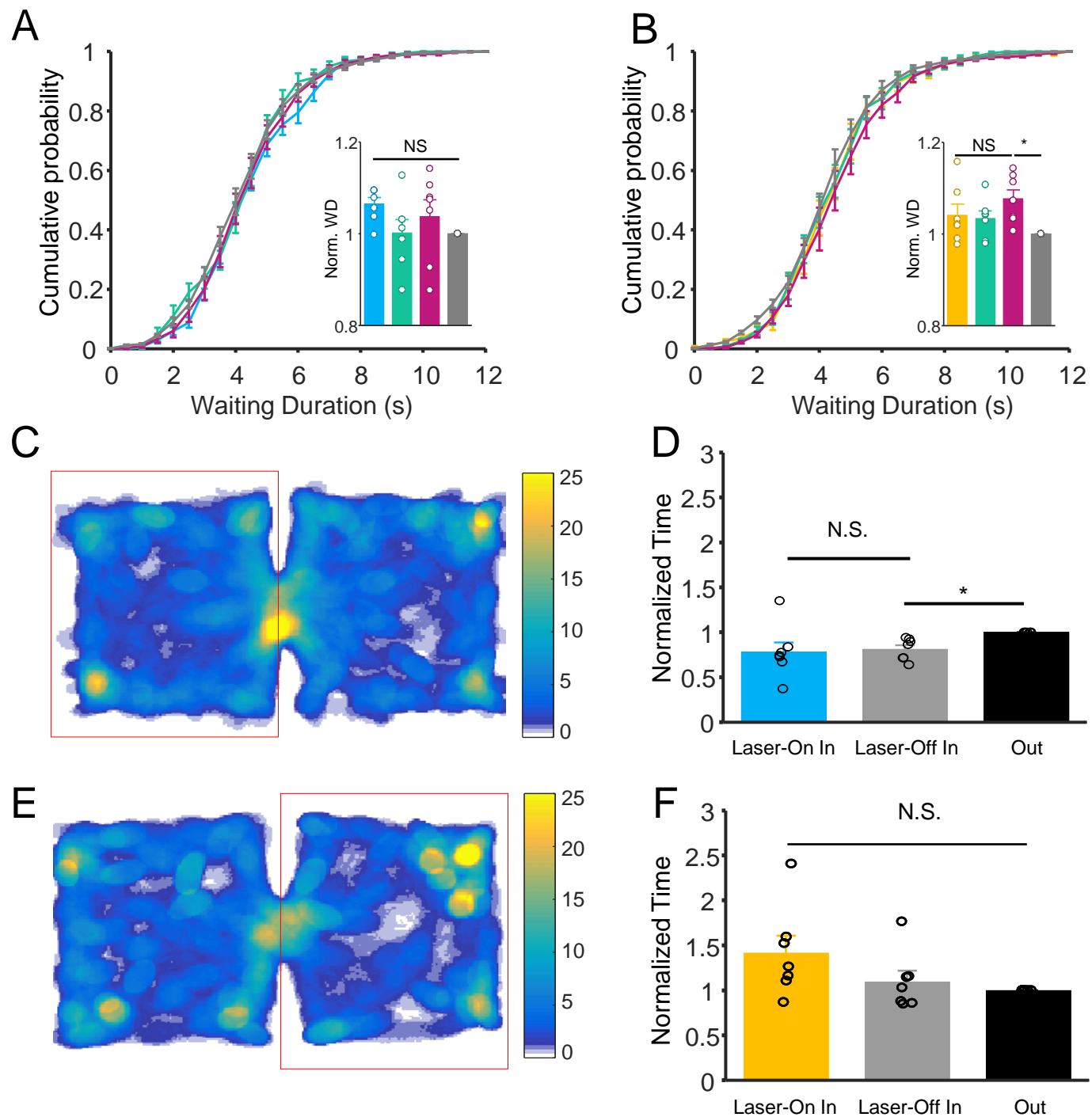


Fig. S6. Optogenetic manipulation of DAT-Cre mouse expressed mCherry in the delayed gratification tasks and RPPT.

(A) Waiting durations in 473nm laser delivered trials (blue) are not different compared with those of all other trials ($p=0.17$, Friedman test, $n=7$). (B) Waiting durations of 589nm laser un-delivered trials (magenta) slightly increased compared with the waiting duration of the previous day ($p=0.02$, Friedman test, $n=7$). (C) Heat-map of mouse traces in RPPT in which the VTA of the mouse was delivered 473nm laser pseudo-randomly in 20% probability while the mouse entered into a randomly chosen box (red rectangle). (D) Mean durations that the mouse stayed in chosen box while the laser delivered (Laser-On In), laser off (Laser-Off In) and the other box. There is no significant difference in waiting duration between Laser-On In and Laser-Off In ($F=3.54$, $p=0.09$, one-way ANOVA, $n=7$). (E) Heatmap of mouse traces same as shown in (c) while the mouse was delivered 589nm laser in a randomly chosen box (red rectangle). (F) Mean durations that the mouse stayed in chosen box while 589nm laser delivered (Laser-On In), laser off (Laser-Off In), and in the other box. 589nm laser delivering to mCherry mouse didn't alter waiting duration mouse stayed in any boxes under all experimental conditions ($F=2.64$, $p=0.14$, one-way ANOVA, $n=7$). All error bars represent the s.e.m..

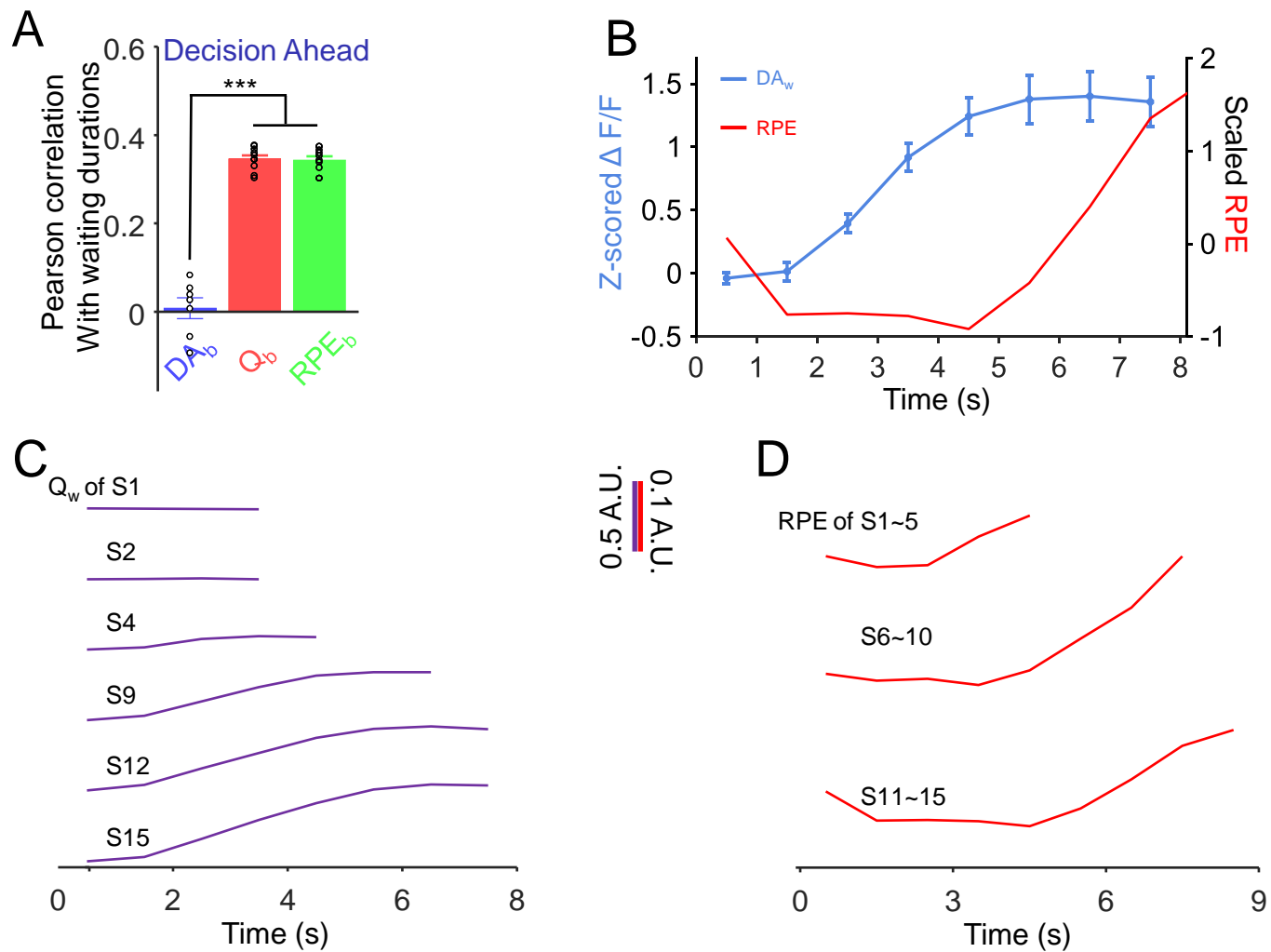


Fig. S7. The Data from RL models. (A) The correlation coefficient of mean DA activity 1 s before waiting (DA_b), the value of action (Q_b), and RPE of action (RPE_b) before waiting in the Decision Ahead model, with waiting durations. The correlation of V_b ($r = 0.35 \pm 0.01$, $p = 0.001$, $n = 10$, Pearson correlation) and RPE_b ($r = 0.34 \pm 0.01$, $p = 0.001$, $n = 10$, Pearson correlation) with waiting duration were significantly ($p < 0.001$, Kruskal-Wallis test) higher than the CC of DA_b ($r = 0.01 \pm 0.02$, $p = 0.36$, $n = 7$) with waiting durations. (B) Plots of Z-scored $\Delta F/F$ values (DA_w , light blue) at 0.5s before the waiting ended and RPE of waiting (RPE_w). There were no significant correlation between DA_w and RPE_w ($r = 0.34$, $p = 0.41$, Pearson correlation). (C-D) Value of waiting (Q_w) (C) and RPE (D) changed in the Continuous Deliberation model. All error bars represent the s.e.m..

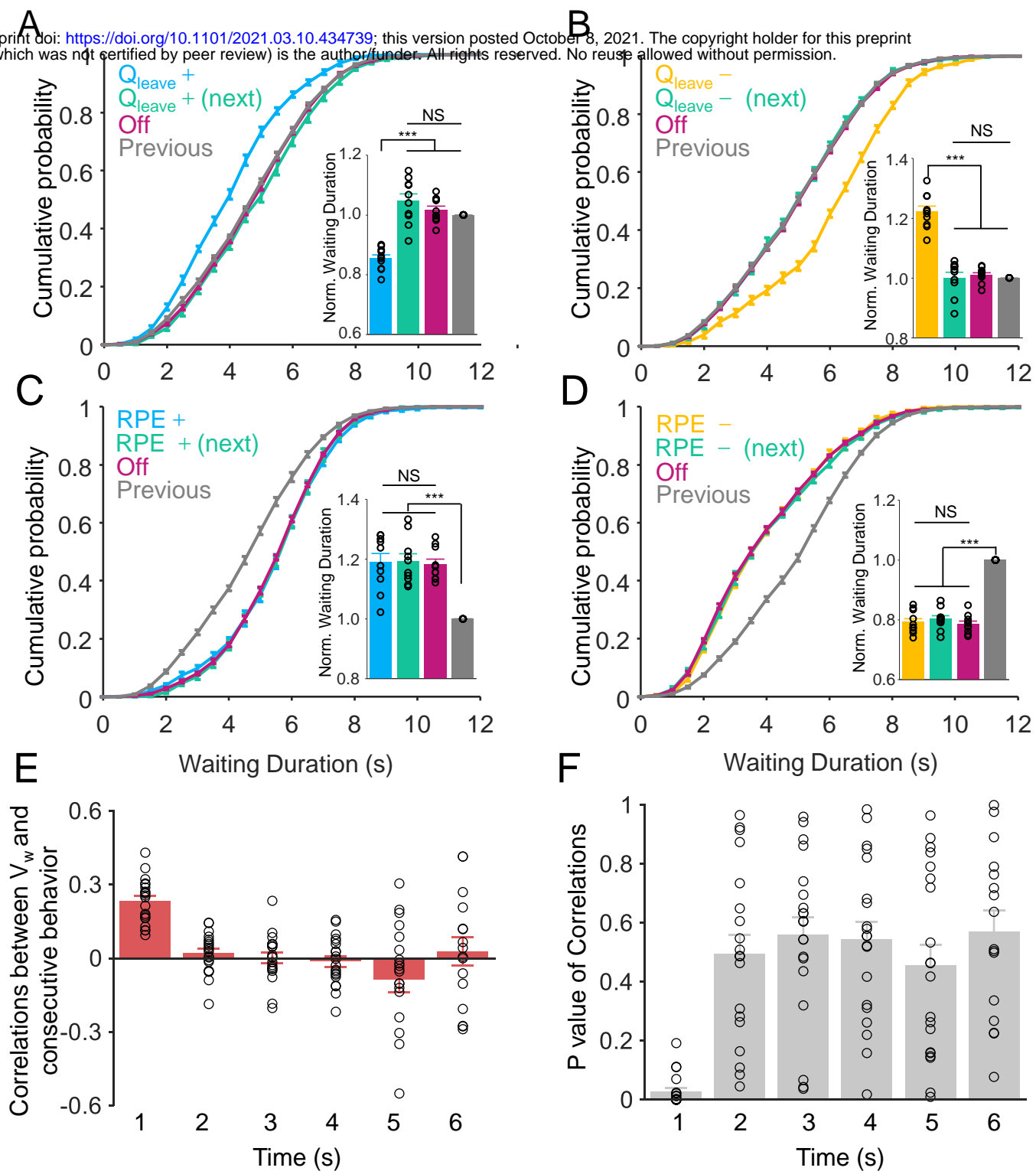


Fig. S8. Manipulation of value of leaving and RPE in RL model. (A-B) Either increasing or decreasing the value of leaving (Q_{leave}) in the continuous deliberation model as with the manipulation of Q_{wait} induced the opposite results compared with the optogenetics manipulating DAergic activity (increasing Q_{leave} in **A**: $p < 0.001$, Friedman test, $n = 10$; decreasing Q_{leave} in **B**: $p < 0.001$, Friedman test, $n = 10$) and had no influences on other trials (A-B, $p > 0.999$, Friedman test, $n = 10$). (C-D) Either increasing (C) or decreasing (D) the RPE in the continuous deliberation model as with the experimental data, alters the waiting durations in the same direction in all trials, whether or not the RPEs-manipulation (increasing RPEs in **C**: $p < 0.001$, Friedman test, $n = 10$; decreasing RPE in (D): $p < 0.001$, Friedman test, $n = 10$). (E) The value of waiting is only positively correlated (0.23 ± 0.02 , $p = 0.03 \pm 0.01$, $n = 20$) with the adjacent behavior in the *Continuous Deliberation* RL model. (F) The p values of correlation coefficients in E. All error bars represent the s.e.m..