# Emergence of rapid value inference through meta-reinforcement learning

Jaeeon Lee[1,2*], Jay A. Hennig[2,3,4,5], Vanessa Frelih[1,2], Samuel J. Gershman[2,3], Naoshige Uchida[1,2*]

[1]Department of Molecular and Cellular Biology, Harvard University, Cambridge, MA, USA

[2]Center for Brain Science, Harvard University, Cambridge, MA, USA

[3]Department of Psychology, Harvard University, Cambridge, MA, USA

[4]Department of Neuroscience, Baylor College of Medicine, Houston, TX, USA

[5]Neuroengineering Initiative, Rice University, Houston, TX, USA

[*]Correspondence: uchida@mcb.harvard.edu, jaylee@g.harard.edu

The ability to estimate the value associated with a specific stimulus or action is essential for adaptive behavior. Value can be updated either incrementally through experience or rapidly by inference based on latent environmental structure. Yet, how the brain implements and transitions between these modes of value computation remains unclear. To address this question, we examined the neuronal mechanisms underlying reversal learning. Mice were trained in an odor-outcome association task either with stable or dynamically changing contingencies. Mice trained on stable contingencies formed long-term value representations that depended on synaptic plasticity in the basolateral amygdala (BLA). In contrast, mice exposed to repeated reversals acquired the ability to infer values, independent from plasticity in BLA, enabling faster learning but with more rapid memory decay. Recurrent neural network models (RNNs) trained with continuous weight updates recapitulated this transition, shifting from plasticity-based to dynamics-based value computation. Neural activity in the BLA encoded both value and contextual information necessary for computing value based on latent task structure, similar to those found in the RNNs. Disrupting BLA activity before cue delivery preferentially impaired dynamics-based value updating. Furthermore, mice could learn distinct correlation structures that enabled structure-specific value inference. Together, these findings provide a mechanistic framework for fast value updates via inference, a core feature of intelligent behavior.

30

**Main**

Animals must continually estimate the value of sensory cues and actions to guide adaptive behavior. Reinforcement learning provides an algorithmic framework for this process[1–3], yet the neural mechanisms by which value is learned, stored, and flexibly updated remain unresolved. Classical theories emphasize incremental trial-and-error learning, in which changes in synaptic strength encode long-term value memories[4–11]. Indeed, plasticity within circuits such as the striatum and amygdala has been linked to the formation of stable appetitive and aversive associations.

In contrast, accumulating evidence suggests that animals can update value through inference by exploiting structural knowledge of the environment[12–21]. In RNNs, meta-reinforcement learning — the emergence of a fast reinforcement learning algorithm through a slower reinforcement learning algorithm — gives rise to inference-like behavior[22–24]. In these pre-trained models, value can be updated without online synaptic changes, through recurrent dynamics that encode hidden task states. Although both incremental learning and inference are observed behaviorally, the neural mechanisms supporting each — and how the brain transitions between them — remain elusive. A major challenge is that behavioral performance alone often cannot distinguish between plasticity-based and dynamics-based value updating, and the two mechanisms may coexist within the same circuit[25–27].

To address this, we developed two classical conditioning paradigms that differed in outcome stability. Combining behavior, electrophysiology, and computational modeling, we find that rapid value inference emerges through a gradual transition from synaptic plasticity-dependent learning to plasticity-independent value updating mediated by recurrent dynamics that encode task structure.

**Timescale of value update/decay in stable vs dynamic task**

We trained mice to perform a head-fixed classical conditioning task in which an odor cue (conditioned stimulus, CS+ or CS-) was followed by either a water reward or no reward (**Fig. 1a**). We used two versions of the task. In the stable task, the reward contingency was fixed, and mice had to initially learn a fixed value and maintain the value memory in the subsequent sessions (**Fig. 1b**, *top*). In the dynamic task, the reward contingency reversed every session with the reversal

61 happening in the middle of each session (**Fig. 1b**, *bottom*). Each session started with the reward
62 contingency from the previous session so that mice had an incentive to remember previously
63 learned values. We used anticipatory licking during the odor period as a proxy for value to quantify
64 performance on both tasks. In the stable task, mice quickly learned to lick more to the CS+ odor
65 than to the CS- odor from day 1 of training, reaching expert performance by day 3 (**Fig. 1c,**
66 **Extended Data Fig. 1a, c**). In the dynamic task, mice gradually improved their performance over
67 12 days of training, after which expert mice reliably discriminated CS+ and CS- both before
68 reversal (block 1) and after reversal (block 2), achieving comparable performance as in the expert
69 in the stable task (**Fig. 1d,e, Extended Data Fig. 1b, c**).

70      To better understand how value updating at the beginning of the stable task and late stage
71 of the dynamic task differed, we quantified the learning curve for CS+ or CS- lick rate (**Fig. 1f**).
72 For the stable task, we plotted the CS+/CS- lick rate on the 1$^{st}$ day of the stable task or the 1$^{st}$
73 reversal session (1$^{st}$ dynamic task after being trained on the stable task). This was compared to the
74 CS+/CS- lick rate aligned to the reversal point in expert mice in the dynamic task. For both positive
75 value update (+Δvalue) and negative value update (−Δvalue), learning rate was much faster in the
76 dynamic vs stable task (for positive update, $\tau_{stable}$=80.5 trials; $\tau_{dynamic}$=2.4 trials; for negative update,
77 $\tau_{stable}$=52.3 trials; $\tau_{dynamic}$=8.1 trials). Overall, these results suggest that as mice transition from a
78 stable to a dynamic environment, the timescale of value updates, as measured by the conditioned
79 response (e.g. anticipatory licking), gets faster by an order of magnitude.

80      One interesting possibility is that mice transition from a plasticity- to a dynamics-based
81 value updating strategy as they move from the stable to dynamic environment. Dynamics-based
82 value updating may require mice to maintain information about value via persistent activity, which
83 might be prone to temporal degradation akin to working memory[28–30]. We thus reasoned that value
84 memory might degrade at a distinct timescale in the stable vs. dynamic tasks. To test this
85 hypothesis, we first quantified how fast value memory degrades over time in the stable or dynamic
86 task. In both tasks, the reward contingency at the beginning of each session was the same as that
87 of the end of the previous session (**Fig. 1g**). Thus, we asked if mice correctly discriminated CS+
88 and CS- odors on the very first trial in each session. Mice were fully trained on either the stable
89 task (5 days of training) or dynamic task (12 days of training) and then tested for value memory
90 by quantifying the number of licks during the first trial of each cue. In the stable task, mice
91 performed correctly from the first trial by licking to the first CS+ and not the CS- (**Fig. 1g**, stable

92    task). However, in the dynamic task, mice started the session by licking to both CS+ and CS- even

93    though they had discriminated CS+ and CS- at the end of the previous session (**Fig. 1g**, dynamic

94    task). Mice then quickly learned to suppress licking to the CS-, allowing them to still maintain

95    high discriminability between CS+ and CS- throughout the block (**Extended Data Fig. 1d**). To

96    quantify this "forgetting" effect more systematically, we computed a discrimination index for each

97    session, with index=1 being perfect discrimination between first CS+ and first CS-, index=0 being

98    no discrimination, and index=-1 being the flipped discrimination (see Methods). We also varied

99    the duration of the break between each session of the stable task by pausing behavioral training for

100    up to 8 days. In the stable task, the discrimination index was close to 1 even with 8 days of break

101    whereas in the dynamic task, the index was close to zero even after a 1-day break (**Fig. 1h**). Overall,

102    these results suggest that value memory in the stable task persists over 8 days, whereas value

103    memory degrades to chance level after only 1 day in the dynamic task.

104          To further quantify the timescale at which dynamic value memory degrades, we introduced

105    an inter-trial interval (ITI) that was longer than the average duration that mice were initially trained

106    on (ITI=37.5-300 sec) in the middle of each session in both the stable and dynamic task (two long-

107    ITI trials per session) (**Fig. 1i**, *left*). These extended ITIs caused a similar phenotype as the break

108    between sessions: value memory decayed to chance level in the dynamic task after 300 seconds

109    but not in the stable task (**Fig. 1i**, *right,* **Extended Data Fig. 1e**). The discrimination index was

110    significantly lower for dynamic vs stable tasks for duration of 150 and 300 seconds. As an

111    alternative method to measure value, we expressed a dopamine sensor ($GRAB^{DA3m}$) in the ventral

112    striatum and quantified the dopamine cue response as a measure for value (**Extended Data Fig.**

113    **1f-h**). The discrimination index for the dopamine response was similar to the discrimination index

114    for the CS lick rate, suggesting that both behavioral and neural readout of value display a similar

115    forgetting timescale between the two tasks.

116

117    **RNN model with continual plasticity**

118    The above results suggest that mice might initially store value information in synaptic weights

119    when the environment is stable, but repeated exposure to value reversals might cause a transition

120    from a plasticity-based value update to a dynamics-based value update, allowing faster value

121    updating at the expense of being more forgetful. To better understand how this transition might

122    occur mechanistically, we took a computational modelling approach. Previous works using the

123    temporal difference (TD) learning algorithm have been successful at explaining various features

124    of dopamine responses and value learning[4,31–33]. However, most previous applications of TD to

125    animal learning *a priori* assumed a specific state representation that is fixed and thus does not

126    model the emergence of the state representation itself. On the other hand, more recent work has

127    shown that TD learning models equipped with RNNs can learn useful representations directly (e.g.

128    beliefs about hidden states), mirroring activity seen in actual neural recordings[22,25,34,35]. However,

129    many of these models are trained offline. This implies that plasticity in these models cannot play

130    a causal role in updating value online—in stark contrast to classical TD modeling approaches,

131    where value is updated online but exclusively via plasticity. Thus, to create a biologically realistic

132    model of learning that is both agnostic about state representations and also displays online

133    plasticity updates, we trained RNNs with TD learning in an online fashion (**Fig. 2a**, see Methods),

134    where at every timestep, the RNN's weights were updated based on the recent history of inputs.

135    Unlike previous methods using RNNs to model value learning, weights in the RNN were never

136    frozen, allowing for a continual interaction between recurrent dynamics and plasticity.

137         We trained RNNs on either the stable task or dynamic task (**Fig. 2b-c, Extended Data Fig.**

138    **2a-c-d**; see Methods). The RNNs' value readouts could reliably discriminate between cue A and

139    cue B (**Fig. 2b**). Interestingly, in the dynamic task, the learning rate for updating value displayed

140    an abrupt transition from a slow update regime to a fast update regime (**Fig. 2b**, right). We

141    quantified the learning rate for positive or negative value updates in the stable and dynamic tasks

142    (**Fig. 2c**). Similar to the mouse behavioral data (**Fig. 1f**), RNNs displayed faster value updates in

143    the dynamic task compared to the first reversal in the stable task (**Fig. 2c**). To better understand

144    the mechanisms driving value updates in these two conditions, we manipulated plasticity in the

145    RNNs by setting the learning rate of the RNN weight update to zero (**Fig. 2d**, see Methods).

146    Without plasticity, RNNs were completely impaired at updating value in the stable task, whereas

147    RNNs were still able to update value in the dynamic task (**Fig. 2d**, *right*). We quantified the

148    difference between the value readouts of CS+ and CS- cues in RNNs with plasticity ($\Delta w \neq 0$) or

149    without plasticity ($\Delta w = 0$), for RNNs trained on either the stable or dynamic task (**Fig. 2e**). Without

150    plasticity, the difference between predicted values for CS+ and CS- cues decreased in the stable

151    task but not in the dynamic task. Overall, these findings suggest that RNNs with online weight

152    updating initially use plasticity to learn value in the stable task, but transition to a plasticity-

153    independent mechanism in the dynamic task.

154    What could be the mechanism driving value updating in the dynamic task? To answer this
155    question, we applied principal component analysis (PCA) to the activity in the RNNs (see Methods)
156    and plotted the neural-state space trajectories with an additional axis representing the value
157    predicted by the RNNs (**Fig. 2f, g**). In the stable task, neural trajectories for CS+ gradually changed
158    so that the RNN's response to the cue (green point) moved upwards towards higher value,
159    representing plasticity-driven value updating (cue A, **Fig. 2f**). In the dynamic task, neural
160    trajectories for an expert RNN were segregated by cue type and block type, with PC1 encoding
161    cue type and PC2 encoding block type (**Fig. 2g**, *left*). Information about block identity, or context
162    information, could potentially be used in the RNN to decode value information for each block (**Fig.
163    2g**, *right*). To see if the contextual information present in the RNN activity was driving value
164    computation in the RNNs, we computed the context axis, which was defined as the linear
165    discriminant axis that best separated context information in an expert RNN (see Methods). We
166    then projected the hidden units' activity during the ITI onto the context axis (context$^{proj}$), and
167    plotted it along the value readout for cue A and cue B (**Fig. 2h**). When value was updated slowly
168    during the initial phase of learning, context$^{proj}$ did not discriminate block identity very well.
169    However, as the RNN transitioned into a faster value update regime, context$^{proj}$ started to
170    discriminate block identity. We calculated the Spearman correlation between the difference in
171    value (value diff.) on the last trial to cue A and cue B, and context$^{proj}$ in either the naïve or expert
172    RNN (**Fig. 2i**). The correlation was larger in expert RNNs compared to naïve RNNs (**Fig. 2j**, *left*).
173    Furthermore, the activity in different context (context A or B) projected onto the discriminant axis,
174    became more distinct in expert mice, suggesting that context information became more discrete
175    and separable, and potentially indicating the emergence of fixed points corresponding to each
176    context. Overall, these results suggest that in the dynamic task, RNNs gradually develop a
177    representation of block identity (e.g., using fixed points) which allow the RNNs to rapidly update
178    the value of each cue across blocks through neural dynamics (i.e. without plasticity).

179    Lastly, we simulated the effect of introducing a long ITI during the stable or dynamic task.
180    Given that context information encoded in the hidden units' activity is potentially important for
181    computing value in the dynamic task, we reasoned that a long ITI might cause a drift in the activity
182    in the RNN, resulting in an ITI duration-dependent change in value discrimination. Consistent with
183    this prediction, a long ITI caused a change in context$^{proj}$ and a subsequent change in value readout
184    in the RNN (**Fig. 2k**). We quantified the discrimination index similarly to experimental data in

185 **Fig. 1h, j**. We found that a long ITI caused a duration-dependent change in discrimination index
186 only in the dynamic task, recapitulating experimental data (**Fig. 1j**). This effect could be explained
187 by the magnitude of the drift along the context axis, with larger drift causing a larger change in
188 value difference (**Fig. 2m**). Overall, these results suggest that contextual information, represented
189 by the population activity, is prone to temporal degradation, thus causing a time-dependent
190 degradation of value memory in the dynamic task.

191

192 **The role of plasticity and activity in the BLA**
193 Computational modelling with RNNs suggested that a transition from plasticity to dynamics-based
194 value update could explain the experimental data. This implies that blocking synaptic plasticity in
195 the region of the brain responsible for value updating should impair performance in the stable task
196 on day 1, but not performance in the dynamic task after being fully trained. We tested this
197 prediction by targeting the BLA, a brain region that has been previously shown to undergo synaptic
198 plasticity in both appetitive and aversive tasks as well as showing activity correlated with
199 conditioned responding[7,11,36–39]. To block synaptic plasticity acutely, we injected a CaMKII
200 blocker (KN-93) locally in the BLA and tested the performance on either day 1 of the stable task
201 or on expert stage of the dynamic task (**Fig. 3a, Extended Data Fig. 1a-e**). KN-93 has been shown
202 to effectively block synaptic plasticity in slice, and local infusion *in vivo* has been shown to cause
203 behavioral effects consistent with an impairment of synaptic plasticity[40–44]. KN-93 infusion in the
204 BLA significantly decreased the difference in anticipatory licking between CS+ and CS- trials in
205 the stable task compared to saline infusion (**Fig. 3b**, *left*). In stark contrast, the same manipulation
206 had no effect in the dynamic task (**Fig. 3b**, *right*). These results suggest that BLA plasticity is
207 necessary to initially update value in the stable task, but becomes dispensable in the dynamic task.
208 One alternative explanation for the dissociable role of BLA plasticity in the stable vs
209 dynamic tasks is that BLA might become disengaged in the dynamic task, with another brain
210 region taking over the role of BLA. Thus, plasticity in another brain region other than BLA might
211 still be responsible for updating value in the dynamic task (**Extended Data Fig. 3f**). This
212 alternative model would predict that BLA activity is only necessary to perform the stable task and
213 not the dynamic task. To test if this was true, we acutely inactivated BLA activity. We generated
214 an emx1-Cre × gtACR1 mice, in which the inhibitory opsin is expressed in a Cre-dependent
215 manner in the excitatory neurons throughout the brain including BLA. BLA specificity was

216　achieved by implanting an optical fiber just above BLA in emx1-cre × gtACR1 mice. Mice were

217　trained in either the stable or dynamic task, after which BLA was inactivated on 15% of all trials

218　during the cue period for 3 seconds (**Fig. 3c**). Inactivating BLA neurons impaired performance in

219　both stable and dynamic tasks, by increasing the CS- licks, resulting in poorer discrimination of

220　CS+/CS- (**Fig. 3d**). Thus, BLA activity is still necessary for performing the dynamic task, despite

221　plasticity in BLA becoming dispensable.

222

223　**Value coding in the amygdala**

224　To better understand the nature of value coding in the BLA, we performed acute high-density

225　electrophysiological recording using Neuropixels probes in the BLA (**Fig. 4a, b**). We modified the

226　original task so that both stable and dynamic value could be measured in the same recording

227　session (**Fig. 4a**). The hybrid task consisted of 3 odor cues, two of which were stable odor cues

228　(odor A=reward, odor B=no reward), and one of which was a dynamic odor cue (odor C=reward

229　or no reward depending on block). The value of odor C changed throughout the session across 4

230　blocks. In all sessions, reward contingency for odor C started with the same condition as the

231　previous session to ensure continuity of value (**Fig. 4a**, *bottom*). Consistent with the timescale of

232　value memory decay described previously (**Fig. 1h**), mice forgot the value of odor C between

233　sessions, as indicated by mice always licking to odor C at the beginning of the session (**Extended**

234　**Data Fig. 4**). This behavior was not present in odor B, suggesting that value memory for stable

235　and dynamic odors had distinct timescales of memory decay in the hybrid task.

236　　　　　Analysis of neural recording data revealed that BLA units encoded both stable and dynamic

237　values, with many units encoding both (**Fig. 4c**). An example stable value coding (SV) unit in the

238　BLA consistently differentiated odor A and odor B across blocks (**Fig. 4c**, top). An example

239　dynamic value coding (DV) unit in the BLA responded to odor C differentially depending on its

240　value in that block (**Fig. 4c**, middle). An example unit in the BLA encoded both SV and DV

241　(SVDV, **Fig. 4c**, bottom). Given that Neuropixels recording allowed wide sampling of brain

242　regions surrounding BLA, we analyzed the fraction of SV, DV and SVDV units across all recorded

243　brain regions (**Fig. 4d**). SV, DV and SVDV units were enriched in the amygdala, especially in the

244　BLA and BMA (BLA: SV fraction=35%, DV fraction=13%, SVDV fraction=8%; BMA: SV

245　fraction=34%, DV fraction=9%, SVDV fraction=6%). We defined the polarity of stable value or

246　dynamic value by whether firing rate was higher for the rewarded cue or not (see Methods). We

247    plotted the mean firing rate across all positive SV, negative SV, positive DV or negative DV in

248    the amygdala (**Fig. 4e**). Interestingly, positive SV and DV units tended to fire phasically to both

249    odors (**Fig. 4e**, left panels). In contrast, negative SV and DV units tended to be bidirectionally

250    modulated relative to baseline (**Fig. 4e**, right panels). To see if the polarity of SV and DV was

251    congruent in SVDV units, we plotted the SV selectivity vs DV selectivity for all SVDV units (**Fig.**

252    **4f**, left). Selectivity was defined as the difference in firing rate between the rewarded cue/block

253    and non-rewarded cue/block (SV selectivity=odor A- odor B; DV selectivity=odorC$^{reward}$ –

254    odorC$^{noreward}$). We found a strong positive correlation between SV selectivity and DV selectivity

255    (**Fig. 4f**). The fraction of units that had congruent polarity for SV and DV was higher than units

256    that had incongruent polarity (**Fig. 4f**, right). Overall, these results suggest that amygdala contains

257    neurons that encode both SV and DV in a congruent manner.

258

259    **Context coding in the amygdala**

260    BLA can compute and store stable value using plasticity (**Fig. 3a, b**) but requires context

261    information to compute dynamic value (**Fig. 2h-j**). We considered two hypotheses for how BLA

262    could acquire dynamic value. In one scenario, dynamic value is computed outside BLA and

263    inherited by BLA (model 1, **Extended Data Fig. 5a**). In another scenario, BLA locally computes

264    dynamic value using context information within BLA (model 2, **Extended Data Fig. 5a**). To

265    distinguish between these two scenarios, we looked for cells in the BLA that could differentiate

266    block types during the ITI. To our surprise, we found many units in the BLA that contained context

267    information during the ITI period (**Fig. 5, Extended Data Fig5. b)**. An example unit in the BLA

268    differentiated context by firing higher during the ITI period of non-rewarded blocks (**Fig. 5a-b**).

269    This unit's firing rate was negatively correlated with the odor C CS licks, suggesting that ITI firing

270    rate was predictive of upcoming anticipatory licking or value for odor C (negative context unit).

271    Out of all context units in the amygdala, we found that a large fraction of them also encoded stable

272    and dynamic value, suggesting that context and value information might be multiplexed in the

273    amygdala (**Extended Data Fig. 5c**). We analyzed all the brain region recorded and found that

274    BLA and BMA were enriched with context coding units, with negative context units being more

275    dominant than positive context units (BLA: 5.7%; BMA: 4.2%) (**Fig. 5c**). To test if the ITI firing

276    rate was predictive of the anticipatory CS licks to odor C, we computed the Spearman's correlation

277    coefficient between odor C CS licks and the ITI firing rate before cue onset for each neuron (**Fig.**

278   **5d**). Importantly, we restricted our analysis to contain only one block type to avoid circular

279   analysis (see Methods).  We found that positive or negative context coding units had a statistically

280   significant correlation coefficient, suggesting that the ITI activity of context units predicts

281   upcoming behavioral choice for dynamic odor C (**Fig. 5e**). Lastly, given that context coding is

282   confounded by reward rate coding in this task, we asked if ITI firing rate was updated in a manner

283   consistent with pure context coding or pure reward rate coding. We reasoned that if context coding

284   is prevalent, then ITI firing rate should only be updated after dynamic odor C, and not after stable

285   odors A or B. However, if reward rate coding is prevalent, then ITI firing rate might be updated

286   similarly to all outcomes regardless of cue type (**Extended Data Fig. 5d**). Context update analysis

287   revealed that change in ITI firing rate was more consistent with context coding than reward rate

288   coding (**Fig. 5f**). This was especially true for negative context units, which was most of the context

289   coding units. To see if this context coding was causal to value computed using dynamics, we

290   trained emx1-Cre × gtACR1 mice on either stable or dynamic mice, after which BLA was

291   bilaterally inactivated during the ITI (**Fig. 5g**). We reasoned that if context information during the

292   ITI is important for computing dynamic value, disrupting activity in this period should impair

293   performance during the dynamic task, but not during the stable task where value is not computed

294   using dynamics. BLA inactivation during the ITI did not change the performance in the stable task

295   but increased CS- licks in the dynamic task (**Fig. 5h**). When we computed the difference between

296   the mean CS+ licks and CS- licks, this difference was significantly reduced in the dynamic task

297   but not in the stable task (**Fig. 5i**). Overall, these results suggest that BLA contains context

298   information which not only predicts upcoming behavior, but is also necessary for performing in

299   the dynamic task.

300

**Dynamics allow value inference**

301

302   We have shown so far that value updating via dynamics can be fast but forgetful compared to value

303   updating via plasticity. Another advantage of using dynamics over plasticity is that dynamics-

304   based value updates support value inference, which is the ability to update value without direct

305   experience. This is because RNNs naturally learn the structure of the task and encode the hidden

306   states of the task (**Fig. 2**). Given that in the dynamic task, there are essentially two hidden states

307   corresponding to each block, RNNs using dynamics can infer value for one cue without direct

308   experience (counterfactual learning; **Fig. 6a**).  To look for evidence that our RNNs can perform

309    value inference, we designed a paradigm in which, after reversal, we only presented one cue type

310    for up to 20 trials, after which we presented the opposite cue (probe cue) for the first time (F**ig. 6b**)

311    An agent that understands the structure of the task (i.e. anti-correlation in the values predicted by

312    odor A and B) would be able to infer the change in value for the probe cue, whereas an agent that

313    does not understand the structure of the task might not infer any change and start from the pre-

314    reversal point. We tested this in both RNNs and mice (**Fig. 6c-g, EDF 6a)**. In RNNs, value

315    inference only emerged in expert RNNs trained on the dynamic task (**Fig. 6c**). This is consistent

316    with the idea that a naïve RNN uses plasticity, whereas an expert RNN that uses dynamics alone

317    can infer value using its learned dynamics. Similarly in mice, we found that naïve mice undergoing

318    reversal for the first time failed to infer a change in value, whereas expert mice fully trained on the

319    dynamic task could infer a change in value (**Fig. 6d**). This effect was consistent across many RNNs

320    and mice as measured by CS licks (**Fig. 6e, f**, **EDF 6c, d**). If cue-evoked dopamine reports value,

321    then we should also expect dopamine signals to reflect inferred value change. Consistent with

322    previous studies[14,16,17], we also found that dopamine signals reflected inference (**Fig. 6f**, *right,*

323    **EDF 6b**). We also found that the magnitude of the inferred change in value ($\Delta$CS licks)

324    parametrically varied as a function of the number of opposite trials, with a larger number of trials

325    leading to bigger change in $\Delta$CS licks (**Fig. 6g, EDF 6c, d**).

326        Inference has been mostly studied using tasks implementing anti-correlated values between

327    two options[12–14,19]. To test if mice can learn distinct structures to guide structure-specific inference,

328    we implemented three correlation structures of odor values. We trained separate cohorts of mice

329    in environments in which a pair of odor-value associations are either anti-correlated, correlated, or

330    independent. We then tested for the presence of inference in a similar fashion as above (**Fig. 6g**).

331    Indeed, mice displayed signatures of inference consistent with the pre-exposed correlation

332    structure as measured based on anticipatory licks and cue-evoked dopamine responses (**Fig. 6h, i**).

333    For instance, when the value of odor A changes from positive to zero, mice trained in the anti-

334    correlated structure increased the value of odor B without directly experiencing the new

335    contingency for odor B, whereas mice trained in the correlated structure decreased the value of

336    odor B. Mice trained in the independent structure exhibited a negligible sign of inference. Overall,

337    these results suggest that a transition from plasticity- to dynamics-based value update allows rapid

338    value inference to emerge and is specific to the learned correlated structure.

339

**Discussion**

340 

341 Through a combination of behavioral, circuit-level, and computational approaches, we show that

342 value computation can transition from a synaptic plasticity-based to a recurrent dynamics-based

343 mechanism allowing value inference to emerge in a dynamic environment. Because contextual

344 representations were maintained during the ITI, the outcome of one cue updated the inferred value

345 of both cues simultaneously. Mice could also learn distinct correlation structures, suggesting

346 inference can be adaptive. Overall, our work provides a mechanistic framework for understanding

347 how rapid inference emerges in the brain.

348 Computationally, our results extend prior models showing that recurrent network dynamics

349 can encode value without synaptic change[23,25,26,45]. Previous studies typically used an offline-

350 learning rule in which the weights are only updated outside the task and then fixed after

351 convergence. Thus, these networks assume that plasticity does not play a role for incremental

352 online improvement in performance. Our model incorporates continuous online plasticity during

353 task performance, allowing a much more biologically realistic interaction between plasticity and

354 dynamics. Moreover, we show that a single learning rule—truncated backpropagation through

355 time (TBPTT)—can produce both plasticity-based and dynamics-based value updates, depending

356 on the timescale of the learning window. Shorter windows favor value update based on plasticity,

357 while longer windows enable meta-reinforcement learning through emergent recurrent dynamics

358 (**Extended Data Fig. 2f, g**). One interesting possibility is that distinct brain areas may implement

359 the same underlying rule with different effective timescales, enabling diverse computational

360 functions to emerge in different brain regions[46–48].

361 Behaviorally, we demonstrate that mice exploit environmental regularities in the task to

362 infer value, consistent with prior work showing across species that animals display inference

363 behavior and that regions like the hippocampus and OFC contribute to inference-based decision-

364 making[12–14,17,31,32,49,50]. Our findings highlight the amygdala's contribution to context-dependent

365 inference[51], raising the question of whether contextual representations are inherited from inputs

366 such as the ventral hippocampus and OFC, or computed locally within the amygdala.

367 Disentangling these possibilities will be essential for understanding how inference unfolds across

368 distributed neural circuits.

369 Finally, our results reveal a fundamental tradeoff between stability and flexibility in value

370 computation. Dynamics-based value representations enable rapid value update using inference at

371      the cost of degrading over time, whereas plasticity-based representations provide stable long-term
372      storage at the cost of being slow and inflexible. Given the further cost of having to maintain
373      information in the persistent activity, dynamics-based value update might only emerge when the
374      benefits of faster value update (e.g. environment is dynamic) is high. Elucidating the exact
375      conditions under which these two modes of value update dominate will be crucial in the future.

## References

1. Sutton, R. S. & Barto, A. G. *Reinforcement Learning: An Introduction*. (A Bradford Book, Cambridge, Mass, 1998).

2. Botvinick, M., Wang, J. X., Dabney, W., Miller, K. J. & Kurth-Nelson, Z. Deep Reinforcement Learning and Its Neuroscientific Implications. *Neuron* **107**, 603–616 (2020).

3. Lee, D., Seo, H. & Jung, M. W. Neural Basis of Reinforcement Learning and Decision Making. *Annual Review of Neuroscience* **35**, 287–308 (2012).

4. Schultz, W., Dayan, P. & Montague, P. R. A Neural Substrate of Prediction and Reward. *Science* **275**, 1593–1599 (1997).

5. Reynolds, J. N. J., Hyland, B. I. & Wickens, J. R. A cellular mechanism of reward-related learning. *Nature* **413**, 67–70 (2001).

6. Hong, S. & Hikosaka, O. Dopamine-Mediated Learning and Switching in Cortico-Striatal Circuit Explain Behavioral Changes in Reinforcement Learning. *Front. Behav. Neurosci.* **5**, (2011).

7. Tye, K. M., Stuber, G. D., de Ridder, B., Bonci, A. & Janak, P. H. Rapid strengthening of thalamo-amygdala synapses mediates cue–reward learning. *Nature* **453**, 1253–1257 (2008).

8. Yamaguchi, K. *et al.* The minimal behavioral time window for reward conditioning in the nucleus accumbens of mice. 43.

9. Xiong, Q., Znamenskiy, P. & Zador, A. M. Selective corticostriatal plasticity during acquisition of an auditory discrimination task. *Nature* **521**, 348–351 (2015).

10. Tsutsui-Kimura, I. *et al.* Dopamine in the tail of the striatum facilitates avoidance in threat–reward conflicts. *Nat Neurosci* **28**, 795–810 (2025).

11. Rogan, M. T., Stäubli, U. V. & LeDoux, J. E. Fear conditioning induces associative long-term potentiation in the amygdala. *Nature* **390**, 604–607 (1997).

12. Mishchanchuk, K. *et al.* Hidden state inference requires abstract contextual representations in ventral hippocampus. Preprint at https://doi.org/10.1101/2024.05.17.594673 (2024).

13. Vertechi, P. *et al.* Inference-Based Decisions in a Hidden State Foraging Task: Differential Contributions of Prefrontal Cortical Areas. *Neuron* **106**, 166-176.e6 (2020).

14. Bromberg-Martin, E. S., Matsumoto, M., Hong, S. & Hikosaka, O. A Pallidus-Habenula-Dopamine Pathway Signals Inferred Stimulus Values. *Journal of Neurophysiology* **104**, 1068–1076 (2010).

407    15. Qü, A. J. *et al.* Nucleus accumbens dopamine release reflects Bayesian inference during
408        instrumental learning. *PLOS Computational Biology* **21**, e1013226 (2025).

409    16. Blanco-Pozo, M., Akam, T. & Walton, M. E. Dopamine-independent effect of rewards on
410        choices through hidden-state inference. *Nat Neurosci* **27**, 286–297 (2024).

411    17. Takahashi, Y. K., Stalnaker, T. A., Roesch, M. R. & Schoenbaum, G. Effects of inference on
412        dopaminergic prediction errors depend on orbitofrontal processing. *Behavioral Neuroscience*
413        **131**, 127–134 (2017).

414    18. Barron, H. C. *et al.* Neuronal Computation Underlying Inferential Reasoning in Humans and
415        Mice. *Cell* **183**, 228-243.e21 (2020).

416    19. Courellis, H. S. *et al.* Abstract representations emerge in human hippocampal neurons during
417        inference. *Nature* **632**, 841–849 (2024).

418    20. Jang, A. I. *et al.* The Role of Frontal Cortical and Medial-Temporal Lobe Brain Areas in
419        Learning a Bayesian Prior Belief on Reversals. *J. Neurosci.* **35**, 11751–11760 (2015).

420    21. Costa, V. D., Tran, V. L., Turchi, J. & Averbeck, B. B. Reversal Learning and Dopamine: A
421        Bayesian Perspective. *J. Neurosci.* **35**, 2407–2416 (2015).

422    22. Hattori, R. *et al.* Meta-reinforcement learning via orbitofrontal cortex. *Nat Neurosci* **26**,
423        2182–2191 (2023).

424    23. Wang, J. X. *et al.* Prefrontal cortex as a meta-reinforcement learning system. *Nat Neurosci*
425        **21**, 860–868 (2018).

426    24. Kim, C. M., Chow, C. C. & Averbeck, B. B. Neural dynamics of reversal learning in the
427        prefrontal cortex and recurrent neural networks. *eLife* **13**, RP103660 (2025).

428    25. Parker, N. F. *et al.* Choice-selective sequences dominate in cortical relative to thalamic
429        inputs to NAc to support reinforcement learning. *Cell Reports* **39**, (2022).

430    26. Pereira-Obilinovic, U., Hou, H., Svoboda, K. & Wang, X.-J. Brain mechanism of foraging:
431        Reward-dependent synaptic plasticity versus neural integration of values. *Proceedings of the*
432        *National Academy of Sciences* **121**, e2318521121 (2024).

433    27. Durstewitz, D., Averbeck, B. & Koppe, G. What neuroscience can tell AI about learning in
434        continuously changing environments. *Nat Mach Intell* **7**, 1897–1912 (2025).

435    28. Bae, J. W. *et al.* Parallel processing of working memory and temporal information by distinct
436        types of cortical projection neurons. *Nat Commun* **12**, 4352 (2021).

437  29. Zhang, X. *et al.* Active information maintenance in working memory by a sensory cortex.
438      *eLife* **8**, e43191 (2019).

439  30. Bolkan, S. S. *et al.* Thalamic projections sustain prefrontal activity during working memory
440      maintenance. *Nat Neurosci* **20**, 987–996 (2017).

441  31. Babayan, B. M., Uchida, N. & Gershman, S. J. Belief state representation in the dopamine
442      system. *Nat Commun* **9**, 1891 (2018).

443  32. Starkweather, C. K., Babayan, B. M., Uchida, N. & Gershman, S. J. Dopamine reward
444      prediction errors reflect hidden-state inference across time. *Nat Neurosci* **20**, 581–589
445      (2017).

446  33. Gershman, S. J. *et al.* Explaining dopamine through prediction errors and beyond. *Nat*
447      *Neurosci* **27**, 1645–1655 (2024).

448  34. Hennig, J. A. *et al.* Emergence of belief-like representations through reinforcement learning.
449      *PLOS Computational Biology* **19**, e1011067 (2023).

450  35. Qian, L. *et al.* Prospective contingency explains behavior and dopamine signals during
451      associative learning. *Nat Neurosci* **28**, 1280–1292 (2025).

452  36. Beyeler, A. *et al.* Organization of Valence-Encoding and Projection-Defined Neurons in the
453      Basolateral Amygdala. *Cell Reports* **22**, 905–918 (2018).

454  37. Paton, J. J., Belova, M. A., Morrison, S. E. & Salzman, C. D. The primate amygdala
455      represents the positive and negative value of visual stimuli during learning. *Nature* **439**, 865–
456      870 (2006).

457  38. Zhang, X. & Li, B. Population coding of valence in the basolateral amygdala. *Nat Commun*
458      **9**, 5195 (2018).

459  39. Zhang, X. *et al.* Genetically identified amygdala–striatal circuits for valence-specific
460      behaviors. *Nat Neurosci* **24**, 1586–1600 (2021).

461  40. Yagishita, S. *et al.* A critical time window for dopamine actions on the structural plasticity of
462      dendritic spines. *Science* https://doi.org/10.1126/science.1255514 (2014)
463      doi:10.1126/science.1255514.

464  41. Radiske, A., de Castro, C. M., Rossato, J. I., Gonzalez, M. C. & Cammarota, M.
465      Hippocampal CaMKII inhibition induces reactivation-dependent amnesia for extinction
466      memory and causes fear relapse. *Sci Rep* **13**, 21712 (2023).

467   42. Adler, A., Zhao, R., Shin, M. E., Yasuda, R. & Gan, W.-B. Somatostatin-Expressing
468         Interneurons Enable and Maintain Learning-Dependent Sequential Activation of Pyramidal
469         Neurons. *Neuron* **102**, 202-216.e7 (2019).

470   43. Cichon, J. & Gan, W.-B. Branch-specific dendritic Ca2+ spikes cause persistent synaptic
471         plasticity. *Nature* **520**, 180–185 (2015).

472   44. Radiske, A. *et al.* Avoidance memory requires CaMKII activity to persist after recall.
473         *Molecular Brain* **14**, 167 (2021).

474   45. Kim, C. M., Chow, C. C. & Averbeck, B. B. Neural dynamics of reversal learning in the
475         prefrontal cortex and recurrent neural networks. *eLife* **13**, (2025).

476   46. Mohebi, A., Wei, W., Pelattini, L., Kim, K. & Berke, J. D. Dopamine transients follow a
477         striatal gradient of reward time horizons. *Nat Neurosci* **27**, 737–746 (2024).

478   47. Murray, J. D. *et al.* A hierarchy of intrinsic timescales across primate cortex. *Nat Neurosci*
479         **17**, 1661–1663 (2014).

480   48. Song, M. *et al.* Hierarchical gradients of multiple timescales in the mammalian forebrain.
481         *Proceedings of the National Academy of Sciences* **121**, e2415695121 (2024).

482   49. Hampton, A. N., Bossaerts, P. & O'Doherty, J. P. The Role of the Ventromedial Prefrontal
483         Cortex in Abstract State-Based Inference during Decision Making in Humans. *J. Neurosci.*
484         **26**, 8360–8367 (2006).

485   50. Baram, A. B., Muller, T. H., Nili, H., Garvert, M. M. & Behrens, T. E. J. Entorhinal and
486         ventromedial prefrontal cortices abstract and generalize the structure of reinforcement
487         learning problems. *Neuron* **109**, 713-723.e7 (2021).

488   51. Prévost, C., McNamee, D., Jessup, R. K., Bossaerts, P. & O'Doherty, J. P. Evidence for
489         Model-based Computations in the Human Amygdala during Pavlovian Conditioning. *PLOS*
490         *Computational Biology* **9**, e1002918 (2013).

491   52. Lee, J. & Sabatini, B. L. Striatal indirect pathway mediates action switching via modulation
492         of collicular dynamics. *bioRxiv* 2020.10.01.319574 (2020) doi:10.1101/2020.10.01.319574.

493   53. Sumi, M. *et al.* The newly synthesized selective Ca2+calmodulin dependent protein kinase II
494         inhibitor KN-93 reduces dopamine contents in PC12h cells. *Biochemical and Biophysical*
495         *Research Communications* **181**, 968–975 (1991).

496   54. Werbos, P. J. Backpropagation through time: what it does and how to do it. *Proceedings of*
497         *the IEEE* **78**, 1550–1560 (1990).

498    55. Williams, R. J. & Zipser, D. Gradient-based learning algorithms for recurrent networks and
499        their computational complexity. in *Backpropagation: theory, architectures, and applications*
500        433–486 (L. Erlbaum Associates Inc., USA, 1995).
501

502    **Methods**

503

504    **Experimental procedures**

505    **Animals.** A total of 48 wild type (WT) C57BL/6J mice (Jackson Laboratory, male and female)

506    were used in the experiments. For optogenetic inhibition of BLA, we crossed a cre-dependent

507    gtACR1 reporter mouse (JAX:033089) with emx1-Cre mice (JAX: 005628) (n=5 mice) to label

508    the excitatory populations within BLA. We used mice heterozygous for both alleles.

509         Animals were housed on a 12-h dark–12-h light cycle and performed the task at the same

510    time each day ($\pm 1$ h), during the dark period. Ambient temperature was kept at $75 \pm 5$ °F, and

511    humidity was kept below 50%. Animals were group-housed (2–5 animals per cage) until surgery,

512    then individually housed throughout the experiment. All procedures were performed in accordance

513    with the US National Institutes of Health Guide for the Care and Use of Laboratory Animals and

514    approved by the Harvard Institutional Animal Care and Use Committee.

515

516    **Surgeries.** All surgical procedures were conducted under aseptic conditions. Mice (older than 8

517    weeks) were anesthetized with isoflurane (3.5% for induction, followed by 1–2% for maintenance

518    at $1$ L min$^{-1}$). A local anesthetic (2% lidocaine) was administered subcutaneously at the incision

519    site. Analgesia was provided with buprenorphine ($0.1$ mg kg$^{-1}$, i.p.) pre-operatively and ketoprofen

520    ($5$ mg kg$^{-1}$, i.p.) for two days post-operatively. After leveling, cleaning, and drying the skull, a

521    custom-made titanium head plate was attached using adhesive cement (C&B Metabond, Parkell).

522    For all injections, the solution (virus or KN-93) was loaded into a pulled glass pipette (5-000-1001-

523    X, Drummond) backfilled with mineral oil and fitted with a plunger. A small craniotomy (<1 mm

524    diameter) was made using a dental drill, and the pipette assembly was mounted on a stereotaxic

525    holder, lowered to the target coordinates, and injected slowly (~100 nL min$^{-1}$) to minimize tissue

526    damage (MO-10, Narishige). After each injection, the pipette was left in place for at least 5 min to

527    allow diffusion before being raised to the next site or withdrawn from the brain. Target coordinates

528    for the brain regions were: ventral striatum (VS): 1.0/1.5/3.8mm, BLA: -1.65/3.3/4.2mm,

529    (anterior-posterior/medial-lateral/dorsal-ventral; coordinates are relative to bregma, and DV

530    relative to surface of the brain). Post surgery, mice were allowed to recover for at least 2 weeks

531    for beginning the experiment.

532

**Viruses.** To record dopamine using fiber photometry, we expressed AAVs encoding the green dopamine sensor GRAB[DA3m] (AAV-hsyn-DA3m(h-D05), WZ Biosciences, 300nl concentration=$5\times10^{12}$ gc/mL) in the left hemisphere VS of WT mice.

**Behavioral setup.** All behavioral experiments took place inside custom-built enclosed behavioral box in which head-fixed mice could stand on a fixed running wheel. The frame of the behavioral box was built using aluminum frames (McMaster) and the walls were made using black hardboard (Thorlabs, TB4). Mice had access to a water spout which delivered artificially sweetened water (Acesulfame Potassium powder dissolved in water at 6g/L; Prescribed For Life) for water reward, and an odor spout which delivered the odor cures for classical conditioning. Behavioral events were controlled (and licking was monitored) using custom-written software in MATLAB (Mathworks) and the Bpod library (Sanworks) interfacing with the Bpod state machine (1024 and 1027, Sanworks), valve module (1015, Sanworks) and port interface board (1020, Sanworks)/water valve (LHDA1233115H, Lee Company) assembly. Odors were delivered using a custom olfactometer, which directed air through one of eight solenoid valves (LHDA1221111H, Lee Company) mounted on a manifold (LFMX0510528B, Lee Company). Each odor was dissolved in mineral oil at 10% dilution, and 30 μL of diluted odor solution was applied to a syringe filter (2.7 μm pore, 13 mm diameter; 6823-1327, Whatman). Wall air was passed through a hydrocarbon filter (HT200-4, Agilent Technologies) and split into a $100\,\text{mL min}^{-1}$ odor stream and $900\,\text{mL min}^{-1}$ carrier stream using analogue flowmeters (MFLX32460- 40 and MFLX32460-42, Cole-Parmer), which were recombined at the odor manifold before being delivered to the nose of the mouse. During photometry experiments, licking was monitored using closed-loop circuit similar to previously described method[52]. During *in vivo* electrophysiology, an infrared emitter–photodiode. Infrared method is less prone to electrical artifact, and thus more suitable for *in vivo* electrophysiology.

**Behavioral tasks.** We used a classical conditioning paradigm in which odor cues predict specific outcomes. An odor was chosen pseudo-randomly (out of two or three odors) and was delivered for 1.5 second after which an outcome (1.5 μL of artificially sweetened water or no reward) was delivered via the waterspout. An inter-trial interval (ITI) separated each trial. The duration of the ITI followed a truncated exponential distribution with a mean of 8 seconds, minimum duration of

564  5 seconds, and maximum duration of 12 seconds. Each session consisted of 200 trials, which lasted

565  about ~35 minutes on average.

566  In the stable task, the odor A (S)-(-)-limonene was always rewarded and odor B (1-heptanol)

567  was always not rewarded. In the dynamic task, the reward contingency changed every session,

568  with the reversal happening at 101th trial (middle point). In the hybrid task (for *in vivo*

569  electrophysiology), we used 3 odors, two of which were stable (odor A and B) and one was

570  dynamic (odor C; 1-hexanol). The outcome of the $3^{rd}$ dynamic odor changed every 60 trials

571  (instead of 100 trials), and the whole session lasted for 240 trials (total of 4 blocks). The stable

572  cues (odor A and B) were presented 25% of the time and dynamic cue (odor C) were presented

573  50% of the time to counter-balance the ratio of stable and dynamic cues.

574

575  **Odors.** We used (S)-(-)-limonene (odor A), 1-heptanol (odor B), and 1-hexanol (odor C). For

576  stable or dynamic task, we used odor A and B. For the hybrid task, we used odor A, B and C. All

577  odors were diluted in mineral oil at 10%.

578

579  **Behavioral training.** Mice were first handled while undergoing water deprivation. This lasted for

580  at least a week until mice reached around 85% of their baseline weight. We confirmed that mice

581  were comfortable licking to a syringe that delivered water while handling, indicating a reduction

582  in overall stress and a willingness to seek water. After handling, mice were habituated on the rig,

583  for at least 3 days. We head-fixed the mice and immediately dispensed water reward. The $1^{st}$

584  habituation session lasted for 15 minutes, and the subsequent habitation session lasted for 35

585  minutes (similar time as the actual length of a behavioral session). We confirmed on the third

586  session that mice were comfortably licking to the spout. If mice never licked to the spout on the

587  $3^{rd}$ day, we continued the habitation until mice became comfortable licking.

588  For training in the stable task, we trained mice for at least 5 days. Mice were able to perform

589  well on the first day and steadily improved (**Extended Data Fig. 1a, c**). Value memory for stable

590  task was tested after at least 5 days of training. For training on the dynamic task, mice were first

591  trained on the stable task for at least 3 days and then trained on the stable task for another 3 days

592  with reversed contingency (odor A=no reward, odor B=reward). After this, mice were then trained

593  on the dynamic task where reversal happened every session at trial 101, for at least 12 days. We

594    made sure that every dynamic task session started with the contingency that the mouse had last
595    seen from the previous session, to preserve the continuity of reward contingency across days.

596         For training in the hybrid task, mice were initially trained on the stable task for 5 days.
597    Mice were then trained on the hybrid task by introducing the third cue (odor C). Initially, each cue
598    was presented 33% of the time, and there were two blocks (one reversal) for the dynamic cue.
599    After 8 days of training, we transitioned to the final version of the hybrid task with cue presentation
600    ratio of 25/25/50% for cue A/B/C respectively, and with 4 blocks. After about 10 sessions in this
601    final version of the hybrid task, Neuropixels recording was performed in the following sessions.

602

603    **Testing value memory.** We tested the stability of the value memory in two ways. In the first way,
604    after mice were fully trained on the stable or dynamic task, we introduced break between sessions
605    (1, 2, 4 or 8 days for stable task; 1 day for dynamic task). The break duration was randomized
606    across mice. In the second way, after mice were fully trained on the stable or dynamic task, we
607    increased the ITI length (duration=0.625, 1.25, 2.5, 5 min) on trial $50^{th}$ and $150^{th}$. To quantify the
608    mice's ability to remember previously learned values, we computed the discrimination index as
609    follows:

610

611    $$discrimination\ index = \frac{CS^+lick(1st\ trial) - CS^-lick(1st\ trial)}{mean\ CS^+lick(last\ 5\ trials) - mean\ CS^-lick(last\ 5\ trials)}$$

612

613    A similar metric was used for dopamine photometry signal (Extended Data Fig. 1h) or for RNNs
614    value (Main Fig. 2i).

615

616    **Testing value inference.** To test if mice could use structural knowledge of the task to infer a
617    change in value based on value update for the opposite cue, we first trained mice fully on the
618    dynamic task (anticorrelated). Next, at the reversal point, we presented one cue type for 5, 10 or
619    20 trials, followed by the other cue (probe trial). We quantified the change in the value of the probe
620    cue relative to level before reversal. Any change in the value of the probe cue would indicate an
621    inferred value update based on structural knowledge of the task (without direct experience). For
622    mice, we computed the inference change in value by computing the change in CS licks (ΔCS licks)
623    or in dopamine (DA) signal (ΔDA (norm. dF/F)). For RNNs, we quantified the change in inferred

624    value ($\Delta$inferred value). For both mice and RNNs, the change was computed by subtracting the

625    pre-reversal baseline level, which was computed by taking the mean value (CS licks, dopamine

626    signal RNN value readout) over the last 5 trials before reversal.

627         We conducted additional experiments to test if mice could learn distinct correlation

628    structures. Two different cohorts of mice were trained as described before but instead of the

629    dynamic task having anti-correlated structure, the values of the odors were either positively

630    correlated or independent. In the positively correlated structure, each session consisted of three

631    70-trial blocks: both odors rewarded, both unrewarded, then both rewarded again. In the

632    independent structure, one odor was constantly rewarded whereas another odor was rewarded and

633    then not rewarded, and vice versa in the next session. In order to test inference, one cue was

634    presented 5-20 times before presenting the other cue. For the independent structure, the stable

635    value odor was always the cue being tested for inference.

636

637    **Fiber photometry.** We used a commercially available bundle fiber photometry system (BFMC,

638    Doric) to record photometry signals from multiple animals simultaneously. A low-

639    autofluorescence Branching Bundle Patchcord (400μm, 0.57 NA) was connected to the

640    photometry system. Each end of the fiber went to a single behavioral rig, allowing us to

641    independent track photometry signals from up to 3 animals simultaneously. A blue excitation LED

642    (470 nm, 10 μW power at tip of the 400um patchcord ferrule) was used to collect GRAB[DA3m]

643    signal. A purple excitation LED (415 nm, 9 μW power) was used to collect control signal for

644    correcting movement artifact. The following parameters were used for imaging in the Doric

645    Neuroscience Studio V6 software: power=10%, framerate=30Hz, exposure=0.012 seconds,

646    gain=9.9dB.

647

648    **KN-93 infusion experiment.** For inhibiting synaptic plasticity in BLA, we used KN-93, a CaMKII

649    inhibitor known to disrupt synaptic plasticity[40,53]. Once mice were fully habituated to the rig (for

650    testing in stable task) or fully trained on the dynamic task (for testing in dynamic task), the day

651    before the infusion, two small craniotomies were made above BLA on each hemisphere, and sealed

652    with silicone elastomers (Kwik-Cast, World Precision Instruments). On the day of infusion, mice

653    were lightly anaesthetized with isofluorane (1%). KN-93 (water soluble KN-93 dissolved in saline

654    at 100 μM concentration) (422711-1MG, Millipore Sigma) was injected bilaterally into BLA

655   (volume: 300 nl per site). Mice were left to recover for an hour before behavioral testing. To verify

656   the injection site, KN-93 was mixed with a lipophilic tracer DiI (less than 5% of total volume)

657   (V22889, ThermoFisher Scientific). Mice were perfused (see Histology) and the slices imaged

658   under a wide-field microscope (see Extended Data Fig. 3). For saline infusion, the same procedure

659   was performed with saline+DiI.

660

661   *In vivo* **electrophysiology.** We performed acute Neuropixels (1.0, single shank) recording in mice

662   fully trained the hybrid task. A total of 6 recording sessions was performed per mouse, 3 sessions

663   per hemisphere. All recordings were performed in SpikeGLX software

664   (https://github.com/billkarsh/SpikeGLX), with a sampling rate of 30 kHz, local field potential gain

665   of 250 and action potential gain of 500, and we analyzed only the action potential channel (which

666   was high-pass filtered in hardware with a cut-off frequency of 300 Hz). Behavioral and neural

667   recordings were synchronized using a transistor–transistor logic (TTL) pulse sent from the Bpod

668   to the PXIe acquisition module SMA input at the start of every trial.

669          The day before the first recording session, a small craniotomy was made bilaterally above

670   BLA. Slicone elastomers (Kwik-Cast) was used to cover the craniotomies. On recording sessions,

671   the silicone gel was gently removed to expose the brain. A Neuropixels probe was soaked into

672   either DiO, DiI or DiD solutions (Vybran Multicolor Cell-Labeling Kit, ThermoFisher

673   Scientific) for tracking the probe location, then mounted on a vertical manipulator (KDC101 and

674   Z825B, Thorlabs). The probe was slowly lowered until it touched the brain surface. Saline was

675   applied around the craniotomy to prevent drying and helping the insertion through the dura. The

676   probe was slowly inserted into the dura at 0.05 mm/s. Once the probe was confirmed to have

677   penetrated the brain, the speed was lowered to 0.01 mm/s. The probe was lowered to 5.3 mm below

678   the brain surface to record in the BLA and regions around it. Once the probe was fully lowered,

679   we waited for 10 minutes for the probe to fully settle inside the brain. Once the recording session

680   was over, the probe was slowly retracted at 0.01 mm/s. Silicone gel was applied on the craniotomy.

681

682   **Optogenetic experiments.** We performed optogenetic experiments to test the role of BLA in

683   stable and dynamic task. Optical fibers (400μm, 0.57 NA, Doric) were implanted bilaterally above

684   BLA (coordinates: -1.3/3.3/4.0) in emx1-Cre × gtACR1 mice. Mice were allowed to recover for

685   at least 2 weeks before starting the handling procedure. Mice were first trained on the stable task

686    for 5 sessions, and then underwent 2 stimulation sessions. The same mice underwent 2 weeks of

687    training on the dynamic task and then underwent 2 stimulation sessions. 470nm LED was used to

688    inhibit BLA neurons. For testing the role of BLA during cue period, LED was on for a duration of

689    3 seconds aligned to odor onset, and randomly in 15% of all trials. for testing the role of BLA

690    during ITI period, LED was on for a duration of 3 seconds, 0.5 seconds prior to the odor onset of

691    the upcoming trial. We performed 2 stimulation sessions per mouse. We made sure to interleave a

692    non-stimulation session to avoid chronic effect from repeated stimulation of BLA neurons during

693    the dynamic task. The effect of stimulation was quantified by comparing the effect of stimulation

694    on CS+ licks and CS- licks (Fig. 3d, Fig. 5h) or by quantifying the difference between CS+ licks

695    and CS- licks (Fig. 5i).

696

697    **RNN modelling.**

698    **Recurrent neural network implementation.** We implemented RNNs as previously described[34,35].

699    Briefly, we trained recurrent neural networks composed of GRU units (N=20) to estimate value at

700    each timestep. The hidden unit activity is given by the following equation:

701

702    $$z_t = f_\phi(o_t, z_{t-1})$$

703

704    $z_t$ is hidden units' activity at time t

705    $\phi$ is the parameter vector of the RNN

706    $o_t$ is the observation produced by environment at time t

707    $z_{t-1}$ is the hidden units' activity at time t-1

708

709    $o_t$ is a one-hot vector defined as follows:

710

711    $\mathbf{o_t} = [odorA, odorB, outcome], where\ odorA \in \{0,1\}, odorB \in \{0,1\}, outcomeB \in \{0,1\}$

712

713    Thus, at each timestep, the RNN had access to a vector indicating whether odor A was present,

714    odor B was present, and whether reward was delivered or not. Each trial began with an intertrial

715    interval (ITI) of duration that followed a geometrical distribution with parameter p=0.8. After the

716    ITI, a cue was presented ($\mathbf{o_t} = [1,0,0]\ or\ \mathbf{o_t} = [0,1,0]$ for odor A or B respectively) followed by

717     the outcome ($\mathbf{o_t} = [0,0,1]$ *or* $[0,0,0]$) for reward or no reward respectively).$V_t = w^\top z_t + w_0$ for

718     $z_t, w \in \mathbb{R}^H$ (where H=20 is the number of hidden units), $w_0, V_t \in \mathbb{R}$. The full parameter

719     vector $\theta = [\boldsymbol{\phi}, \boldsymbol{w}, w_0]$ was learned using TD learning. This involved backpropagating the gradient

720     of the squared error loss $\delta_t^2 = (r_t + \gamma V_{t+1} - V_t)^2$ with respect to $V_t$. The discount factor $\gamma$ was

721     set to $\gamma$=0.2.

722

723     **Task implementation.** To implement the stable task, we used 200 trials where each odor A or B

724     was chosen at random, with odor A being rewarded and odor B not being rewarded. For the

725     dynamic task, reward contingencies flipped every 50 trials (block length = 50 trials). The dynamic

726     task lasted for 18 blocks.

727

728     **Training.** We first initialized the network to PyTorch's default. We used a truncated

729     backpropagation through time (T-BPTT) learning rule[54,55] where at each timestep, the network

730     used recent inputs (defined by the window size W) to compute the gradient and update the weights

731     of the network. We explored the effect of varying window size W (see Extended Data Fig. 2). For

732     the majority of the analyses we fixed the window size to W=720 (~100 trials), which was roughly

733     equal to the length of two blocks in the dynamic task. Given that the network is continuously

734     learning in this setting, it is possible that the learning gets stuck in non-optimal local minima. To

735     avoid such networks that might have sub-optimally learned the task, we first computed mean loss

736     (squared reward prediction error) for the last 20 trials of the last four blocks in the dynamic task

737     for a range of W (Extended Data Fig. 1a-b). Most networks had a loss less than 0.0005. Thus, we

738     excluded networks whose loss exceeded 0.0005. This was to ensure that the final RNNs could

739     solve the task, regardless of the mechanism being used. Learning rate was set to 0.0005, and we

740     used the Adam optimizer with the AMSGrad variant (amsgrad=True in PyTorch), which enforces

741     a non-decreasing second-moment estimate for more stable updates.

742

743     **Plasticity manipulation.** We manipulated the plasticity in RNNs to test the role of plasticity in

744     updating value in either stable or dynamic task. The learning rate was set to zero before training

745     on the stable task to test the role of plasticity in the stable task. For the dynamic task, RNNs were

746     first fully trained on the dynamic task, while excluding RNNs that never converged (see **Training**).

747     We then set the learning rate to zero to test the role of plasticity in updating value in the dynamic

748    task. To quantify the effect of plasticity manipulation, we computed the area under the receiver

749    operating characteristic (AUROC) between value readout of CS+ cue and CS- cue.

750

751    **PCA analysis.** We applied principal component analysis (PCA) to better understand the neural

752    state-space trajectories during stable and dynamic task. PCA was applied to the hidden units'

753    activity over the entire period of training. Example state-space trajectories were plotted in the PC1-

754    PC2 and value readout space. Note that the value readout space is not strictly a function of hidden

755    units' activity, but a combination of hidden units' activity and the readout weights.

756

757    **Context axis analysis.** We defined the context axis as the axis that could best discriminate the

758    context (block type) in the hidden unit activity space in the dynamic task. We first took the ITI

759    activity of the last 4 blocks of the entire 18 blocks training of the dynamic task. The activity was

760    then segregated into two contexts, and then a Fisher linear discriminant analysis was used to

761    compute the line that best separate context. The resulting axis was defined as context axis. Activity

762    during the ITI was projected onto the context axis for further analysis ($\text{context}^{\text{proj}}$).

763          To understand how the context information was being used, we computed the Spearman

764    rank correlation coefficient between the $\text{context}^{\text{proj}}$ and the value difference between the last two

765    cues (value readout of last cue A – value readout of last cue B). The correlation was computed

766    using the first 4 blocks for naïve RNNs or using the last 4 blocks for expert RNNs.

767

768    **Long ITI effect.** To understand whether long ITI could have differential effects on the value

769    memory in RNNs, we simulated the effect of long ITI by increasing the length of ITI to 5, 10, 20

770    or 40 in either stable or dynamic task. Discrimination index was computed similarly for RNNs

771    (see **Testing value memory**).

772          To understand the effect of the long ITI, we projected the activity during the long ITI onto

773    the context axis. We computed the drift by computing the distance moved along the context axis

774    during the long ITI. We then computed the correlation coefficient between the value differential

775    after the long ITI and the drift along the context axis. Value differential was corrected to be positive

776    at the beginning of the long ITI. $\text{Context}^{\text{proj}}$ was also corrected to be positive if drift was moving

777    away from the initial point towards the other fixed point.

778

779 **Analysis**

780 **Photometry preprocessing.** dF/F was first calculated by computing $F_0$, and using the formula

781 dF/F=$(F_{raw} - F_0)/F_0$. $F_0$ was defined as the 10th percentile of $F_{raw}$ in a rolling window of 30 s. dF/F

782 traces were upsampled from 20 Hz to 1000 Hz through linear interpolation (MATLAB function

783 interp1) and then smoothed with a Gaussian filter (SD = 50 ms). We then normalized dF/F by z-

784 scoring (MATLAB zscore function).

785

786 **Spike sorting.** Neuropixels recording data were spiked sorted offline with Kliosort4

787 (https://github.com/MouseLand/Kilosort?tab=readme-ov-file) with default parameters, followed

788 by manual curation of individual units using Phy (https://github. com/cortex-lab/phy).

789

790 **Brain registration for *in vivo* electrophysiology.** To label individual units with their

791 corresponding location in the brain, we registered the histology to the Allen Mouse Brain Atlas.

792 We first used AP_histology to register each histology slices with the tracer to the Allen Mouse

793 Brain Atlas (https://github.com/petersaj/AP_histology/tree/master). Coordinates for each probe

794 track was then converted in the relevant format to be read out by the IBL Ephys Atlas alignment

795 tool (https://github.com/int-brain-lab/iblapps/tree/master/atlaselectrophysiology). We obtained a

796 brain location for each recorded single unit, which was used for brain region specific analysis.

797

798 **Histology.** Mice were deeply anesthetized with an overdose of ketamine/medetomidine,

799 exsanguinated with 0.9% phosphate buffered saline (PBS), and transcardially perfused with cold

800 4% paraformaldehyde (PFA) in PBS. The brain was extracted from the skull and stored in 4% PFA

801 for 24-48 hours at 4°C, after which it was rinsed with PBS, stored in PBS, and cut into 100 $\mu$m

802 sections on a vibratome (VT1000S, Leica). Sections mounted on slides and then imaged using a

803 slide scanner (Zeiss Axioscan 7).

804

805 **Neuropixels recording analysis**

806 **Value coding cells.** We defined a cell as stable value coding if it passed the following criteria:

807 significantly different firing rate between odor A and odor B during the 1.5s odor on period. We

808 tested for significance using student's *t*-test (MATLAB ttest2 function) at α=0.01, for different

809 blocks combination (block1+2, block 3+4, block1+4, block2+3, block1+2+3+4). A cell had to pass

810  all 5 tests to qualify as a stable value coding cell. This was to ensure that the cell was firing

811  consistently across blocks. For dynamic value coding, we similarly defined dynamic value coding

812  cell if it passed the following criteria: significantly different firing rate between odor C during

813  rewarding block and odor C during non-rewarding block. We tested for significance using

814  student's *t*-test at α=0.01for different combinations of blocks (block 1 vs 2, block3 vs 4, block1+3

815  vs 2+4). A cell had to pass all 3 tests to qualify as a stable value coding cell. This was to ensure

816  that the dynamic value was consistently maintained throughout the session and not just in specific

817  blocks. Lastly, a stable and dynamic value coding cell was defined as a single unit that was both

818  stable value coding and dynamic coding following the criteria mentioned above. When quantifying

819  the percentage of cells encoding value, we excluded cells whose baseline firing rate (total

820  spikes/session duration) was below 0.2 Hz.

821

822  **Context coding cells.** We defined a cell as context coding if it passed the following criteria:

823  significantly different firing rate during the ITI (4 seconds before cue onset till cue onset) between

824  rewarding block and non-rewarding block. We tested for significance using a student's *t*-test at

825  α=0.01, for different blocks combination (block 1 vs 2, block3 vs 4, block1+3 vs 2+4). A cell had

826  to pass all 3 tests to qualify as a context coding cell.

827

828  **Correlation between CS licks and context coding.** To quantify the relationship between the

829  context coding during the ITI and value for the dynamic cue, we computed the Spearman

830  correlation between the ITI firing rate during the ITI (last 4 seconds before odor onset) and the

831  number of licks during odor C delivery (1.5 seconds duration). We combined trials from block2

832  and block4 or block 1 and block 3 and computed the correlation for the two sets of trials, and then

833  computed the mean correlation across the block types. This was to avoid introducing spurious

834  correlation given that context coding was initially defined by being able to differentiate block 1+3

835  vs block 2+4, and given that mice's licks to odor C also differentiated block type. For control, we

836  shuffled the CS licks for odor C and ITI firing rate pairing.

837

838  **Context update analysis.** In the hybrid task, average reward rate and context coding are

839  confounded. To determine if context coding reflected block identity based on the outcome of odor

840  C only, we quantified the amount of context coding update for each cue type. If context coding

841    neurons reflect reward rate, then one would expect context to be updated solely based on outcome

842    regardless of cue type. However, if context coding truly reflects odor C specific context, then one

843    would expect context coding to be only updated after odor C. We computed the mean context

844    update at the beginning of each block (Δcontext coding) by first taking the change in ITI firing

845    rate (4 seconds before odor onset till odor onset) after each trial type. We focused on the first six

846    trials for each cue type because context update occurred mostly at the beginning of each block. We

847    restricted our analysis for block 2 and block 4. Thus, for each context coding unit, we obtained a

848    mean context update for each cue type, which was taken from a total of 12 trials for each cue type.

849    Predictions for the reward rate model vs pure context coding model is shown in Extended Data

850    Fig. 5.

**Figure + legends**



**Fig. 1 | Environmental stability sets distinct timescales for value update and forgetting. a,** Trial structure. Following an ITI (5~12 s), an odor cue (1.5s) was presented randomly (odor A or B), followed by an outcome (water reward or no reward). **b,** Task type. In the stable task (top), the outcome was fixed (odor A-reward, odor B-no reward). In the dynamic task (bottom), the outcome flipped once every session (2 blocks). Each session started with the same reward contingency as the last block from the previous session. **c,** Expert performance of mice trained on the stable task (training duration=3days). Lick rate (Hz) for rewarded odor CS+ (blue) and unrewarded odor CS- (red) is shown aligned to odor onset (n=10 mice). **d,** Similar quantification as in **c** for expert mice on the dynamic task, for block 1 (left) and block 2 (right) (n=10 mice). **e,** AUROC between CS+ and CS- in stable (green) and dynamic (purple) task for naïve (empty circle) and expert mice (filled circle) (n=10 mice; ***, $P<0.001$, two-tailed $t$-test). **f,** Learning curve for updating value (+Δvalue:

positive update, –Δvalue: negative update) in stable/1$^{st}$ reversal and dynamic task. *left*, CS+ lick rate for naïve mice first exposed to the stable task, aligned to 1$^{st}$ day (green) and expert mice on the dynamic task, aligned to reversal point (0 trial=reversal point) (purple). *right*, similar quantification for CS- lick rate: mice experiencing reversal for the first time (green) aligned to reversal point and expert mice on the dynamic task decreasing CS- lick rate aligned to reversal point (stable: n=6 mice; dynamic: n=5 mice). **g,** Testing the timescale of value forgetting in dynamic vs stable task with inter-session break. *left,* schematic showing the break (grey box) separating session N from session N+1. *right*, lick raster plot for an example session aligned to odor onset for stable and dynamic task. Grey box is the break (24 hours) between session N and session N+1. Each row indicates a trial with the colored box indicating the odor on period (red, CS+; blue, CS-). **h,** Quantification forgetting for different inter-session break duration (1, 2, 4, or 8 days break). Discrimination index represents how well mice could discriminate CS+ and CS- on the first trial for each CS on session N+1 (normalized by previous session performance; see Methods) for stable (green) or dynamic (purple) task. (stable: n= 4 mice; dynamic: n= 7 mice; ***, $P<0.001$, two-tailed *t*-test). **i,** Testing the timescale of value forgetting in dynamic vs stable tasks with inter-trial break. left, schematic showing two breaks within a session (grey boxes). *right*, lick raster plot for an example session similar to **g**. **j,** Similar quantification as in **h** for inter-trial breaks (37.5, 75, 150 or 300 sec) (stable: n=6 mice; dynamic: n=7 mice; *, $P<0.05$; **, $P<0.01$, two-tailed *t*-test). All data shown are mean ± s.e.m.

**Fig. 2 | RNNs with online weight update recapitulates mouse behavior. a,** *left,* Schematic showing the RNN with plastic synapses (pink). States are represented by the hidden units' activity, and value is readout using RNN representation (see Methods). *right,* Schematic showing the online learning rule using TBPTT. At each timestep, the RNN updates its weight based on the past experience using a sliding window (see Methods). **b,** *left,* Example RNN learning the stable task. Blue and red traces re the RNN value readout for rewarded cue (CS+) and unrewarded cue (CS-) respectively. *right,* Example RNN learning the dynamic task (blue, cue A; red, cue B). Background color denotes the cue type that is being rewarded in that block. **c,** *left,* Example RNNs (n=10) showing positive value update (+Δvalue) in the stable task (green) and in the expert level dynamic

task (purple) For stable task, value starts from naïve to first exposure to the task. For dynamic task, value starts from the start of reversal for the cue that was previously unrewarded. *right,* similar quantification for (–Δvalue) in the first reversal (green) and in expert level dynamic task (purple). **d,** Example RNNs (n=10) showing the effect of freezing the weight in the stable (green) or dynamic task (purple) for positive value update (black=weight update intact; red=weight frozen). **e,** Quantification of the AUROC for CS+ and CS- value readout in the stable (left) and dynamic (task) when weight update was intact (black) or frozen (red). Weight freezing had a significant effect on updating value in the stable task (***, $P<0.001$, two-tailed $t$-test) but no effect in the dynamic task ($P=0.31$, two-tailed $t$-test). **f,** Example neural trajectories in PC1, PC2 and value space during stable task for the first 30 trials. Green dot indicates the time of cue and black dot indicates the time of outcome delivery. Cue A (reward cue) is shown in blue and cue B (non-rewarded cue) is shown in red. **g,** Similar plot as in **f** for an expert RNN trained on the dynamic task. *left*, example neural trajectories are shown for both block. PC1 separate cue type whereas PC2 separates block type. *right*, same neural trajectories plotted in PC1, PC2 and value space. **h,** An Example RNN with the hidden units' activity during the ITI projected onto the context axis (context$^{proj}$) (*top*, see Methods) and the value readout (*bottom*) for cue A (blue) and cue B (red). Colored background represents which cue is being rewarded in that block. **i,** Quantification of the correlation between the value differential (value diff.) and the context projection (context$^{proj}$) for naïve (*left*) and expert (*right*) RNN (n=1). Value differential which was defined as the value difference between the last CS+ and CS- (see Methods). Each dot represents a single trial (red/blue=types of blocks). **j,** *left,* Quantification of Spearman's rank correlation coefficient ($\rho$) between value diff. and context$^{proj}$ in naïve (grey) vs expert (purple) RNNs (n=100). Context projection became more predictive of value differential in expert compared to naïve RNNs (***$P<0.001$, two-tailed $t$-test). *right,* AUROC between block types (cue A-reward block vs cue b-reward) in naïve (grey) vs expert (purple) RNNs (n=100). Context projection became more separable in naïve vs expert mice (***, $P<0.001$, two-tailed $t$-test). **k,** An example expert RNN on the dynamic task, quantified similar to **h**, when long ITI is introduced (black arrow). **l,** Quantification of the effect of introducing long ITI (discrimination index) as a function of the length of the ITI (5-40AU) for stable (green) and dynamic task (purple). Long ITI trials decreased the discrimination index for dynamic but not stable task (***, $P<0.001$, two-tailed $t$-test). **m,** Relationship between value difference (value diff.) between CS+ and CS- after long ITI and the

drift caused by the long ITI along the context axis (see Methods). The amount of drift on the context axis predicted the change in value readout between CS+ and CS- ($R^2=0.62$, two-tailed Pearson's correlation coefficient test, ***, $P<0.001$). All data shown are mean ± s.e.m.
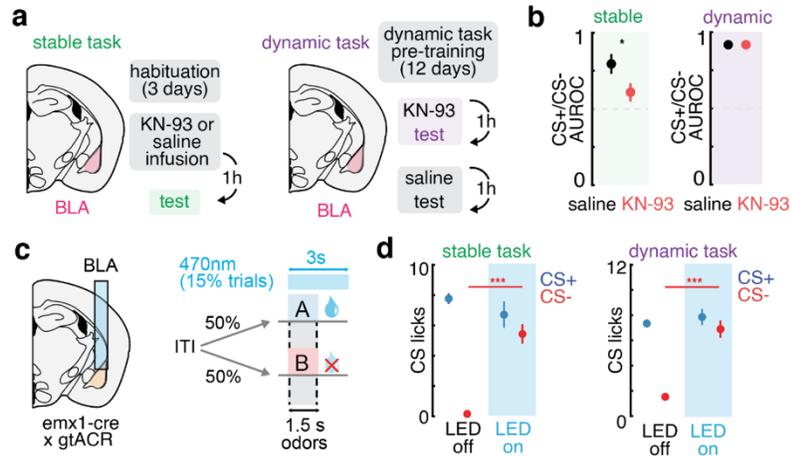
**Fig. 3 | Dissociable roles of BLA plasticity vs activity. a,** Schematic showing experimental flow for testing the role of BLA plasticity in stable vs dynamic task. *left,* Mice were first habituated to the rig for 3 days, after which they underwent infusion of either KN-93 or saline, targeting BLA (pink). Mice were tested on the stable task after 1h post-infusion. *right,* Mice were trained on the dynamic task for 12 days until they reached expert performance, after which they underwent infusion of KN-93 and saline in the BLA on different consecutive sessions (see Methods). **b,** Effect of KN-93 infusion on the performance in stable (left, green) or dynamic (right, purple) task. AUROC between CS+ and CS- licks are quantified for each experiment. (stable/saline, n=11 mice; stable/KN-93, n=8 mice; dynamic, n=7; *, P<0.05, two-tailed *t*-test). **c,** Schematic showing experimental flow for inhibiting BLA activity in the stable vs dynamic task. *left,* emx1-Cre × gtACR1 mice were implanted with bilateral fibers above BLA. *right,* stimulation light (470 nm LED) was on for 3 sec starting from odor onset on 15% of all trials. **d,** Quantification of the effect of inhibiting BLA activity in stable (left) and dynamic (right). CS+ (blue) and CS- (red) licks are shown for LED off (left) and LED on (right) (n=10 sessions form 5 mice; ***, P<0.001, two-tailed *t*-test). All data shown are mean ± s.e.m.

**Fig. 4 | Stable and dynamic value coding in the BLA. a,** Schematic showing trial structure and reward contingency across session. *top*, in the hybrid, task, three odors were presented, odor A or odor B with 25% probability and odor C with 50 % probability. Odor A and B had stable value (A-reward, B-no reward) whereas the value of odor C was dynamic. *bottom*, Each session consisted of four blocks where the value of odor C alternated between reward and no reward. Initial value of odor C was kept the same as the value of odor C in the last block of previous session to preserve continuity. **b,** Neuropixels probe trajectories targeting BLA. *left*, 3d rendering of the probe trajectories of recording sessions. BLA is shown in red. *right*, example histological slice showing probe trajectories. DiO (green), DiI (yellow) and DiD (red) were used probe tracking (see Methods). **c,** Example units coding for stable value (top), dynamic value (middle), or stable and dynamic value (bottom). Each graph shows the smoothed peri-stimulus time histogram (PSTH) averaged across each block (grey, odor C → reward; white, odor C → no reward). Each trace shows firing rate (spikes s$^{-1}$) for CS+ odor (blue), CS- odor (red) or dynamic CS+/- (purple). **d,** *top,* Fraction of units encoding stable value (SV, top), dynamic value (DV, middle), or stable and dynamic value (SVDV, bottom). **e,** *top,* Mean firing rate of all amygdala positive (left, n=780) or negative (right, n=1146) SV neurons, aligned to odor onset (blue/red, CS+/CS-). *Bottom,* Mean firing rate of all positive (left, n=344) or negative (right, n=510) DV neurons aligned to odor onset (purple/grey, CS+/CS-). **f,** *left,* Stable value selectivity (firing rate$^A$ – firing rate$^B$) vs dynamic value selectivity (firing rate$^{C-Reward}$ – firing rate$^{C-No\ reward}$) for all DV neurons ($R^2$=0.62, two-tailed Pearson's correlation coefficient test; ***, $P$<0.001). *right,* Fraction of units categorized by the

polarity of SV and DV selectivity (SV+DV+, dark blue; SV-DV-, light blue; SV+DV-, light pink; SV-DV+, dark pink). All data shown are mean ± s.e.m.
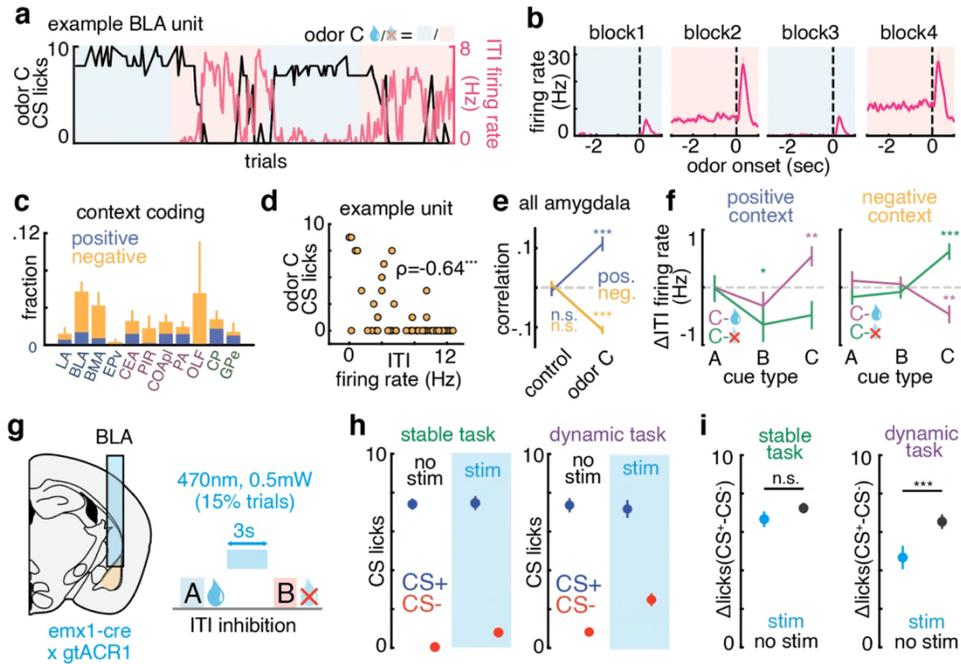
**Fig. 5 | Context coding in the BLA. a,** Example BLA unit that encodes context in the ITI. CS licks to odor (black line) is shown along the mean ITI firing rate during the 3 seconds before cue onset (pink line, see Methods). Background colored rectangles indicate the block type (light blue, rewarded odor C; light pink, not rewarded odor C). **b,** Same example BLA unit as in **a** showing mean firing rate during each block aligned to odor onset (background color indicates block type as in **a**). **c,** Fractions of units context coding units across all brain regions recorded. Positive/negative context (blue/yellow) was defined by the whether ITI firing rate was higher in the odor C reward block vs odor C no reward block (see Methods). **d,** Example negative context encoding BLA unit showing the correlation between CS licks to odor C and ITI firing rate before odor C onset. Data is from block 2+4 combined, showing correlation within same block type (see Methods) ($\rho$=-0.64, two-tailed Spearman's rank correlation test; ***, $P<0.001$). **e,** Quantification of the mean Spearman's rank correlation between odor C CS licks and ITI firing rate for all positive (pos., blue, n=59) and negative (neg., yellow, n=171) in the amygdala (***, $P<0.001$, two-tailed $t$-test). For control, we shuffled the CS licks-ITI firing rate pairing (see Methods). **f,** ITI firing update for all positive (left) and all negative (right) context coding amygdala units. At the beginning of each block, we quantified the change in ITI firing rate followed by each cue type (odor A, B or C) in either odor C reward block (purple) or odor C no reward block (green) (*, $P<0.05$, **, $P<0.005$; ***, $P<0.001$, two-tailed $t$-test). **g,** Schematic showing experimental flow

for inhibiting BLA activity during the ITI period. *left,* emx1-Cre × gtACR1 mice were implanted with bilateral fibers above BLA. *right*, stimulation light (470 nm LED) was on for 3 sec during the ITI on 15% of all trials. **h,** Quantification of CS licks (blue: CS+; red: CS-) for no stim trials (black) or trials where stimulation was on preceding cue (blue area) for stable task (left, n=12 sessions) or dynamic task (right, n=10 sessions). **i,** Quantification of the difference in CS licks ($\Delta$licks ($CS^+$-$CS^-$)) for stim (light blue) and no stim trials (black) in stable (left) or dynamic (right) task (n.s., $P>0.05$; ***, $P<0.001$, two-tailed $t$-test). All data shown are mean $\pm$ s.e.m.

**Fig. 6 | Recurrent dynamics enable structure-specific value inference. a,** Schematic showing distinct predictions for value inference. Plasticity-based value update (left) has no learned structure, thus cannot infer value. Dynamics-based value update (right) learns the hidden states of the environment. The learned structure enables value inference. **b,** Task design to test value inference. Once mice were fully trained on the dynamic task, mice were tested around the reversal point by only presenting one cue type (red, 5-20 trials), after which the other cue was presented (blue). Value (e.g. CS licks) of the other cue could be inferred to be high using structural knowledge (value inference). **c,** Testing value inference in naïve (left) vs expert (right) RNNs trained on the dynamic task (n=4 example RNNs). **d,** Testing value inference in naïve (left) or expert (right) mice on the dynamic task (n=6 mice). **e,** Quantification of change inferred value in RNNs (Δinferred value) in naïve (green, n=100 RNNs) and expert (purple, n=100 RNNs) RNNs (***, $P<0.001$, two-tailed $t$-test). **f,** Quantification of inferred value in mice using CS licks (left, n=6) or dopamine photometry signal (right, n=6 mice) (*, $P<0.05$; ***, $P<0.001$, two-tailed $t$-test). **g,** Relationship between ΔCS licks and the number of trials of the opposite cue, indicating strength of belief that a change of state has occurred (n=6 mice) (*, $P<0.05$, two-tailed $t$-test). **h,** schematic showing three distinct correlations structures of odor values used for dynamic task (see Methods). **i,**

quantification of changed in inferred value for anti-correlated (-corr., orange), correlated (+corr., turquoise), or independent (ind., gray) structures, when the value of the opposite cue decreased. Change in inferred value is quantified using either licks during CS ($\Delta$CS licks) or dopamine signals ($\Delta$DA dF/F) (-corr.=12, +corr.=6, ind.=6). **j,** similar quantification as in **i** when the value of the opposite cue increased (*, $P<0.05$, **, $P<0.005$; ***, $P<0.001$, two-tailed *t*-test). All data shown are mean $\pm$ s.e.m.

**Extended Data Fig. 1 | Further quantification of behavioral training, forgetting after breaks, and dopamine signaling related to forgetting. a,** Lick rate for session 1-3 for the stable task. Mice could learn to discriminate from session 1, with performance gradually improving over 3 days (n=10 mice) **b,** Comparing naïve vs expert performance on the dynamic task. *left*, 1st reversal session for block 1 and block 2. *right,* similar quantification for expert performance. Naïve mice initially had difficulty updating value and adapting to the new reward contingency on block 2 (n=10 mice) **c,** Quantification of the AUROC for CS+ and CS- for stable dynamic task shown as a function of session number (n=10 mice). **d,** Effect of inter-session break (1 day) for stable (left) or dynamic (right). CS lick rate is plotted for each odor (A, red; blue, B) as a function of trial number (last 10 trials or first 10 trials for session N or session N+1 respectively). Mice start licking to the CS- odor after 24 hours break in the dynamic task but gradually learn to suppress licking to the CS-. This effect is not present in the stable task (stable: n=4 mice; dynamic: n=7 mice). **e,** similar plot as in **d** for inter-trial break (300 sec). Grey box indicates the long break (stable: n=6
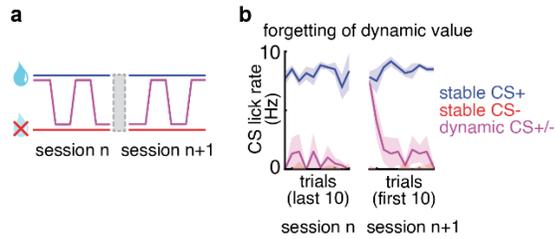
mice; dynamic: n=7 mice). **f,** *left,* Histological slice of am example mouse showing the expression of the dopamine sensor GRAB$^{DA3m}$ (green) in the ventral striatum. *right,* fiber locations (red) aligned to the atlas for all mice (n=12 mice). **g,** Example sessions showing dopamine response to CS after long ITI (300 sec). *top*, Normalized dF/F for CS+ (red) and CS- (blue) aligned to odor onset for stable task. The plot shows the last five trials (column=trial) before the long ITI (grey box). After long ITI, the dopamine response to the CS+/CS- are similar to the level before. *bottom*, similar quantification for dynamic task. After long ITI, the dopamine response to the CS- increases to the level similar to the CS+ response. **h,** Quantification of the forgetting timescale as in **Fig. 1j** for CS dopamine response (mean normalized dF/F during odor period: see Methods). Discrimination index was lower in dynamic task compared to stable task for longer ITI duration (300 sec) (n=7 mice; ***, $P<0.001$, two-tailed $t$-test).
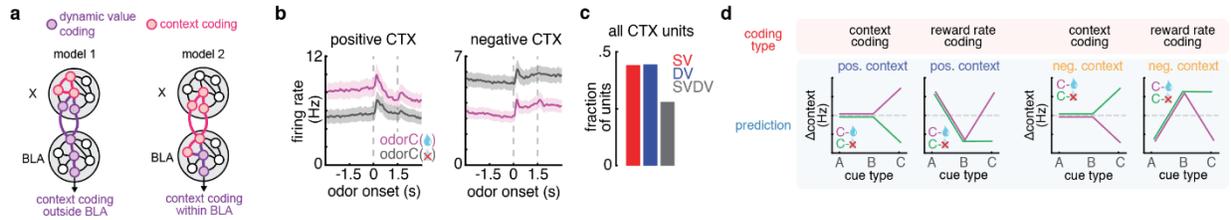
**Extended Data Fig. 2 | RNN modelling with different hyperparameters. a,** Histogram showing the distribution of mean loss (mean $RPE^2$ in the last 4 blocks, only computing the last 20 trials for each block) for different window size (W=100, 300, 500, 700, 900) during the dynamic task. Red line represents the threshold (0.005) for convergence (see Methods). For each window, we trained n=100 RNNs. **b,** Histogram showing the fraction of RNNs (mean $RPE^2$ <0.005) for different window size W. **c,** Value readout of CS+ (blue) and CS-(red) of example RNNs trained on the dynamic task with different window sizes (top: W=90, middle: W=360, bottom: W=720). RNNs with short W=90 do not display meta-reinforcement learning (speeding up of value update). RNNs with medium sized W=360 displays meta-reinforcement learning where the value update transitions from slow to fast process. RNNs with long W=720 displays meta-reinforcement learning but struggles to update value prior to meta-reinforcement learning due to the window encompassing multiple blocks. **d,** Quantification of the mean $RPE^2$ for each block in the dynamic task. Each color indicates a different window size (n=100 RNNs for each window). **e,** Quantification of the AUROC of value readout of CS+ and CS- for each block in the dynamic. Each color indicates a different window size (n=100 RNNs for each window). **f,** Effect of freezing weight update on the value update during the stable task for different window size. AUROC is plotted against window size for RNNs network with plasticity (w. plasticity: grey) and without plasticity (wo. plasticity: pink). Dotted line represents chance level (AUROC=0.5) (**\*\*\***, *P*<0.001, two-tailed *t*-test) **g,** Similar quantification as **f** but for the dynamic task (**\*\*\***, *P*<0.001, two-tailed *t*-test; n.s., not significant).
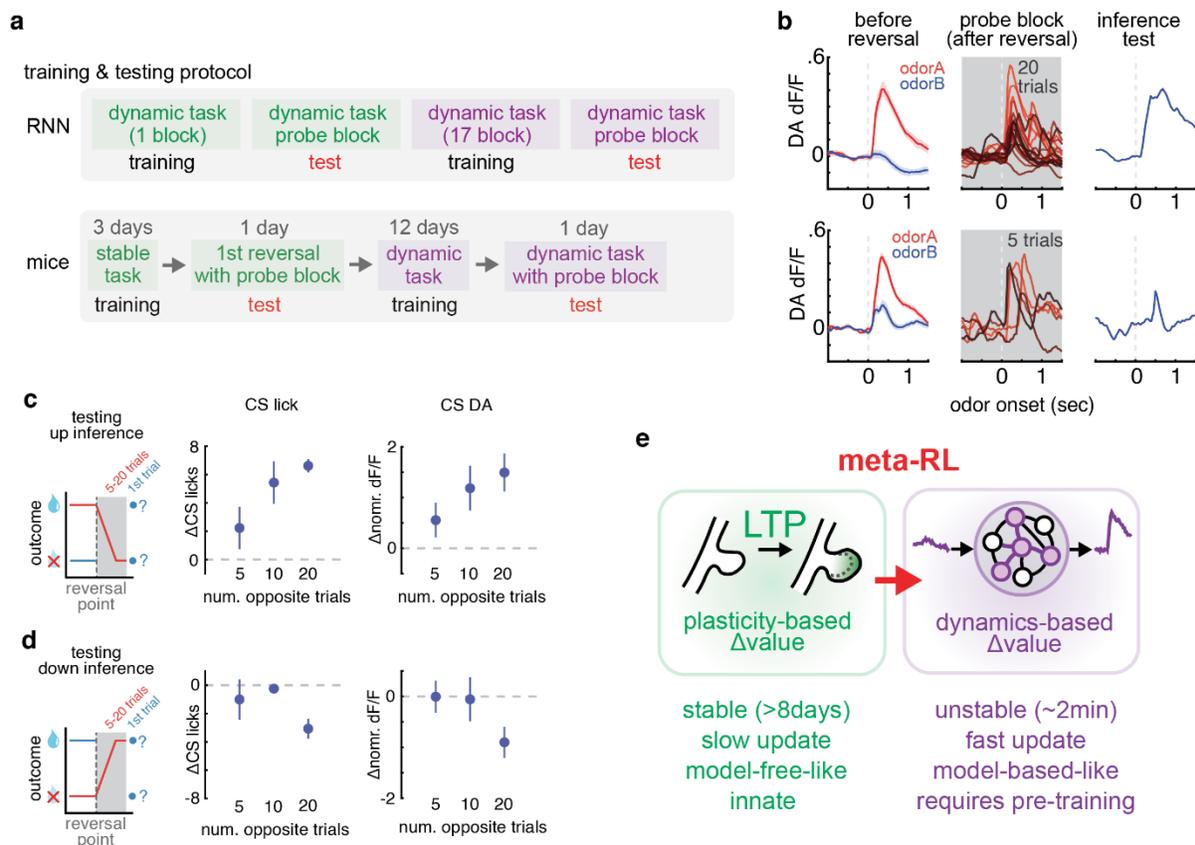
**Extended Data Fig. 3 | Histology, detailed behavioral data and alternative models for BLA manipulation experiments. a,** Example histology showing a coronal slice and the tip of the injection site. Slice shows DAPI (blue) and DiI (orange) which was mixed with saline or KN-93 to locate the injection site (see Methods). **b,** Summary plot showing all the injection sites for the KN-93/saline infusion experiment (see **Fig. 3a-b**) aligned to the reference atlas. Each dot (orange) represents the injection site for one mouse. **c,** Lick rate for CS+ (red) and CS-(blue) aligned to odor onset after saline infusion (left) or KN-93 infusion (right) (stable/saline, n=11 mice; stable/KN-93, n=8 mice). **d,** Similar plot as in **b** for dynamic task. **e,** Similar plot as in c for the dynamic task (n=7).   **f,** Two models consistent with the experimental results of blocking plasticity in the BLA. In model 1, plasticity in the BLA drives value update in the stable task (top, left), but transitions into dynamics-based value update (top, right) after becoming expert in the dynamic task. In model 2, plasticity in BLA drives value update in the stable task (bottom, left), but the locus of plasticity transitions from BLA to another region (region X). Both models can explain why blocking synaptic plasticity in the BLA in the stable task impairs performance but does not impair performance in the dynamic task. A prediction of model 2 is that BLA activity should not be necessary to perform the dynamic task. **g,** Example histology of an emx1-Cre × gtACR1 mouse. Slice shows DAPI (blue) and gtACR1/mCherry (red).

**Extended Data Fig. 4 | Distinct timescale of forgetting of stable vs dynamic value in the hybrid task. a,** Schematic showing two sessions separate by a 1-day break (grey box). Blue/red lines indicate the outcome for stable cues (CS+/CS-) and purple line indicates the outcome for dynamic cue (CS+/-). Session N+1 starts with the same reward contingency as the end of session N. **b,** Quantification of CS lick rate for stable and dynamic cues for the last 10 trials of session N, and first 10 trials of session n+1. CS lick rate for stable cues remained similar to previous day's rate whereas CS lick rate to the dynamic cue increased relative to its previous level.

**Extended Data Fig. 5 | Further quantifications of context coding in the BLA. a,** Two models for dynamic value computation in the BLA. In model 1 (left), dynamic value is computed outside BLA (shown here as unknown region X). BLA inherits dynamic value from region X. Context coding units are not present in BLA. In model 2 (right), context coding is either computed locally or inherited from region X. BLA locally computes dynamic value using context information (purple=dynamic value coding neurons; pink=context coding neurons). **b,** Firing rate of all amygdala positive context units (positive CTX, n=59, left) or negative context units (negative CTX, n=171, right) aligned to odor onset. Firing rate is shown for odor C in reward blocks (purple) or non-rewarded blocks (grey). **c,** Within all CTX units in the amygdala, fraction of units that are SV coding (red), DV coding (blue), or SVDV (grey). **d,** Distinct predictions for context update (Δcontext) when coding is reward coding or reward rate coding. Predictions are shown for positive context coding neurons (pos. context, blue) and negative context coding neurons (neg. context, yellow), in either reward blocks (purple) or non-reward blocks (green).

**EDF 6 | Further quantification of value inference in RNNs and mice. a,** Training and test protocol for value inference related to Main Fig. 6. RNNs were first trained on the 1st block of the dynamic task (equivalent to stable task). On the first reversal point, RNNs were tested for the ability to infer value (probe block). We define this point as naïve RNNs (see main Fig. 6). RNNs were further trained for 17 blocks in the dynamic task, and then tested again on for value inference using the probe block. We define this point as expert RNNs (see main Fig. 6). **b,** Two example sessions showing dopamine photometry signal from a mouse being tested for inference when the opposite trial was presented 20 trials (top) or 5 trials (bottom). *left*, mean traces for odor A and odor B before the reversal (last 10 trials for each odor). *center*, dopamine traces from trials where only one odor was presented consecutively. Color indicates the trial order (light red->dark red=early->late trials). *right*, inference trial where the target cue is presented for the first time (see **Main Fig. 6f, g**). **c,** Up inference. *left,* Schematic showing the protocol for testing upward inference. After reversal, 5-20 trials of one cue were presented serially and then the other

cue was presented the first time. Upward inference was quantified as the change in value (CS licks or CS DA) relative to baseline level, defined as the mean of the last 5 trials before reversal (see Methods). *middle,* Quantification of change in CS licks ($\Delta$CS licks) as a proxy for inferred value for different number of opposite trials (5, 10 and 20). *right,* similar quantification for dopamine response ($\Delta$normalized dF/F) (n=6 mice). **d,** Similar quantification as in **c** but for downward inference (n=6 mice). **e,** Summary diagram showing the transition from plasticity to dynamics-based value update for meta-reinforcement learning.