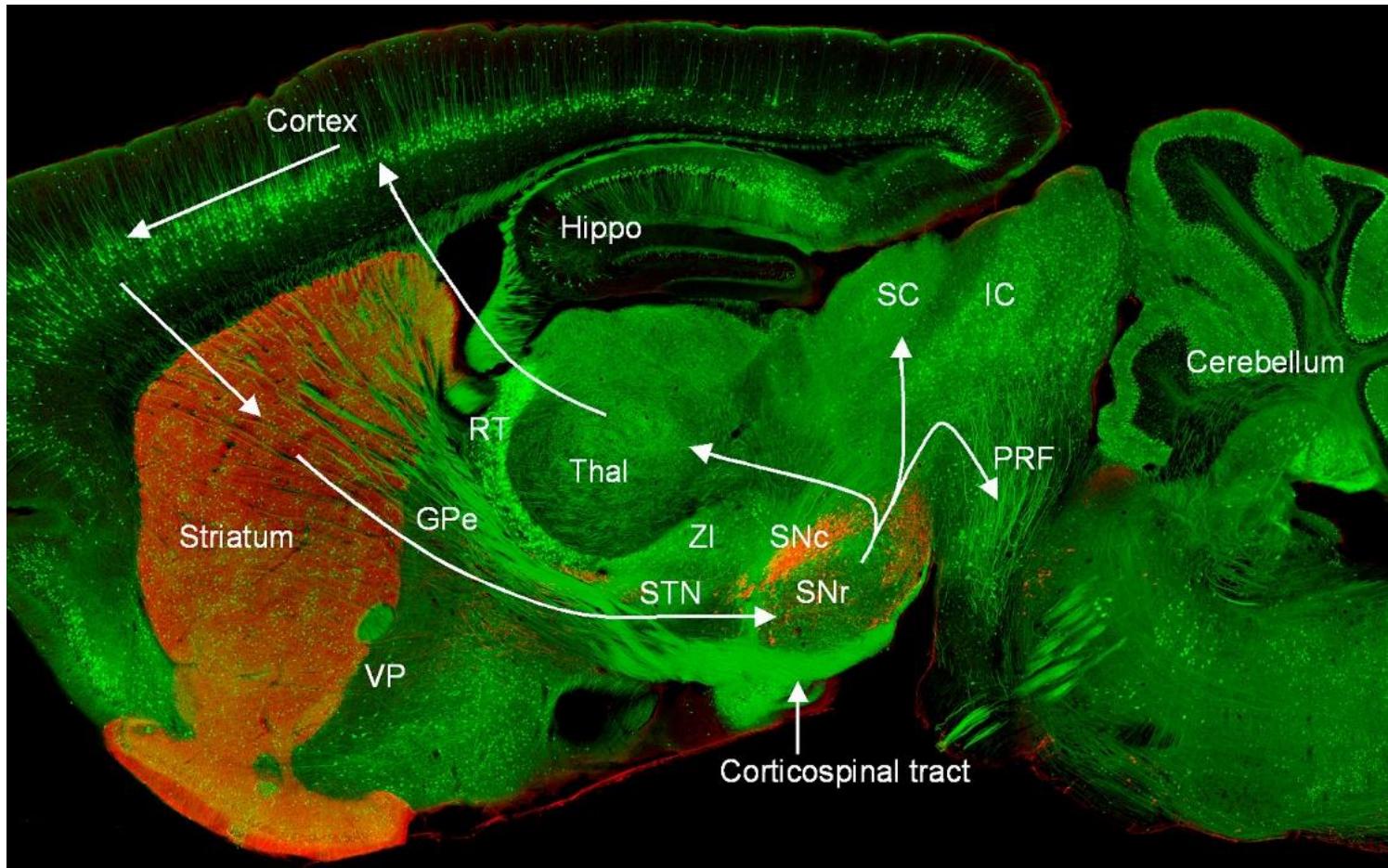


# Brain architecture for adaptive behaviour



Thomas Akam  
RLDM Tutorial  
11/06/2025

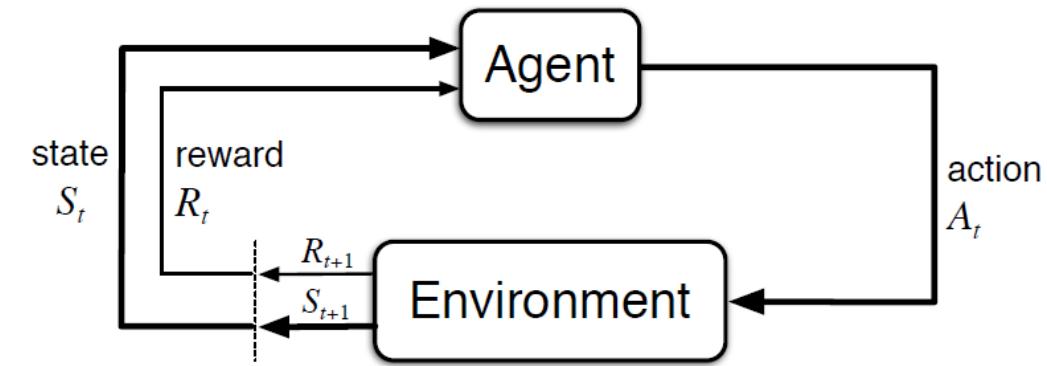
# Behavioural neuroscience and Reinforcement learning

Behavioural neuroscience:

*How do brains generate (evolutionarily) adaptive behaviour?*

Reinforcement learning:

*How can agents learn to maximise long run reward?*



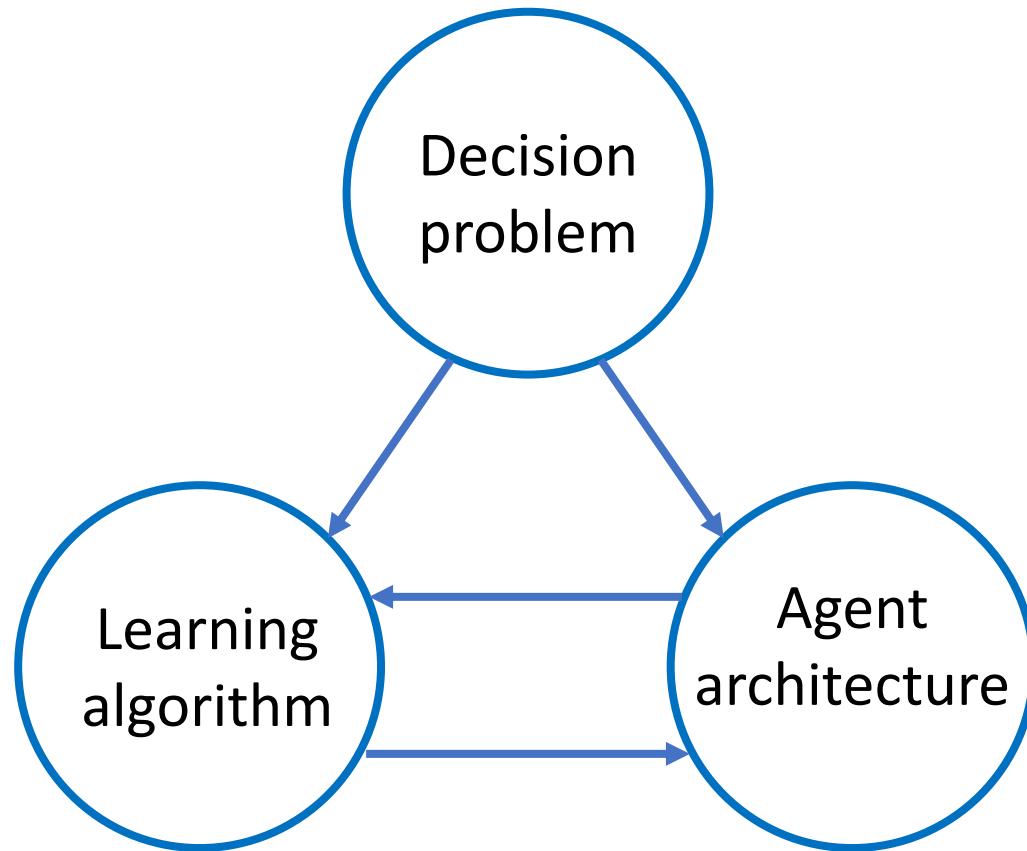
Information content of human genome:

~3 billion base pairs, 2 bit/pair  $\rightarrow$   $\sim 6 \times 10^9$  bits

Information content of human cortex:

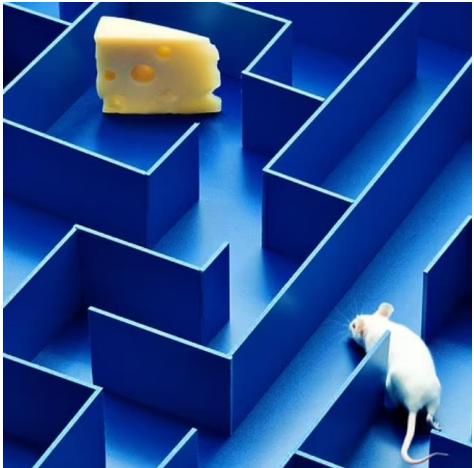
~ $1.6 \times 10^{14}$  synapses, > 1bit/synapse  $\rightarrow$   $> 1.6 \times 10^{14}$  bits

## Behavioural neuroscience and Reinforcement learning



# The brain's control problem

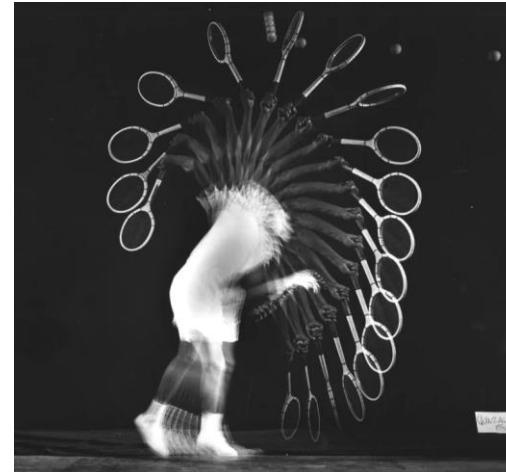
Delayed rewards



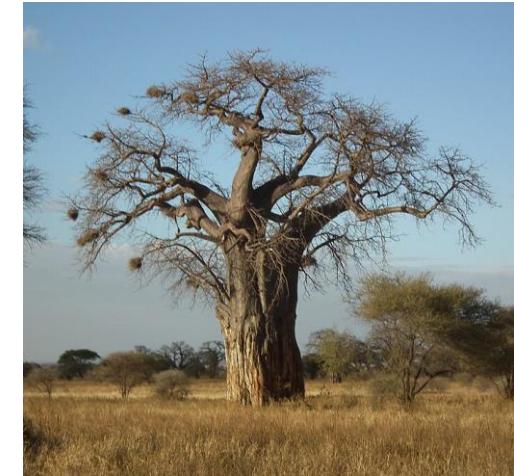
Complex, non-stationary, state space



High-dimensional observation, action space

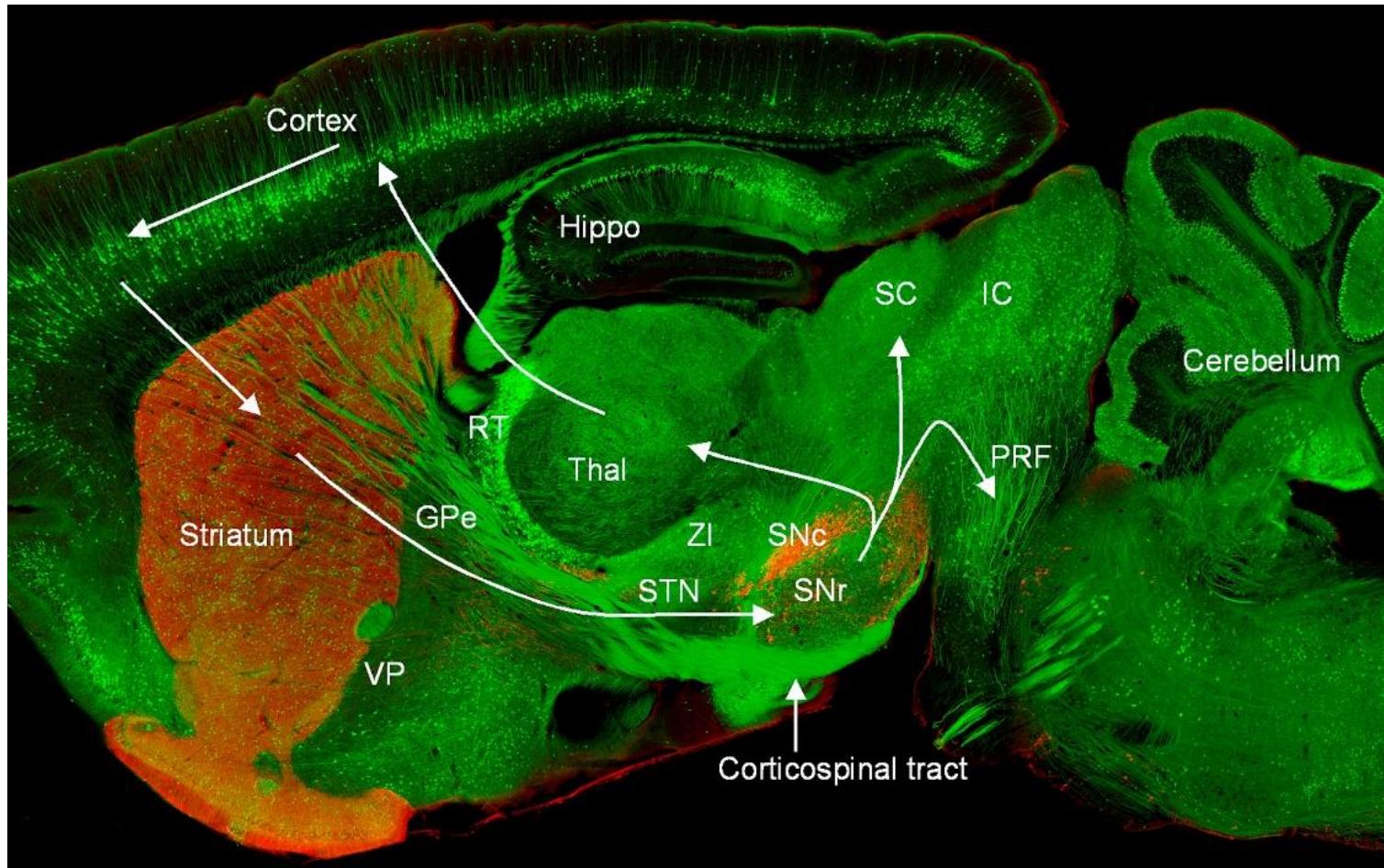


Partial observability



Lion behind tree

## What would a satisfying answer look like?



- What learning algorithm(s) does the brain use?
- How do these map onto the architecture of brain circuits?

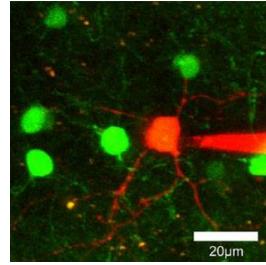
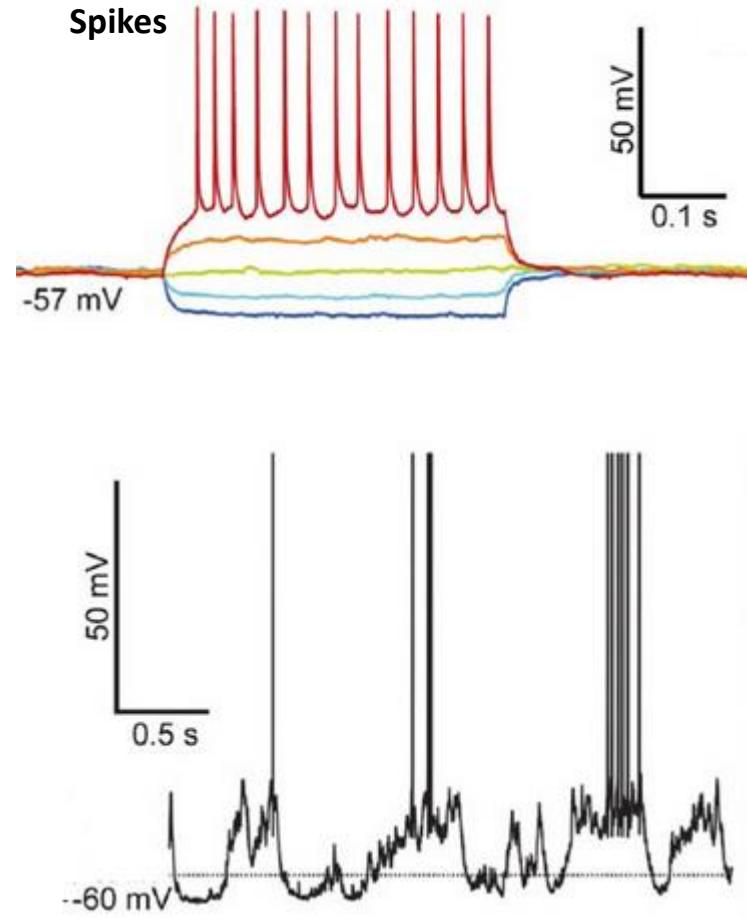
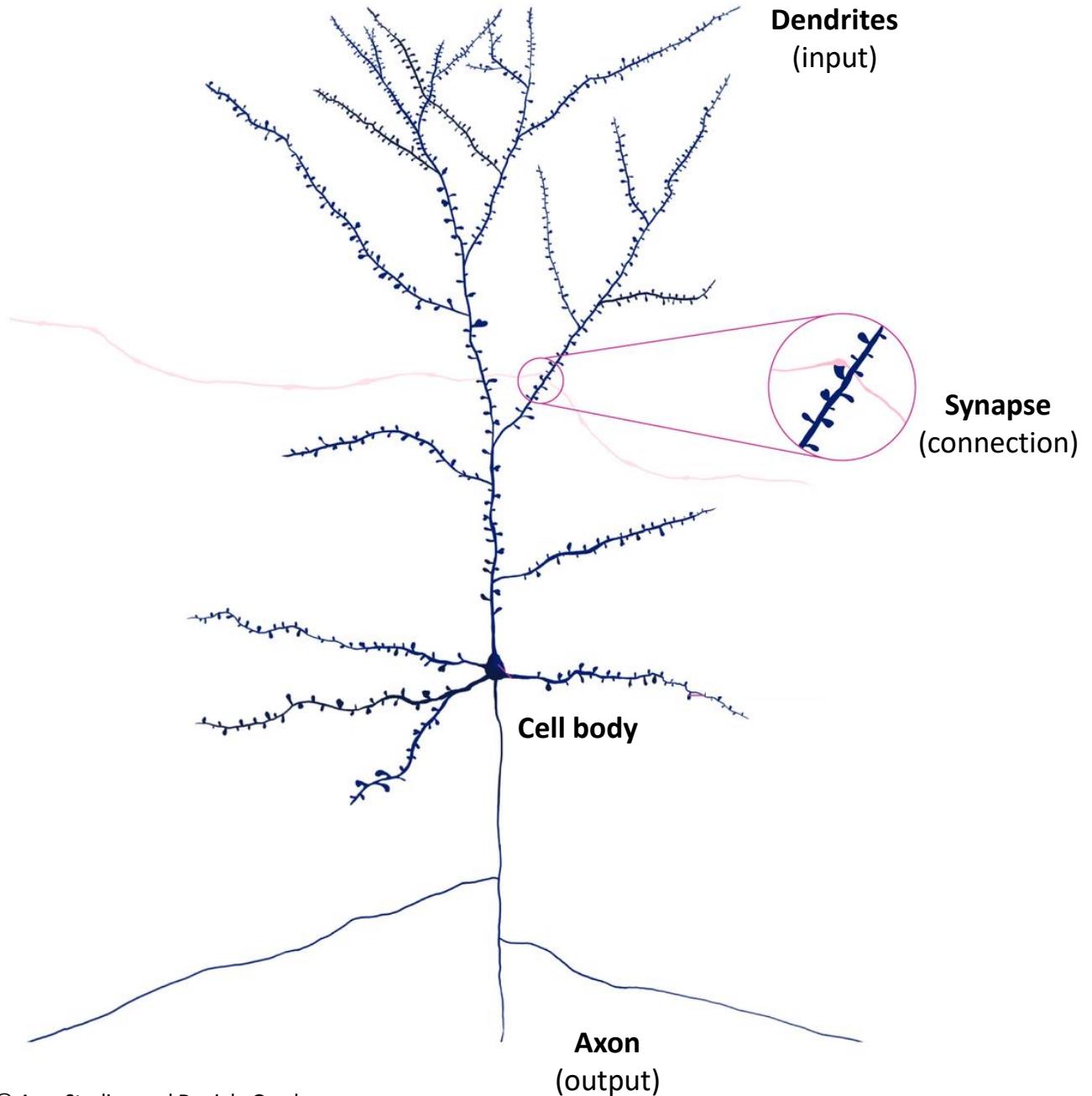
## Talk aim and overview

**Motivating Question:** How do learning algorithms map onto brain structure to solve the biological control problem?

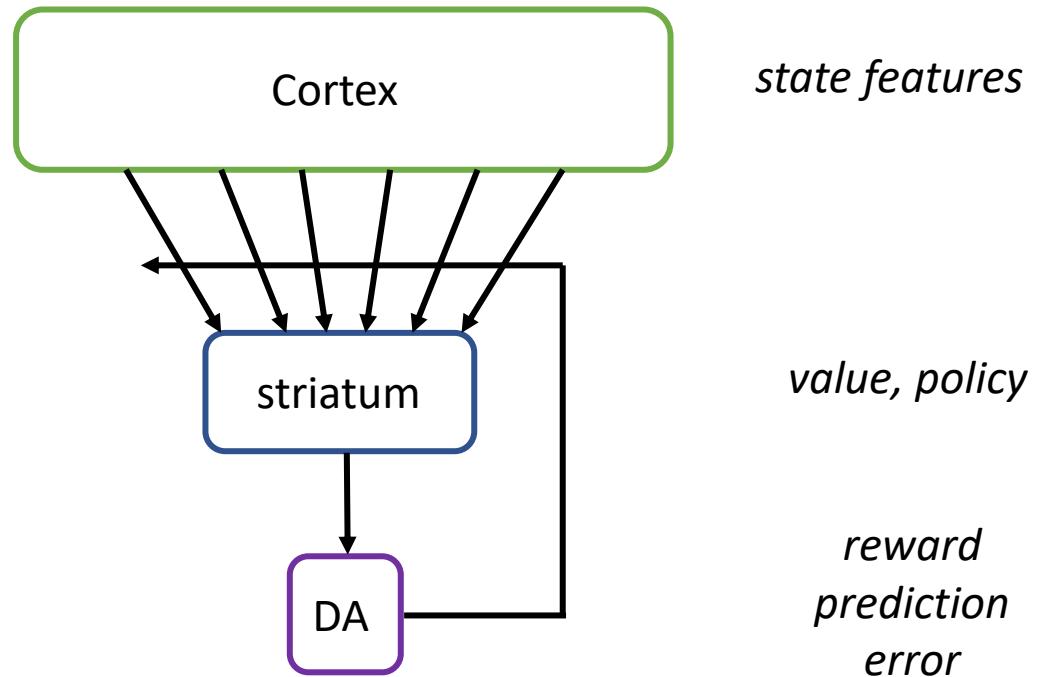
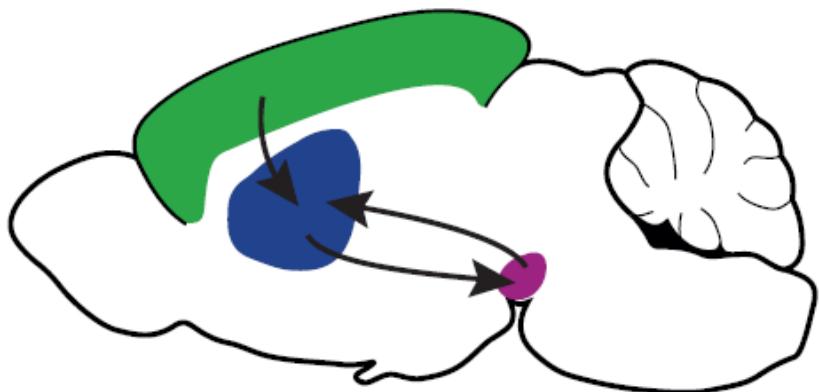
### Talk outline:

1. Striatum and dopamine: The brain's temporal difference reinforcement learning system?
2. Basal ganglia outputs and control of action selection.
3. Cortex: State representation and beyond
4. Hippocampus: Sequence generation and model-based control

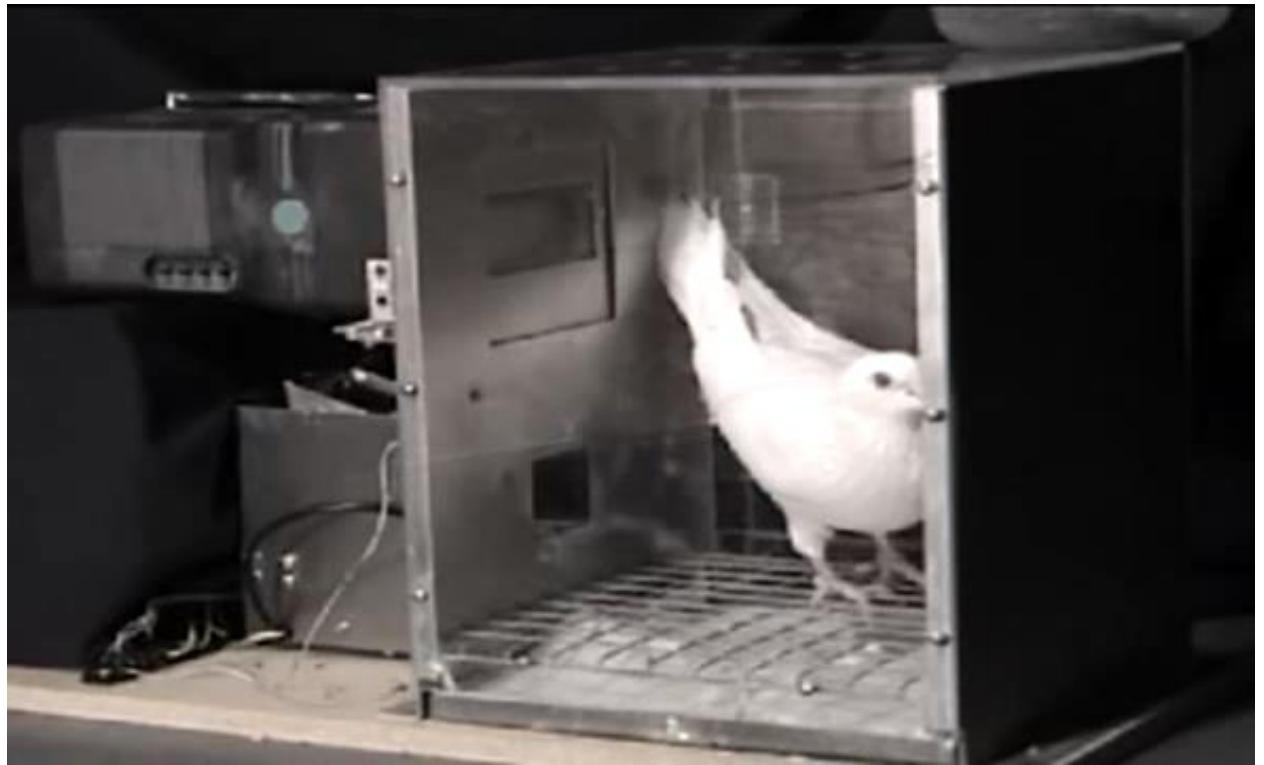
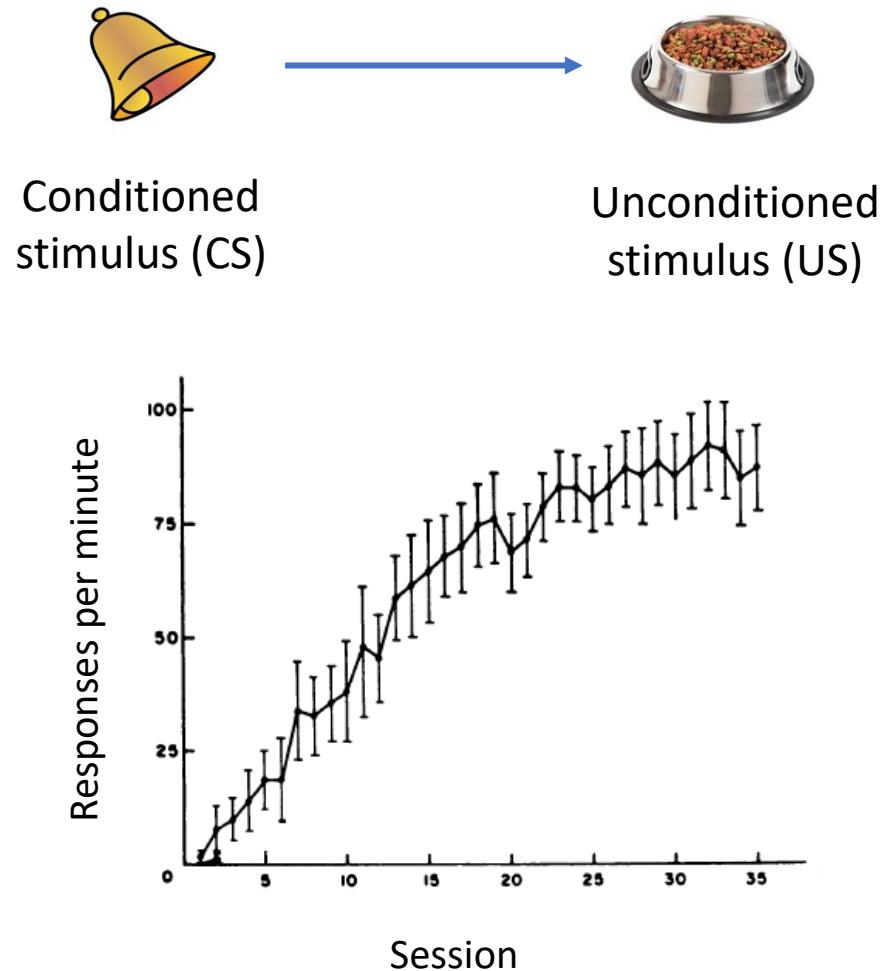
# Neuro 101



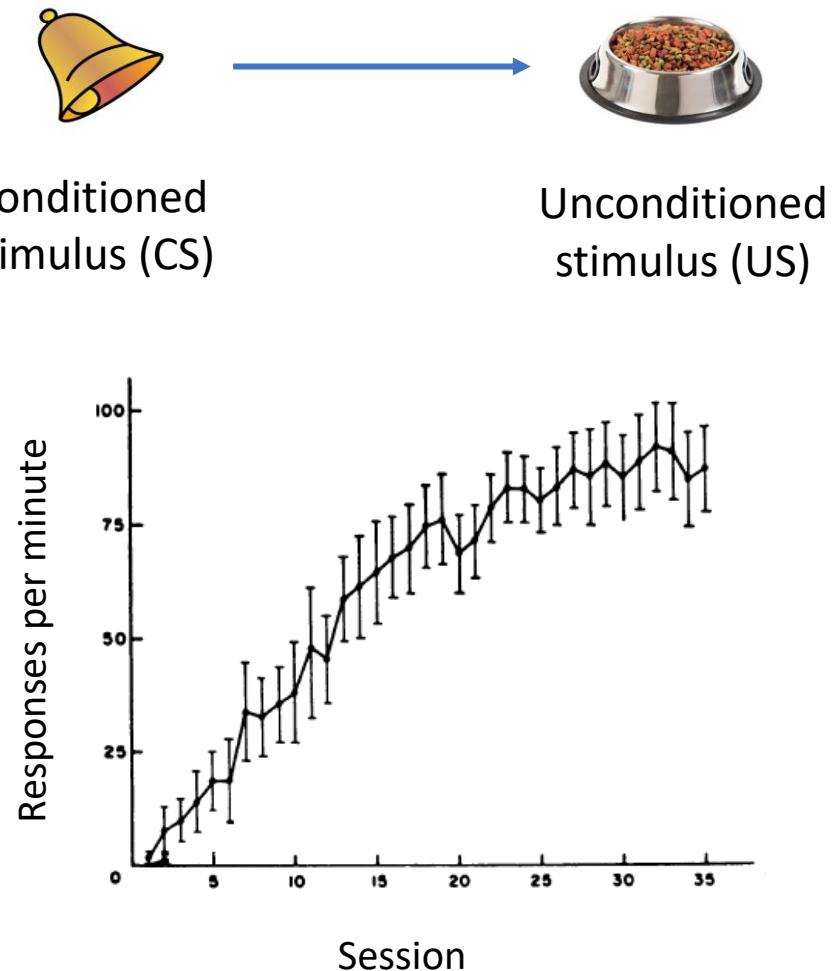
# 1. Striatum and dopamine: The brain's temporal difference (TD) learning system?



## Pavlovian conditioning



## Pavlovian conditioning



## Rescorla – Wagner model (1972)

Prediction:  $V = \sum_i w_i$

Learning:  $\Delta w_i = \alpha(R - \sum_i w_i)$

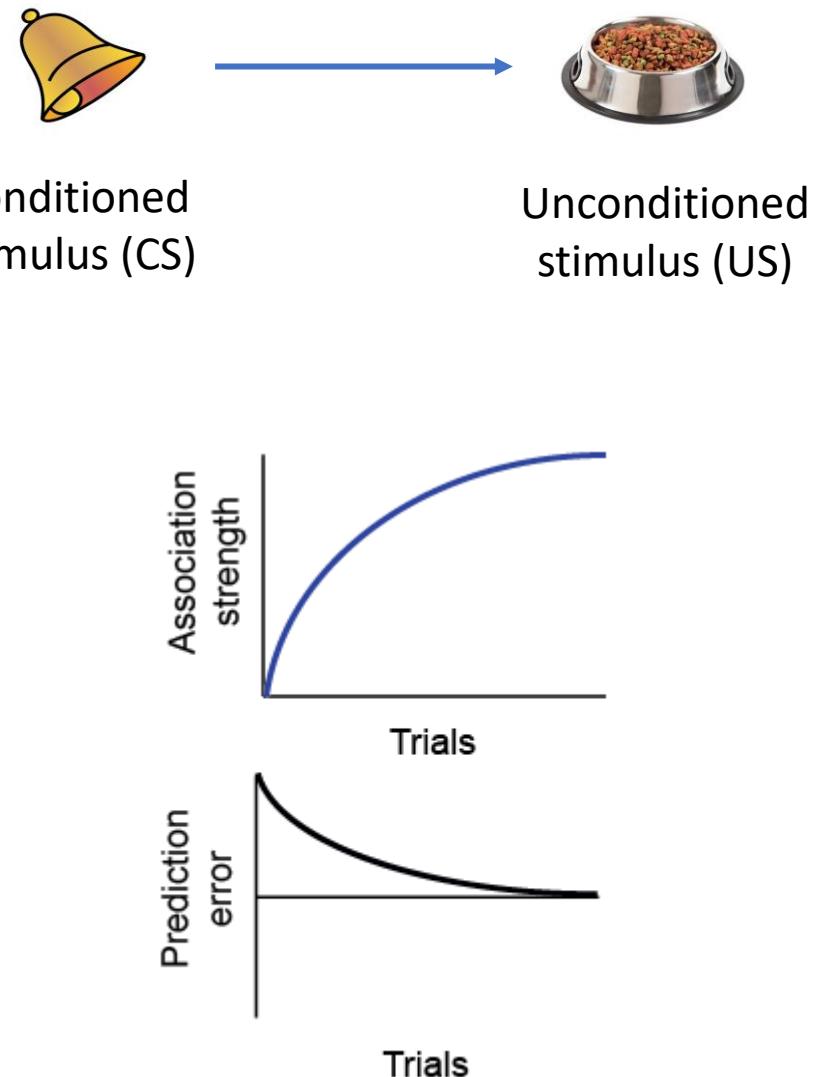
change in association strength

learning rate

strength of US

sum of association strengths for all CS

## Pavlovian conditioning



Rescorla – Wagner model (1972)

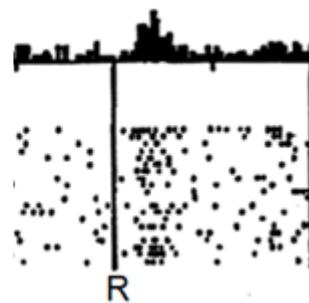
Prediction:  $V = \sum_i w_i$

Learning:  $\Delta w_i = \alpha(R - \sum_i w_i)$

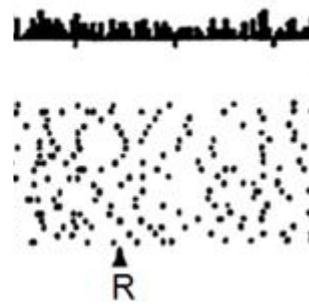
change in association strength  
learning rate  
prediction error

# Dopamine in Pavlovian conditioning

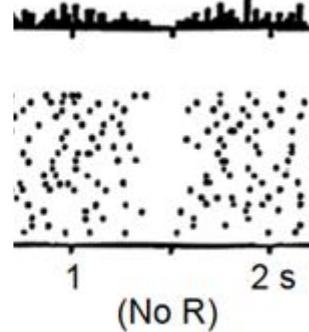
No prediction  
Reward occurs



Reward predicted  
Reward occurs



Reward predicted  
No reward occurs



Activity burst to unexpected rewards

$$\Delta w_i = \alpha(R - \sum_i w_i)$$

Diagram illustrating the change in association strength ( $\Delta w_i$ ) based on the learning rule:

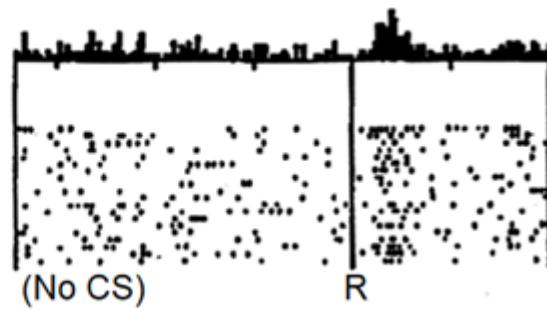
- A blue bracket under the term  $\alpha(R - \sum_i w_i)$  is labeled "learning rate".
- A blue bracket under the term  $R - \sum_i w_i$  is labeled "prediction error".

No response to predicted reward

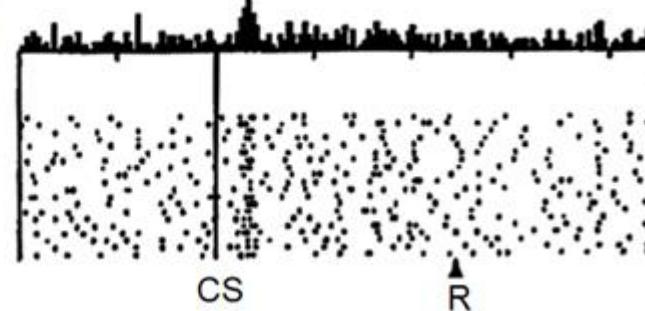
Pause to omission of predicted reward

# Dopamine in Pavlovian conditioning

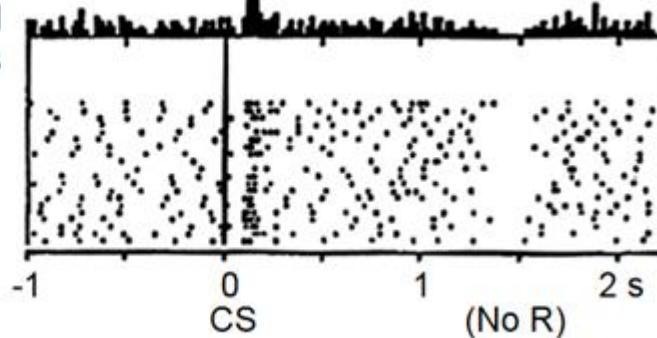
No prediction  
Reward occurs



Reward predicted  
Reward occurs



Reward predicted  
No reward occurs



Activity burst to reward predicting cue

$$\Delta w_i = \alpha(R - \sum_i w_i)$$

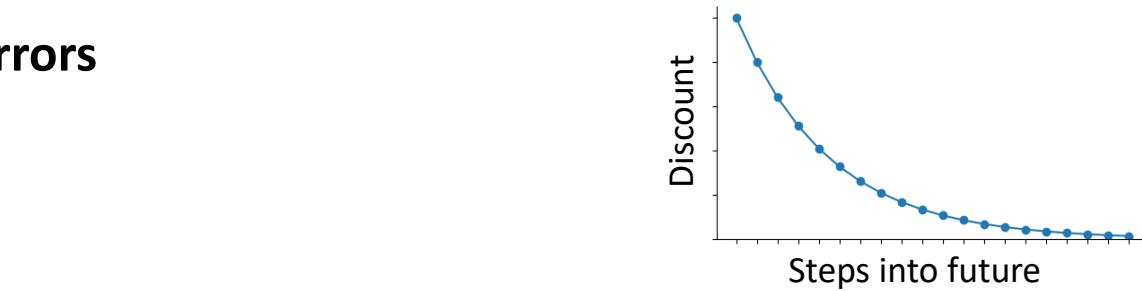
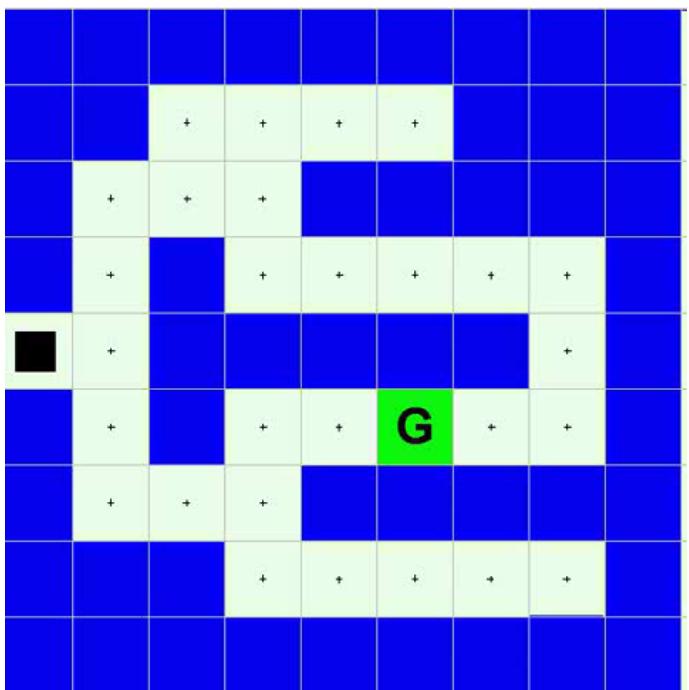
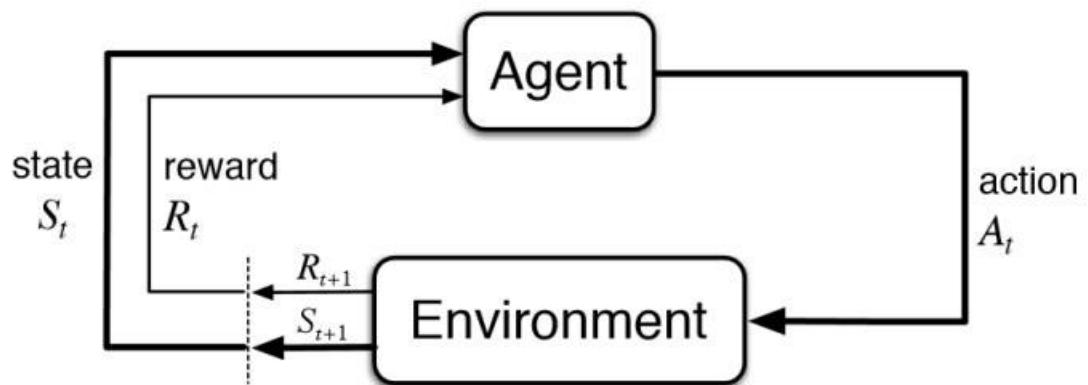
change in association strength

learning rate

prediction error

The diagram illustrates the dopamine learning rule. It shows the equation  $\Delta w_i = \alpha(R - \sum_i w_i)$ . Brackets indicate that the change in association strength ( $\Delta w_i$ ) is determined by the learning rate ( $\alpha$ ) and the prediction error ( $R - \sum_i w_i$ ). The prediction error is the difference between the actual reward ( $R$ ) and the predicted reward ( $\sum_i w_i$ ).

## Delayed rewards & temporal difference reward prediction errors



$$V(s_t) = E \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \right]$$

long run value  
of state  $s$

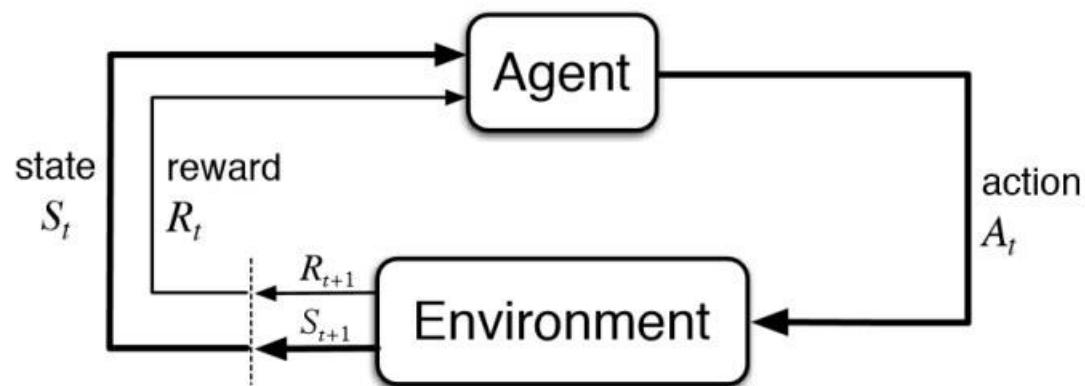
sum of discounted  
future rewards

$$V(s_t) = R_{t+1} + \gamma E[V(s_{t+1})]$$

Immediate  
reward

Discounted value  
of next state

## Delayed rewards & temporal difference reward prediction errors



$$\delta_t = \frac{R_{t+1} + \gamma V(s_{t+1})}{\text{reward prediction error}} - \frac{V(s_t)}{\text{prediction}}$$

Update target

prediction

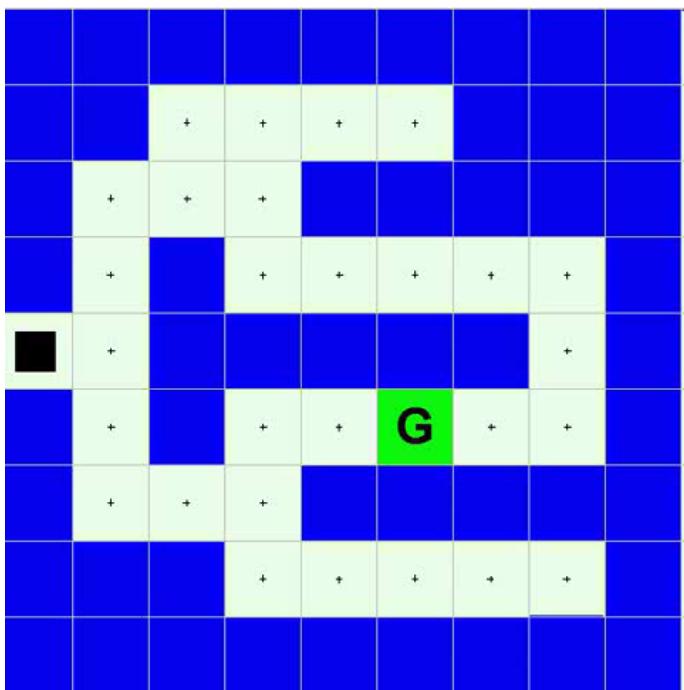
$$V(s_t) \leftarrow \frac{V(s_t)}{\text{old value estimate}} + \frac{\alpha \delta_t}{\text{learning rate}}$$

Updated value estimate

old value estimate

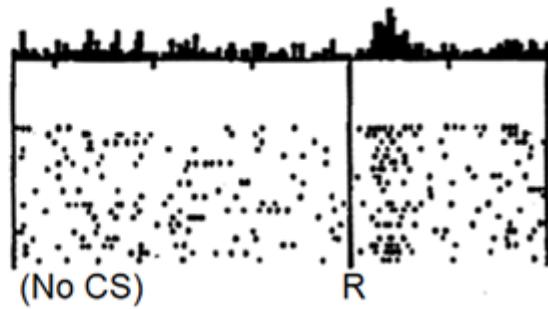
learning rate

reward prediction error

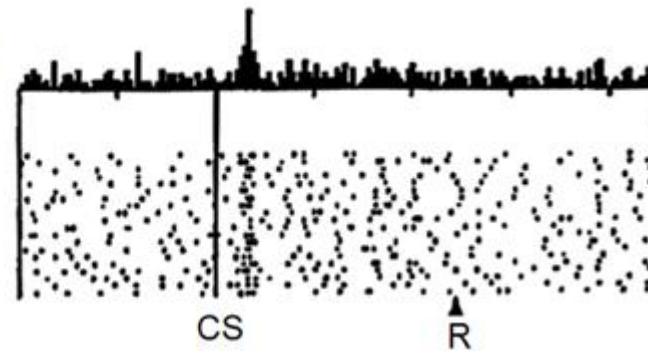


## Dopamine as a reward prediction error signal

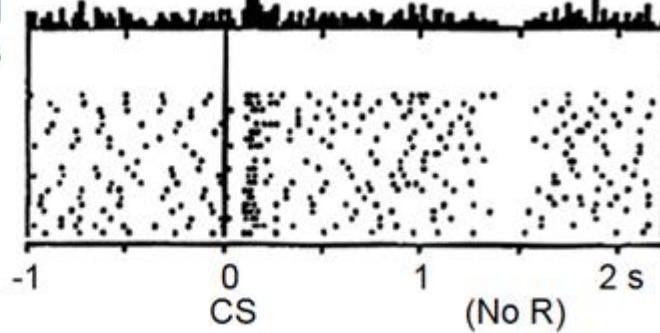
No prediction  
Reward occurs



Reward predicted  
Reward occurs



Reward predicted  
No reward occurs

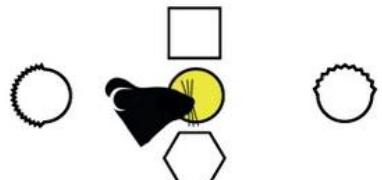


$$\delta_t = \underline{R_{t+1}} + \gamma \underline{V(s_{t+1})} - \underline{\underline{V(s_t)}}$$

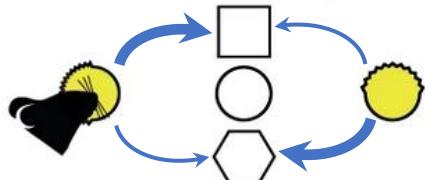
reward prediction error      immediate reward      value of new state      value of previous state

# Dopamine as a reward prediction error signal

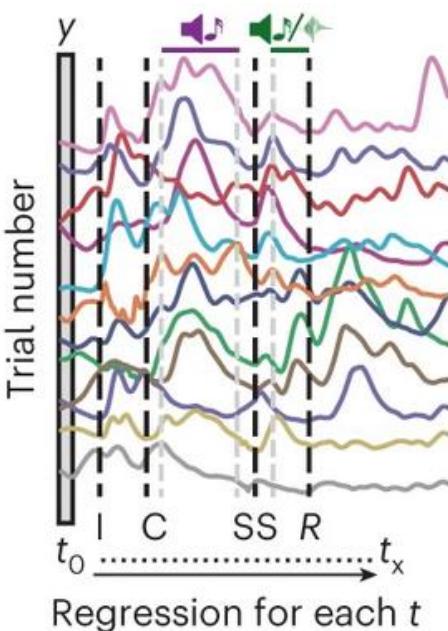
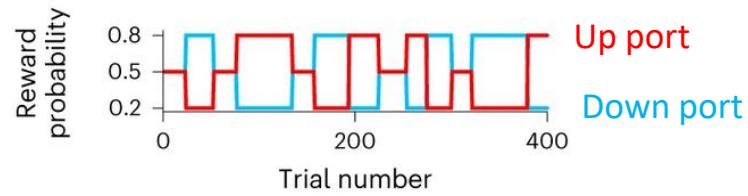
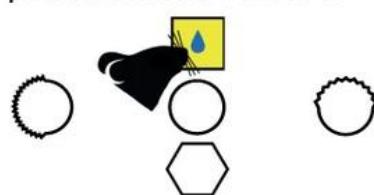
1. Initiate trial in center port



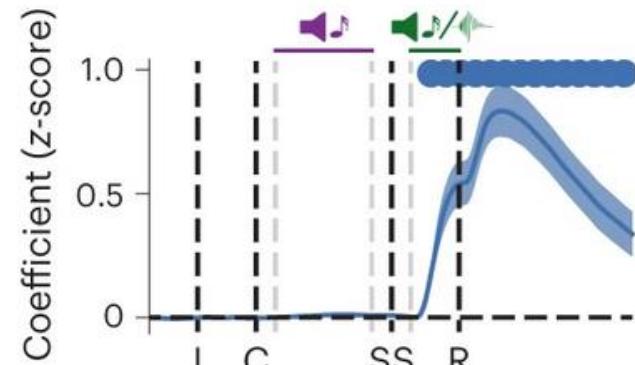
2. Choose between left and right ports



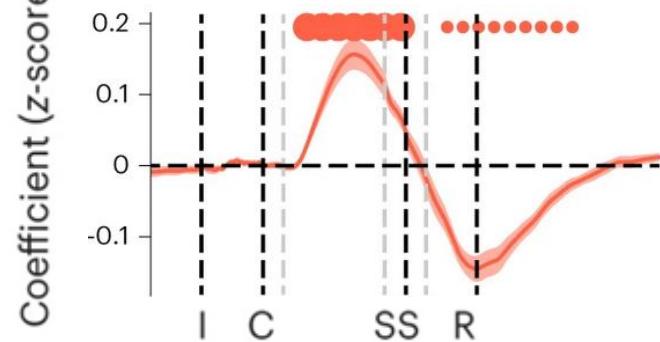
3. Poke active up/down port for probabilistic reward



Trial outcome (rewarded or not)



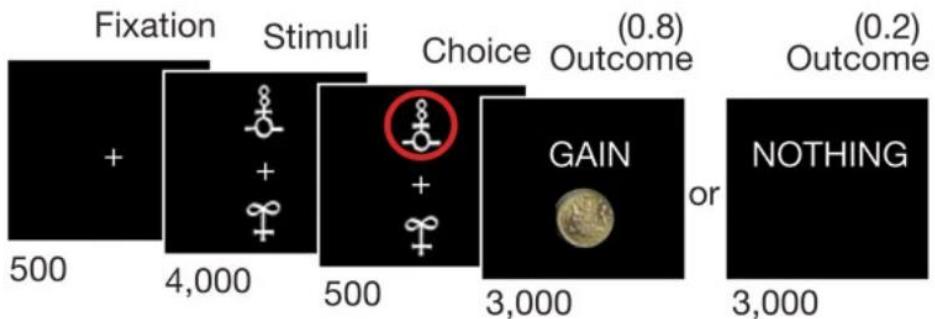
Value of up / down port



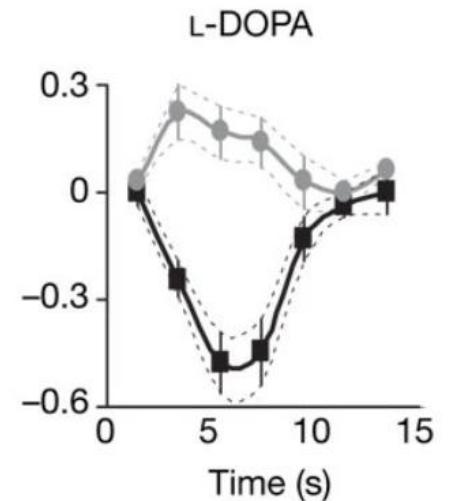
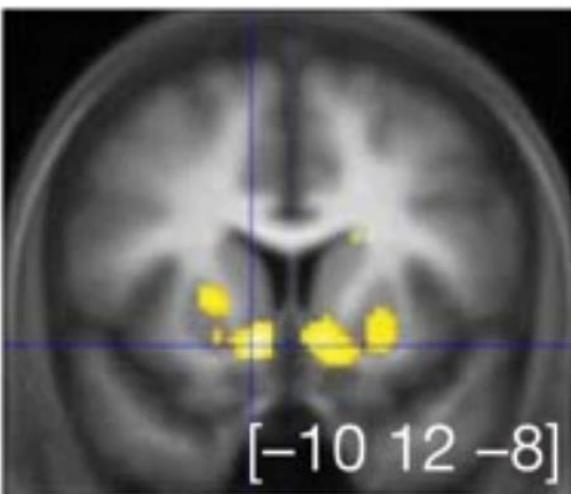
$$\delta_t = R_{t+1} + \gamma V(s_{t+1}) - V(s_t)$$

reward prediction error
immediate reward
value of new state
value of previous state

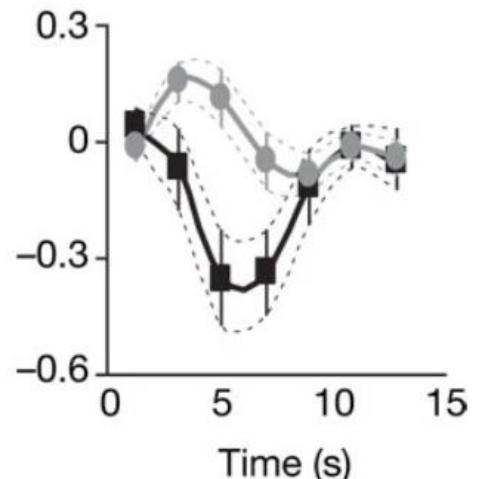
# Dopamine as a reward prediction error signal



## Voxels correlated with RPE

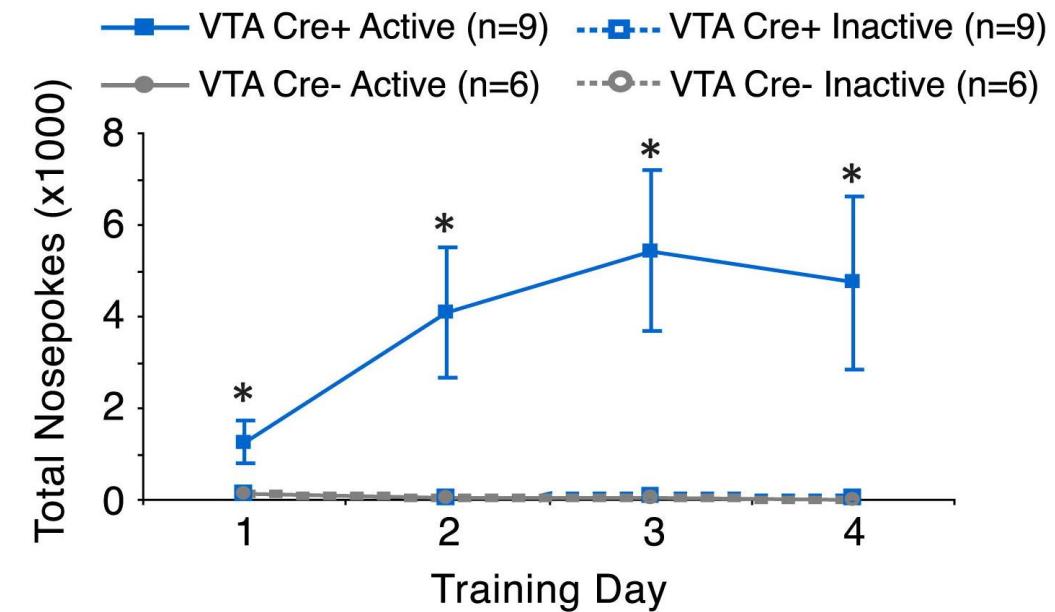
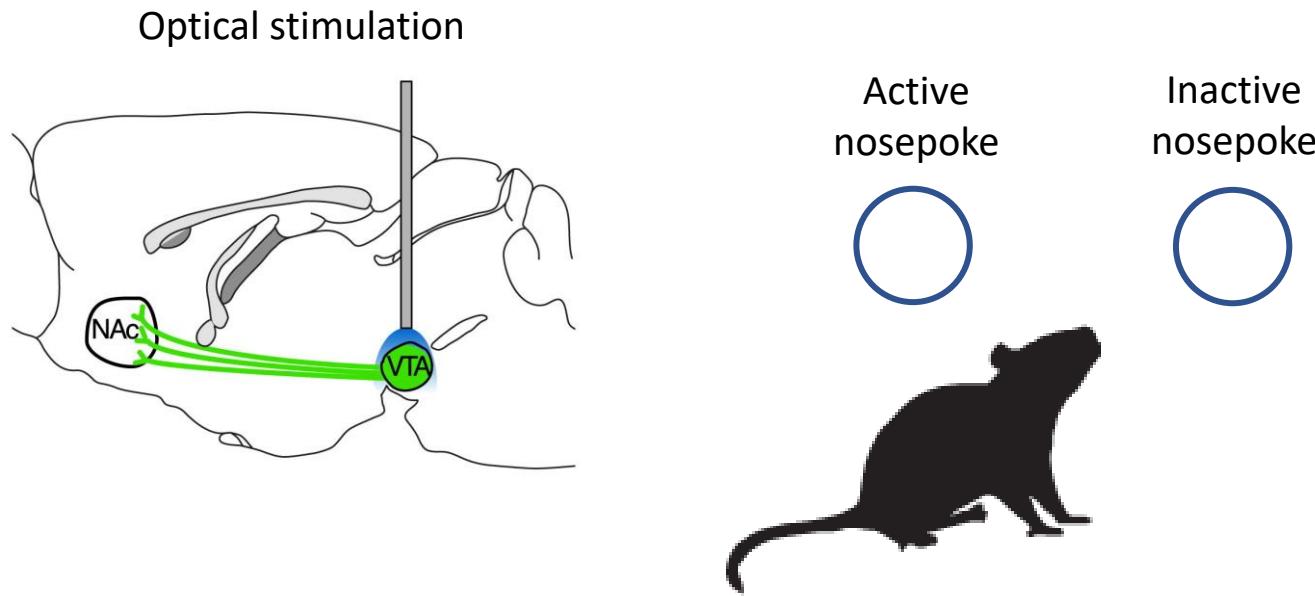


### **Haloperidol**



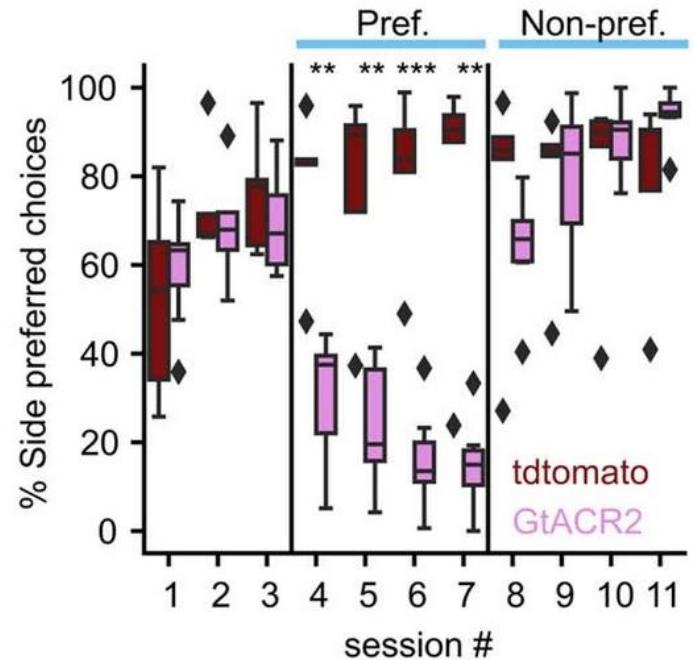
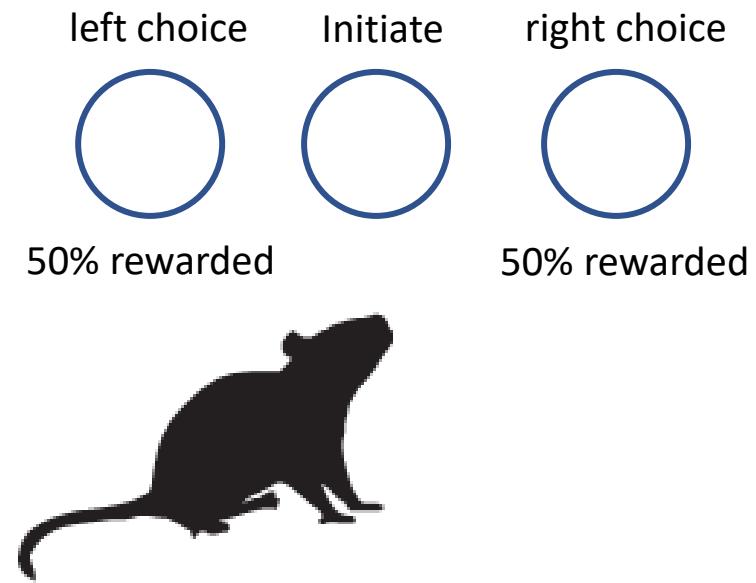
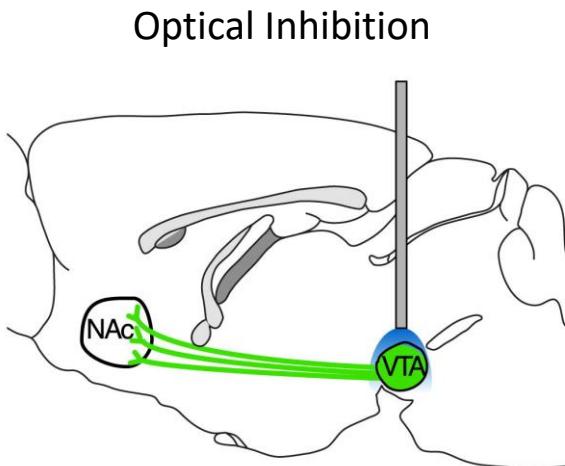
- Brain activity in human ventral striatum correlates with RPE.
  - Boosting / blocking dopamine transmission modulates RPE signal size

## Causal manipulation of dopamine and reinforcement



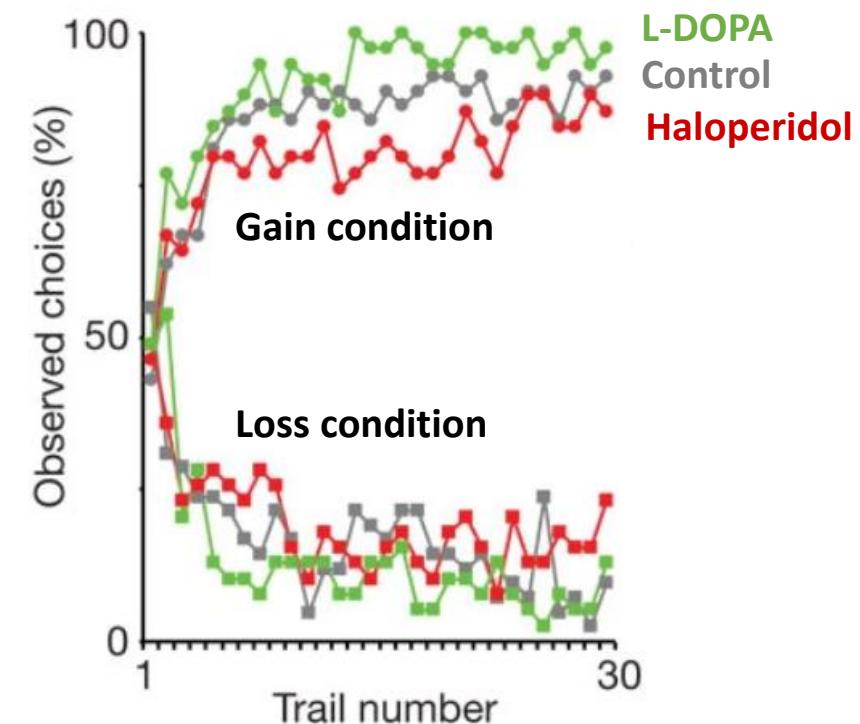
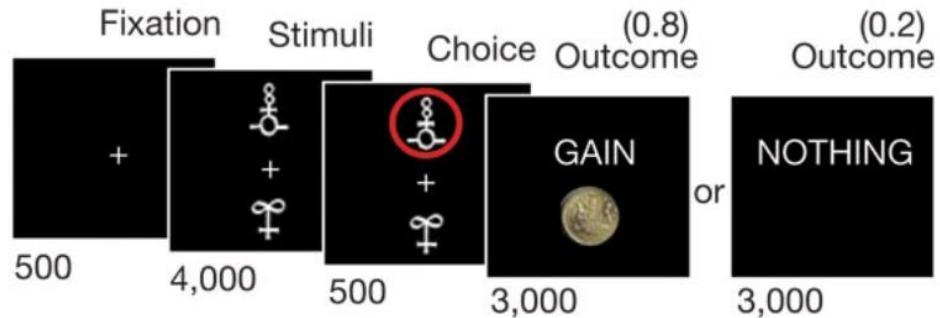
- Optogenetic stimulation of dopamine neurons contingent on action is reinforcing.

# Causal manipulation of dopamine and reinforcement



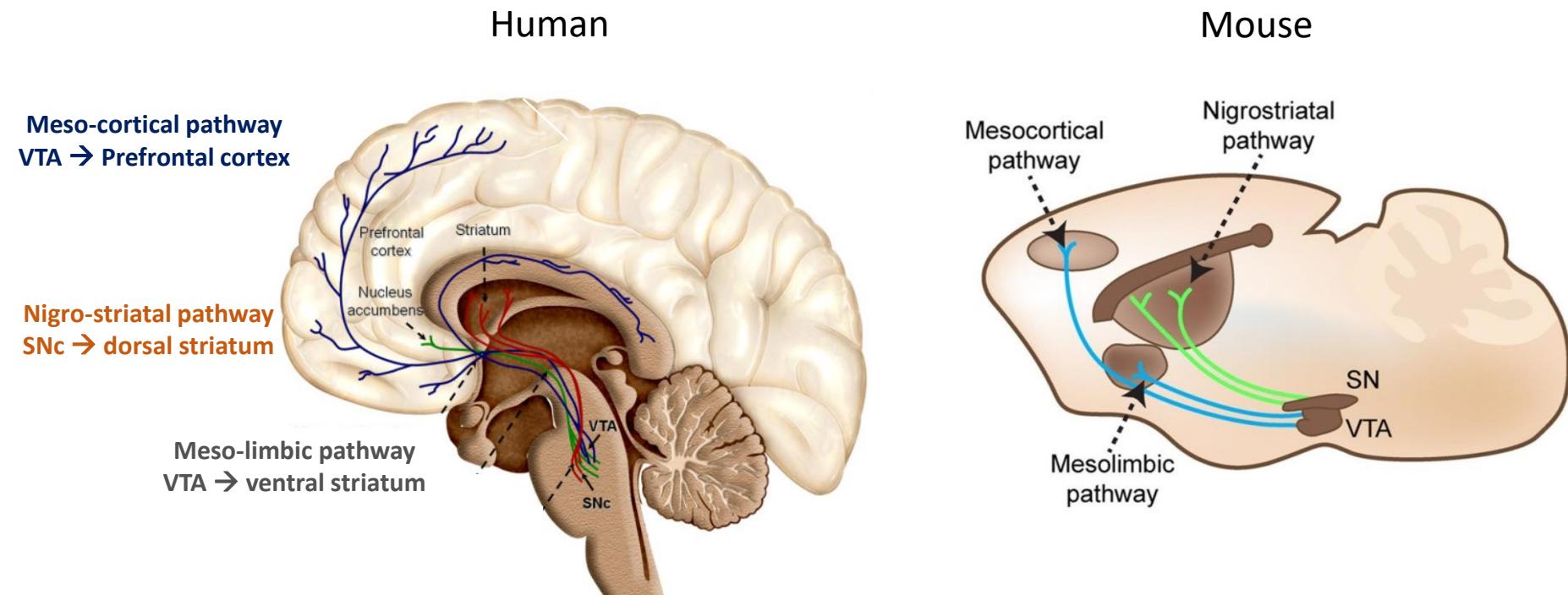
- Optogenetic inhibition of dopamine neurons suppresses preceding actions

## Causal manipulation of dopamine and reinforcement

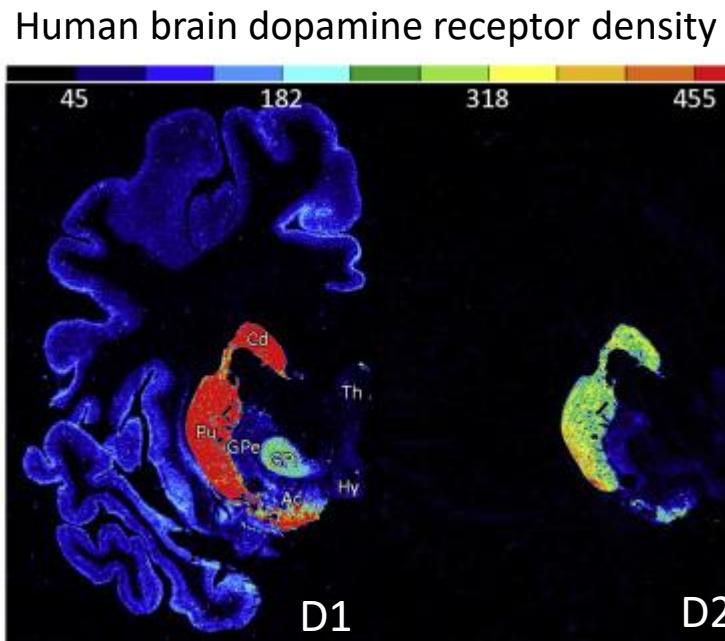


- Boosting / blocking dopamine transmission modulates human learning

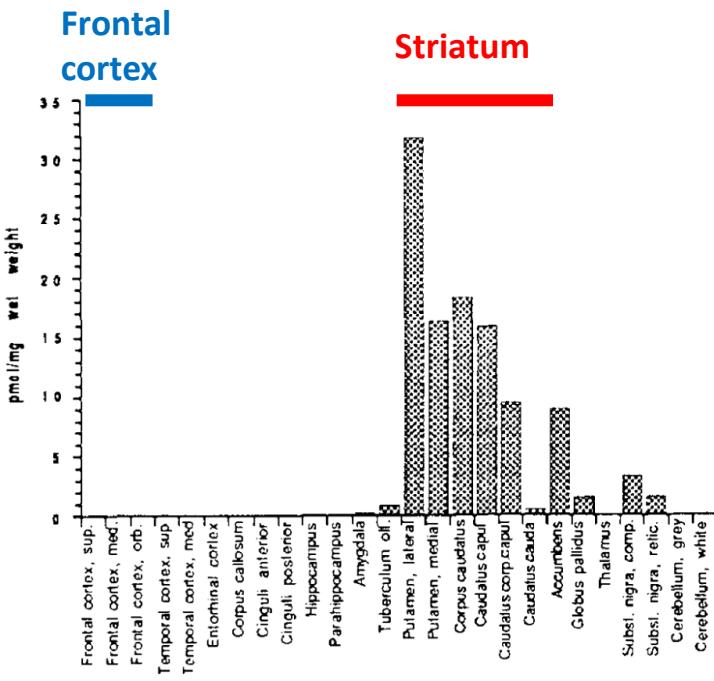
# Dopamine projections



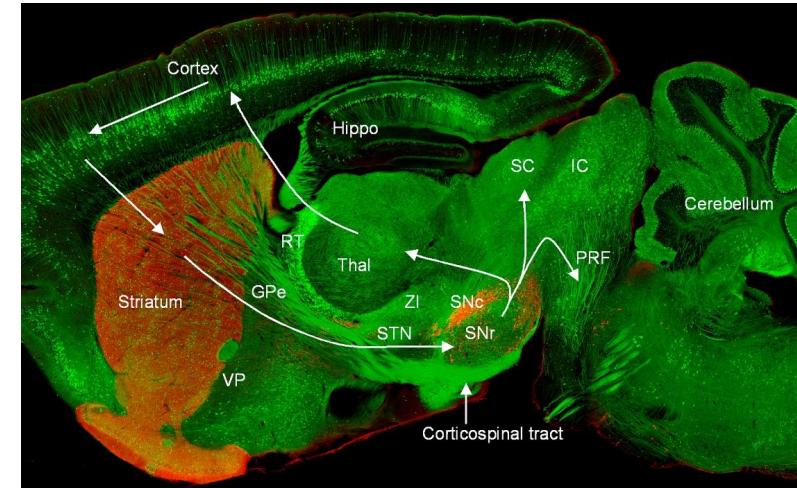
# Dopamine projections



Human brain dopamine concentration

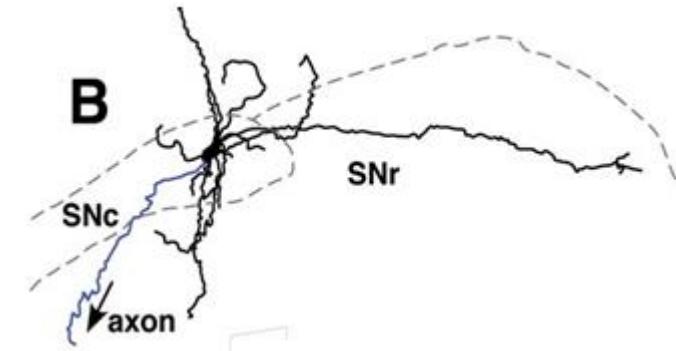
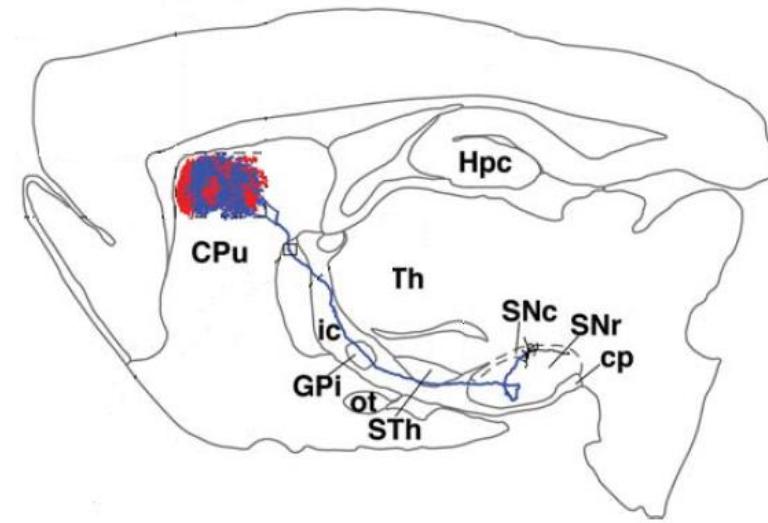
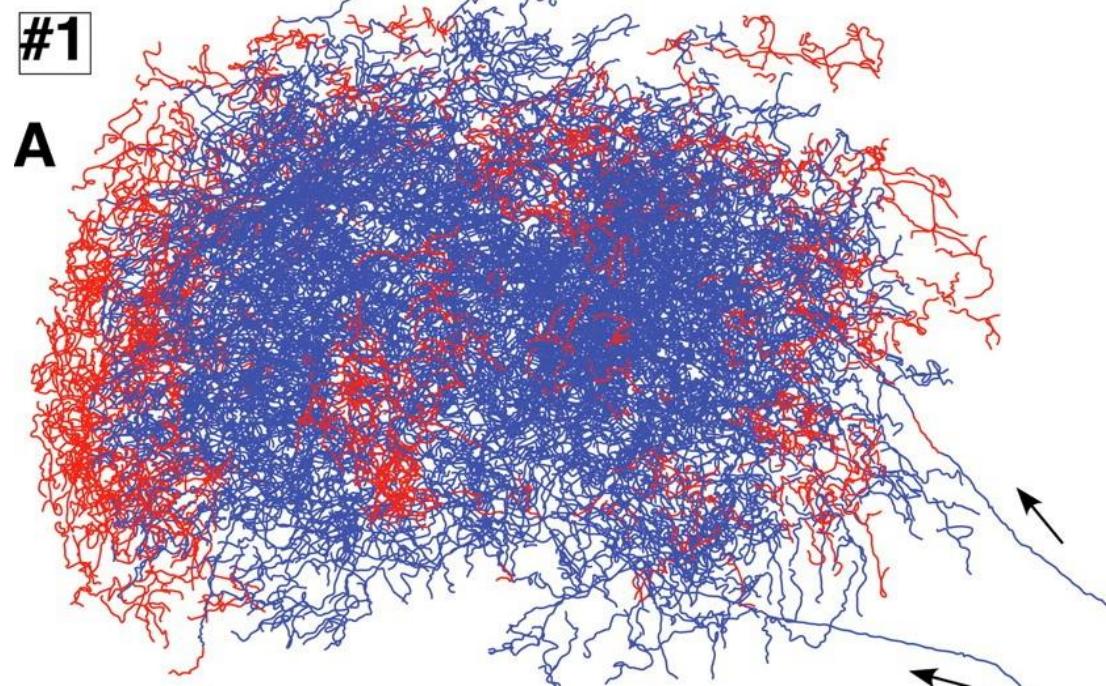


Mouse brain TH stain (dopamine neurons)



- Dopamine receptor density and dopamine neuron innervation are much greater in striatum than frontal cortex

## Dopamine projections



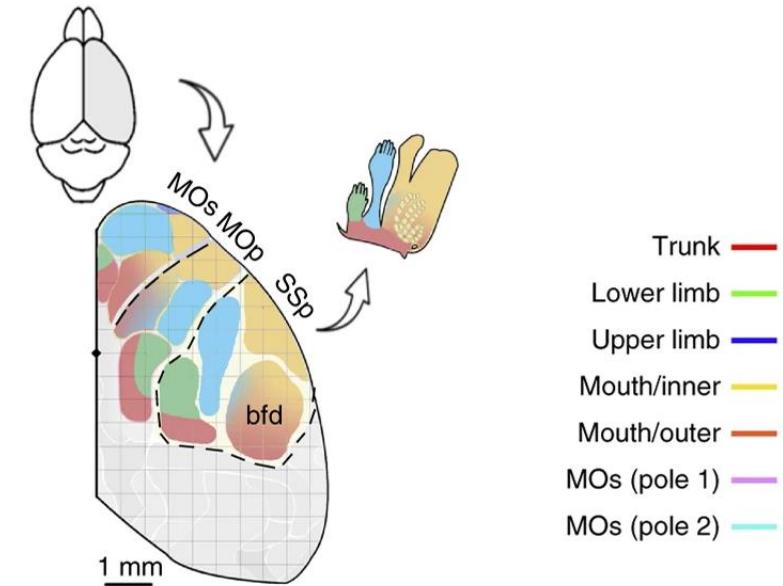
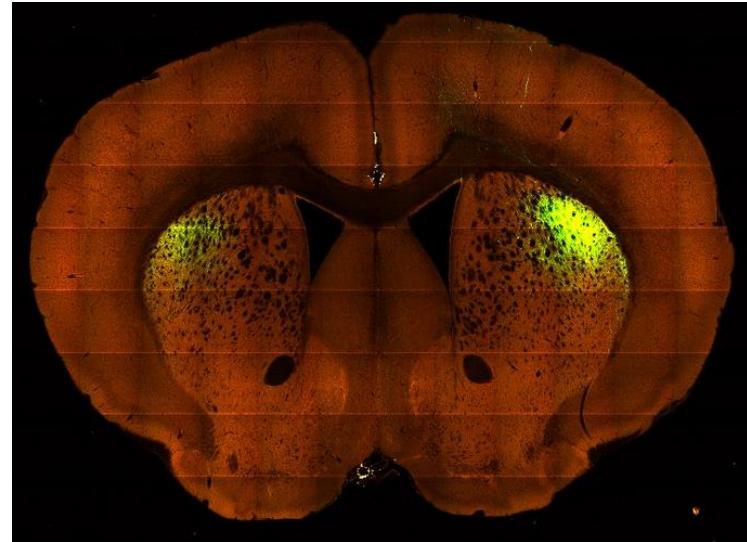
- Individual dopamine neurons densely innervate a large volume of striatum with 250,000 – 400,000 synapses per cell.

# Cortical input to striatum

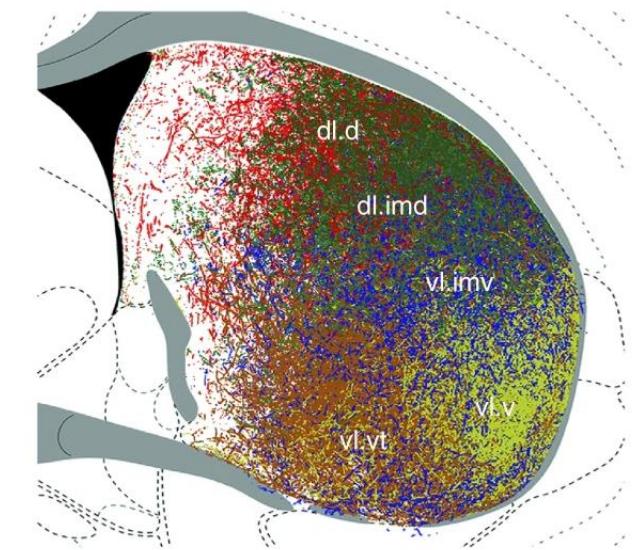
Anterior cingulate cortex →  
dorsomedial striatum



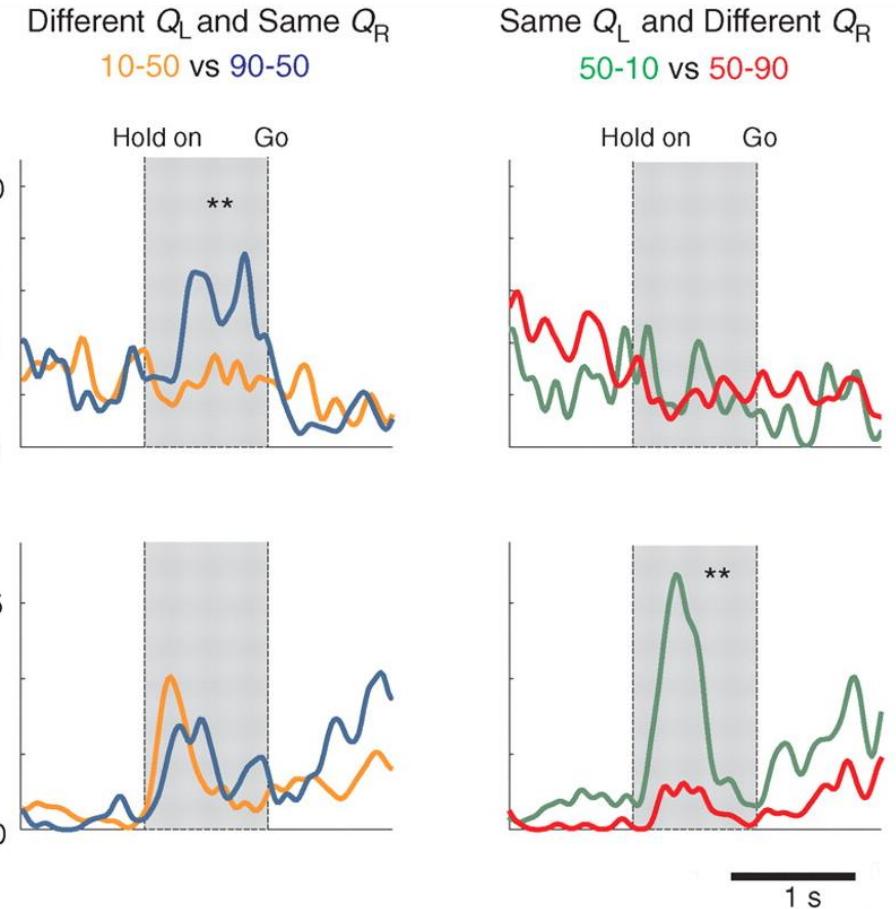
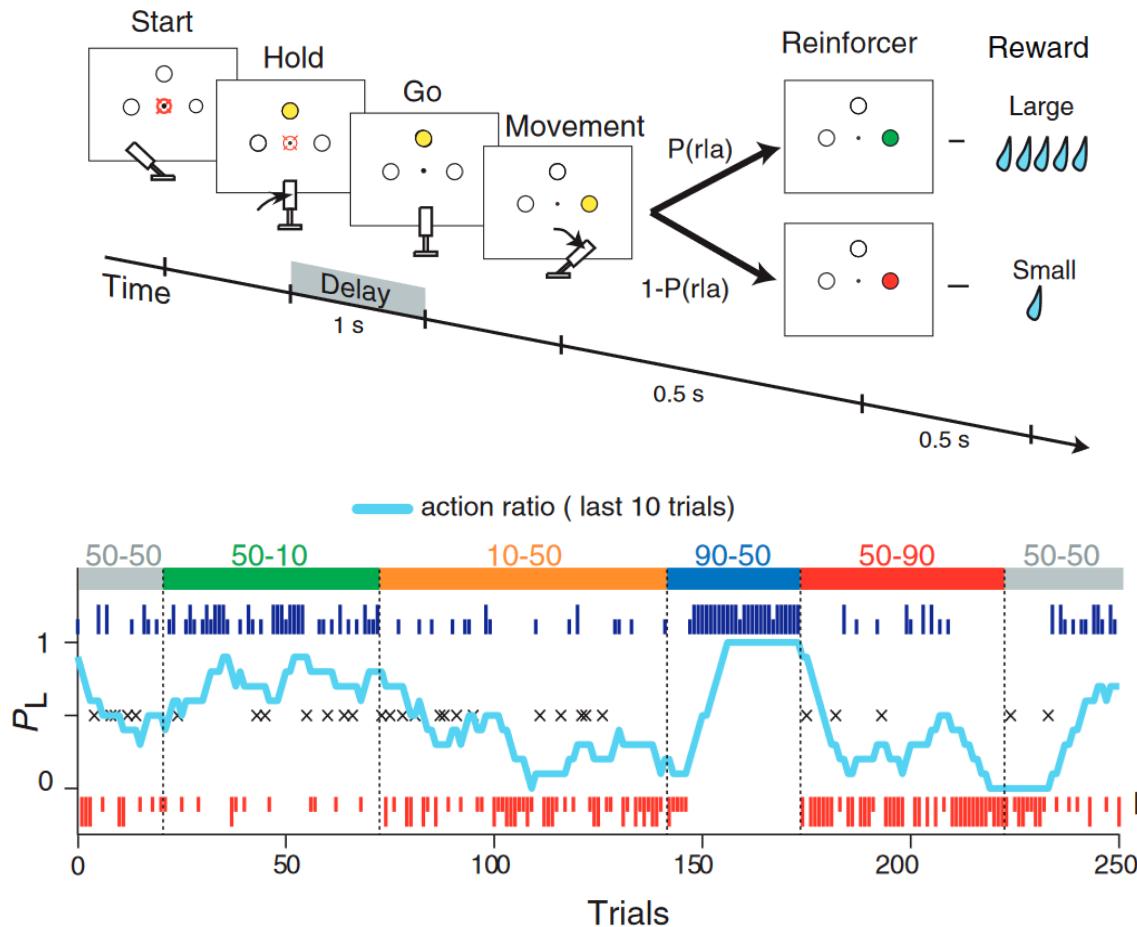
Primary somatosensory cortex →  
dorsolateral striatum



- Massive, highly structured, projection from cortex to striatum

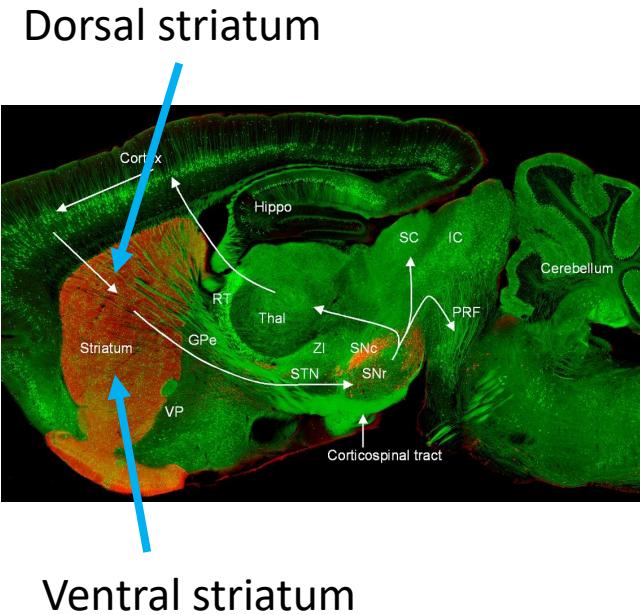
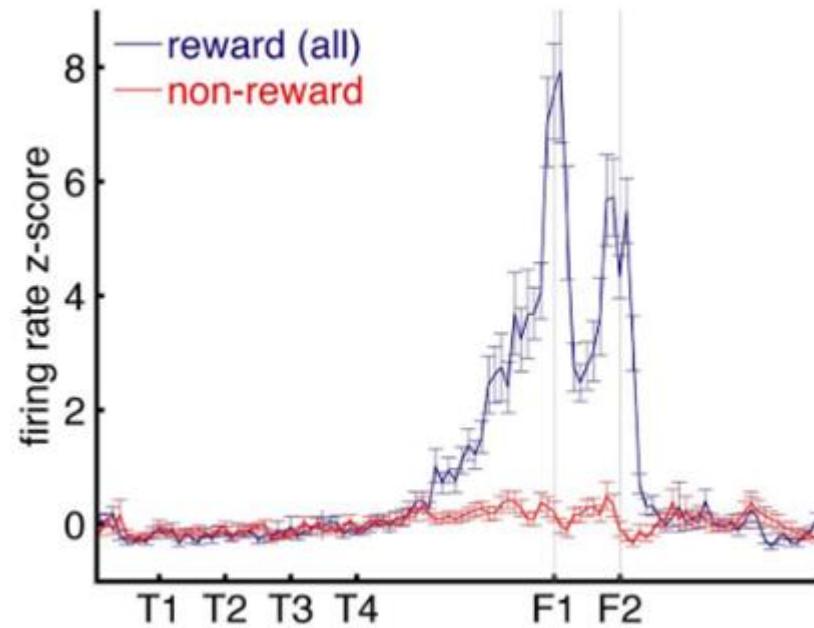
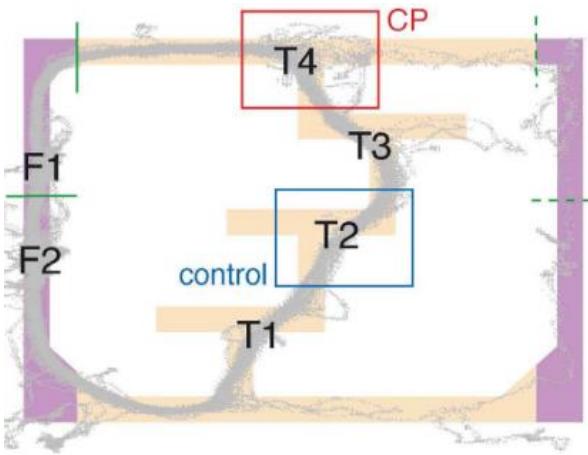


# Striatal action value signals



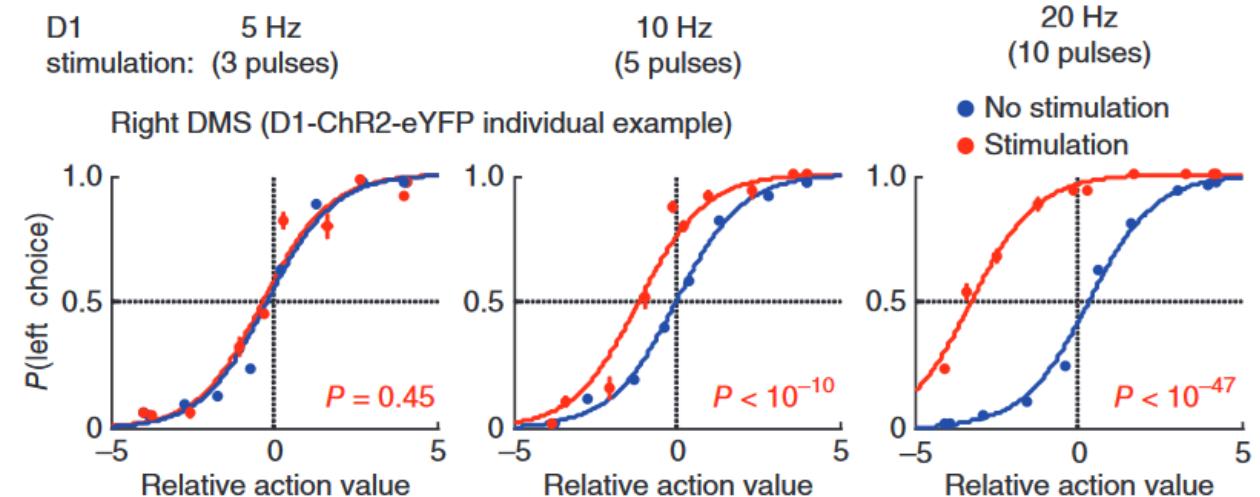
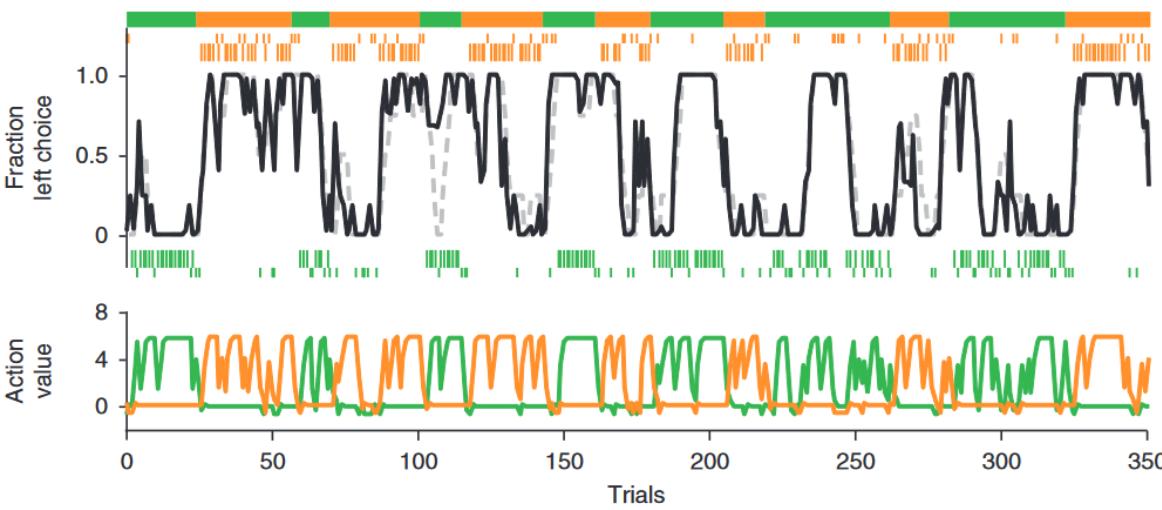
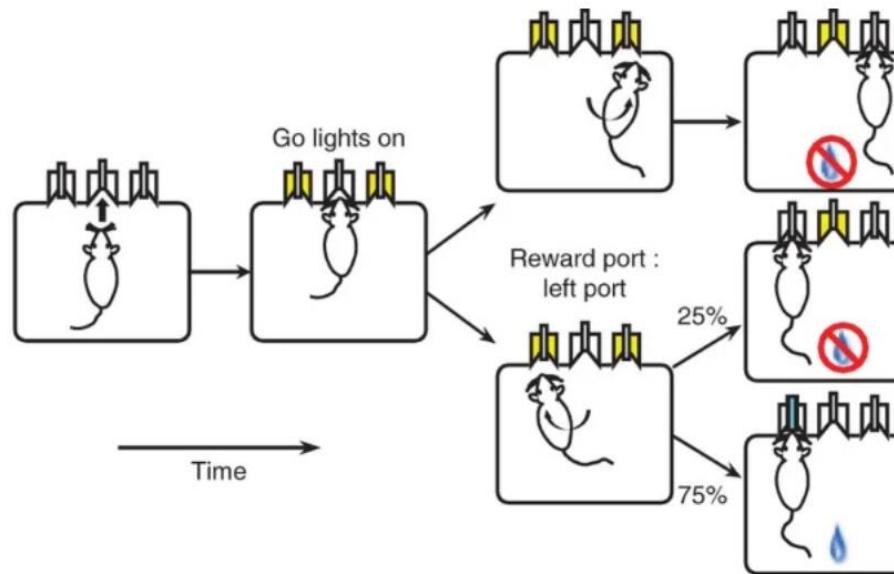
- Striatal neuron activity tracks the value of specific motor actions in a dynamic reward guided decision task

## State-value like signals in ventral striatum



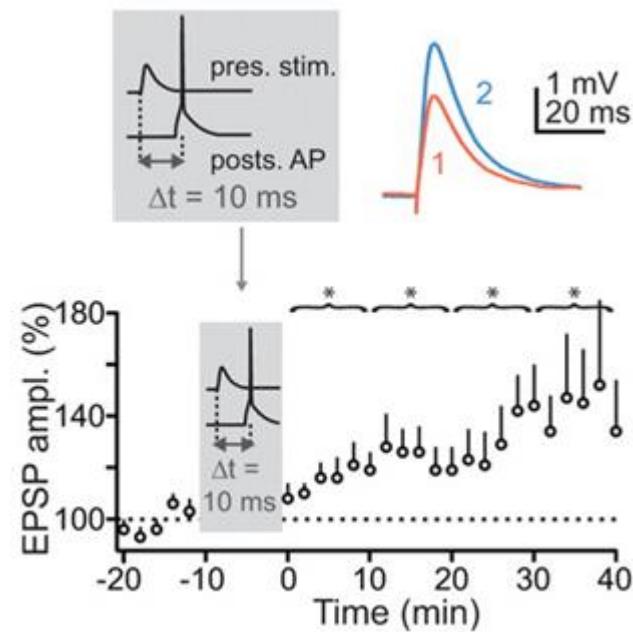
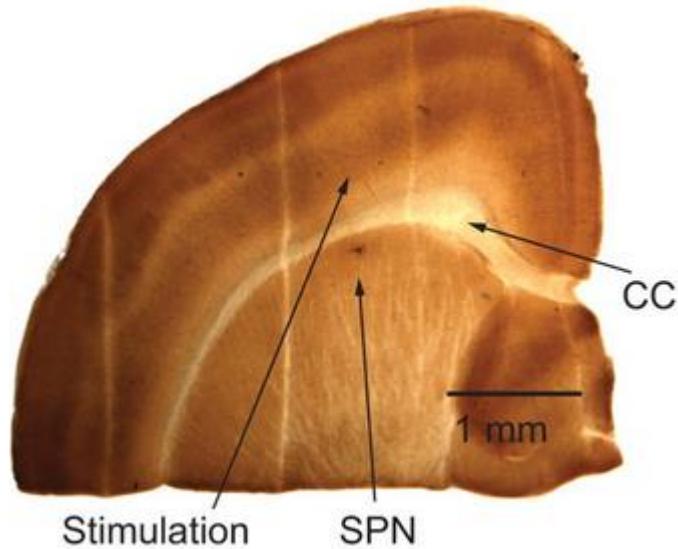
- Ramping activity in ventral striatum during approach to reward resembles state value.
- Differential involvement of ventral/dorsal striatum in state value / action value or policy also seen in human fMRI.

## Striatal Causal manipulations

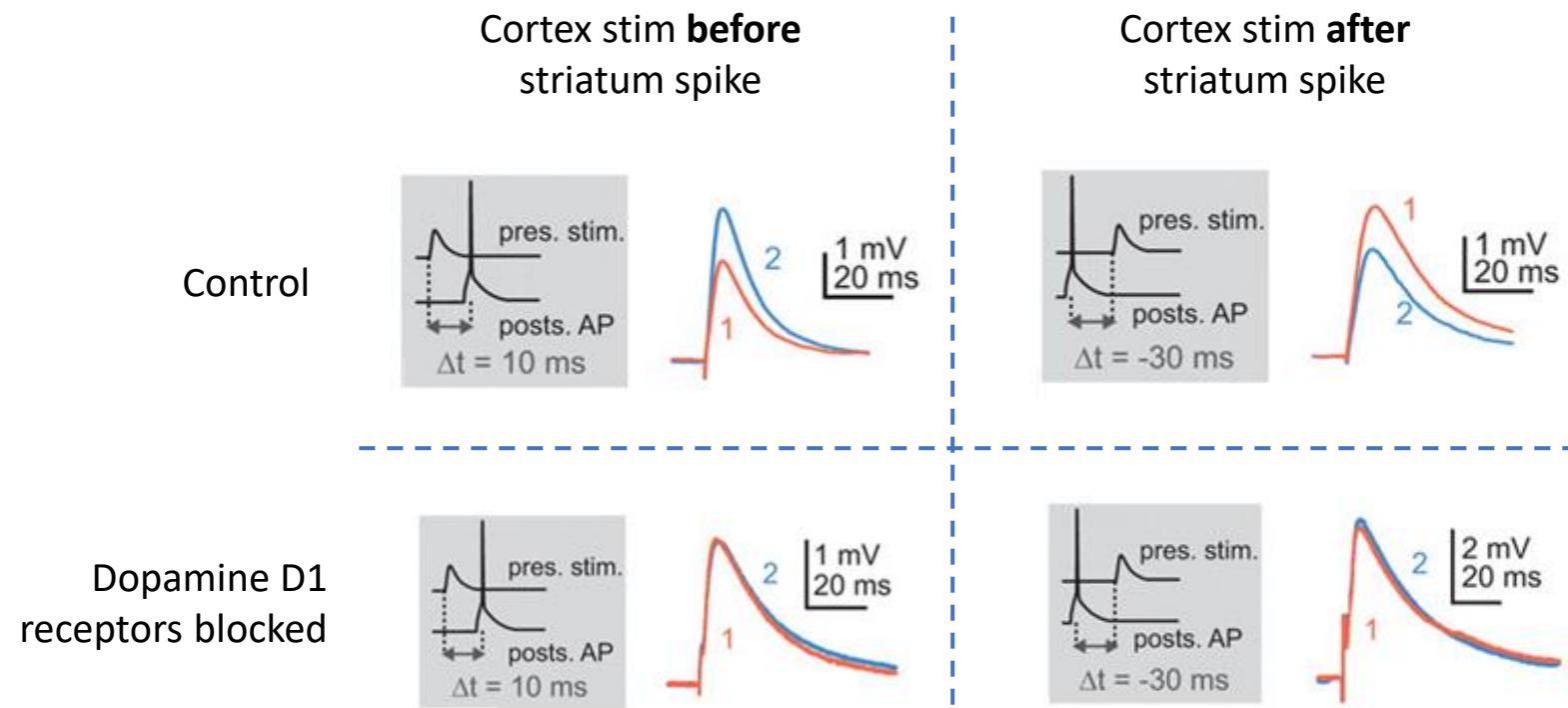
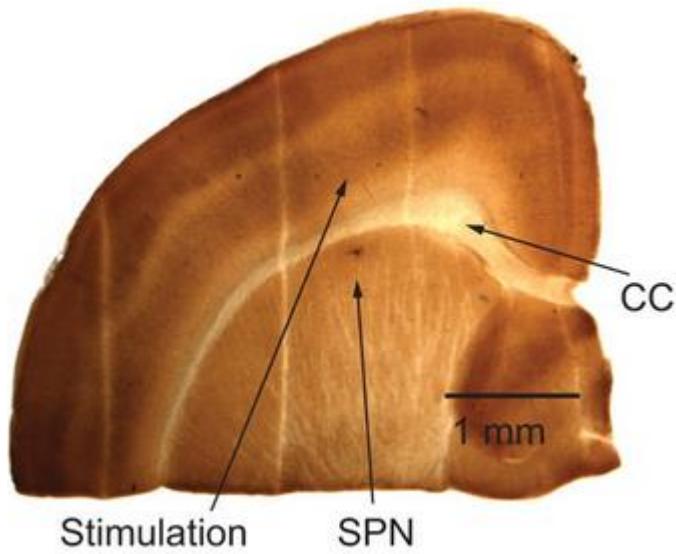


- Stimulation of dorsal striatal neurons biases action selection, consistent with modifying action values

# Dopaminergic modulation of cortico-striatal plasticity



# Dopaminergic modulation of cortico-striatal plasticity



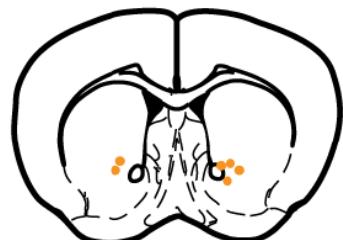
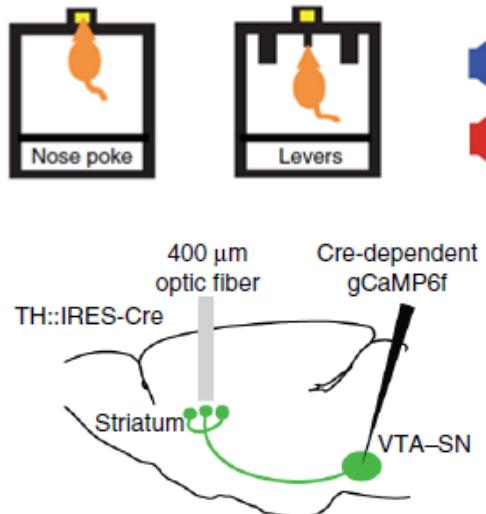
- Plasticity at cortico-striatal synapse depends on pre- and post- synaptic activity + dopamine signalling

## **Interim summary:**

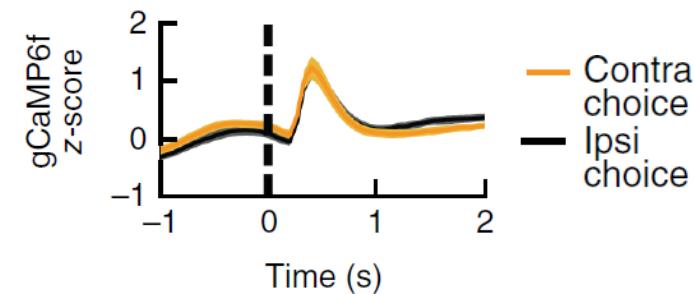
Findings broadly consistent with RPE theory:

- Dopamine activity resembles a TD RPE in many behavioural task and species.
- Dopamine manipulations can reinforce / suppress behaviour.
- Value-like signals are observed in striatum, the primary target of the dopamine system.
- Dopamine potently modulates plasticity at cortico-striatal synapses

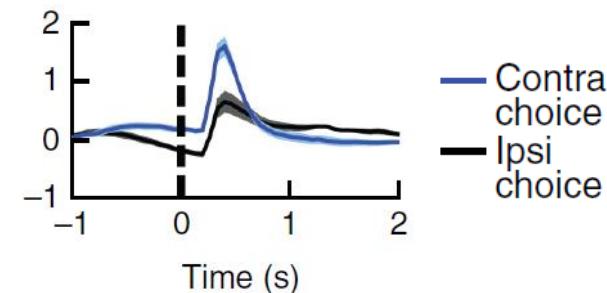
# Heterogeneity of dopamine signals



Ventral striatum

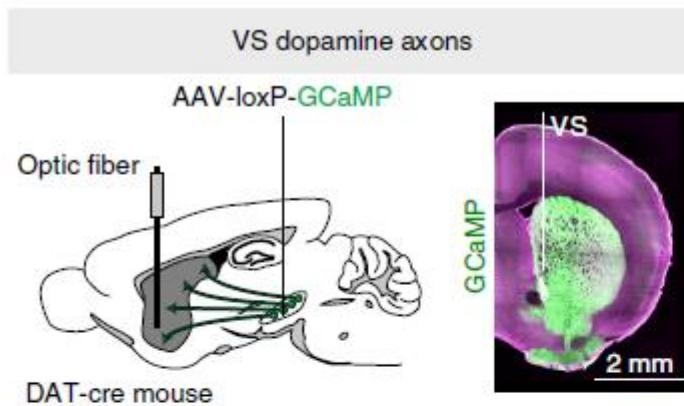


Dorsomedial striatum

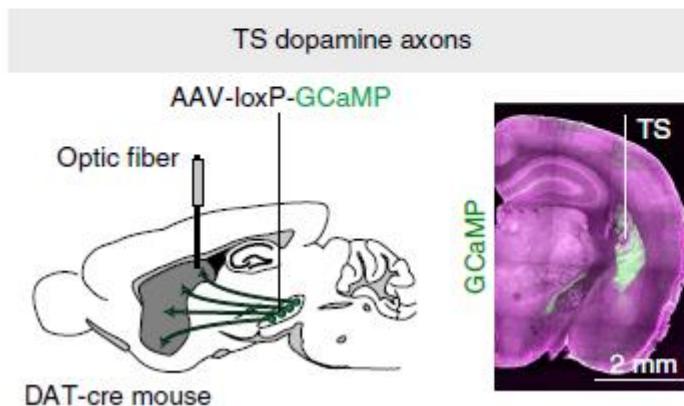
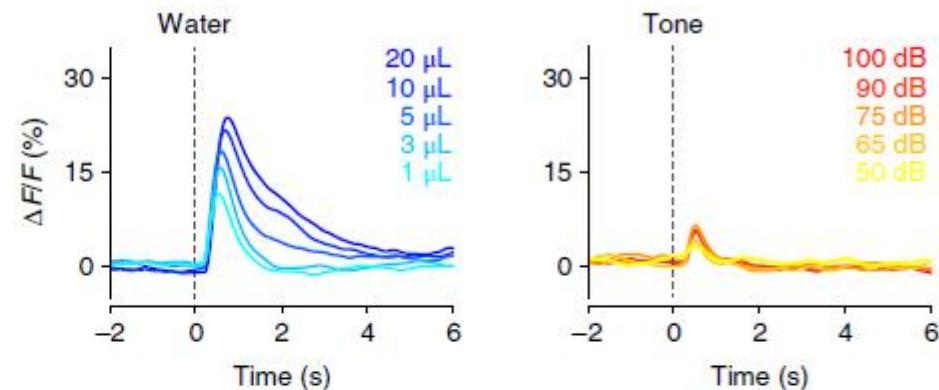


- Dopamine neuron terminals in dorsomedial striatum but not ventral striatum respond to contralateral choices.

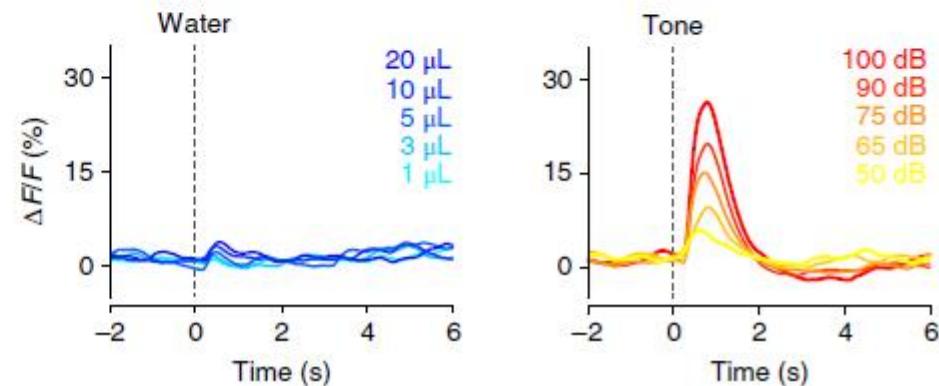
# Heterogeneity of dopamine signals



## Ventral striatum

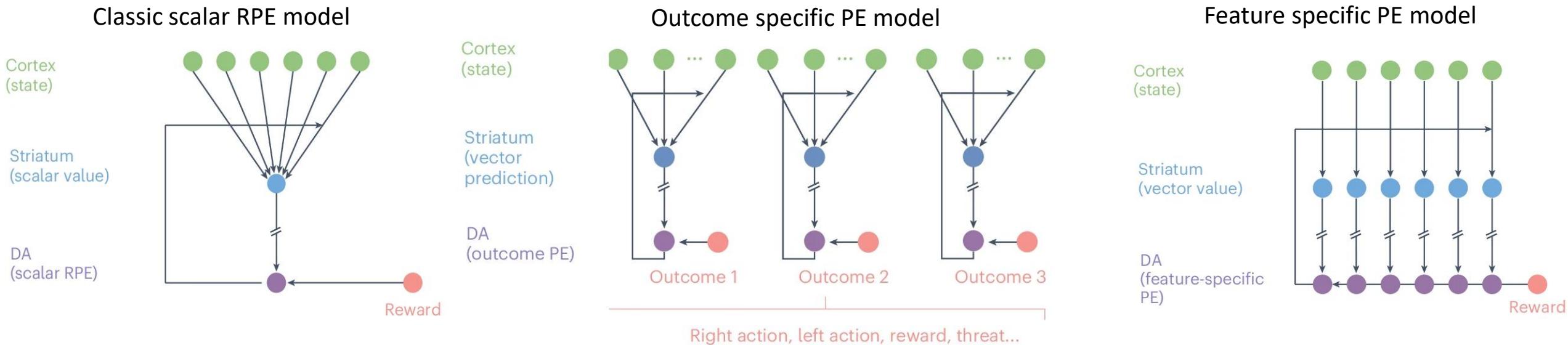


## Posterior-dorsal striatum



- Dopamine neuron axons in the posterior striatum respond to aversive load noise and not to rewards.

# Heterogeneity of dopamine signals

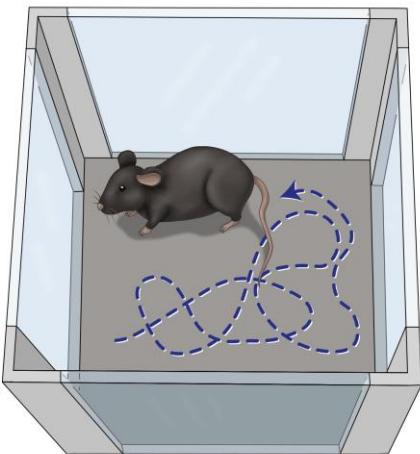


- Heterogeneity is inconsistent with a scalar RPE signal.

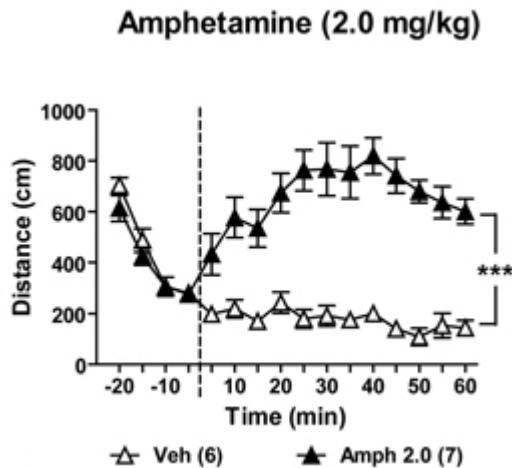
Possible reasons for heterogeneity:

- Outcome specific PE: Prediction of different outcomes (e.g. reward, threat, action) by different cortico-basal ganglia loops.
- Feature specific PE: Different cortico-basal ganglia loops predict reward using different state features.

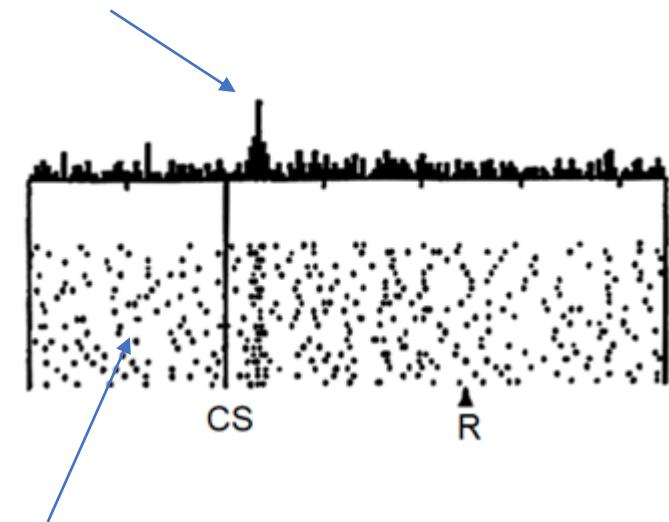
# Dopamine, vigor, and reward rate



© Keri Jones



"phasic" burst of spikes in response to cue

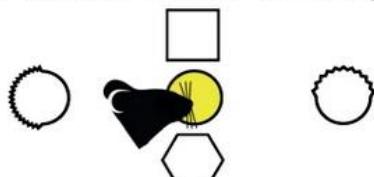


"tonic" activity at low firing rate in absence of external events

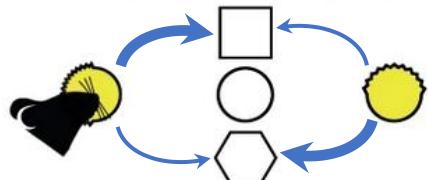
- Pharmacologically boosting dopamine levels induces behavioural activation and increased response vigor.
- Rapid time-course appears inconsistent with RPE-driven learning.
- Niv et al. (2007) propose that "tonic" and "phasic" activity as separate information channels:
  - "Phasic" bursts: RPE → learning
  - "tonic" activity: average reward rate → opportunity cost of time that controls motivation/vigor

## Dopamine, vigor, and reward rate

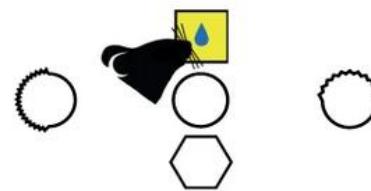
1. Initiate trial in center port



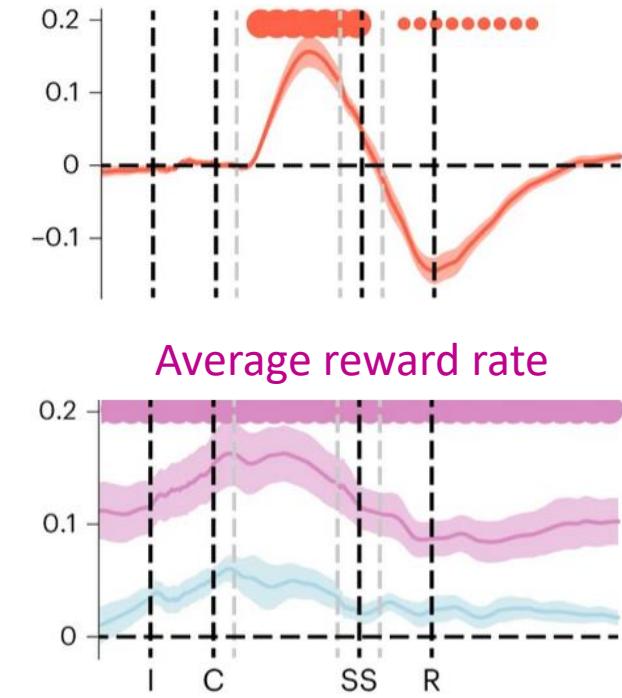
2. Choose between left and right ports



3. Poke active up/down port for probabilistic reward

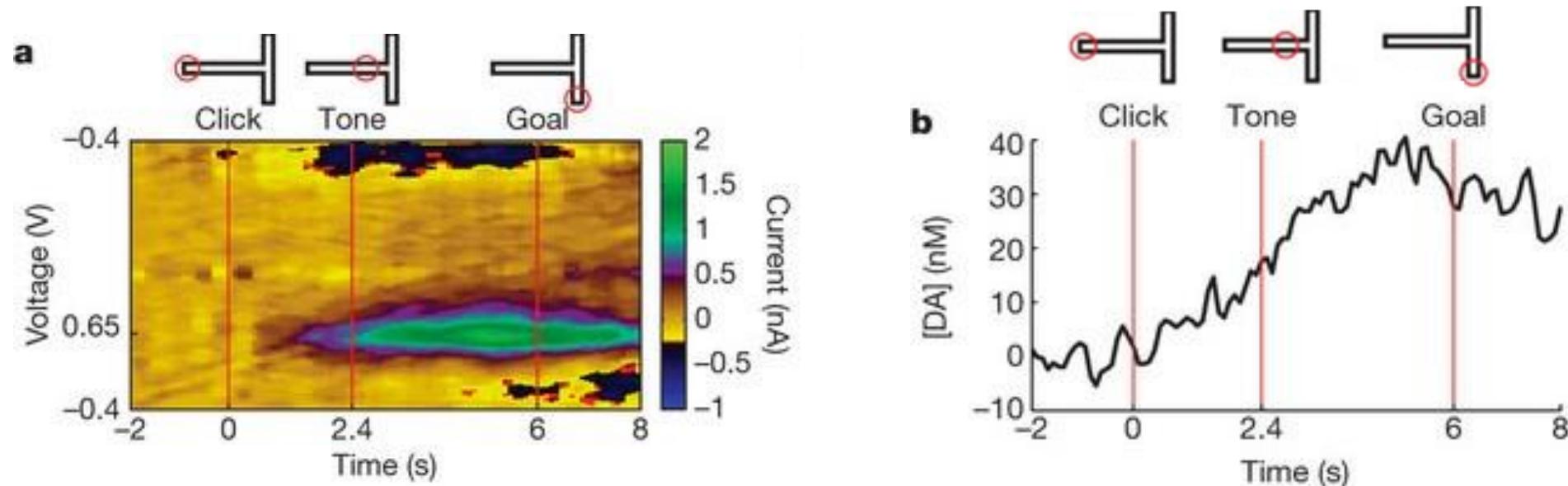


Value of up/down port



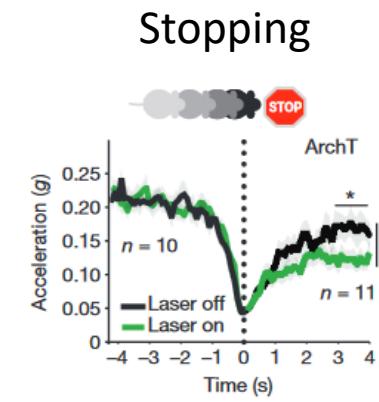
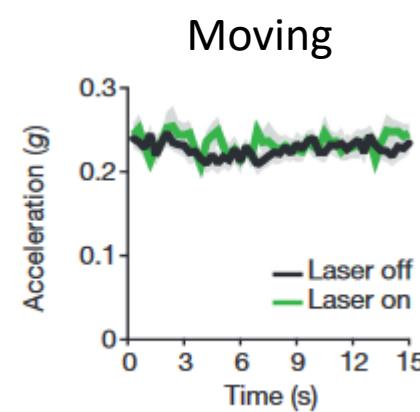
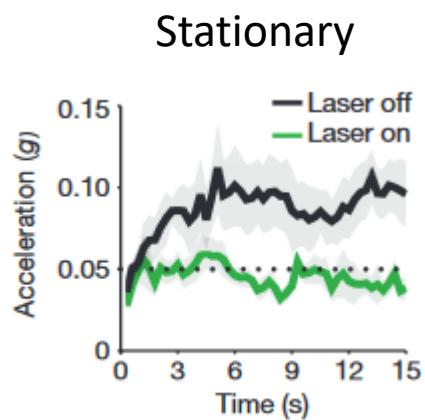
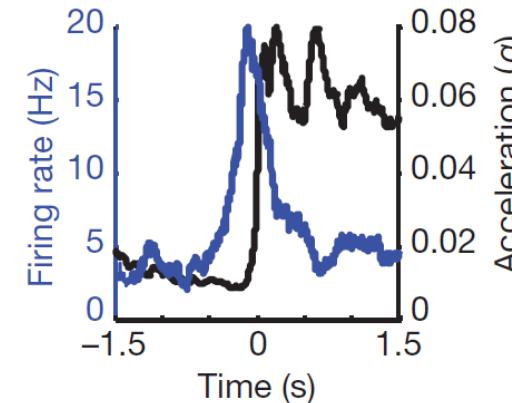
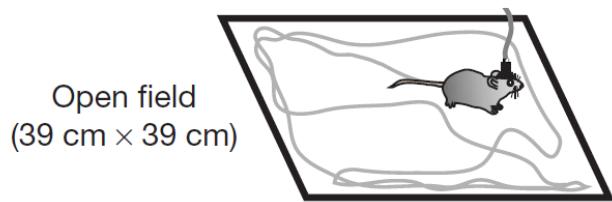
- Average reward rate drives sustained changes in dopamine neuron activity, dissociable from fast phasic RPE responses.

## Dopamine ramps



- Striatal dopamine during navigation to reward show ramping that resembles value more than RPE.
- Many theoretical accounts...

## Dopamine and movement initiation



- SNC dopamine neurons are modulated before spontaneous movement initiation.
- Inhibiting SNC suppresses movement initiation but does not stop stopping ongoing movement.

## Dopamine and RPE summary

Findings broadly consistent with RPE theory:

- Dopamine activity resembles a TD RPE in many behavioural task.
- Dopamine manipulations can reinforce / suppress behaviour.
- Value-like signals are observed in striatum, the primary target of the dopamine system.
- Dopamine potently modulates plasticity at cortico-striatal synapses

Findings not straightforwardly predicted by RPE theory:

- Direct effects on motivation/vigor
- Involvement in action initiation
- Ramping ‘value-like’ activity.

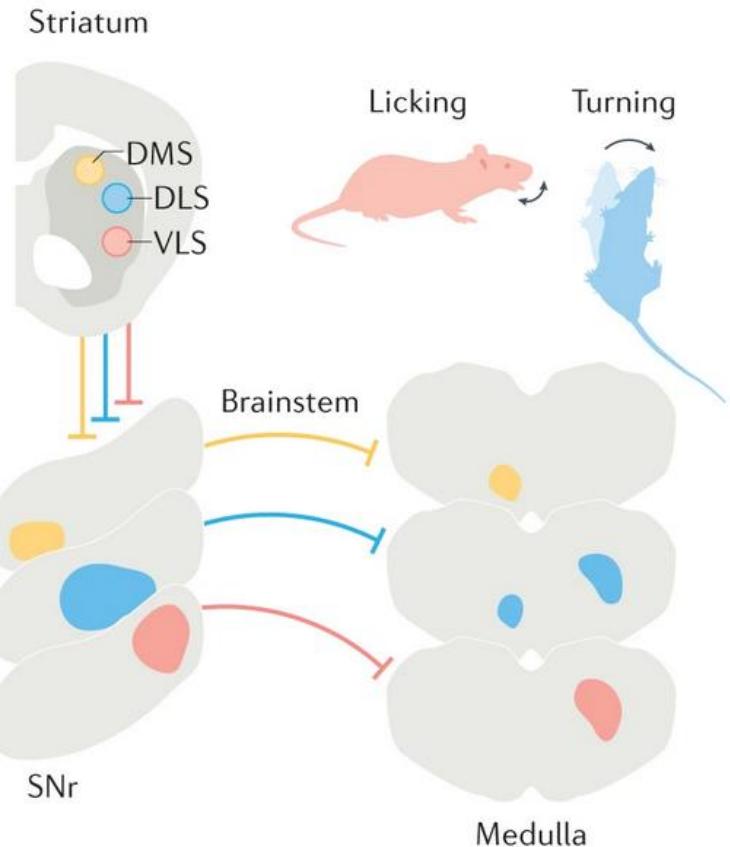
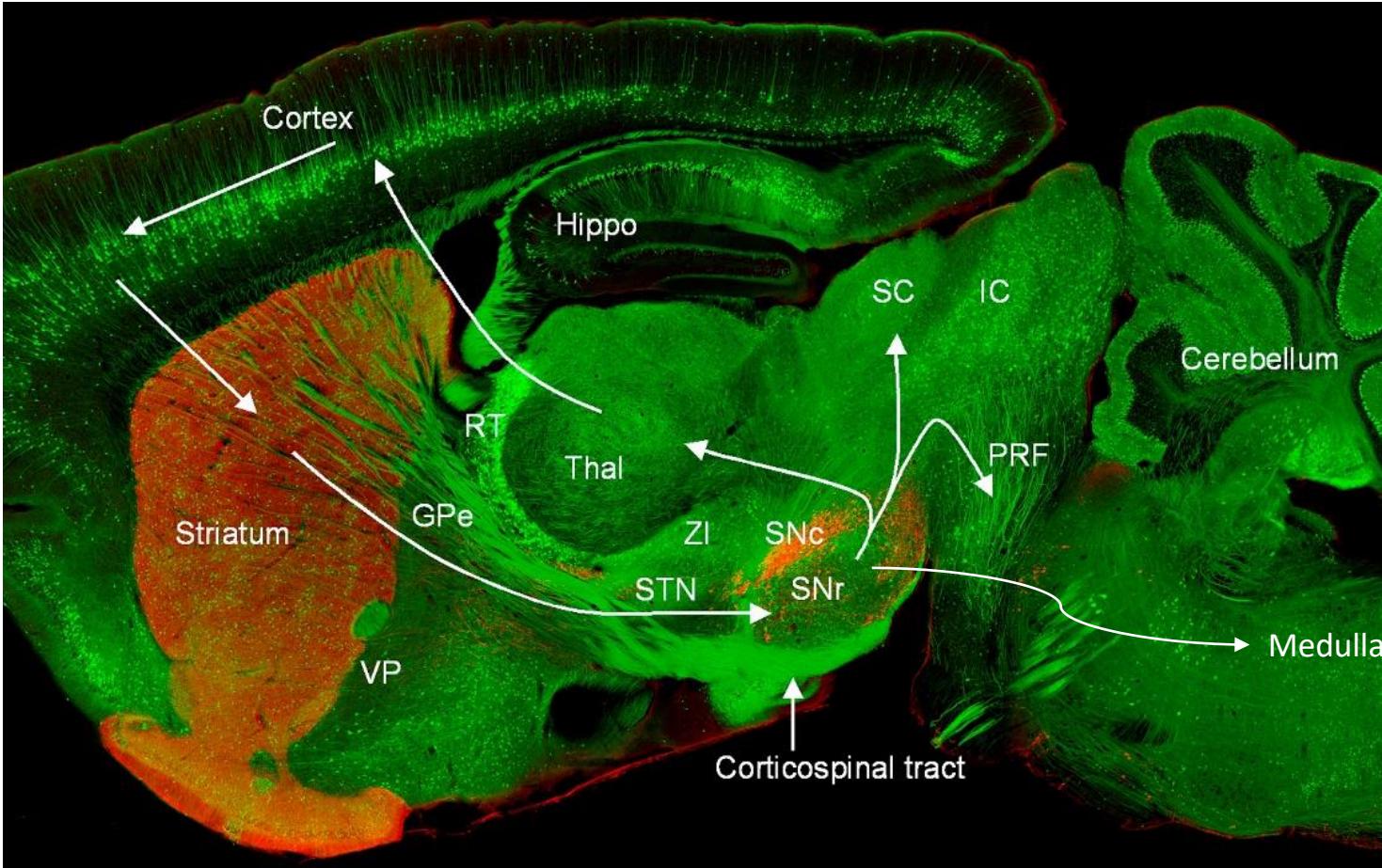
## Talk aim and overview

**Question:** How do learning algorithms map onto brain structure to solve the biological action control problem?

**Talk outline:**

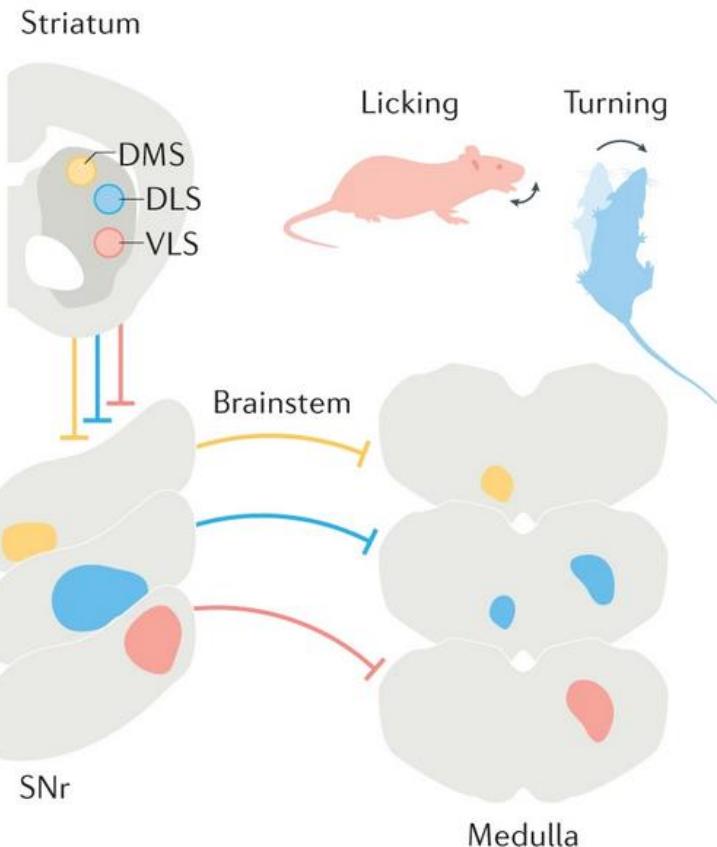
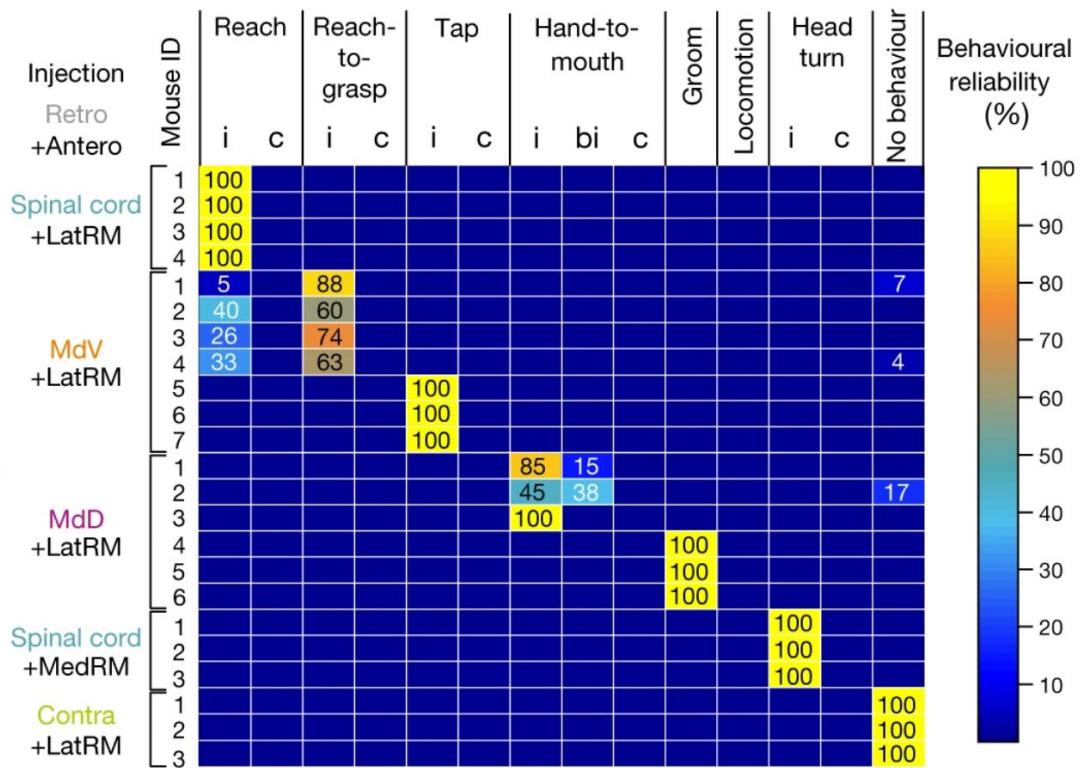
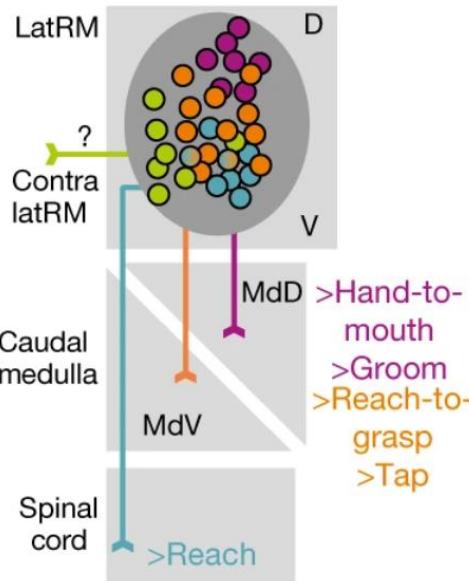
1. Striatum and dopamine: The brain's temporal difference reinforcement learning system?
2. Basal ganglia outputs and control of action selection.
3. Cortex: State representation and beyond
4. Hippocampus: Sequence generation and model-based control

## Basal ganglia outputs to brainstem for control of movement



- Basal ganglia projects to motor control nuclei in the brainstem: Striatum → SNr → Brainstem
- Output to brainstem is topographically organised into parallel pathways

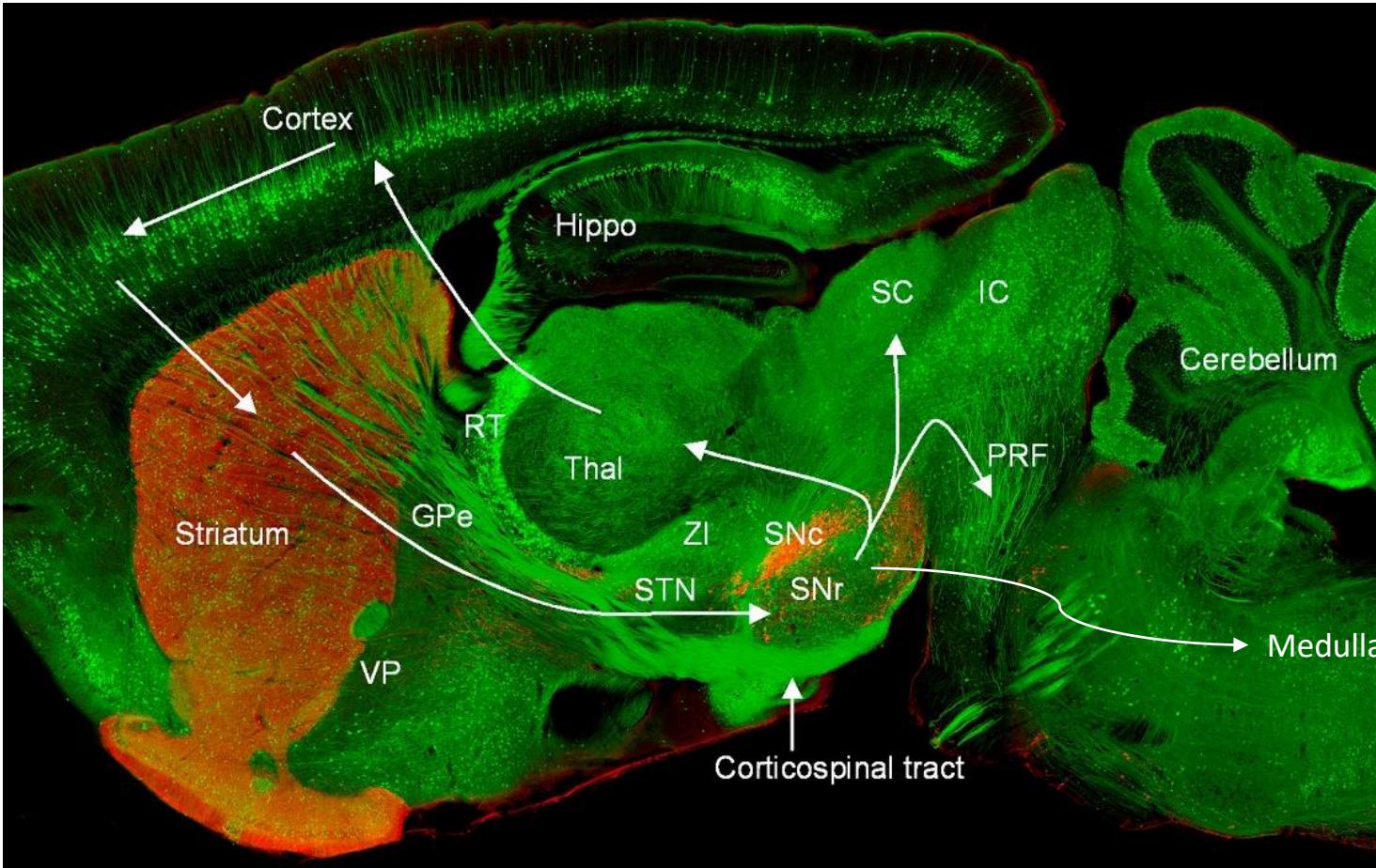
# Basal ganglia outputs to brainstem for control of movement



- Different populations of brainstem neurons control distinct motor actions

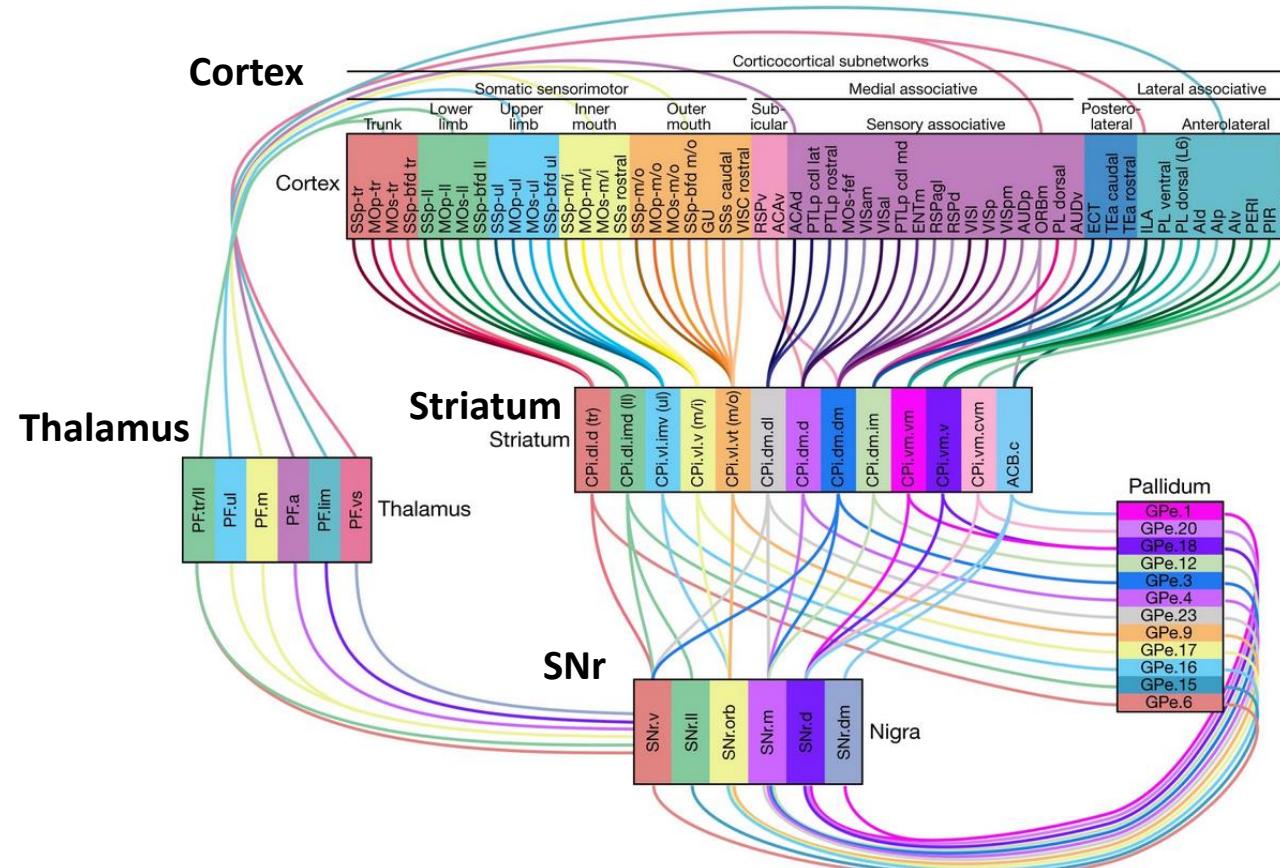
→ Striatum is upstream of motor control nuclei consistent with striatal value / policy signals controlling action

## Basal ganglia outputs: Projections to thalamus and control of cortical activity



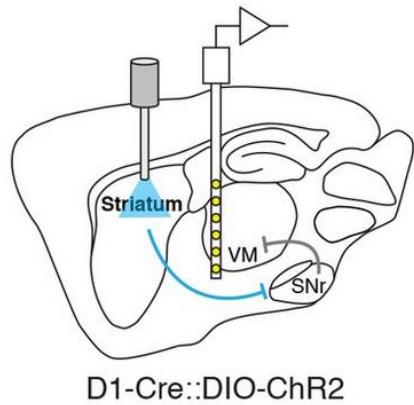
- Basal ganglia output neurons send collaterals to thalamus which has an excitatory projection to cortex

# Basal ganglia outputs: Projections to thalamus and control of cortical activity

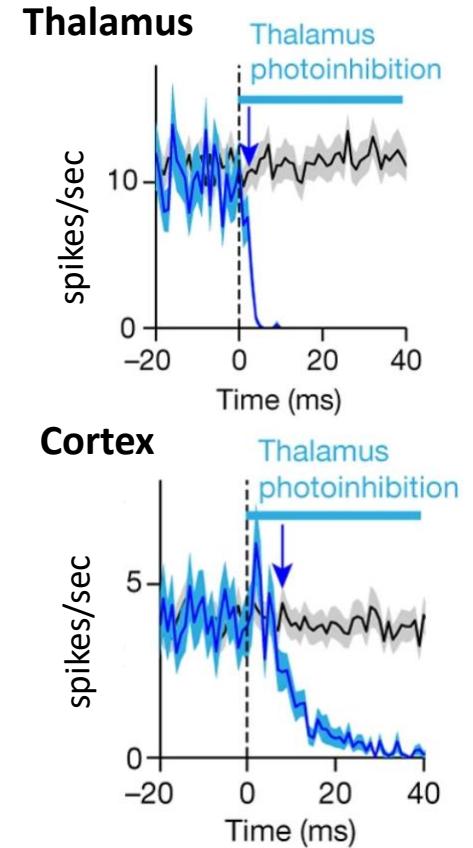
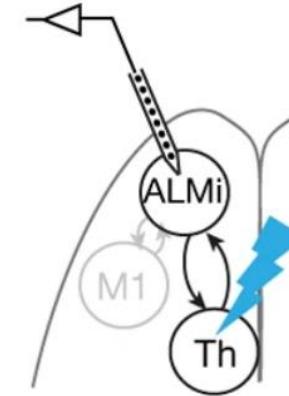
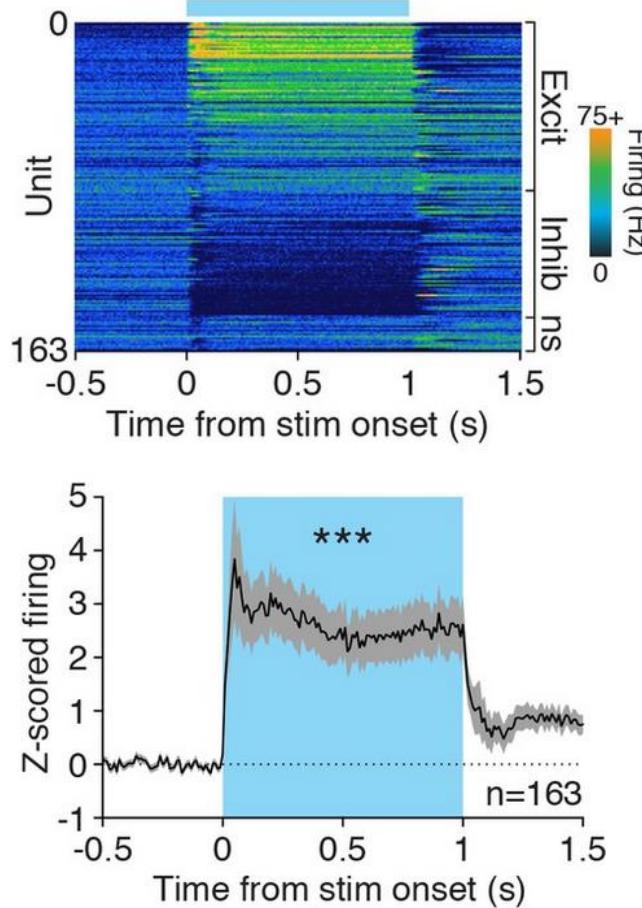


- Basal ganglia output neurons send collaterals to thalamus which has an excitatory projection to cortex
- Parallel loops architecture maps BG outputs channels back to the cortical inputs that drive them

## Basal ganglia outputs: Projections to thalamus and control of cortical activity

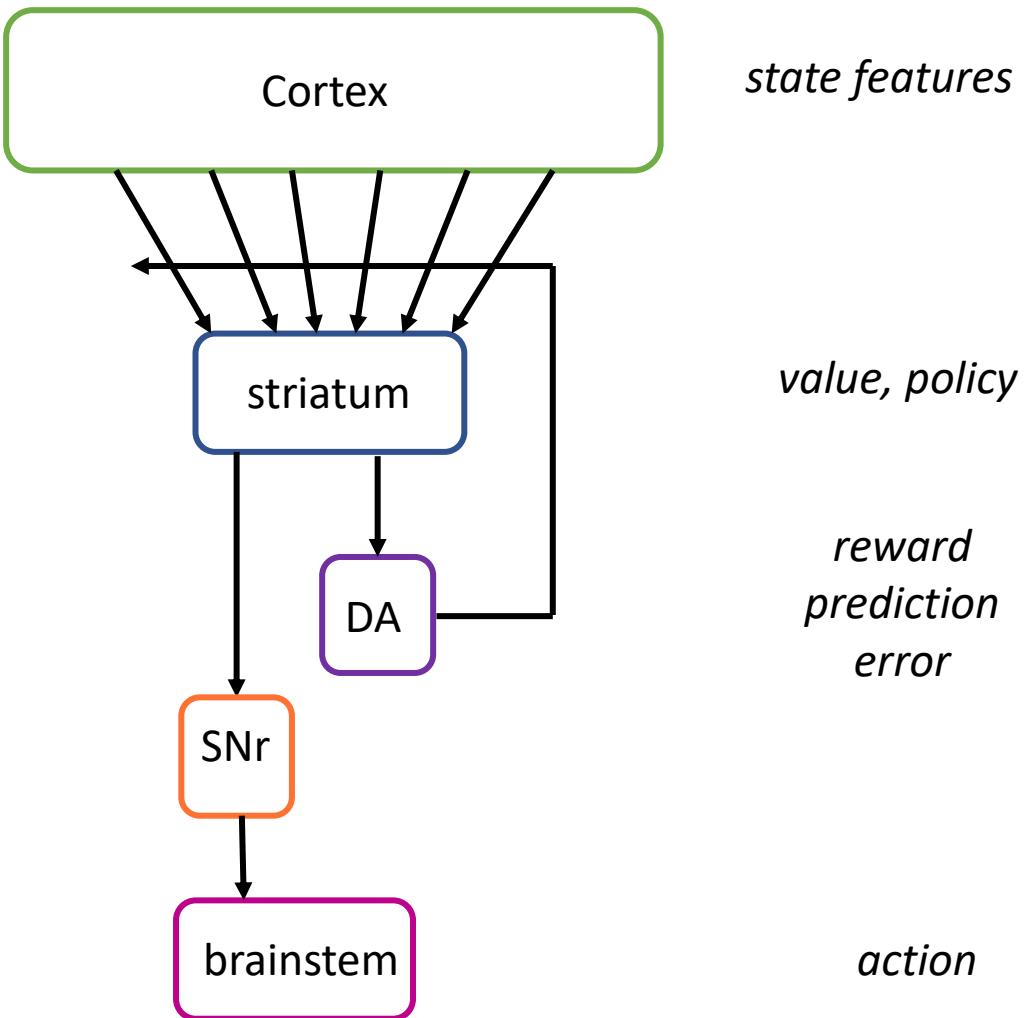
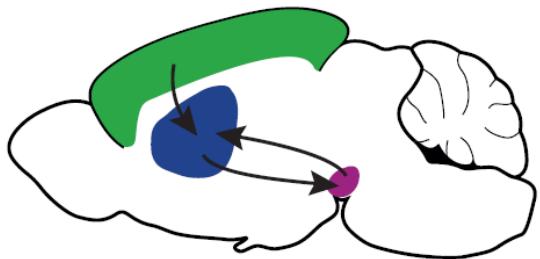


D1-Cre::DIO-ChR2

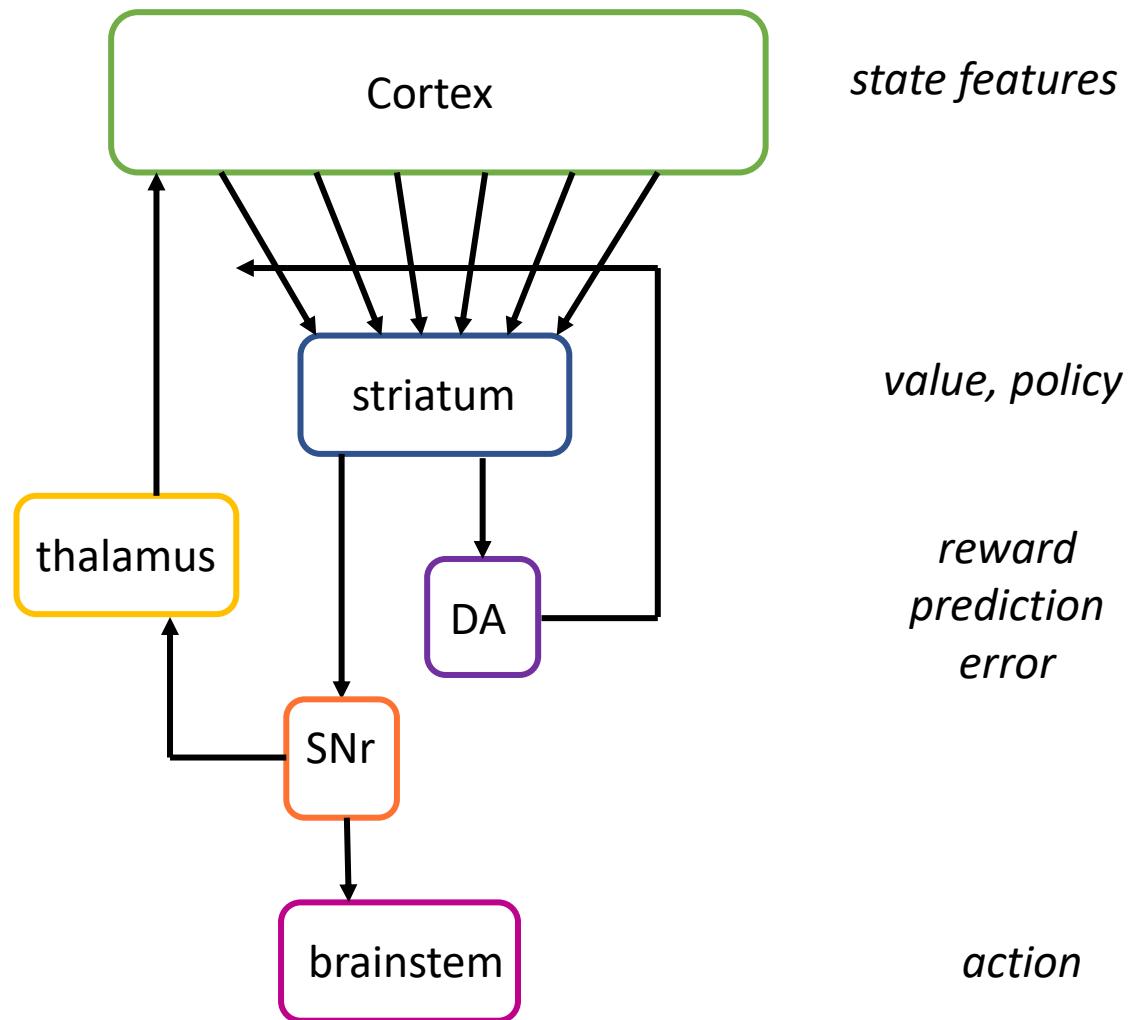
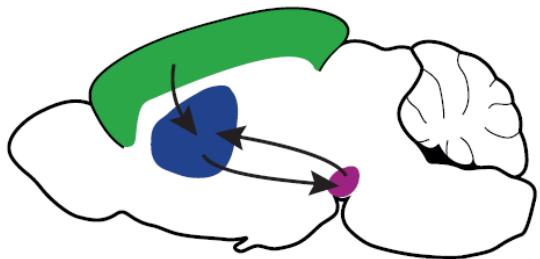


- Basal ganglia output neurons send collaterals to thalamus which has an excitatory projection to cortex
- Parallel loops architecture maps BG outputs channels back to the cortical inputs that drive them
- Striatum can potently modulate thalamic activity and hence cortical activity

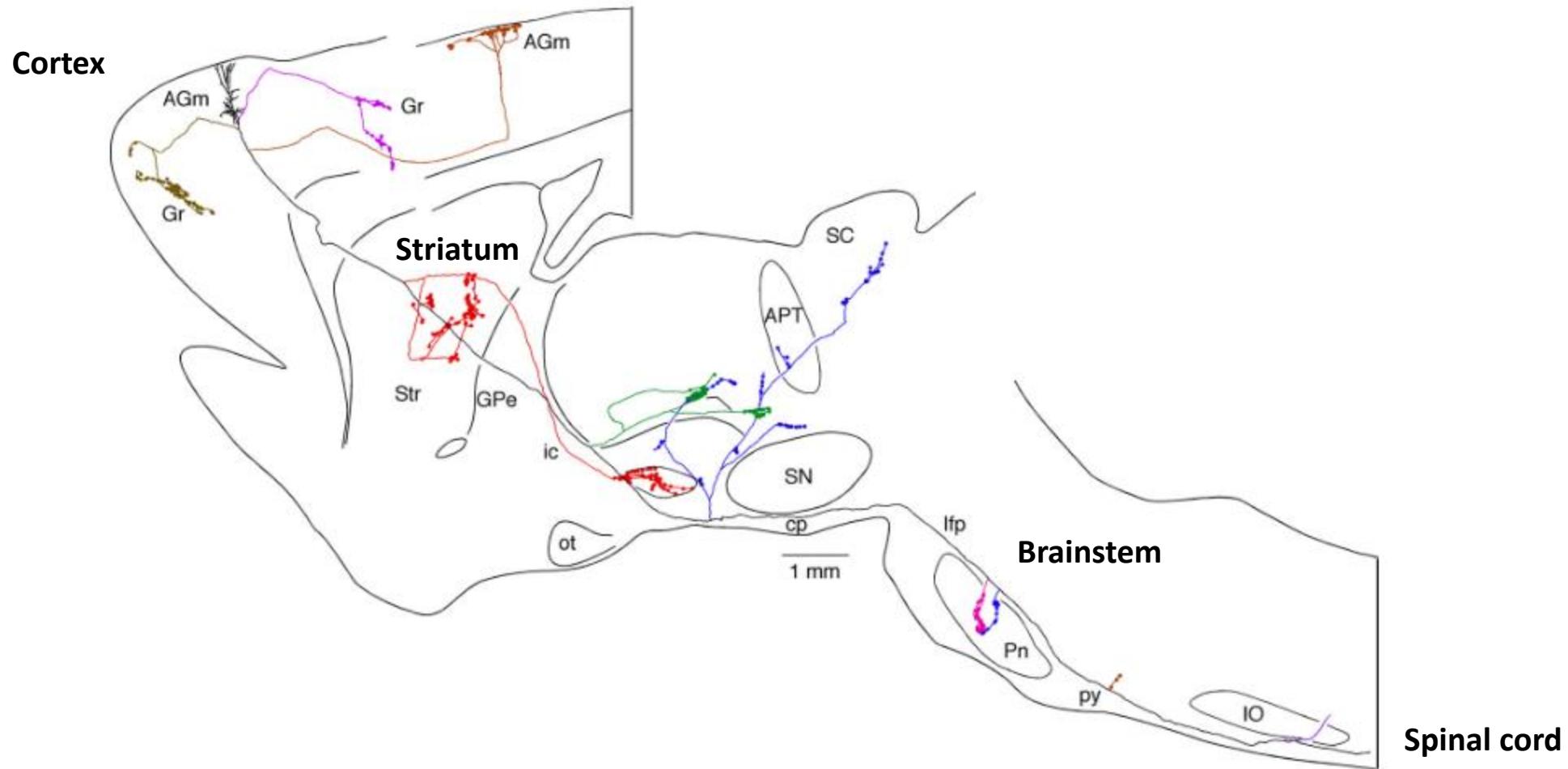
# Where are we at with the boxes and arrows..



**Where are we at with the boxes and arrows..**

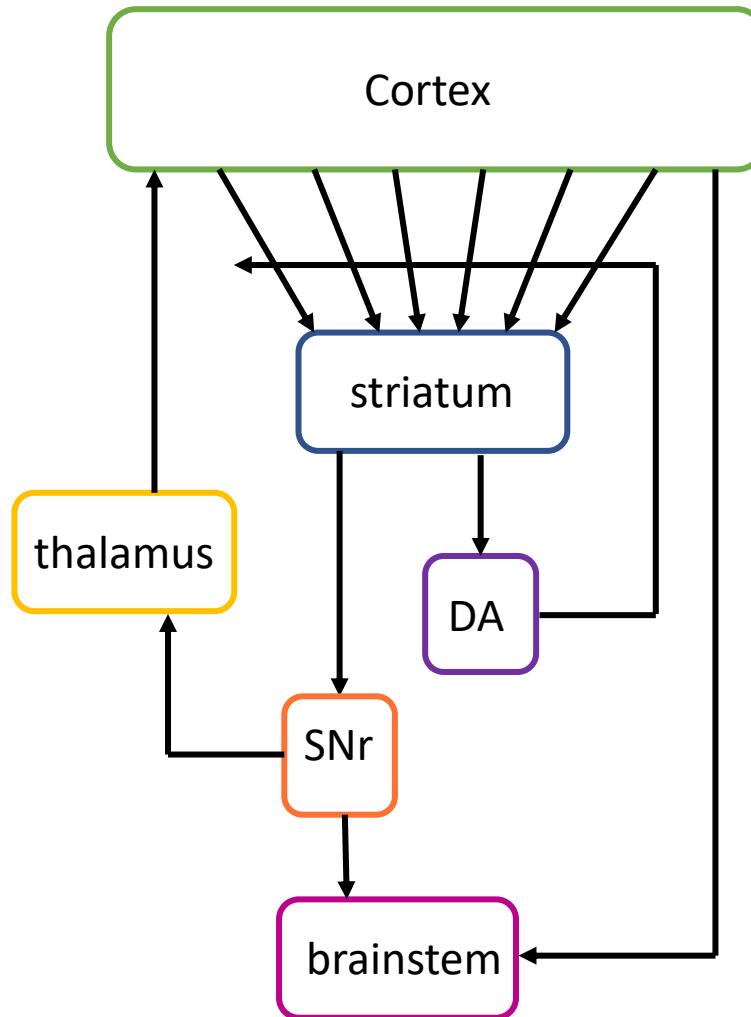
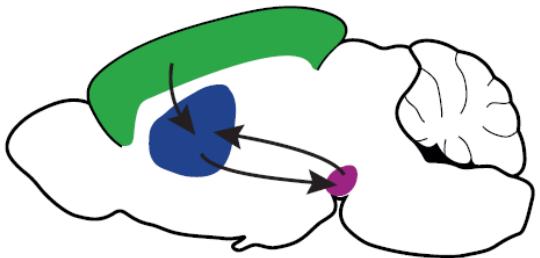


## Further complications



- Frontal and motor cortex make direct projections to brainstem and spinal cord.

It's complicated



*state features/action/value/policy???*

*value, policy*

*reward  
prediction  
error*

*action*

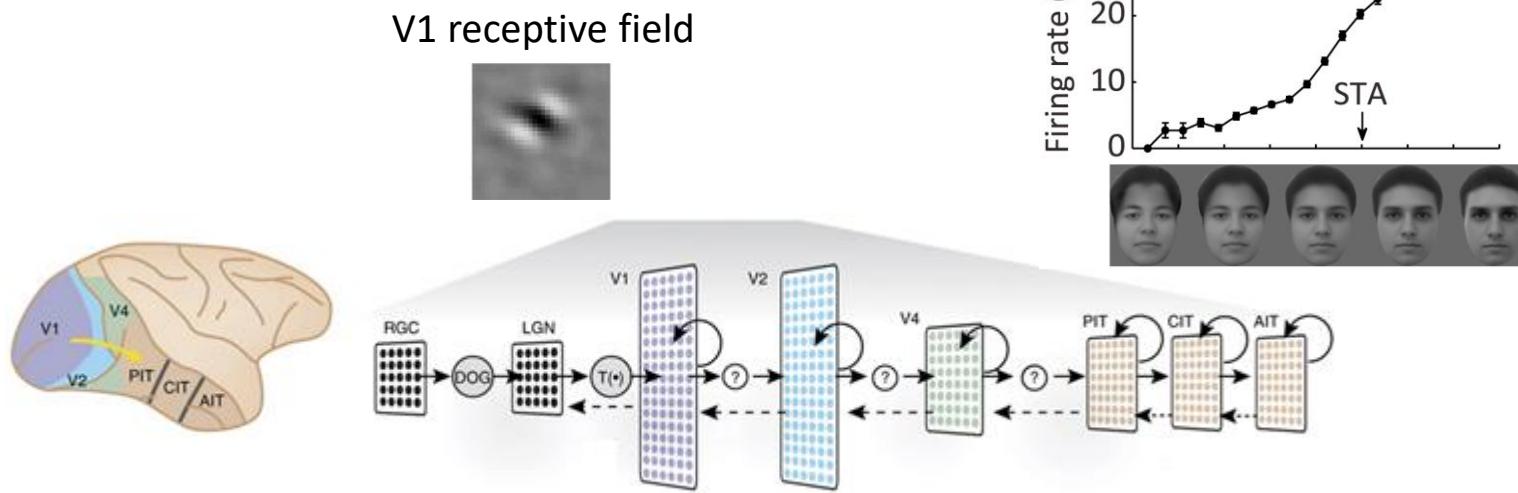
## Talk aim and overview

**Question:** How do learning algorithms map onto brain structure to solve the biological action control problem?

**Talk outline:**

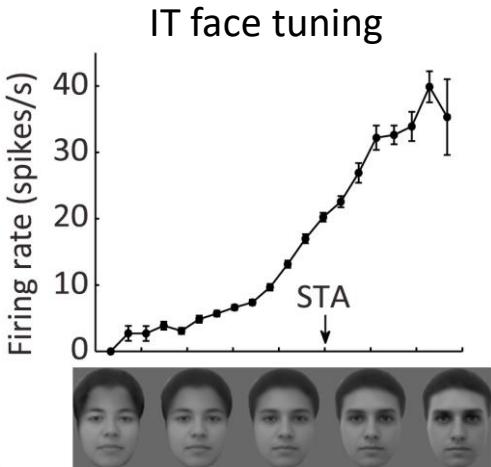
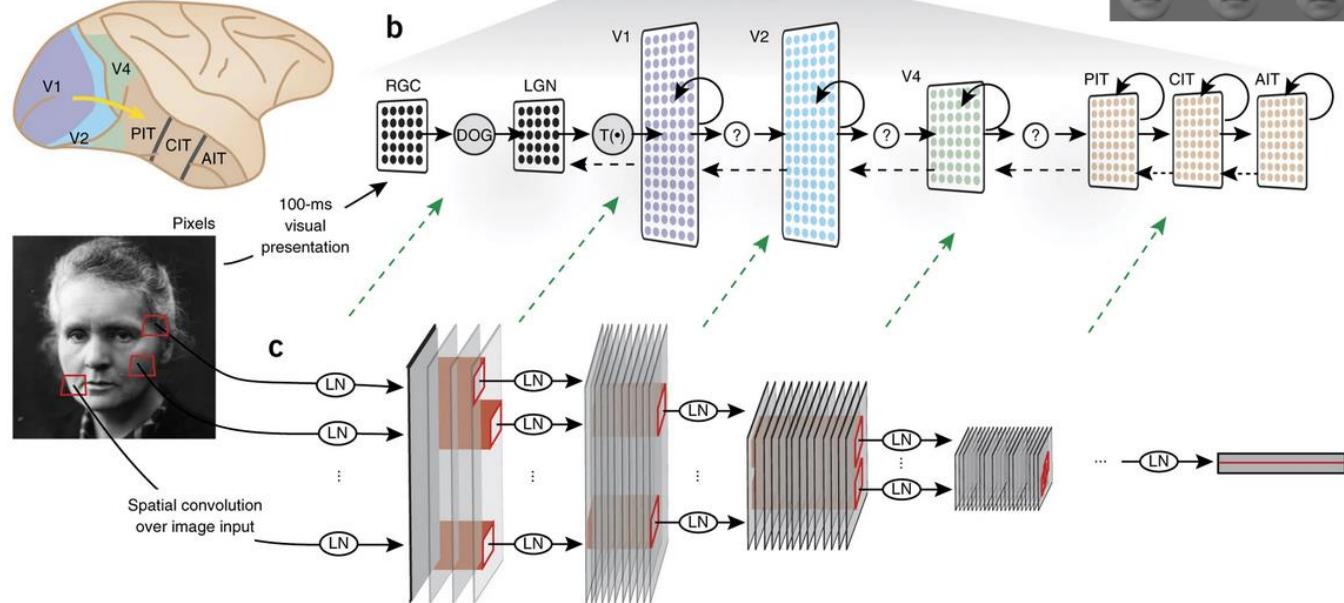
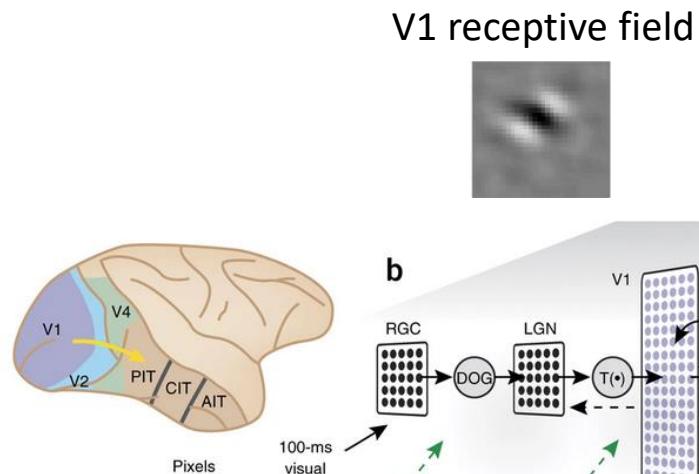
1. Striatum and dopamine: The brain's temporal difference reinforcement learning system?
2. Basal ganglia outputs and control of action selection.
3. Cortex: State representation and beyond
4. Hippocampus: Sequence generation and model-based control

## Cortex and state representation

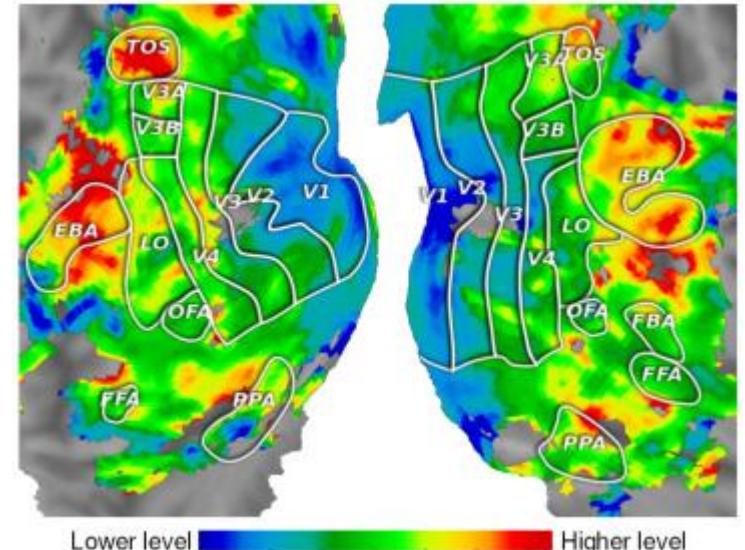


- Sensory cortex extracts behaviourally relevant features (e.g. what and where) from high dimensional sensory input.

# Cortex and state representation

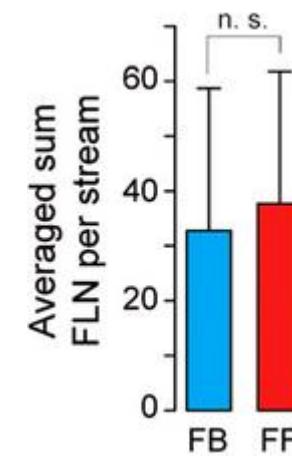
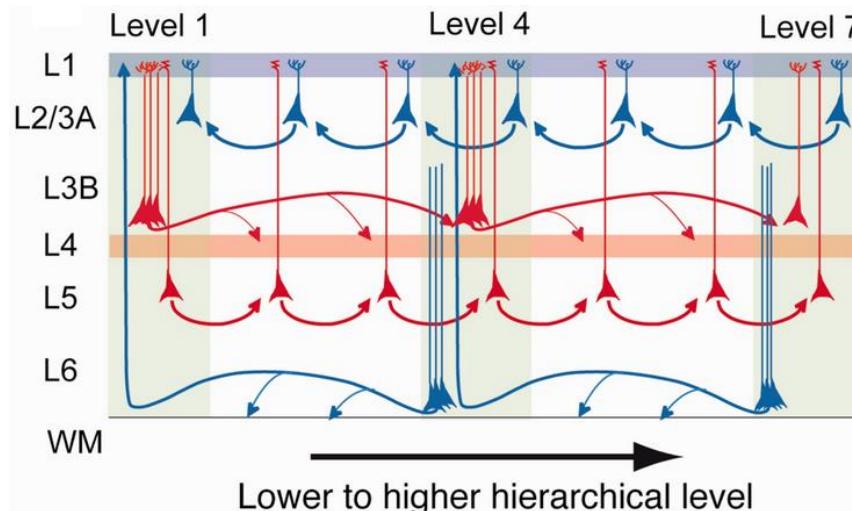
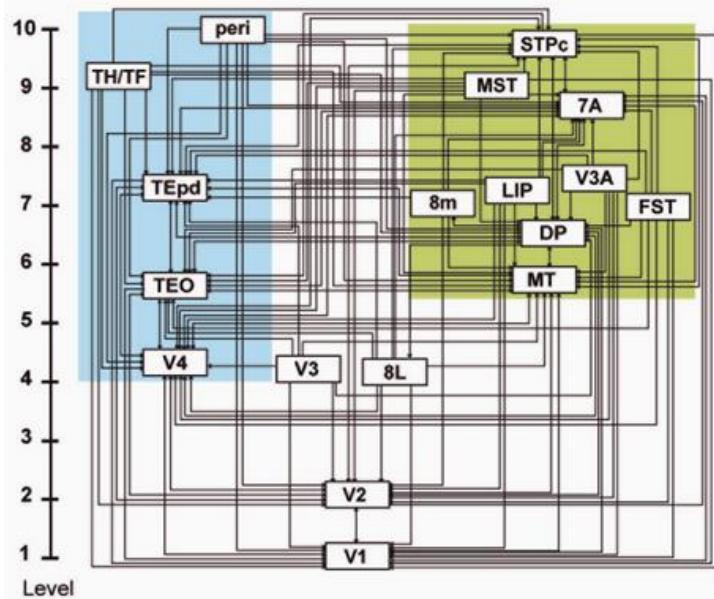


Linear prediction from CNN layer to brain activity



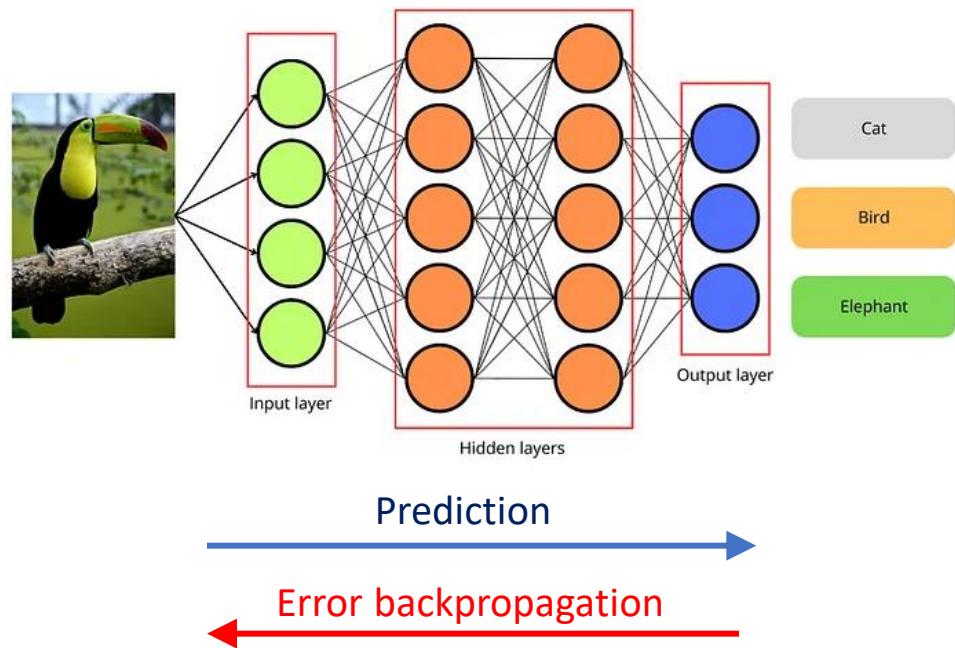
- Sensory cortex extracts behaviourally relevant features (e.g. what and where) from high dimensional sensory input.

# Cortex as a deep, hierarchical, recurrent network

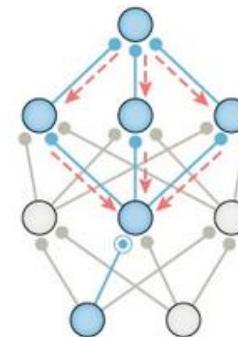


- Cortex is hierarchically organised, with defined feed-forward and feed-back pathways.
- Majority of inter-regional connections are reciprocal.
- Approximately equal numbers of neurons participate in feed-forward and feed-back pathways.

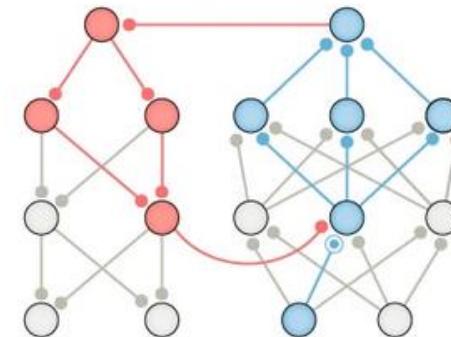
# Reciprocal connectivity and learning in biological networks



Backpropagation

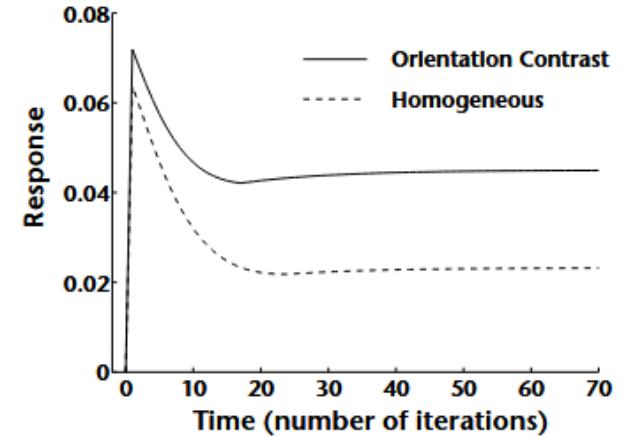
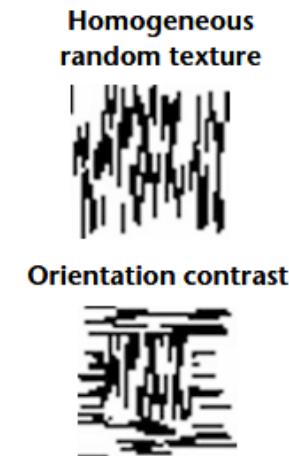
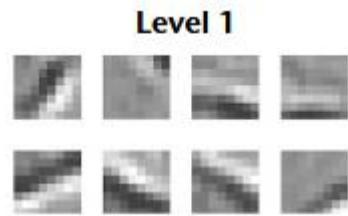
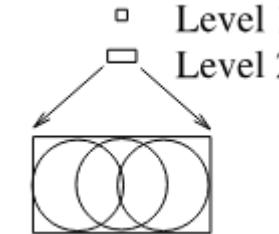
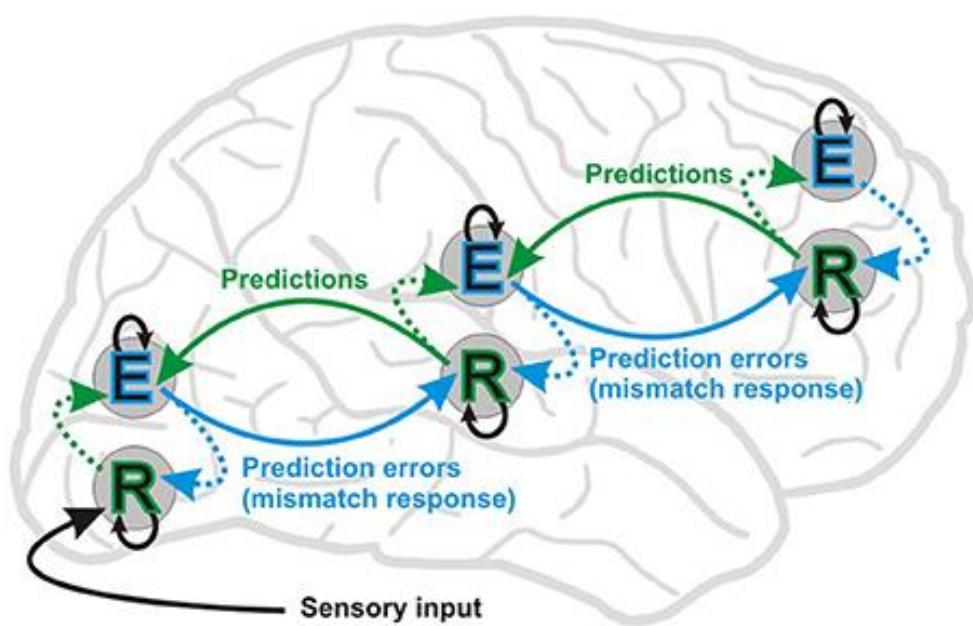


Backprop-like learning  
with feedback network



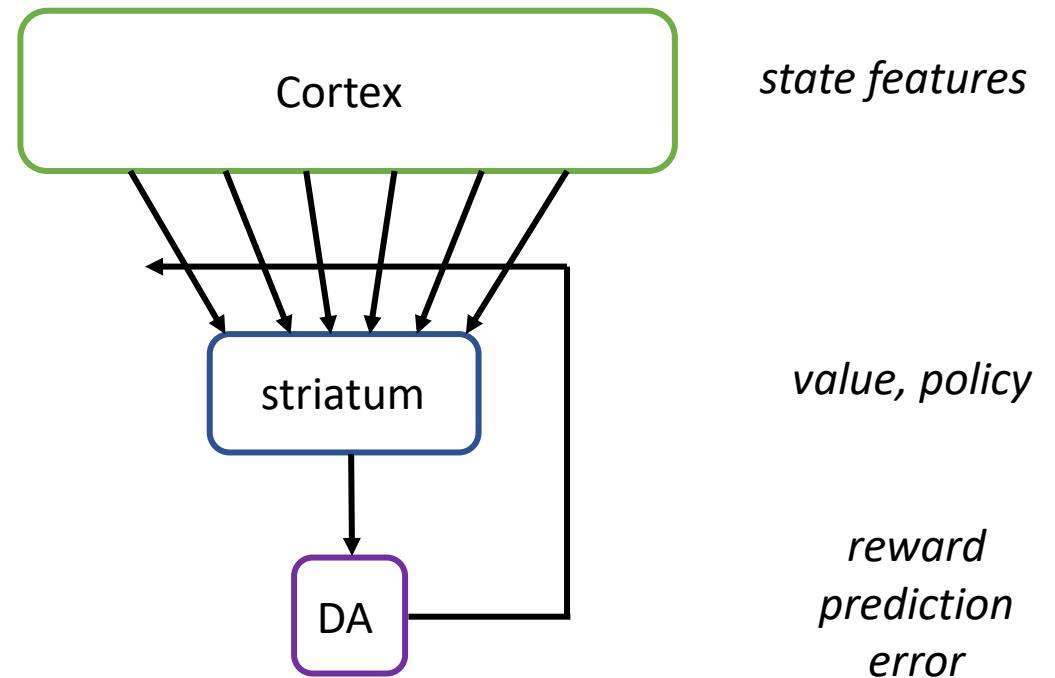
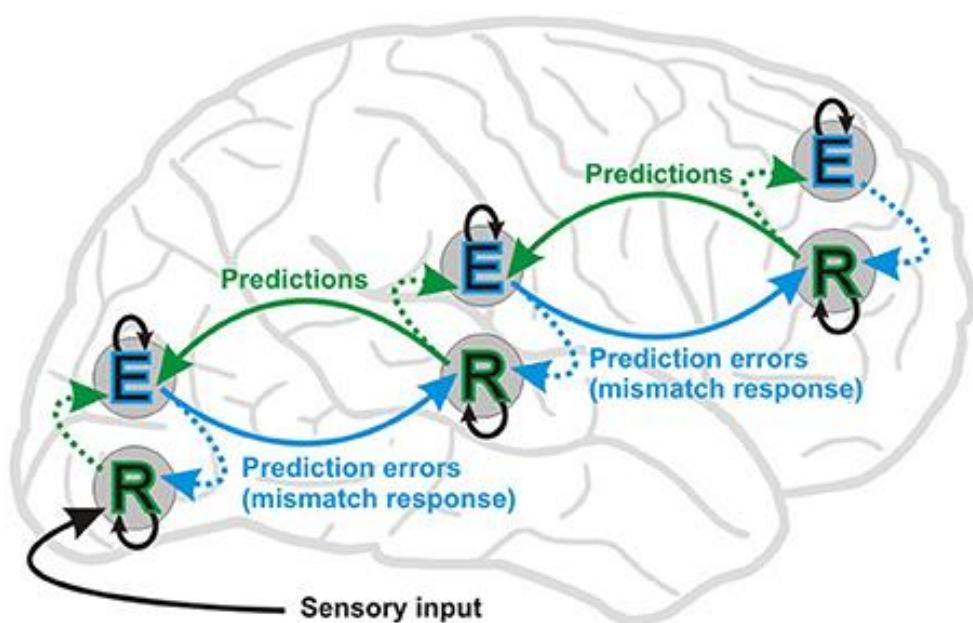
- Deep artificial neural networks learn mappings from input to output by backpropagating errors.
- Backpropagation use non-local information to compute weight updates → not biologically plausible
- Feedback connections can enable backprop-like updates using only local learning rule.

## Hierarchical predictive coding (HPC)



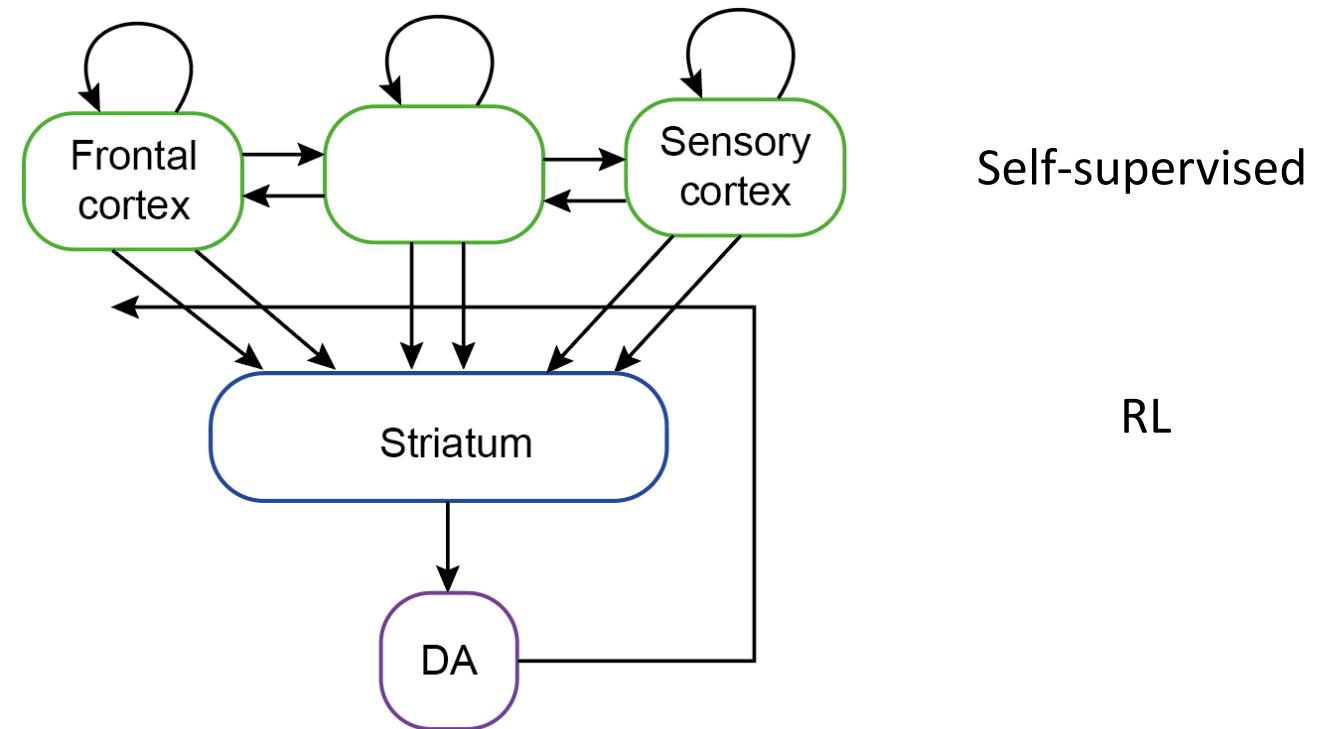
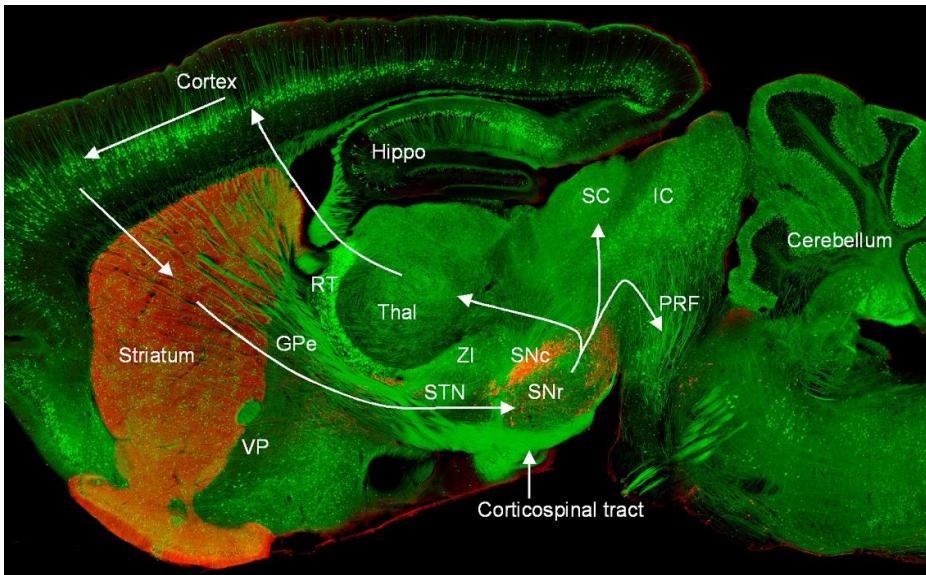
- In HPC, each level of the hierarchy tries to predict activity in the level below, sends prediction errors to the level above.
- HPC learns a generative model of sensory input, using only local learning rules.
- Trained on natural image data, HPC learns neuron-like responses.

## Prediction and error in cortex and basal ganglia



- In HPC models of sensory cortex the prediction target is the next sensory observation → un-/self-supervised learning
- In TD models of basal ganglia function the prediction target is long-run reward → reinforcement learning

# Why this architecture?



- Cortex: Hierarchical, recurrent, largely symmetric projections.
- Striatum: Shallow, no recurrent excitation, no direct projection back to cortex

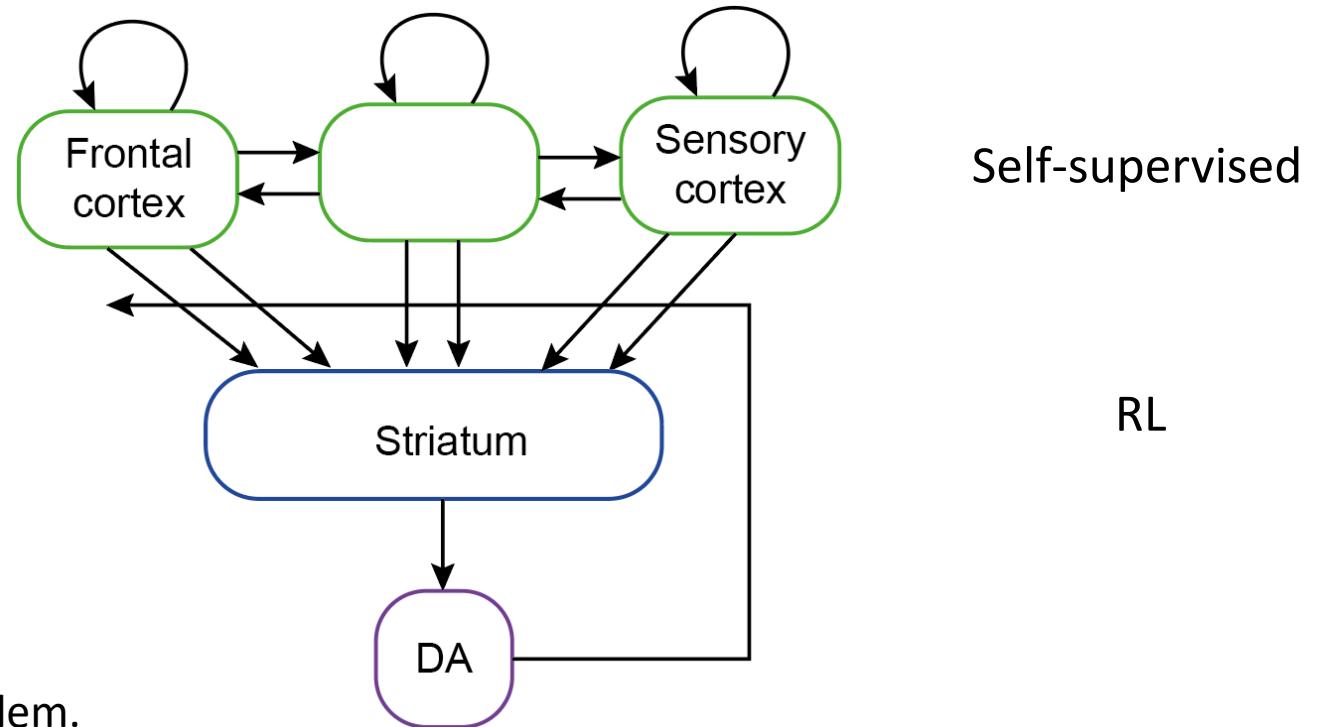
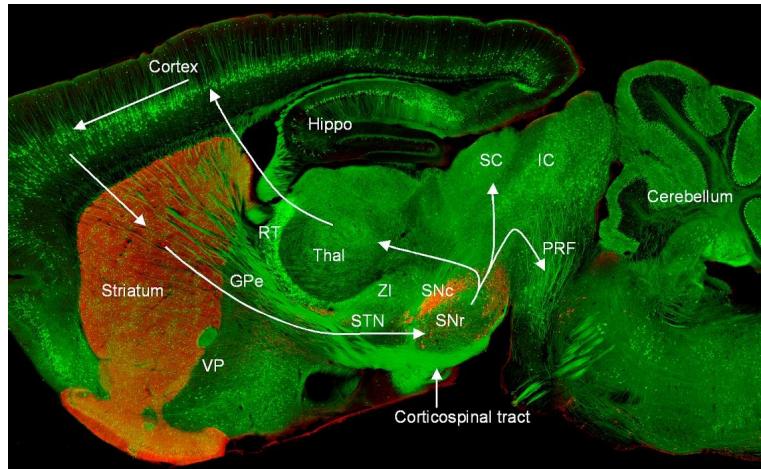
## Bootstrapping and stability in TD learning

$$\delta_t = \frac{R_{t+1} + \gamma V(s_{t+1})}{\text{Update target}} - \frac{V(s_t)}{\text{previous estimate}}$$

reward prediction error

- Estimated value of next state contributes to update target (bootstrapping).
- Bootstrapping + neural network function approximators is notoriously unstable.
- Stabilizing methods from deep RL don't appear biologically plausible:
  - Use of future rewards in computing weight updates (A2C, PPO)
  - Duplication of network weights (Deep Q)
  - Experience replay (Deep Q)

## Why this architecture?



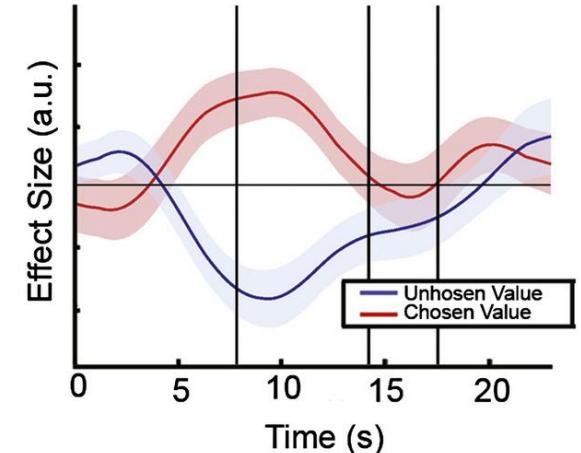
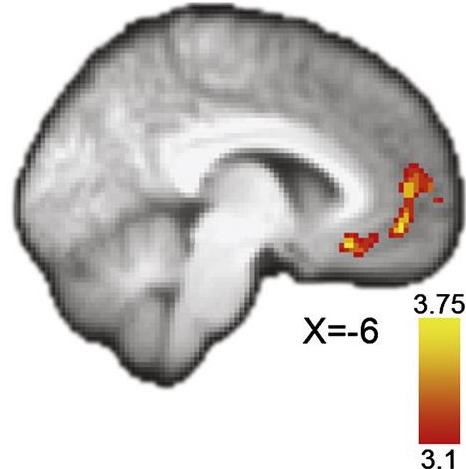
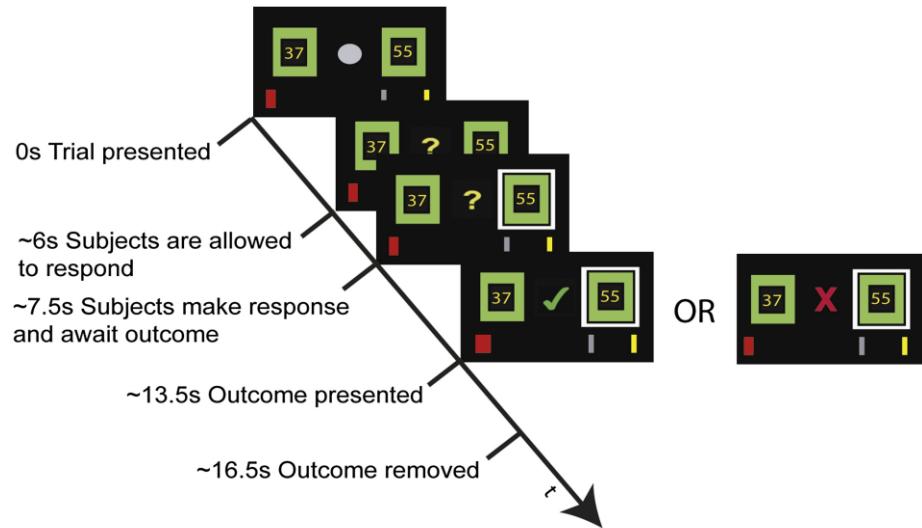
Constraints:

- Need TD learning to solve delayed reward problem.
- Need deep, recurrent, neural networks to handle high dimensional observations and partial observability.

Possible solution:

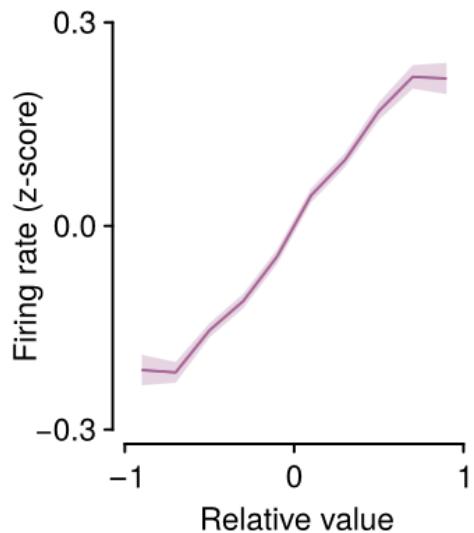
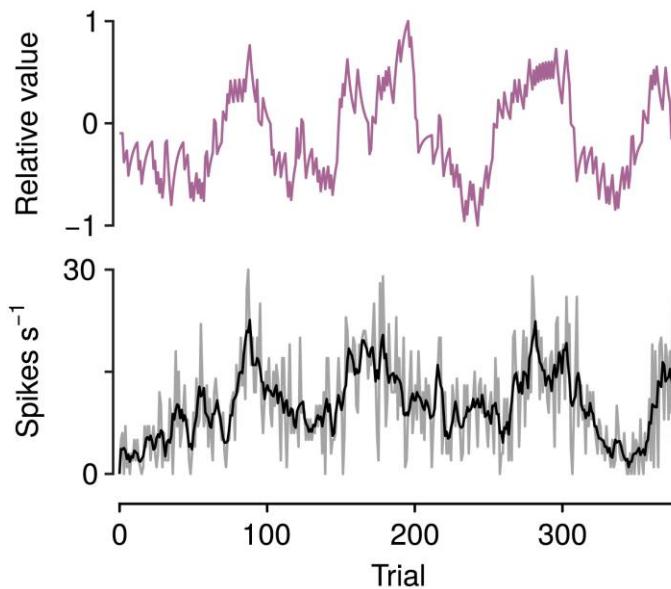
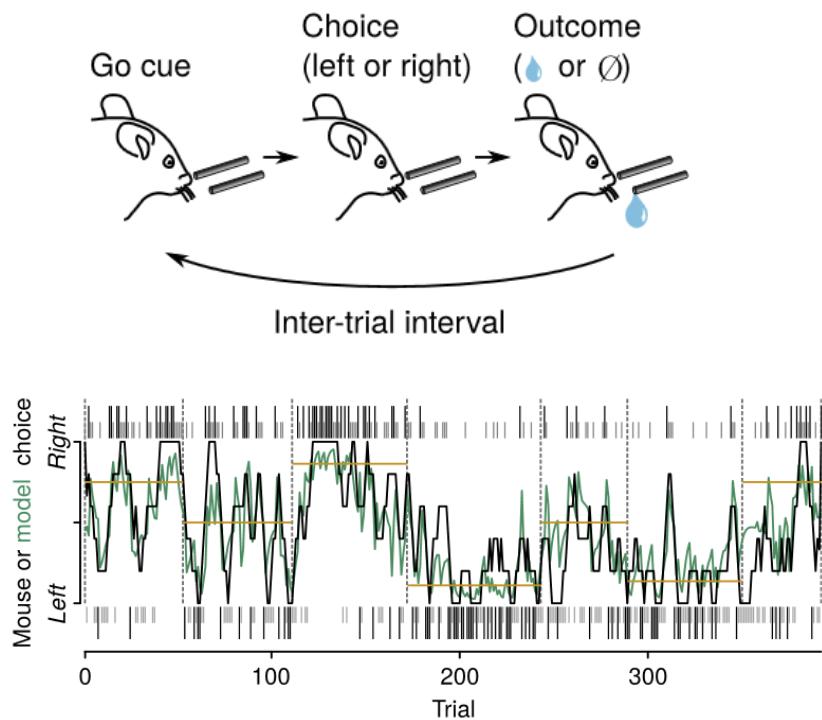
- Separate representation learning from TD learning.
- Learn representations in the deep recurrent cortical network using self-supervised learning
- Learn long-run values in the shallow, feed-forward, cortico-striatal network.

# Cortex as an active player in valuation and action selection



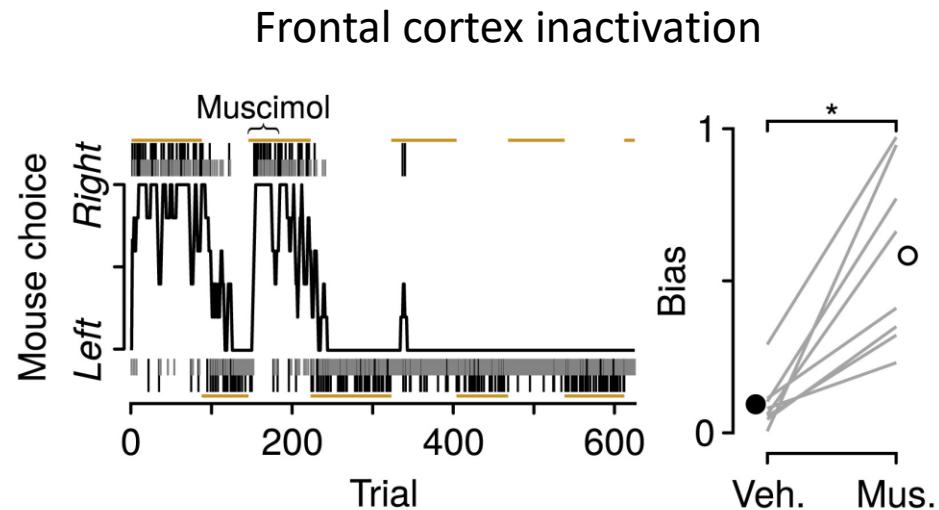
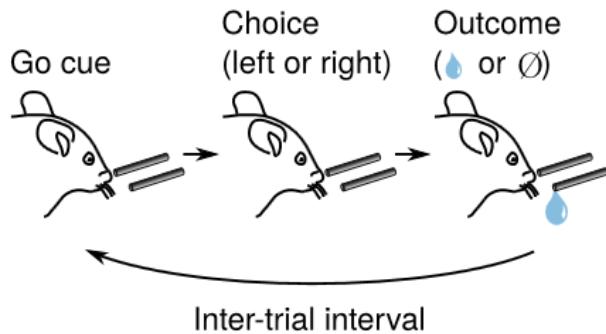
- Value signals are seen in frontal cortex during reward-guided decision making

# Cortex as an active player in valuation and action selection



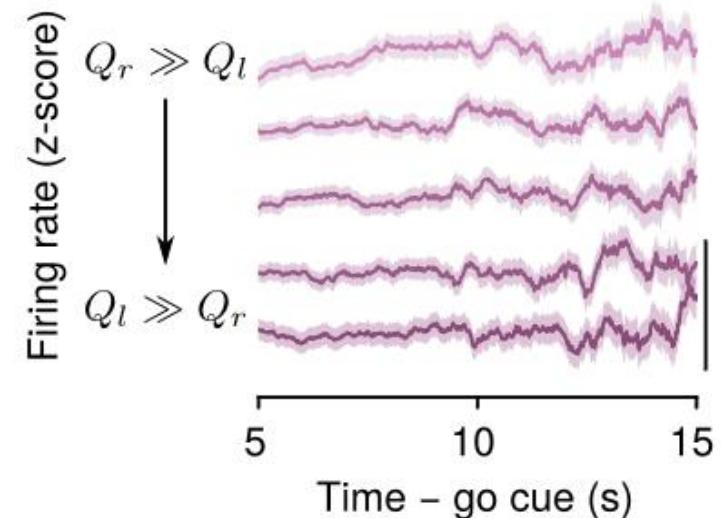
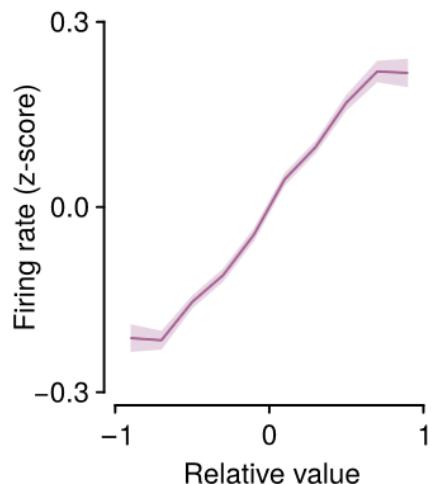
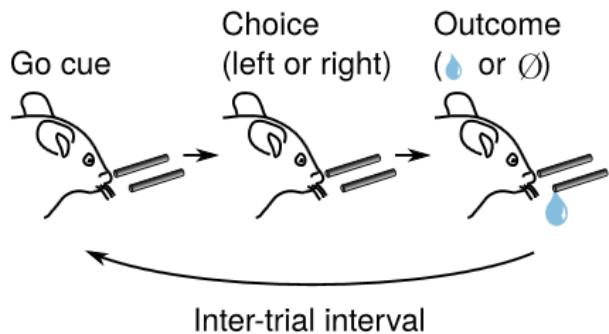
- Value signals are seen in frontal cortex during reward-guided decision making

## Cortex as an active player in valuation and action selection



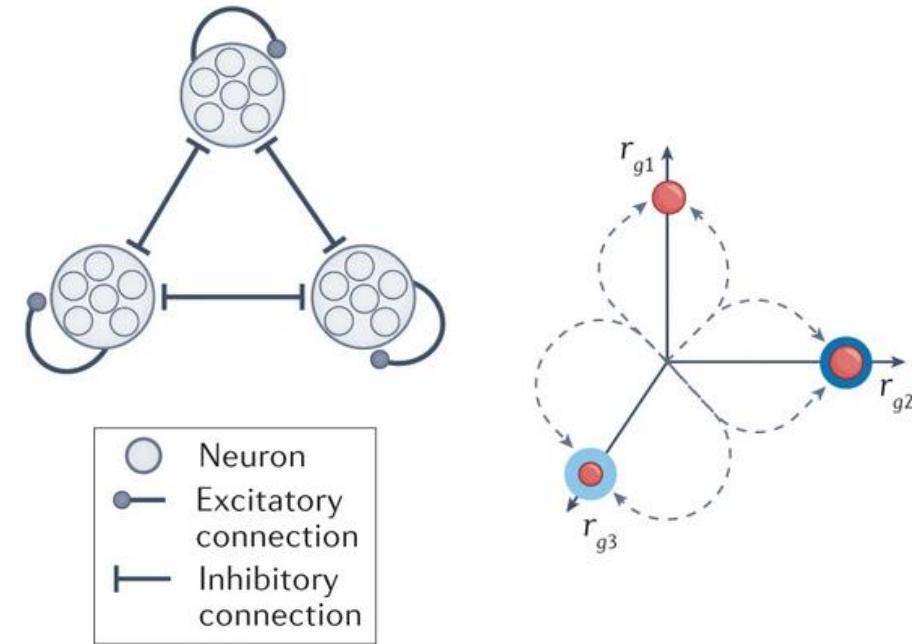
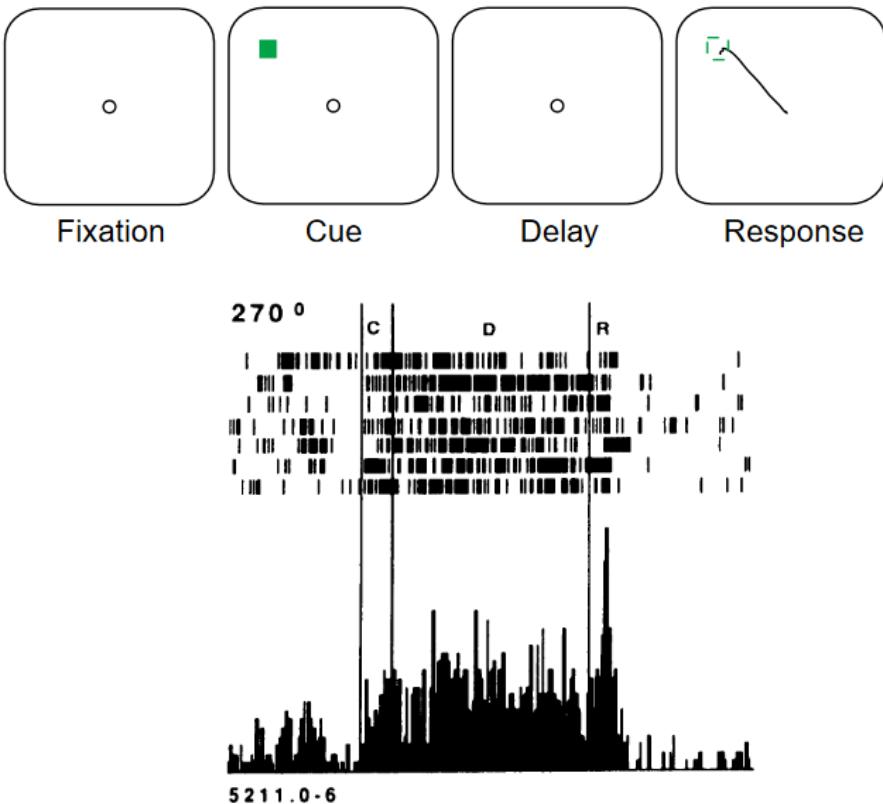
- Value signals are seen in frontal cortex during reward-guided decision making
- Inactivation of frontal cortex profoundly disrupts reward guided decision making

## Cortex as an active player in valuation and action selection



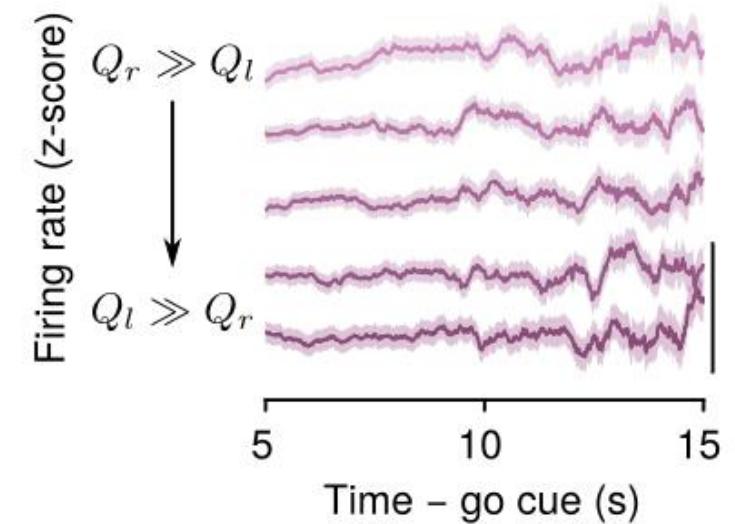
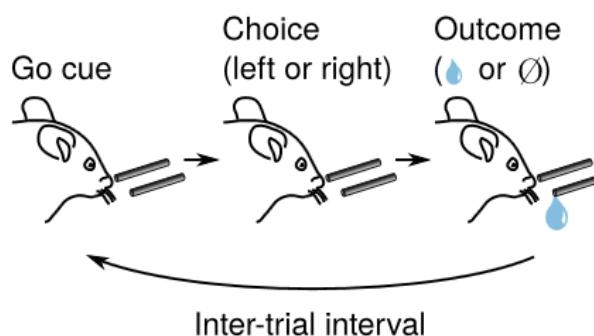
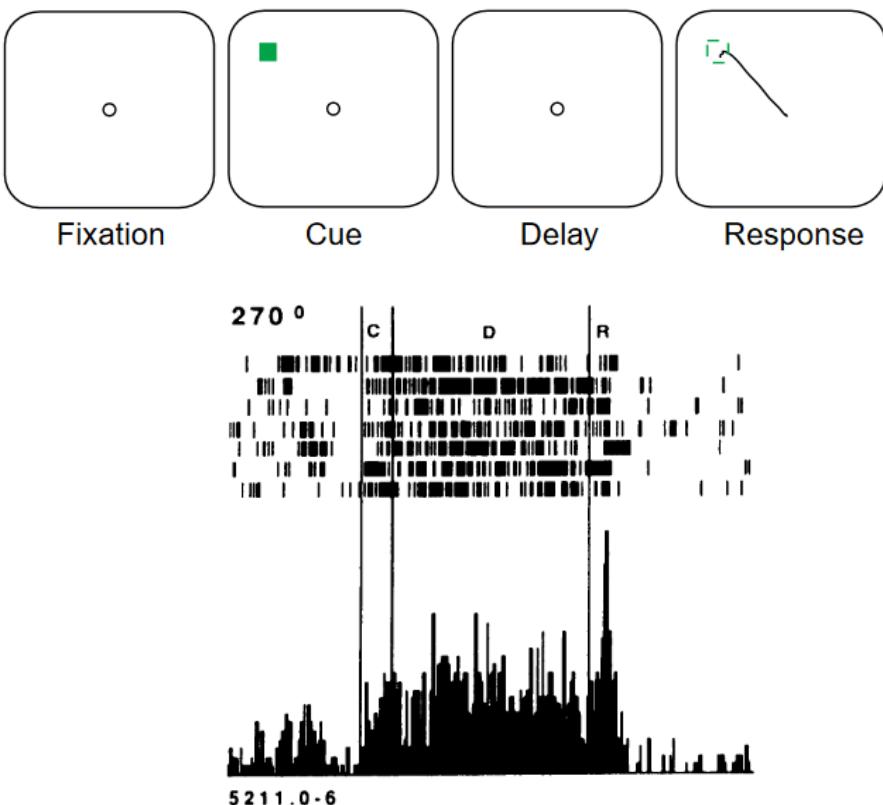
- Value signals are seen in frontal cortex during reward-guided decision making
- Inactivation of frontal cortex profoundly disrupts reward guided decision making
- Value information is present in persistent activity across long inter-trial intervals.

## Recurrent excitation and working memory



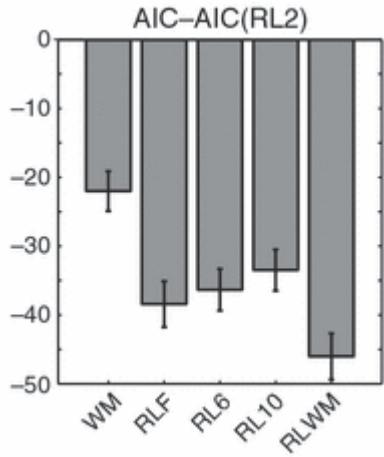
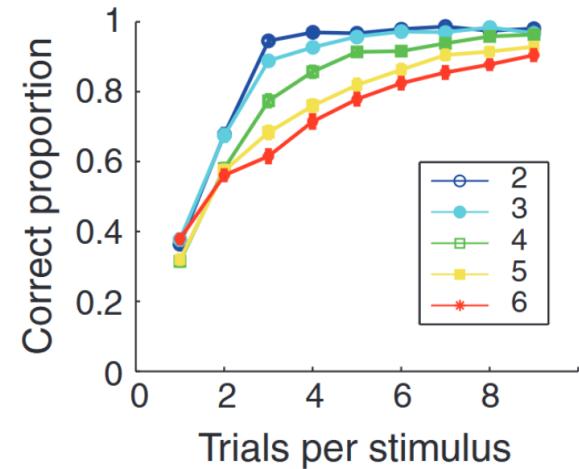
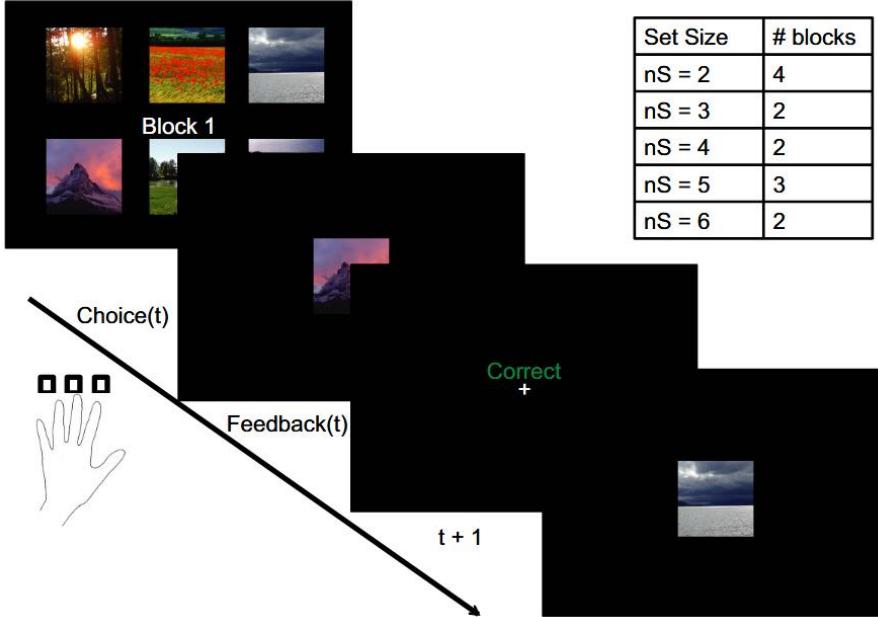
- Frontal cortex neurons famously show sustained persistent activity during the delay period of working memory tasks
- Networks with attractor dynamics have multiple stable fixed-points in their activity space, allowing information to be stored as patterns of recurrent activity.

# Reward guided decision making and working memory



- Both working memory tasks and reward guided decision tasks are partially observable.
- Persistent value coding suggests recent reward history is stored by attractor dynamics not synaptic weight changes.

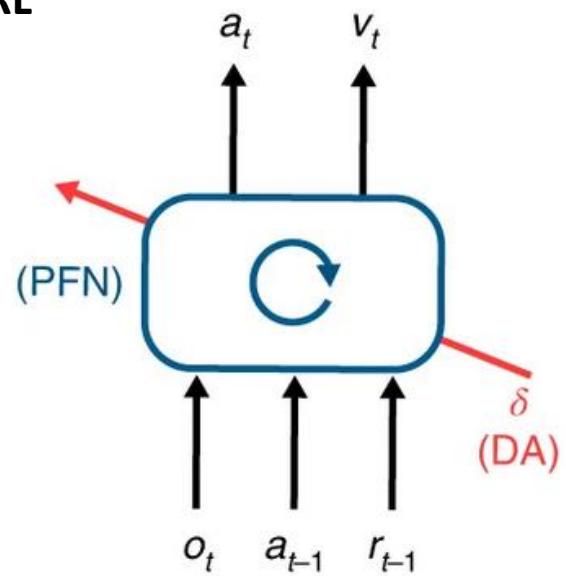
# Working memory in human action learning



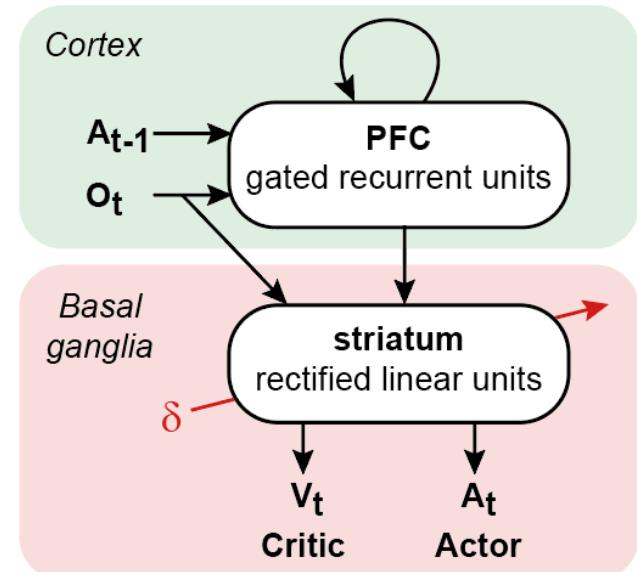
- Human trial-and-error learning shows strong set size effects consistent with limited working memory capacity.
- Data best fit by model combining fast capacity limited working memory with slow capacity unlimited RL.

# Learning recurrent dynamics for reward guided decision making

Meta RL



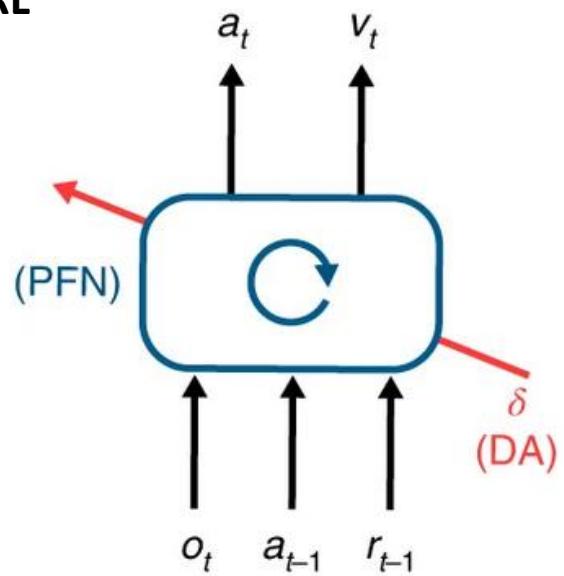
PFC-BG model



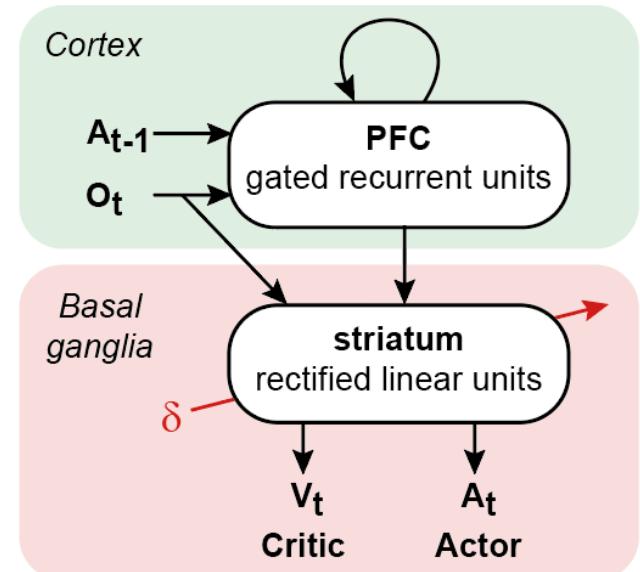
- Meta RL:
  - Recurrent network (LSTM) trained end-to-end using RL (A2C) to predict long run reward and select actions.
- PFC-BG model:
  - Recurrent “frontal cortex” network trained with self-supervised learning to predict the next observation.
  - Feed-forward “basal-ganglia” network trained with RL (A2C) to predict long run reward and select actions.

# Learning recurrent dynamics for reward guided decision making

Meta RL

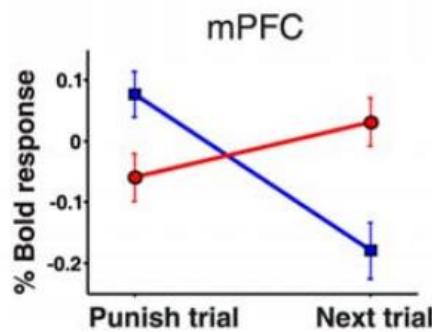
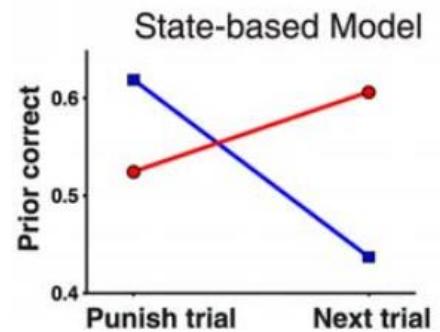
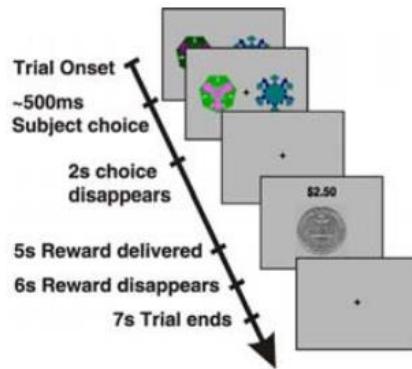


PFC-BG model

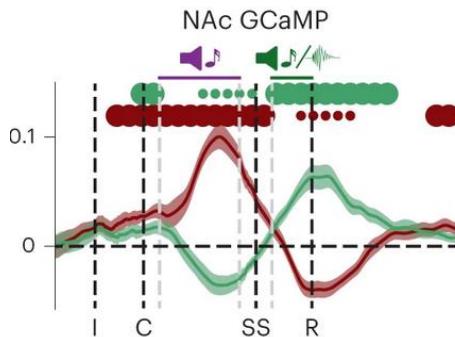
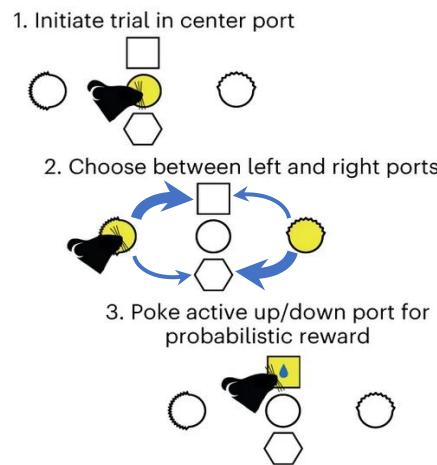
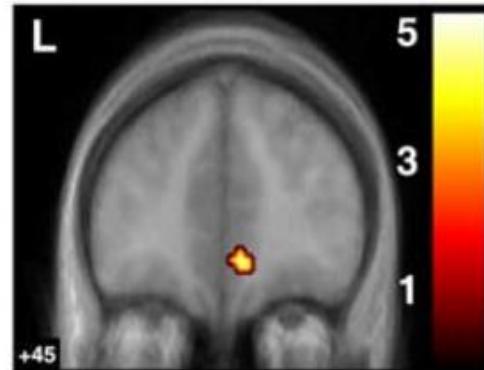


- In both models, synaptic plasticity on slow timescale (task acquisition) sculpts recurrent dynamics that drive behavioural adaptation on a fast timescale (trial-to-trial).

# Recurrent network models and inferred value updates



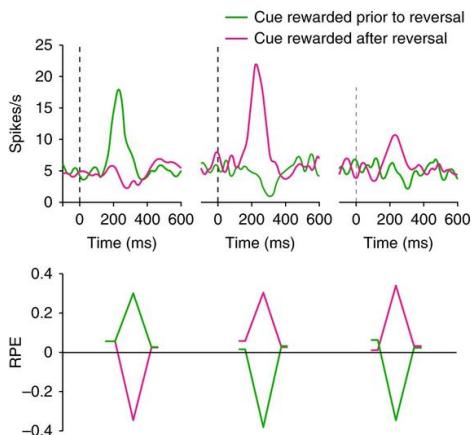
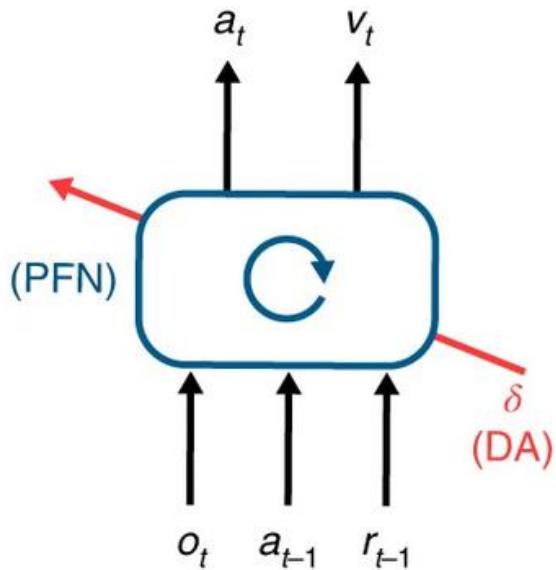
State-based > standard RL



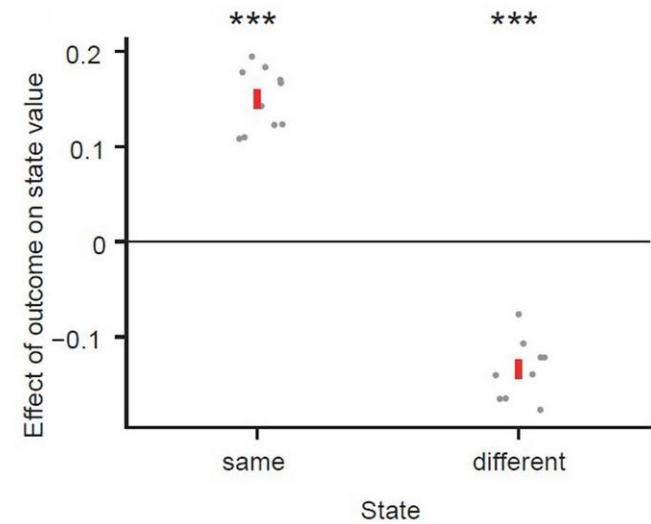
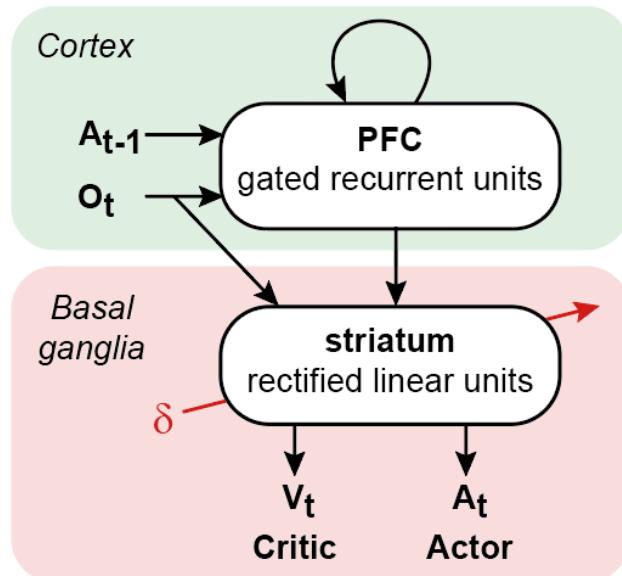
- In reversal learning tasks, evidence one option is good is also evidence the other option is bad.
- Brain signals reflect such inferred value updates, both in humans and mice, frontal cortex and dopamine.

# Recurrent network models and inferred value updates

## Meta RL



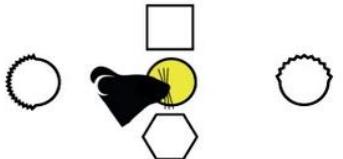
## PFC-BG model



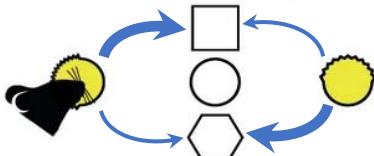
- Both models reproduce inferred value updates when reward probabilities are anticorrelated.

# Frontal value signals or state representations?

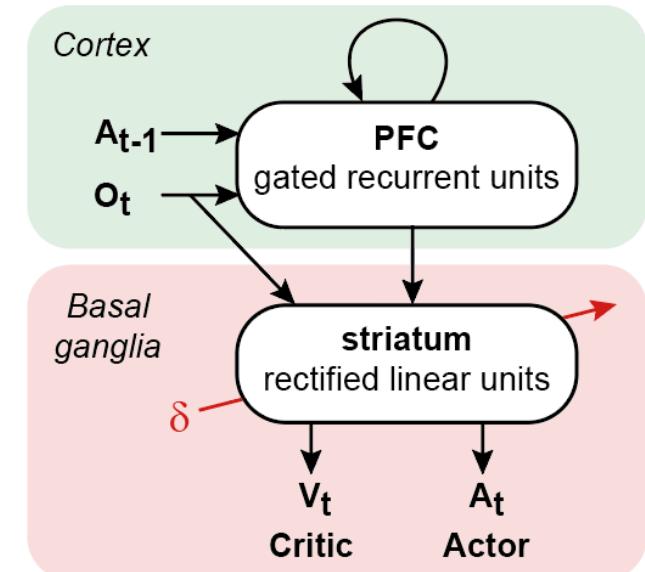
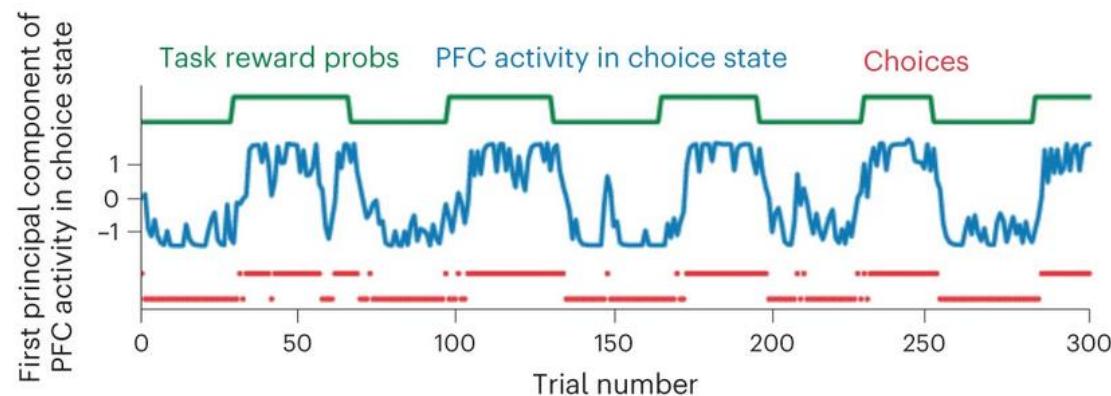
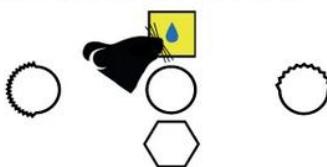
1. Initiate trial in center port



2. Choose between left and right ports

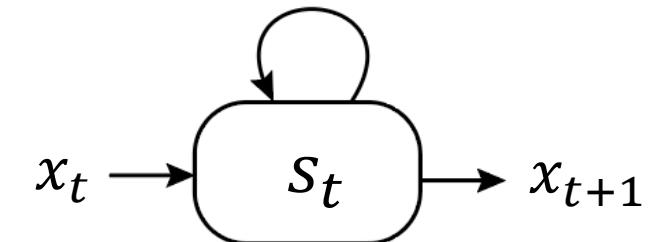
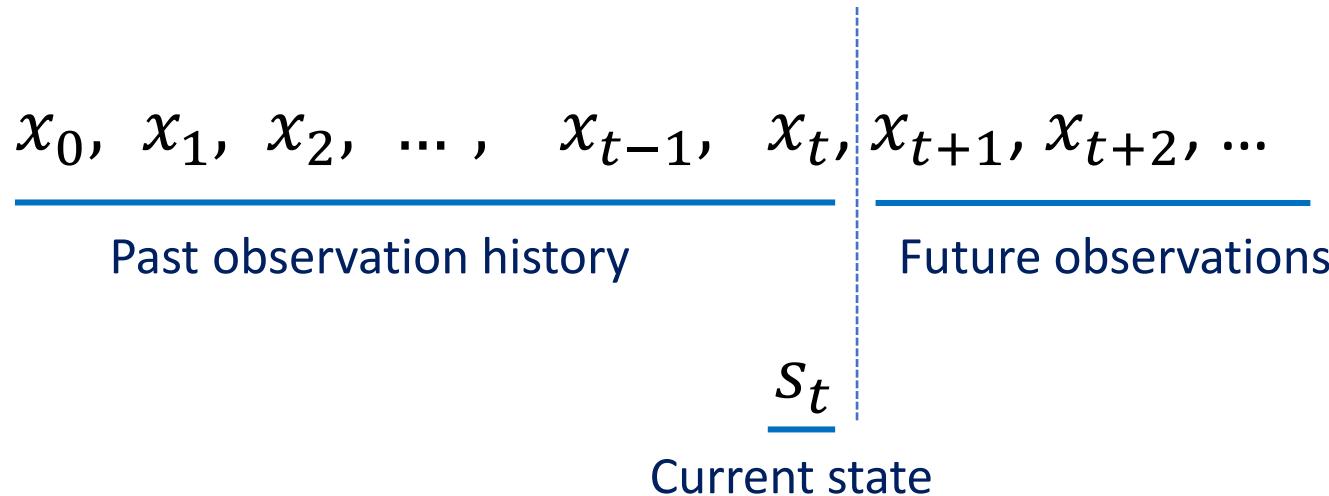


3. Poke active up/down port for probabilistic reward



- Recurrent network activity tracks task reward probabilities despite being trained only to predict the next observation.
- Not RL “value” (long run reward) as recurrent network trained only with self-supervised learning.

## Frontal value signals or state representations?

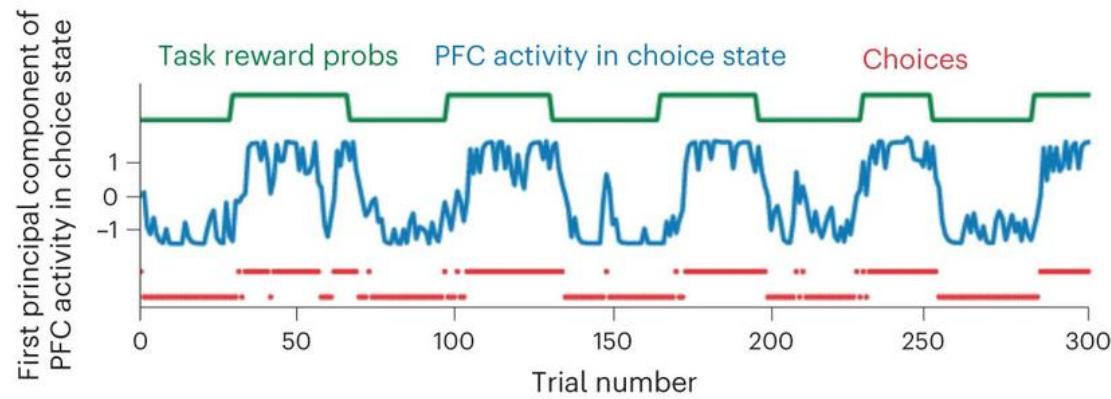
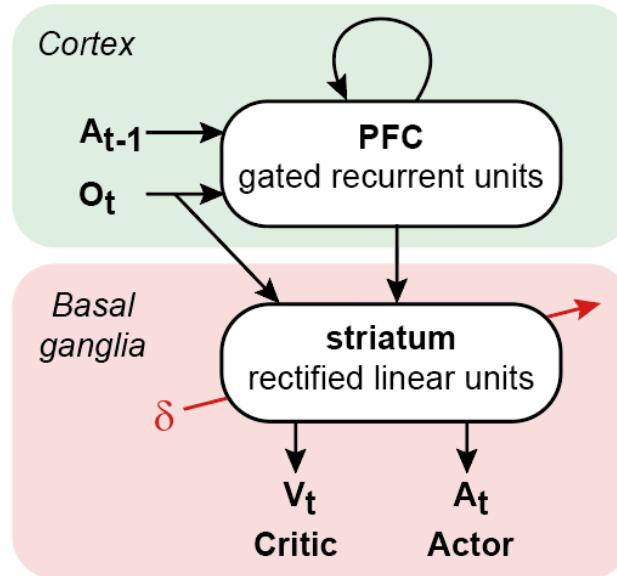


- A “Markov” state summarizes all information in the observation history necessary for making any prediction.

*A common strategy for finding a Markov state is to look for something compact that is recursively updatable and enables accurate short-term predictions. In fact, it is only necessary to make accurate one-step predictions.*

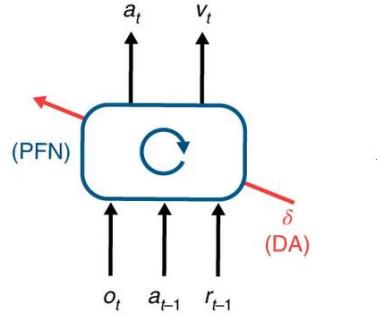
*Sutton and Barto*

# Frontal value signals or state representations?

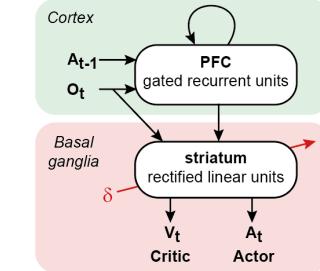


- PFC recurrent network learns a Markov state representation by predicting observations.
- BG network learns values and policy over observable and learned state features.
- “Value” signals in frontal cortex are in fact beliefs about latent task state.

# Frontal value signals or state representations?

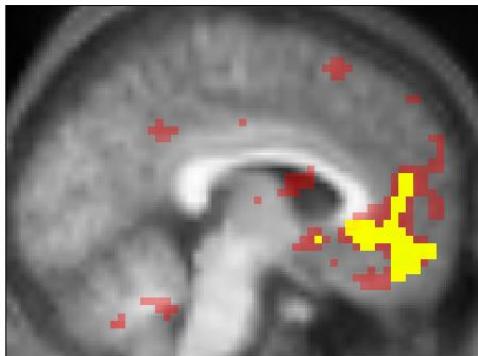


Frontal cortex as  
value/policy

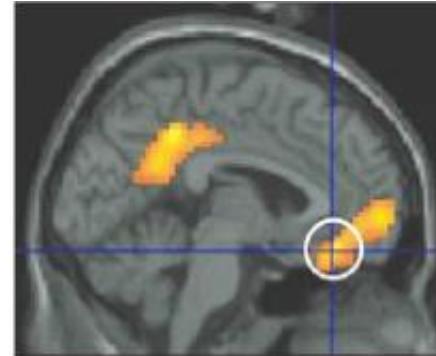


Frontal cortex as  
state representation

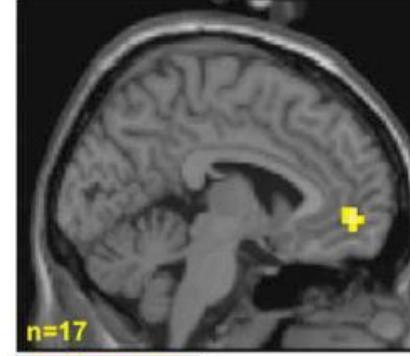
## Frontal value signals or state representations?



Money –  
value predicted

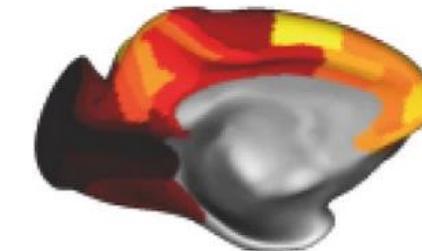


Faces -  
attractiveness



Coke or Pepsi –  
behavioral preference

## Dopamine D1 receptor density



First principal component of 14  
receptor densities



- Ventromedial frontal cortex consistently responds to valued stimuli across diverse domains
  - Gradients in receptor densities and many other neuronal properties across cortical hierarchy
- Reward does appear special in frontal cortex, learning algorithm remain poorly constrained.

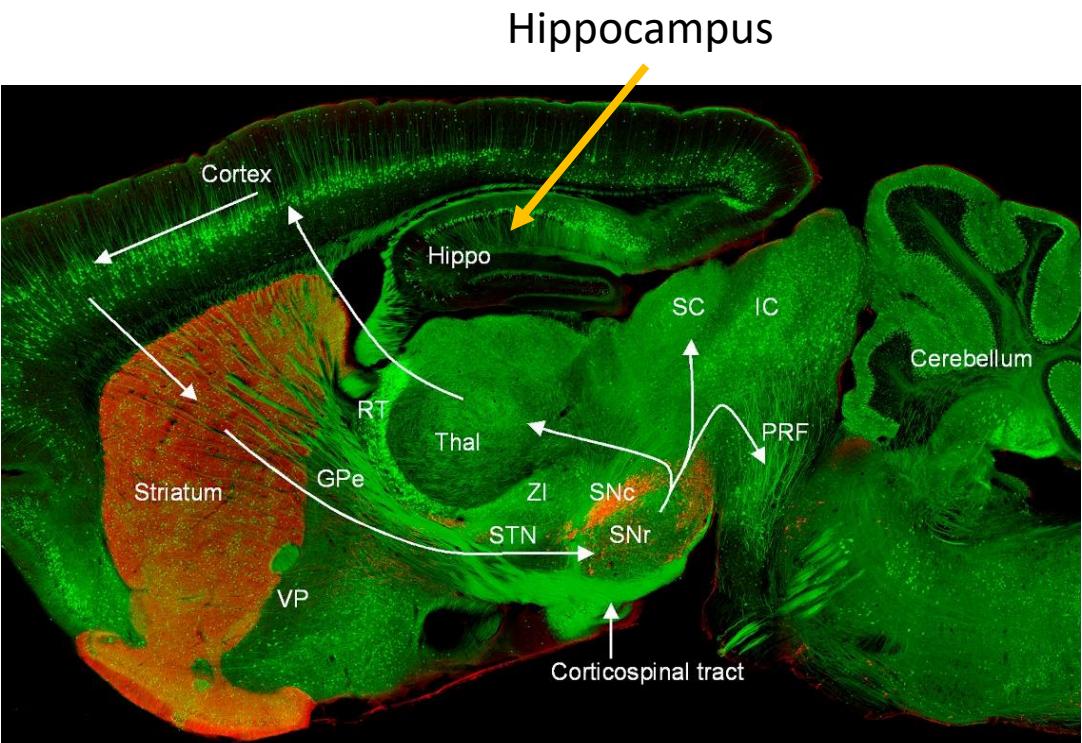
## Talk aim and overview

**Question:** How do learning algorithms map onto brain structure to solve the biological action control problem?

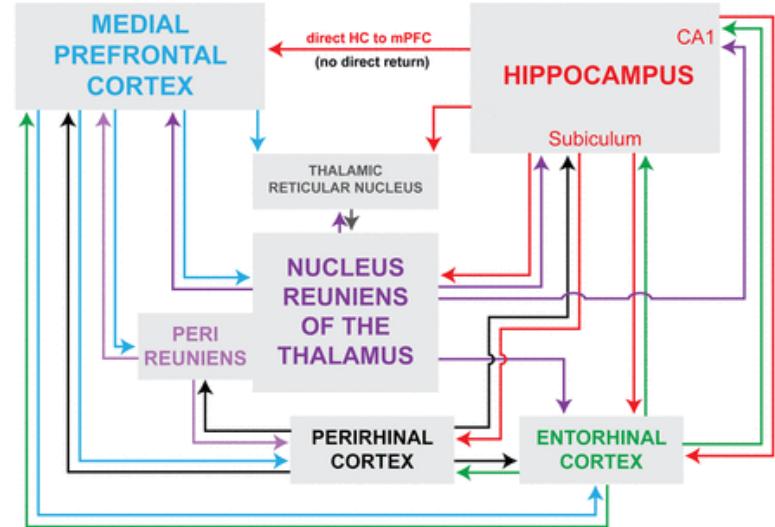
**Talk outline:**

1. Striatum and dopamine: The brain's temporal difference reinforcement learning system?
2. Basal ganglia outputs and control of action selection.
3. Cortex: State representation and beyond
4. Hippocampus: Sequence generation and model-based control

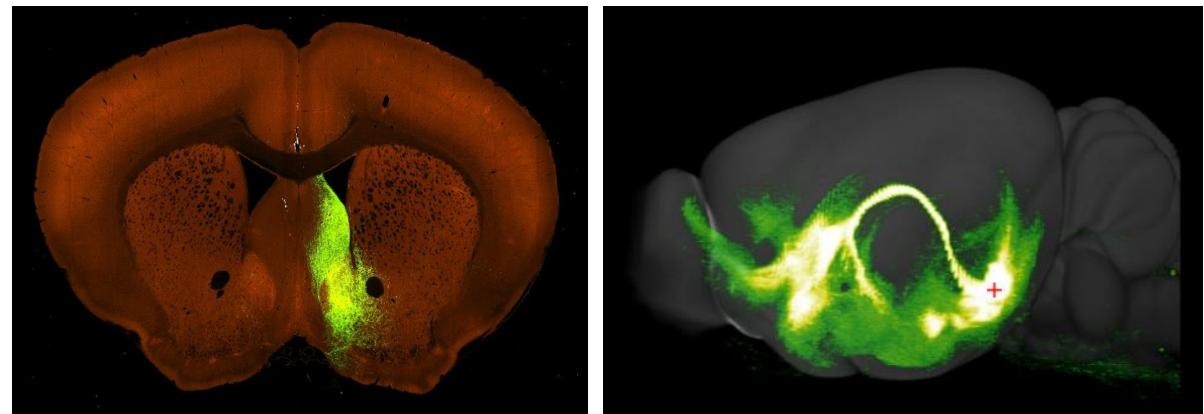
# Hippocampus: Sequence generation and model-based control



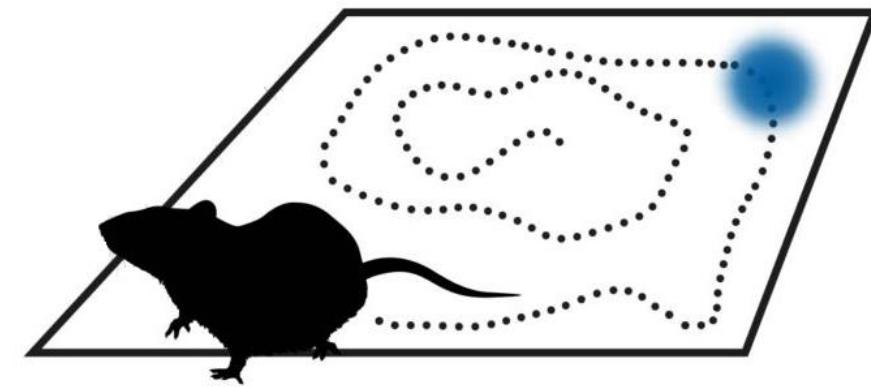
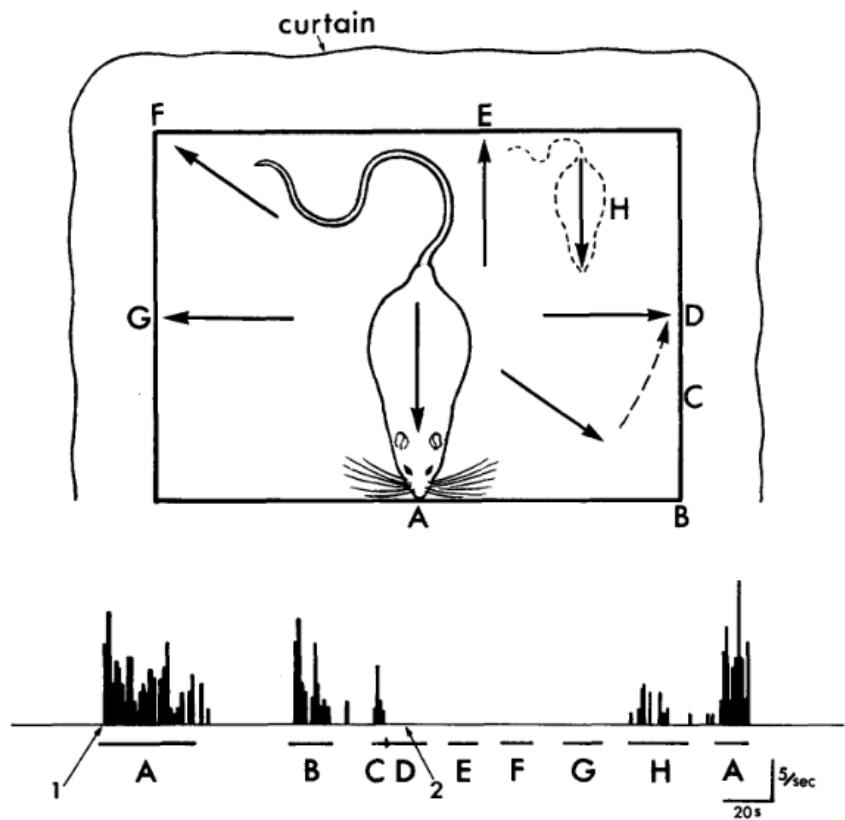
## Connections with frontal cortex



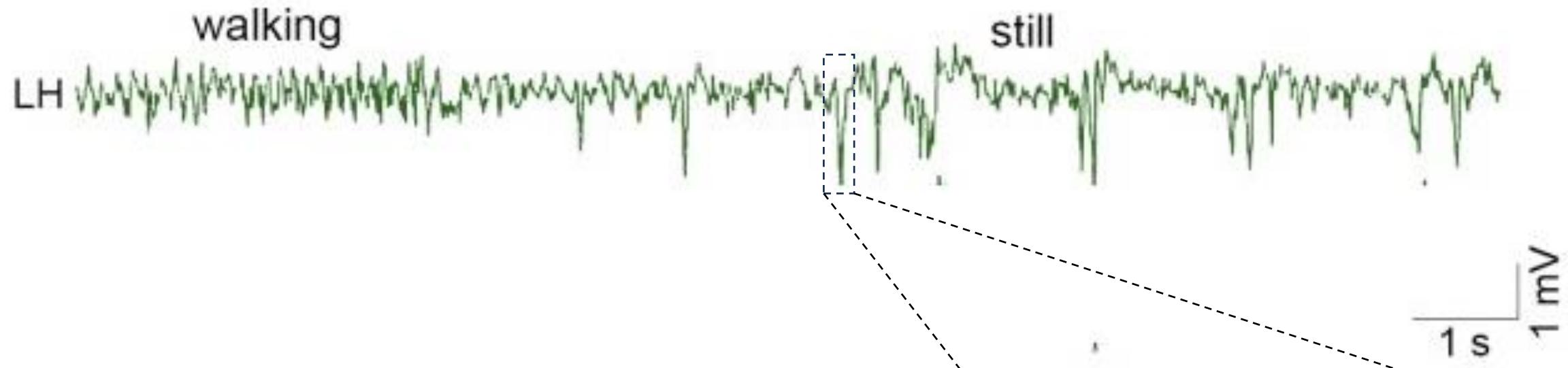
## Projection to ventral striatum



## Hippocampal place cells

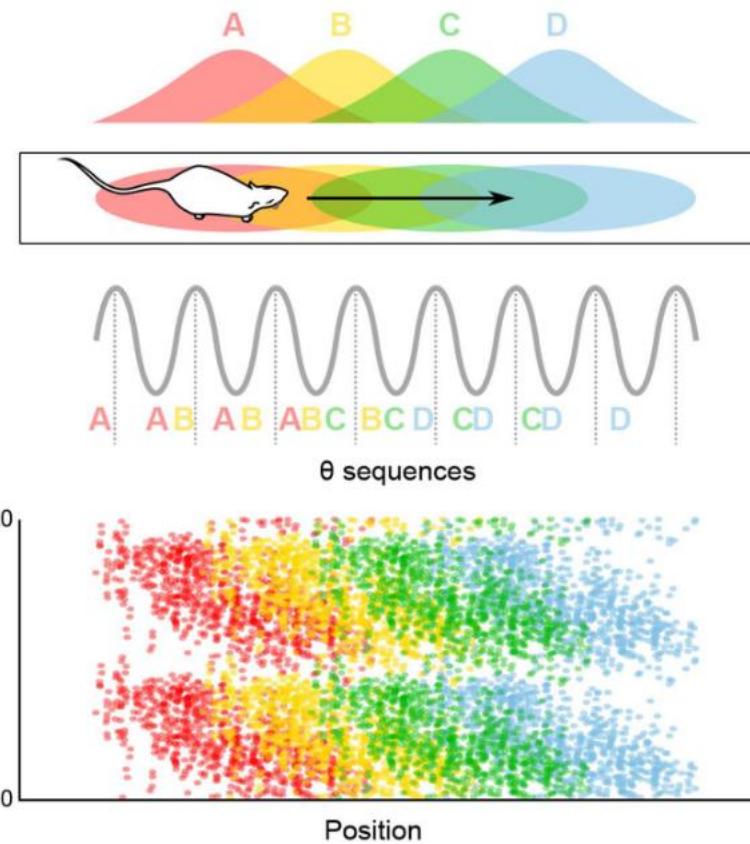
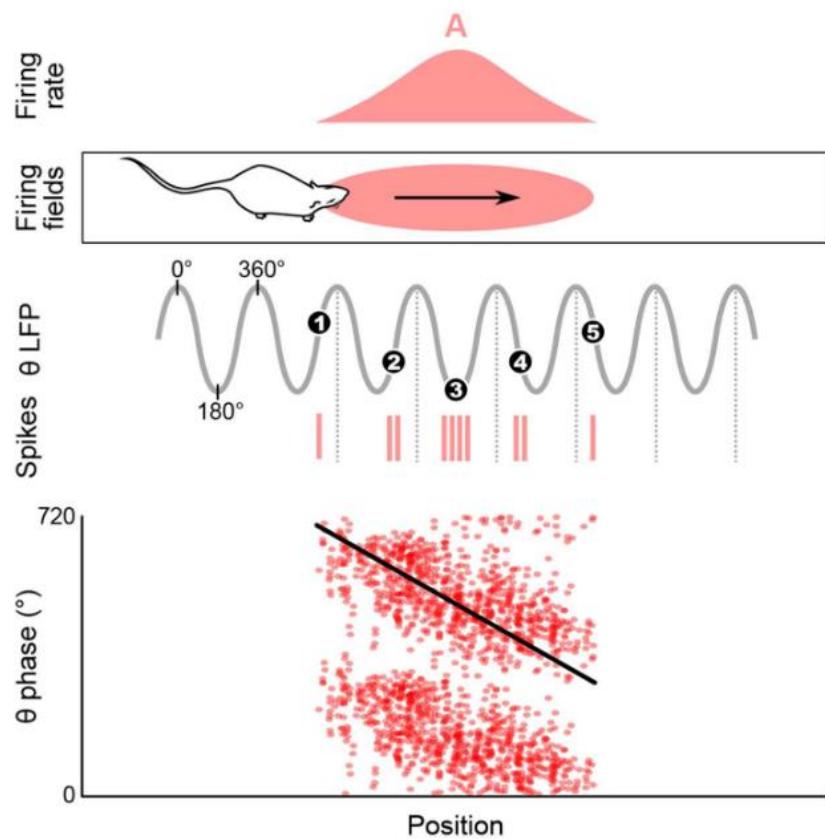


## Hippocampal sequences



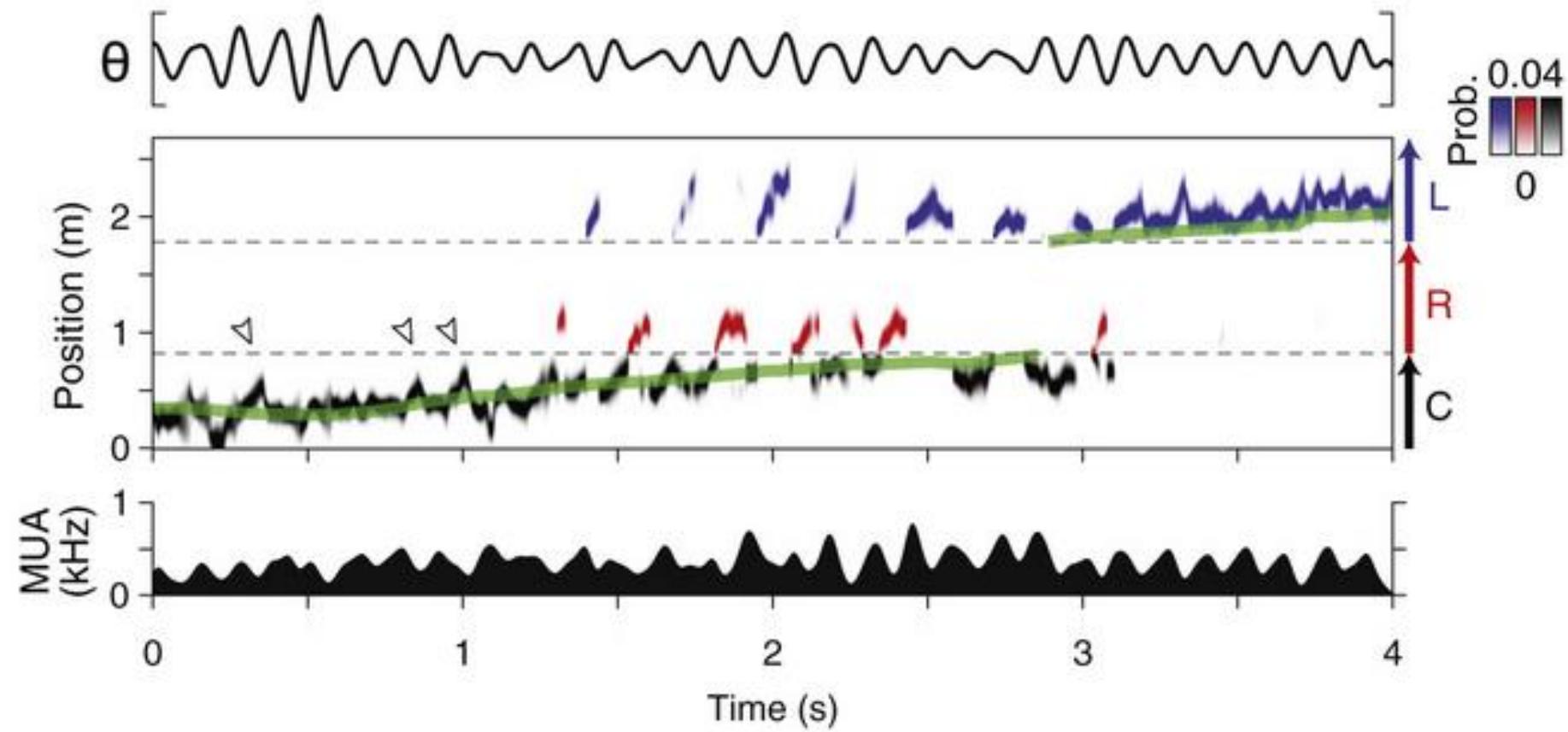
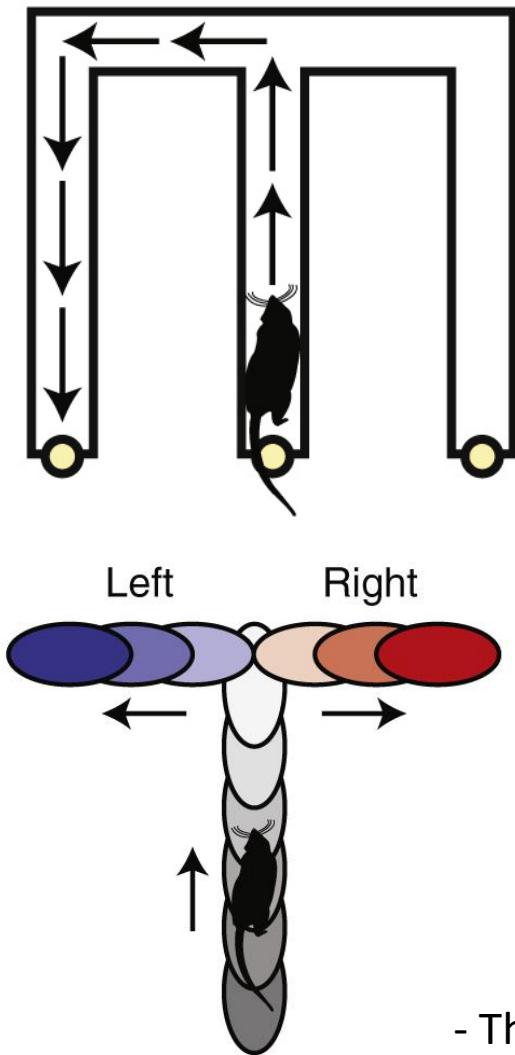
- ~10Hz “Theta” oscillations during movement
- High frequency “sharp wave ripples” (SWR) during immobility, reward consumption, sleep.

## Hippocampal sequences: Theta oscillations



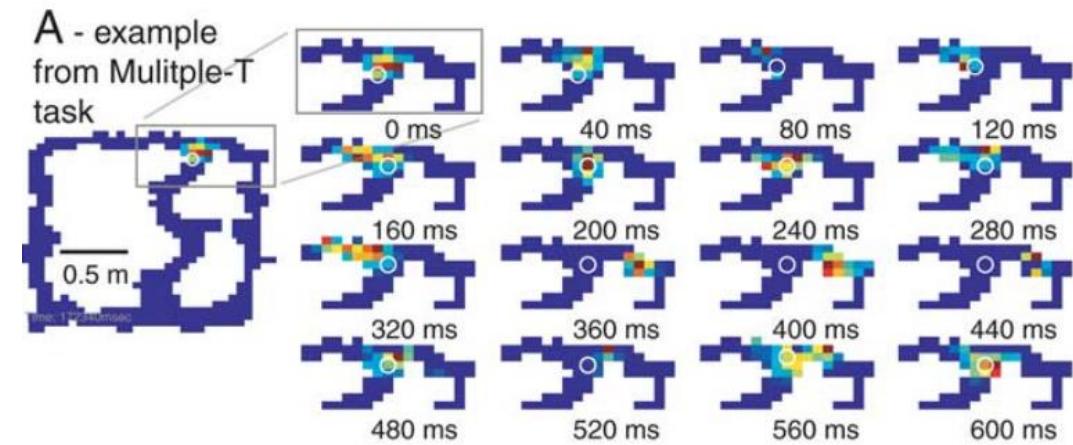
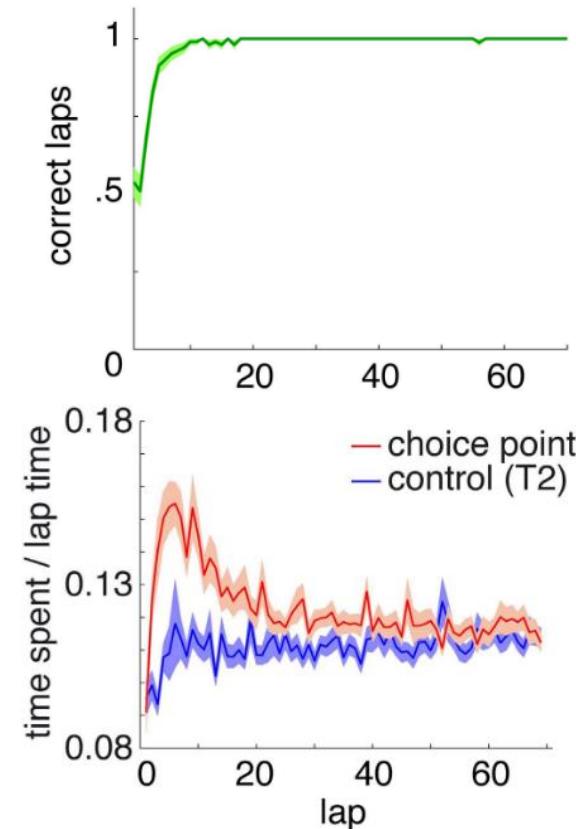
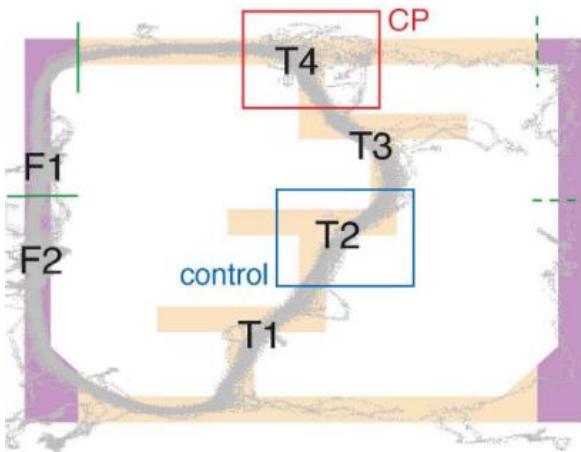
- Firing phase shifts earlier relative to theta in cycle during movement through place field
- Place cells along trajectory fire in a compressed sequence on each theta cycle.

## Hippocampal sequences: Theta oscillations



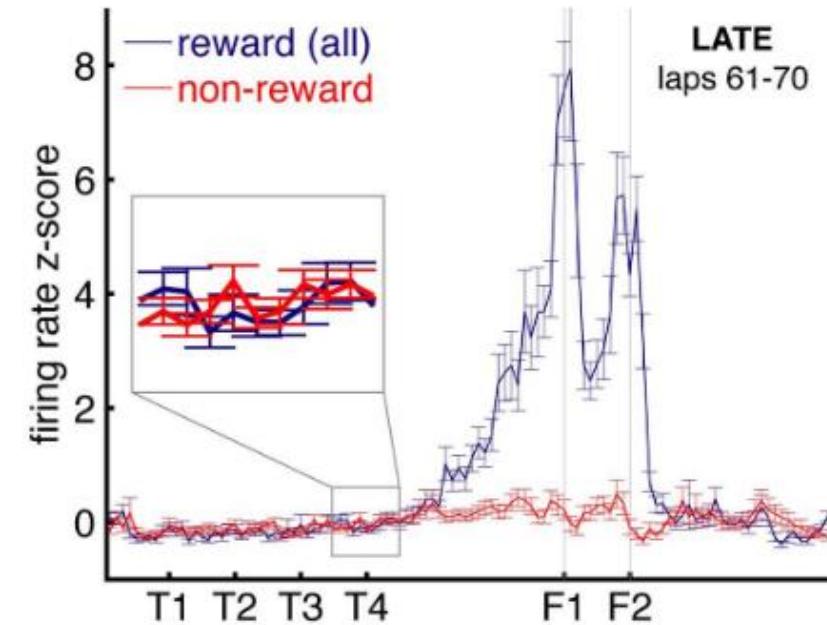
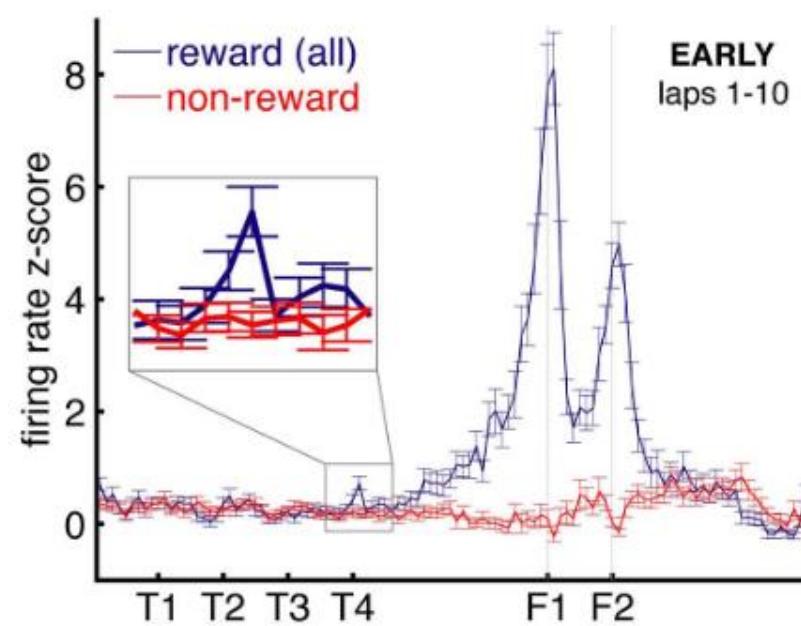
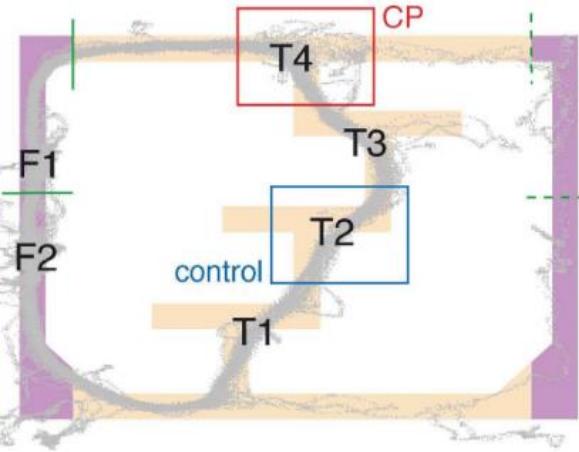
- Theta sweeps alternate between possible future trajectories at decision points.

## Hippocampal sequences: Theta oscillations



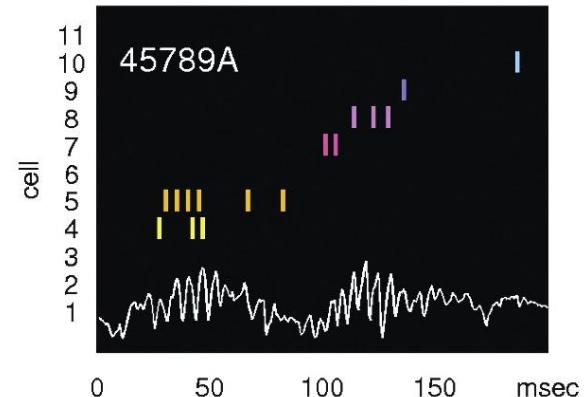
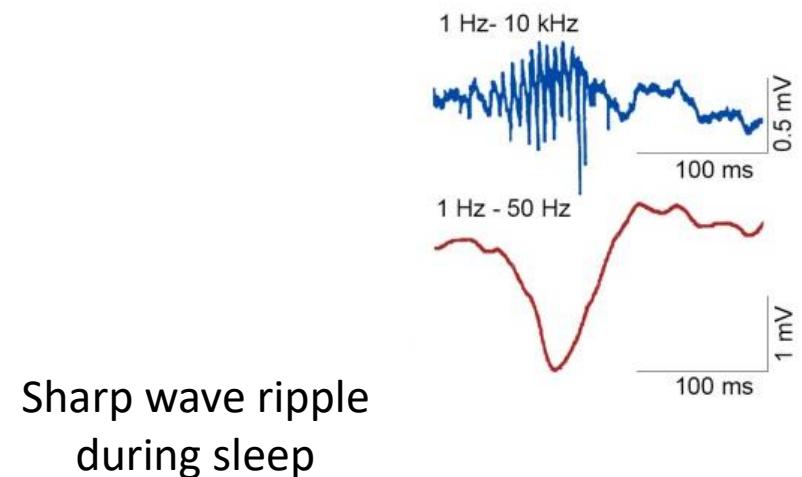
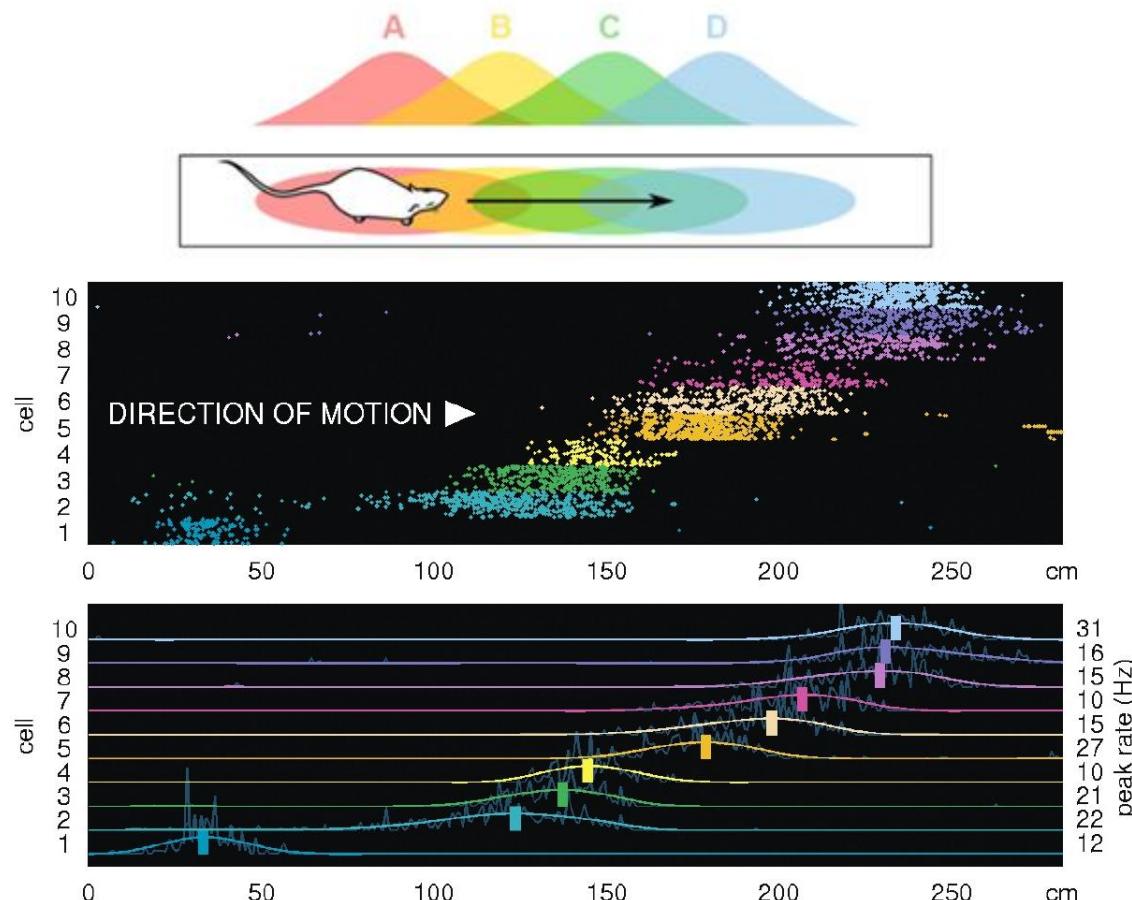
- “Deliberation” at important decision point during initial learning accompanied by theta sweeps down choice arms.

## Hippocampal sequences: Theta oscillations



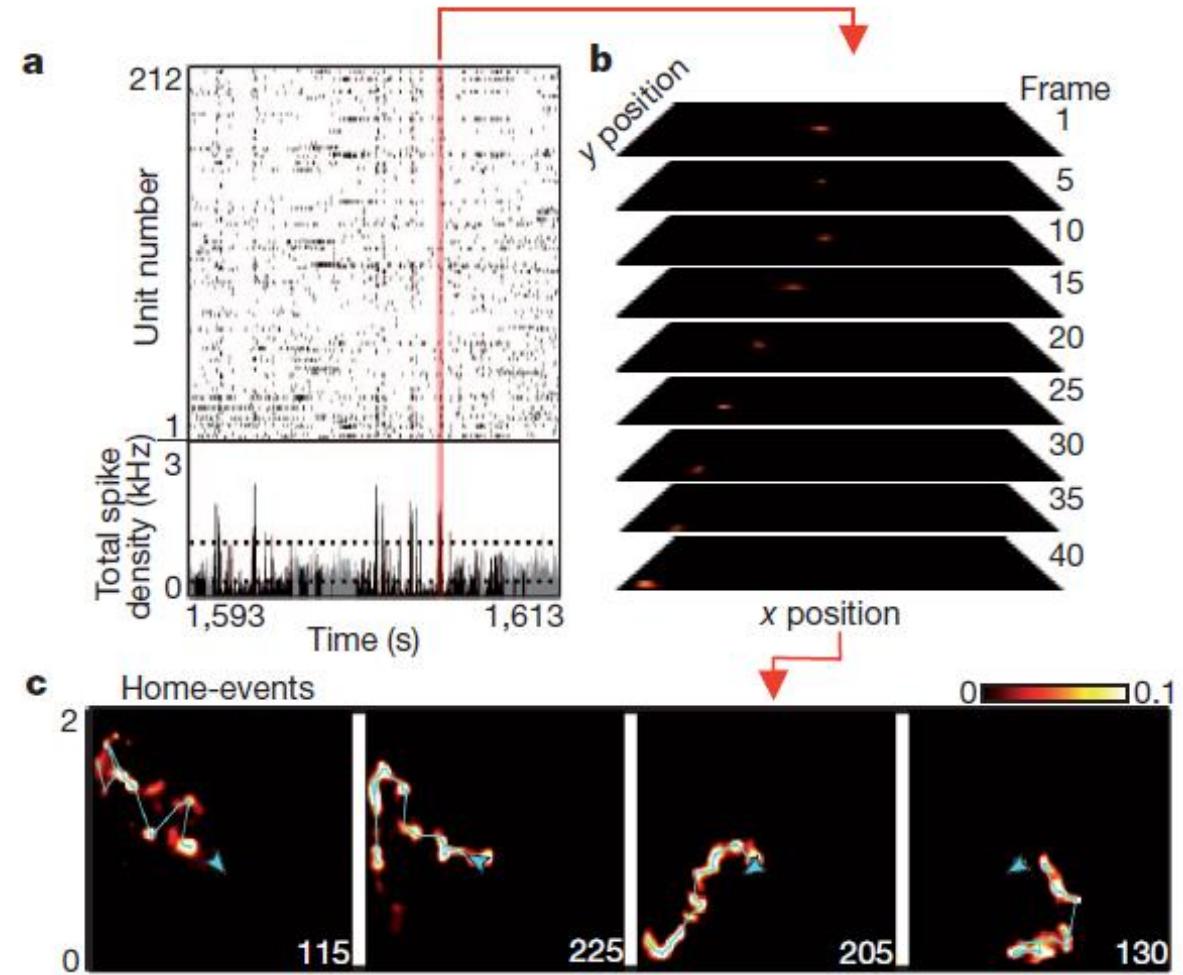
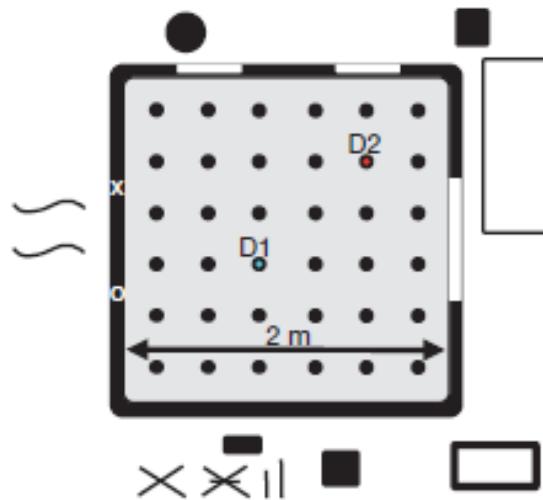
- “Deliberation” at important decision point during initial learning accompanied by theta sweeps down choice arms.
- Activation of “value” neurons in ventral striatum at decision point during early learning.
- Theta sweeps as online planning?

## Hippocampal sequences: Sharp wave ripples



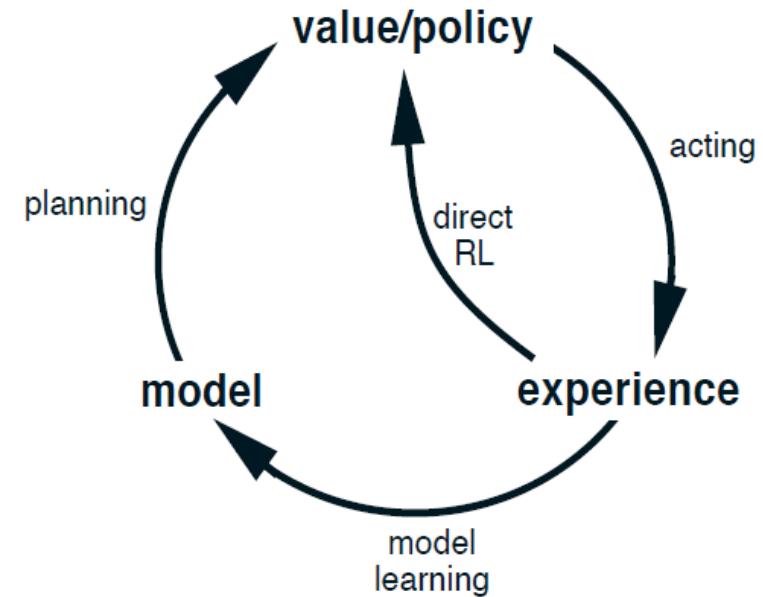
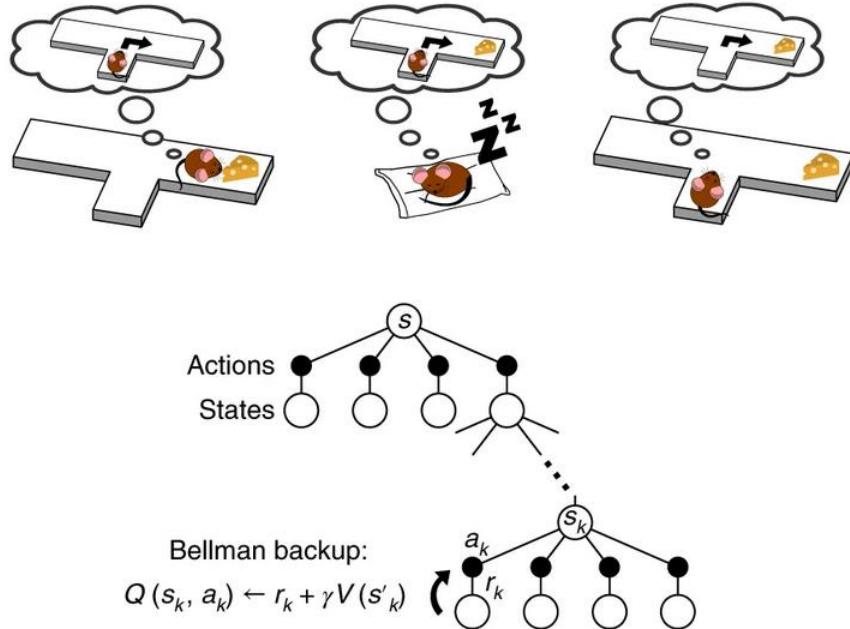
- “Replay” of behavioural sequences during sleep sharp wave ripples

## Hippocampal sequences: Sharp wave ripples



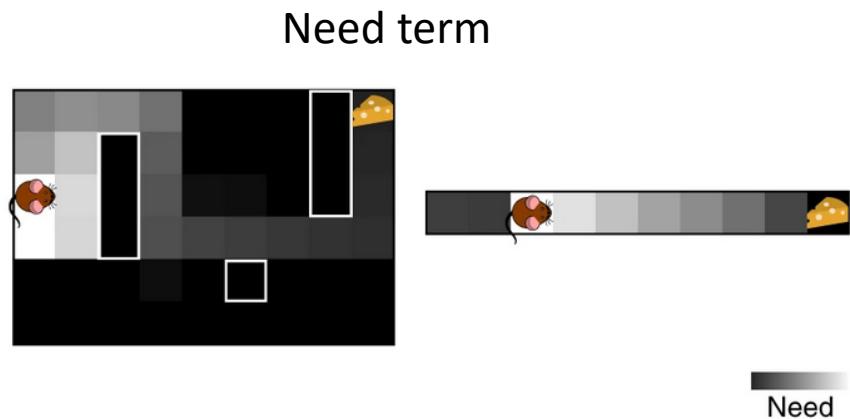
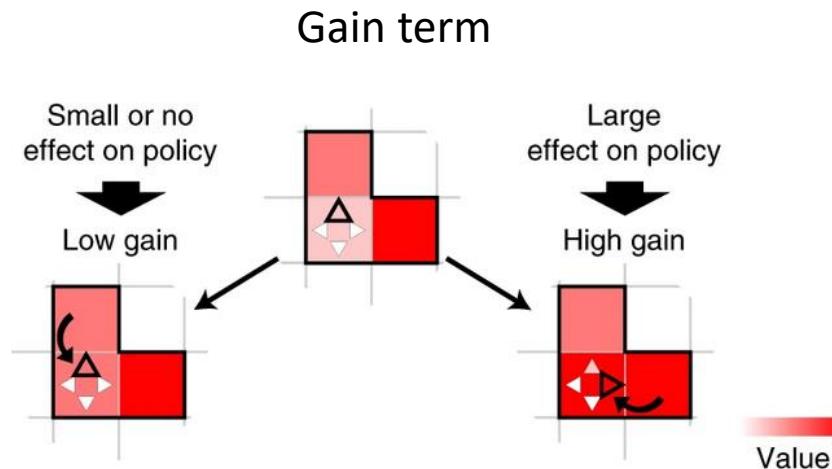
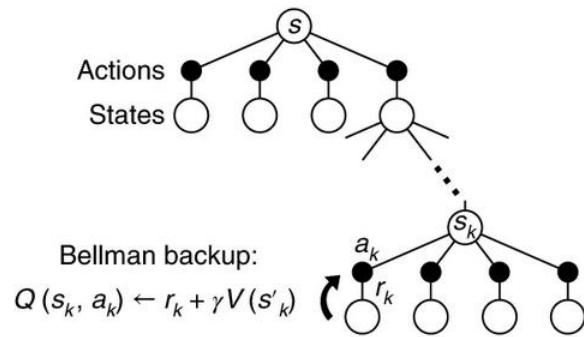
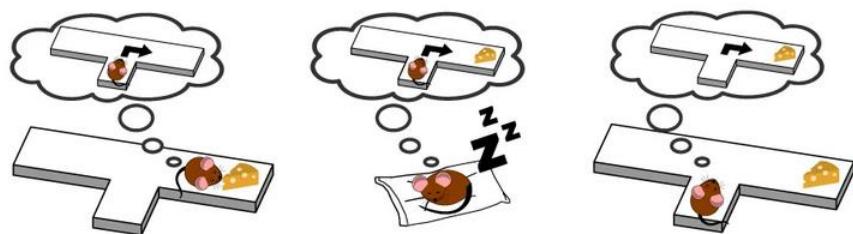
- Awake SWR sequences play our trajectories towards remembered goal locations before navigating to them.

## Sharp wave ripples as DYNA



What state transitions should you replay if you are using them to do DYNA-like action-value updates?

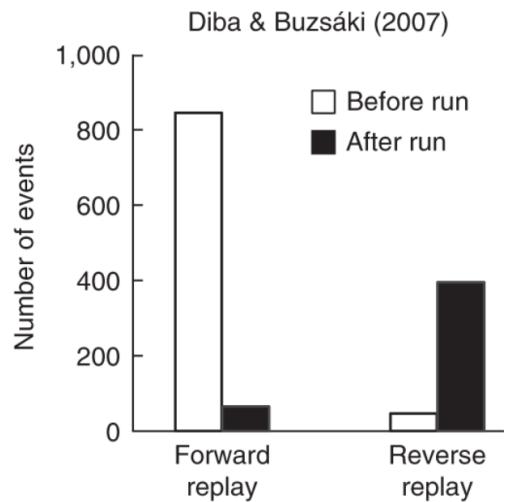
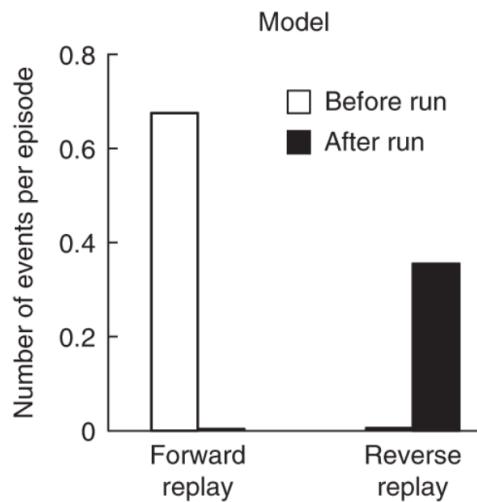
## Sharp wave ripples as DYNA



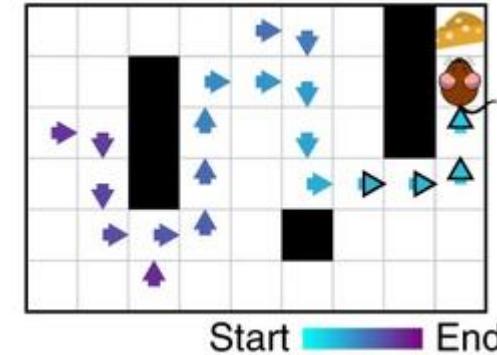
What state transitions should you replay if you are using them to do DYNA-like action-value updates?

- Gain term: Replay transitions that will drive large value updates.
- Need term: Replay transitions that you are likely to take in the near future.

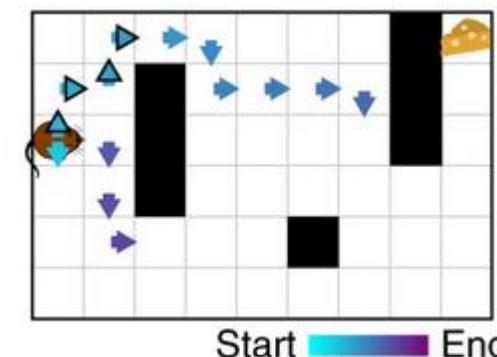
## Sharp wave ripples as DYNA



Reverse replay after reward driven by Gain term



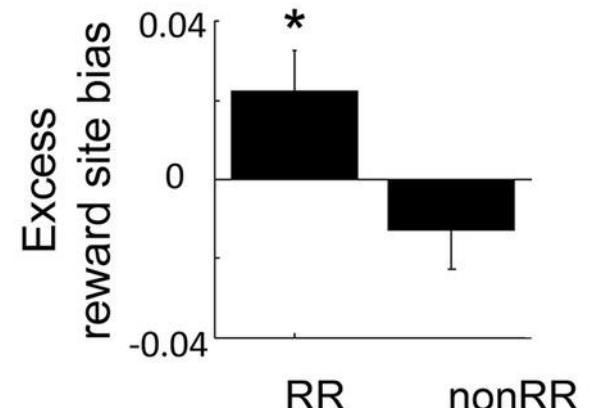
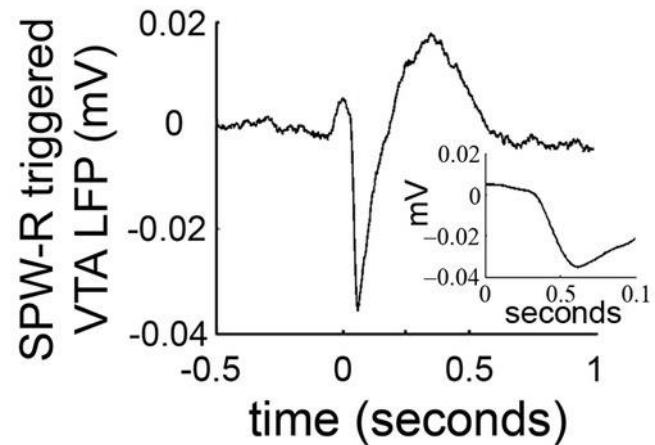
Forward replay before navigation driven by Need term



What state transitions should you replay if you are using them to do DYNA-like action-value updates?

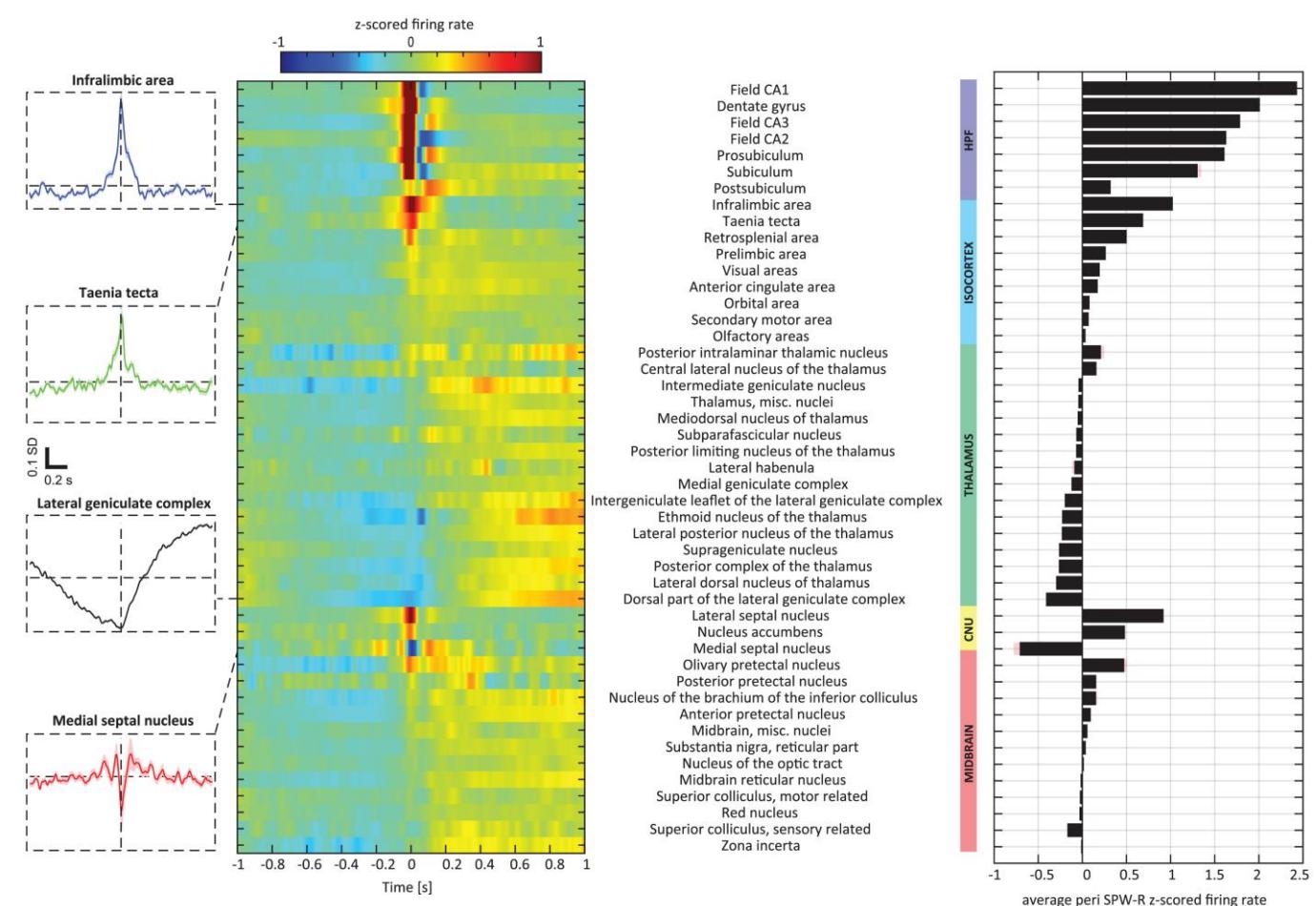
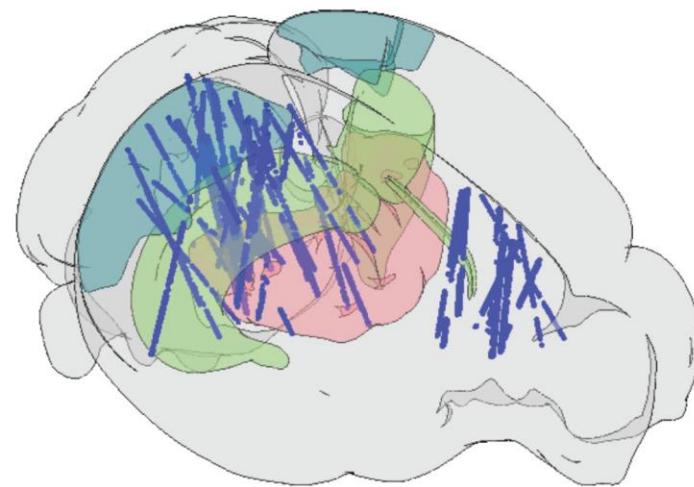
- Gain term: Replay transitions that will drive large value updates.
- Need term: Replay transitions that you are likely to take in the near future.

## Sharp wave ripples as DYNA



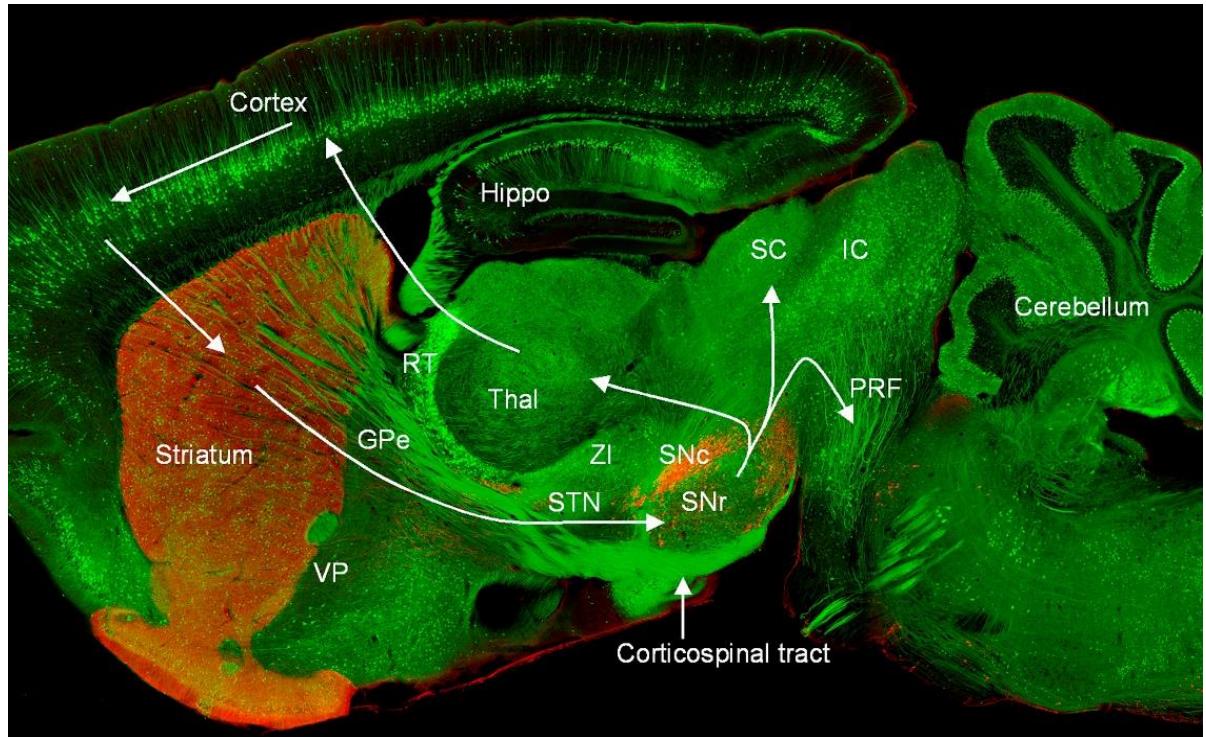
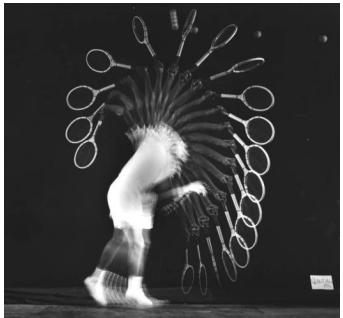
- Reward responsive neurons in VTA (putative dopamine neurons) are activated during SWRs

# Sharp wave ripples engage brain-wide networks



- Brain-wide activity is modulated around sharp-wave ripples.
  - Ripples are thought to be important for consolidation of episodic memories from hippocampus into cortex.
- DYNA-like value learning is likely part of the story but not the whole picture.

# Summary



- Brains evolved to solve a challenging control problem characterised by delayed rewards and a complex state, action and observation space.
- The brain's solution likely combines aspects of temporal difference value learning, self-supervised representation learning, and model-based planning, distributed across a set of specialised brain networks implementing different algorithms.

Thank you for listening.