

# Reinforcement Learning

Quentin Huys

Division of Psychiatry and Max Planck UCL Centre for Computational Psychiatry and Ageing Research, UCL  
Complex Depression, Anxiety and Trauma Service, Camden and Islington NHS Foundation Trust

Systems and Theoretical Neuroscience Course, 4/12/18





$$\{a_t\} \leftarrow \operatorname{argmax}_{\{a_t\}} \sum_{t=1}^{\infty} r_t$$

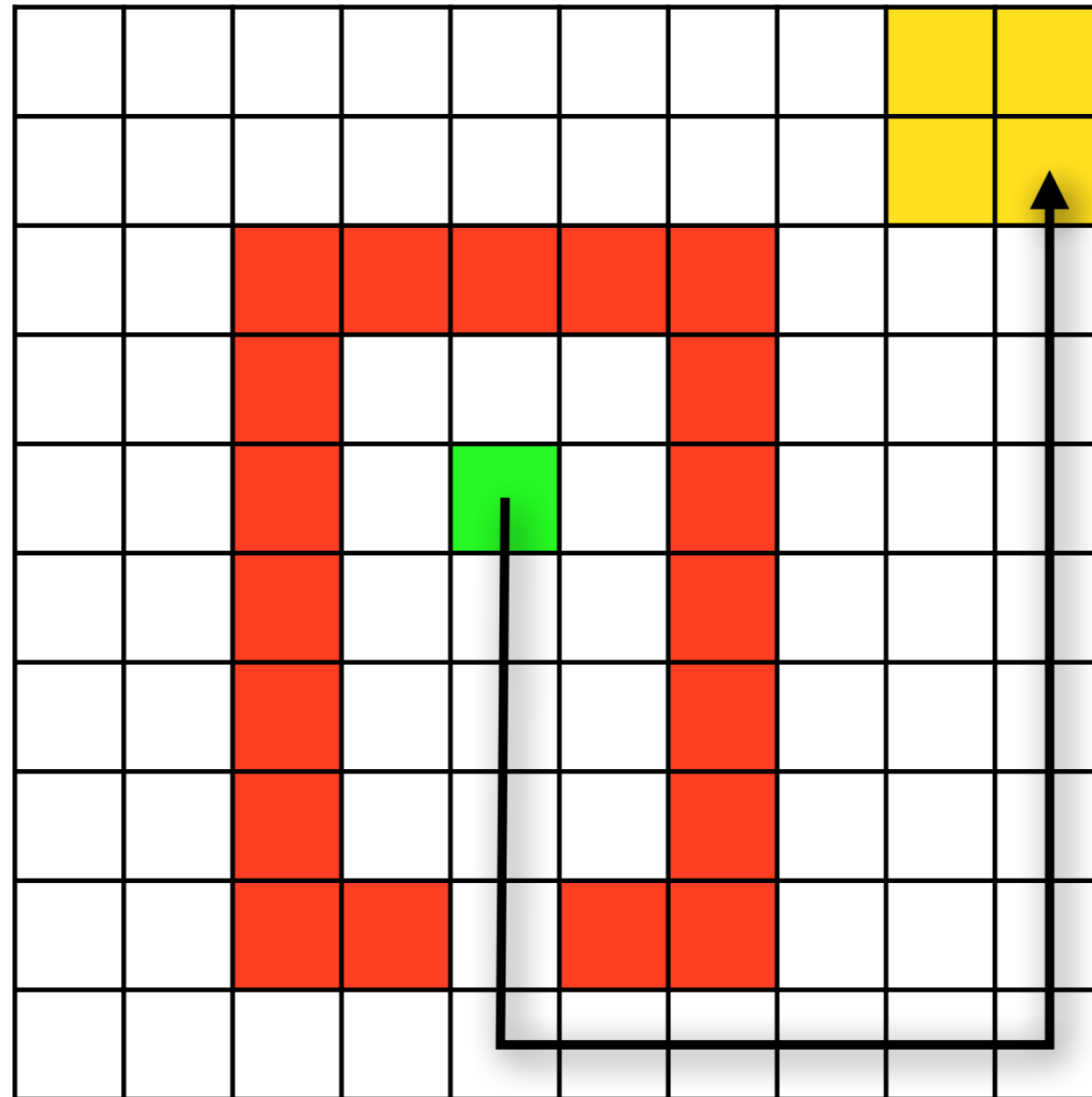


$$\{a_t\} \leftarrow \operatorname{argmax}_{\{a_t\}} \sum_{t=1}^{\infty} r_t$$

Minimizing Loss = Maximizing Reward



Electric shocks  
-1



Gold  
+1

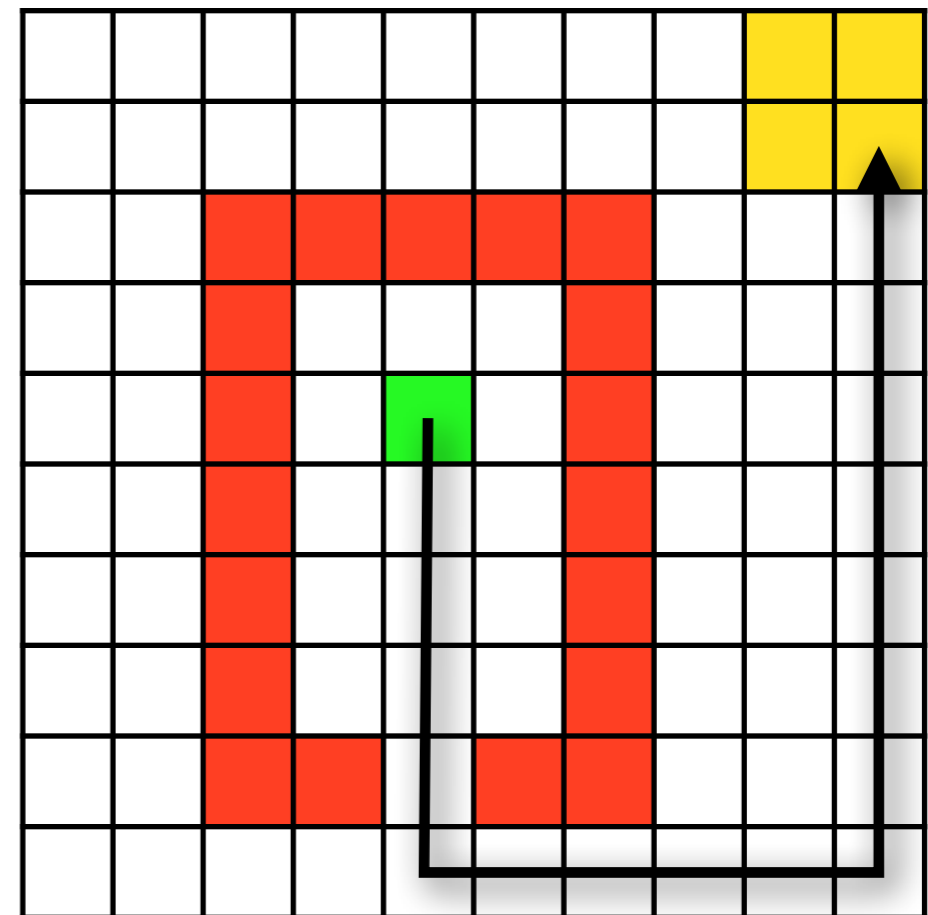
$$s_t \in \mathcal{S}$$

$$a_t \in \mathcal{A}$$

$$\mathcal{T}_{ss'}^a = p(s_{t+1} | s_t, a_t)$$

$$r_t \sim \mathcal{R}(s_{t+1}, a_t, s_t)$$

$$\pi(a|s) = p(a|s)$$



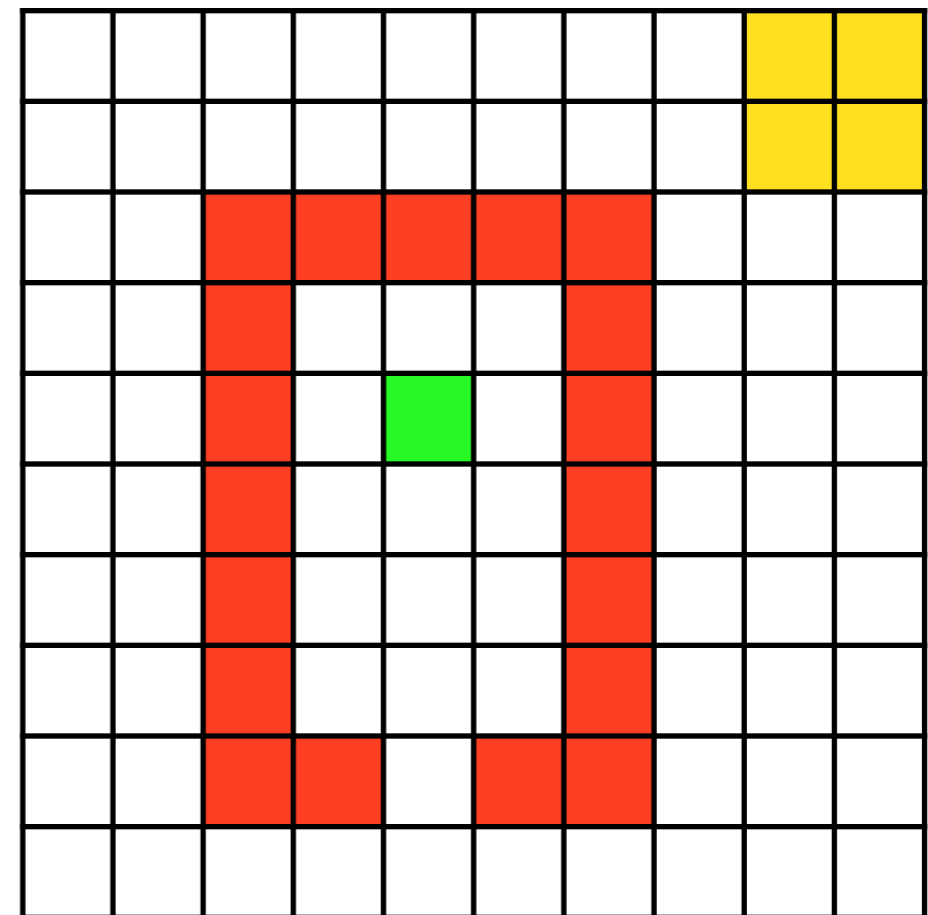
$$s_t \in \mathcal{S}$$

$$a_t \in \mathcal{A}$$

$$\mathcal{T}_{ss'}^a = p(s_{t+1} | s_t, a_t)$$

$$r_t \sim \mathcal{R}(s_{t+1}, a_t, s_t)$$

$$\pi(a|s) = p(a|s)$$



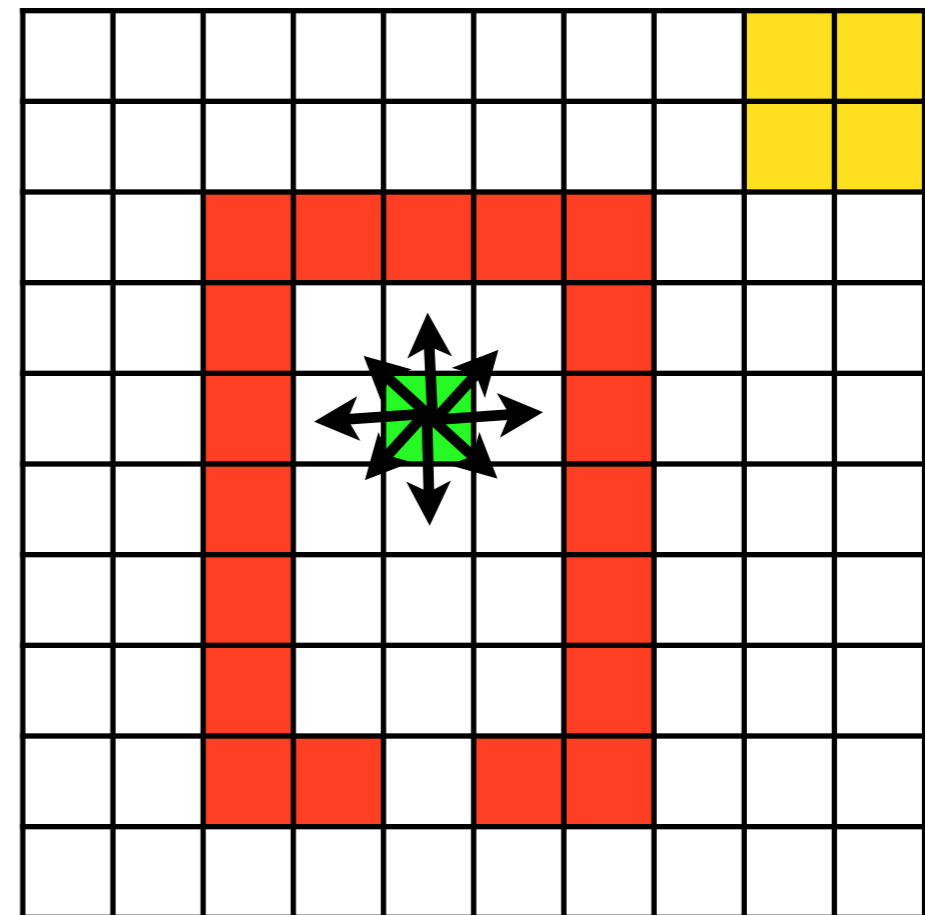
$$s_t \in \mathcal{S}$$

$$a_t \in \mathcal{A}$$

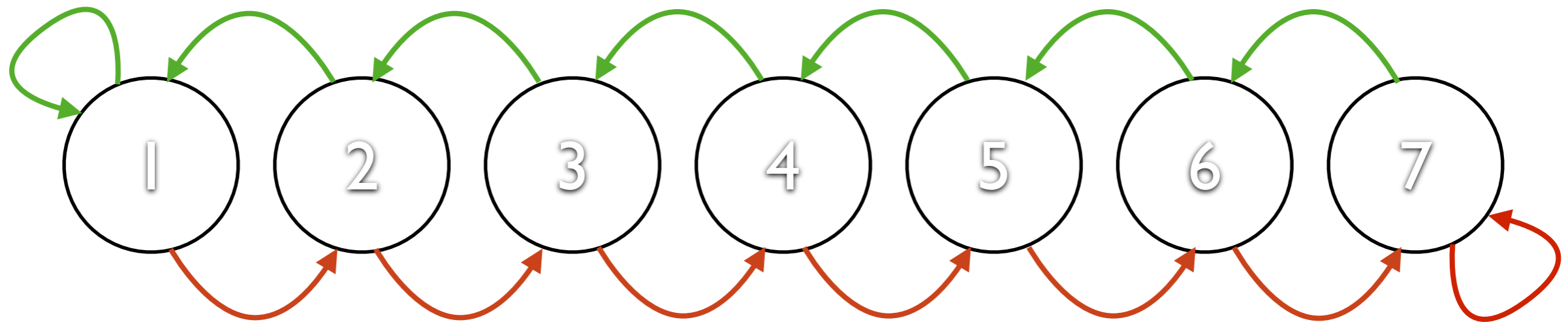
$$\mathcal{T}_{ss'}^a = p(s_{t+1} | s_t, a_t)$$

$$r_t \sim \mathcal{R}(s_{t+1}, a_t, s_t)$$

$$\pi(a|s) = p(a|s)$$



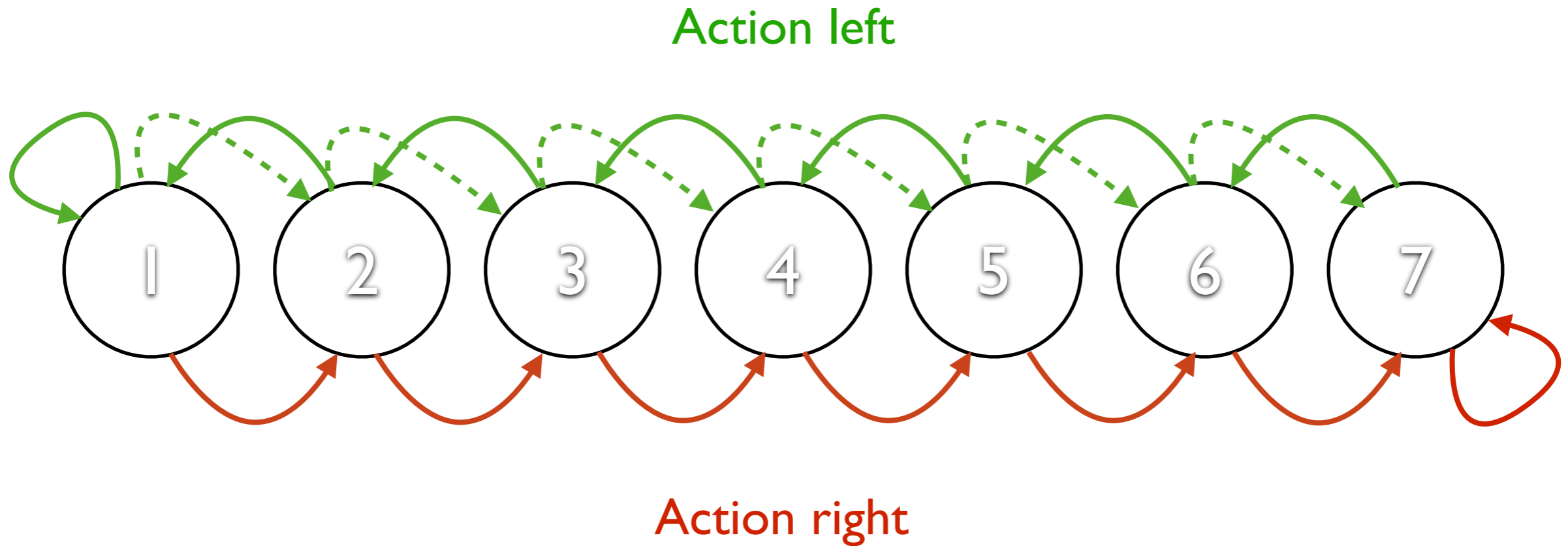
Action left



Action right

$$T^{\text{left}} = \begin{bmatrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

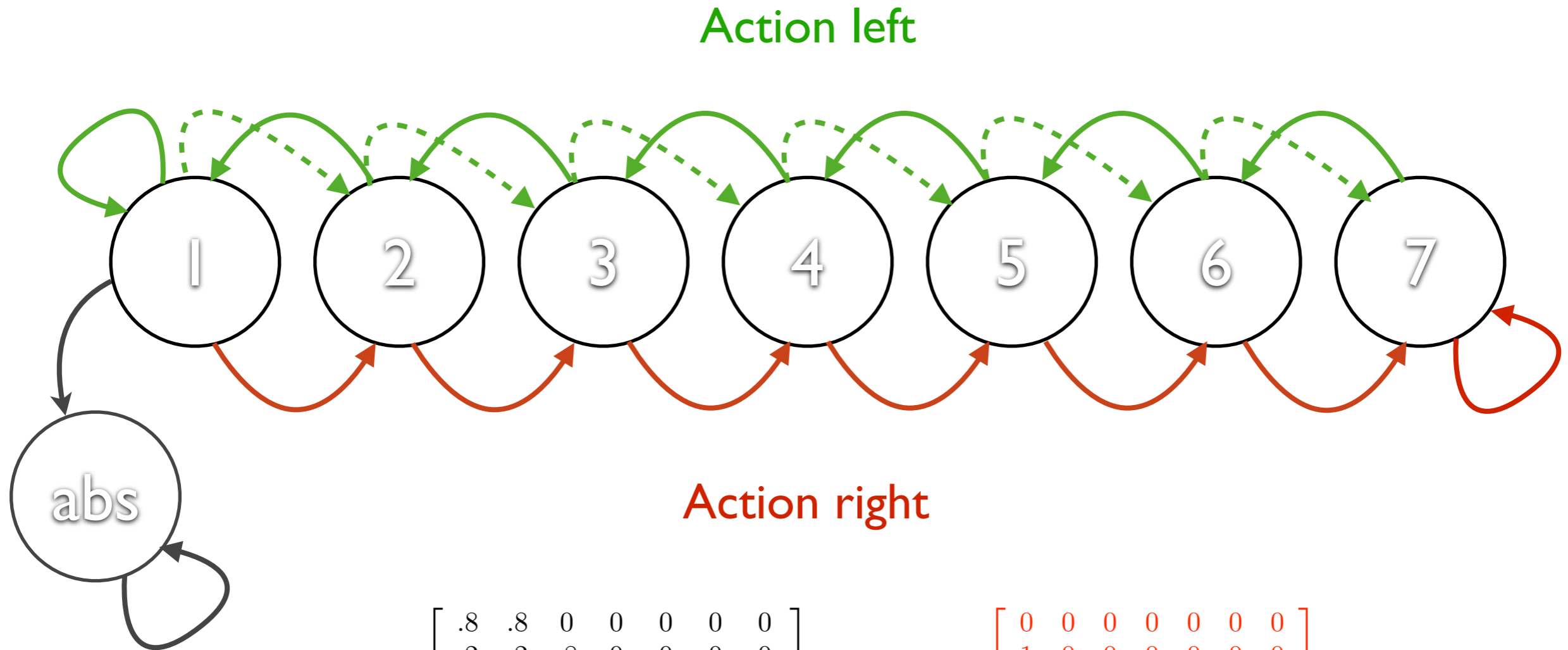
$$T^{\text{right}} = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 1 \end{bmatrix}$$



$$T^{\text{left}} = \begin{bmatrix} .8 & .8 & 0 & 0 & 0 & 0 & 0 \\ .2 & .2 & .8 & 0 & 0 & 0 & 0 \\ 0 & 0 & .2 & .8 & 0 & 0 & 0 \\ 0 & 0 & 0 & .2 & .8 & 0 & 0 \\ 0 & 0 & 0 & 0 & .2 & .8 & 0 \\ 0 & 0 & 0 & 0 & 0 & .2 & .8 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$$T^{\text{right}} = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 1 \end{bmatrix}$$

Noisy: plants, environments, agent



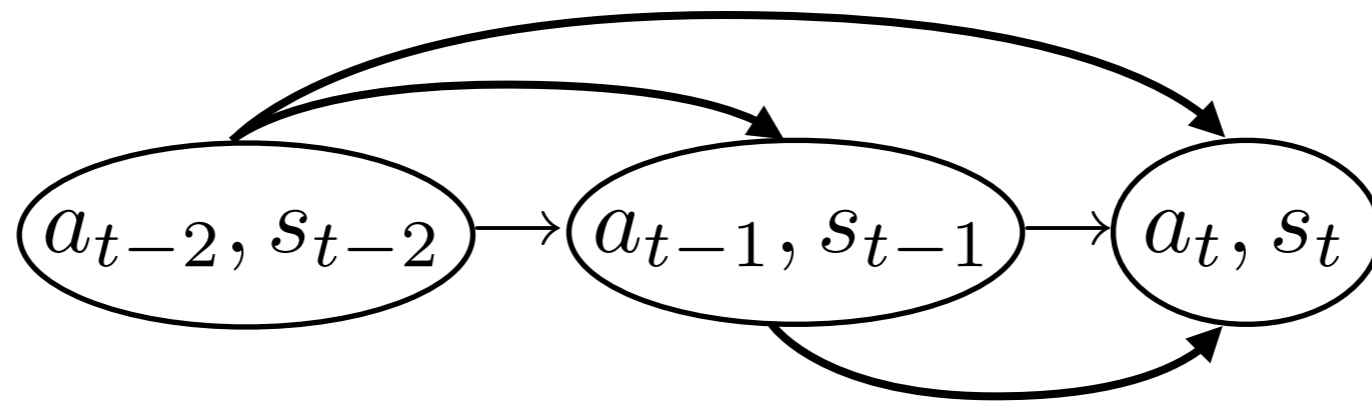
$$T^{\text{left}} = \begin{bmatrix} .8 & .8 & 0 & 0 & 0 & 0 & 0 \\ .2 & .2 & .8 & 0 & 0 & 0 & 0 \\ 0 & 0 & .2 & .8 & 0 & 0 & 0 \\ 0 & 0 & 0 & .2 & .8 & 0 & 0 \\ 0 & 0 & 0 & 0 & .2 & .8 & 0 \\ 0 & 0 & 0 & 0 & 0 & .2 & .8 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$$T^{\text{right}} = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 1 \end{bmatrix}$$

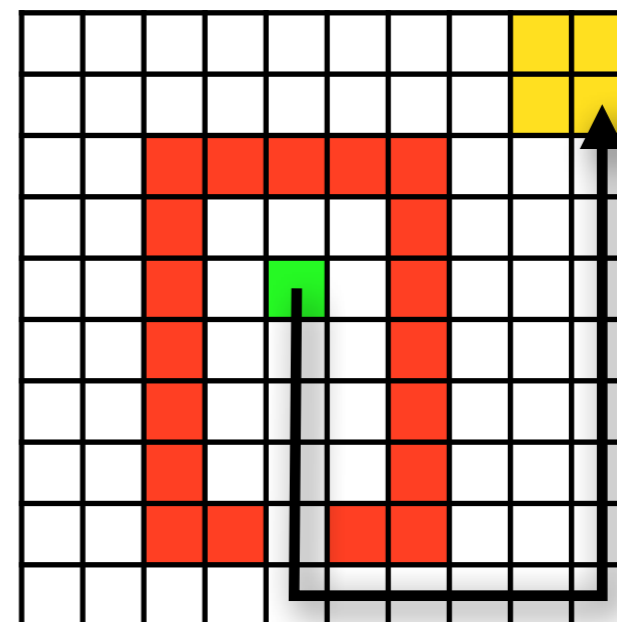
Noisy: plants, environments, agent

Absorbing state  $\rightarrow$  max eigenvalue  $< 1$

$$p(s_{t+1} | a_t, s_t, a_{t-1}, s_{t-1}, a_{t-2}, s_{t-2}, \dots) = p(s_{t+1} | a_t, s_t)$$



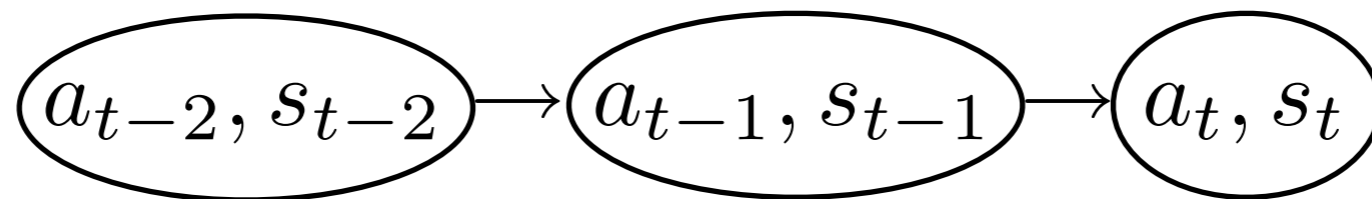
Velocity





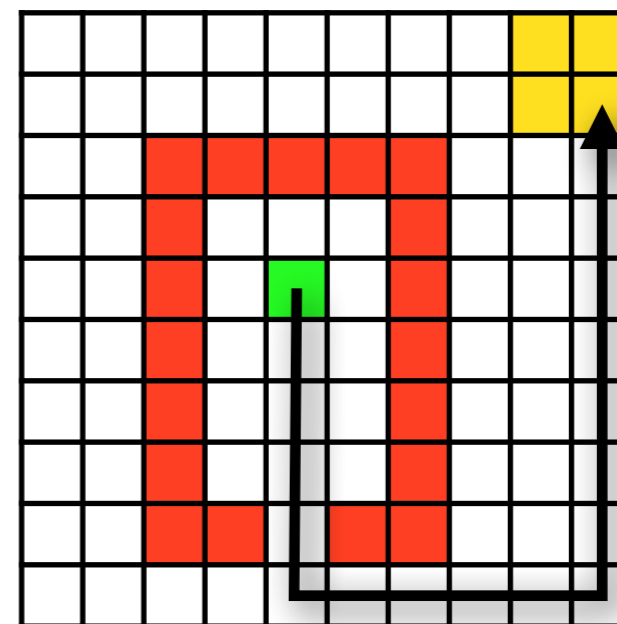


$$p(s_{t+1} | a_t, s_t, a_{t-1}, s_{t-1}, a_{t-2}, s_{t-2}, \dots) = p(s_{t+1} | a_t, s_t)$$



Velocity

$$s' = [\text{position}] \rightarrow s' = \begin{bmatrix} \text{position} \\ \text{velocity} \end{bmatrix}$$



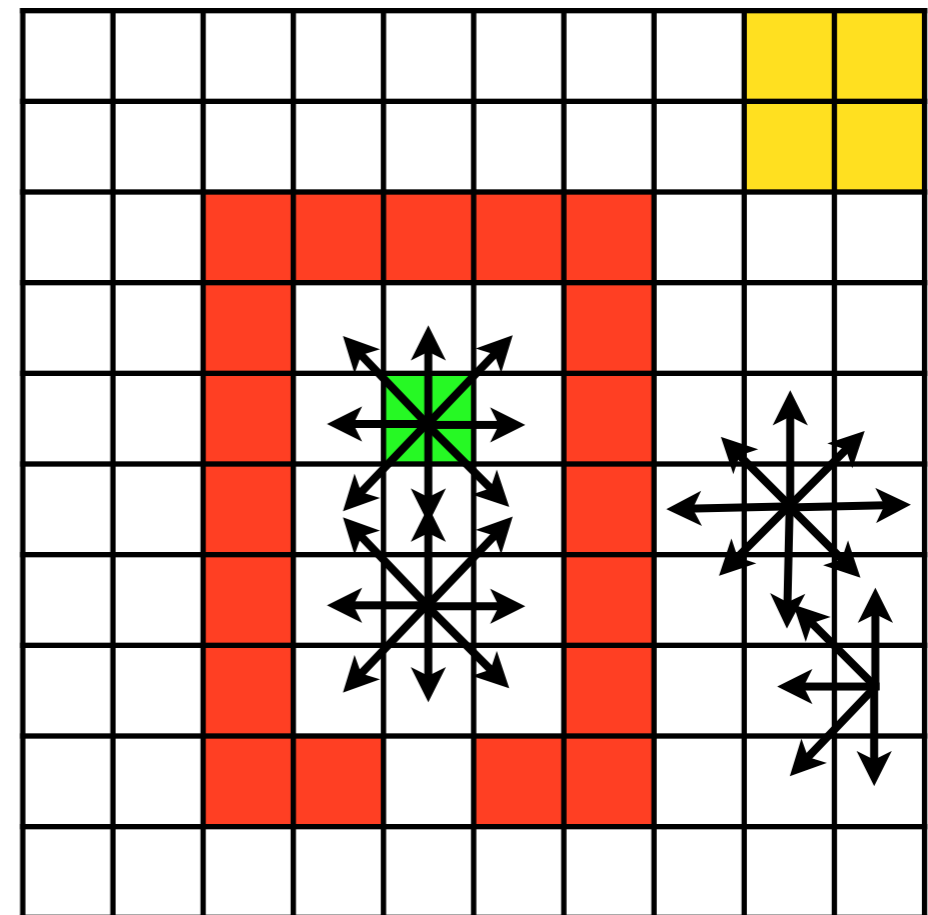
$$s_t \in \mathcal{S}$$

$$a_t \in \mathcal{A}$$

$$\mathcal{T}_{ss'}^a = p(s_{t+1} | s_t, a_t)$$

$$r_t \sim \mathcal{R}(s_{t+1}, a_t, s_t)$$

$$\pi(a|s) = p(a|s)$$



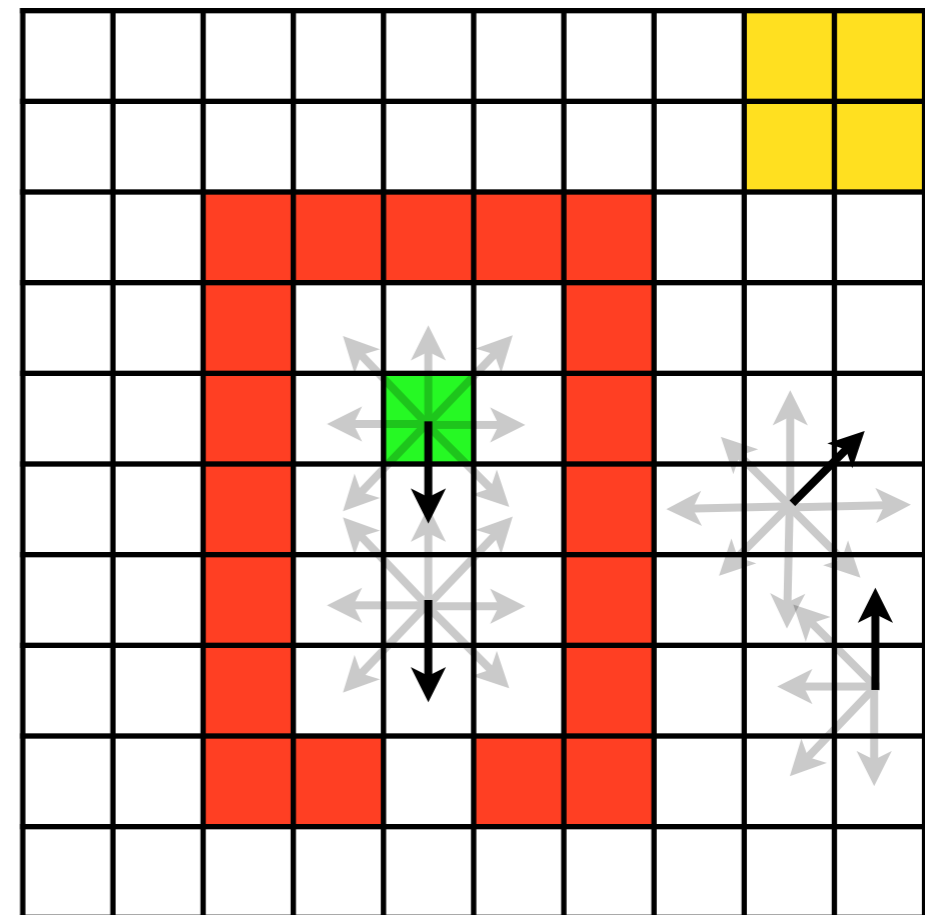
$$s_t \in \mathcal{S}$$

$$a_t \in \mathcal{A}$$

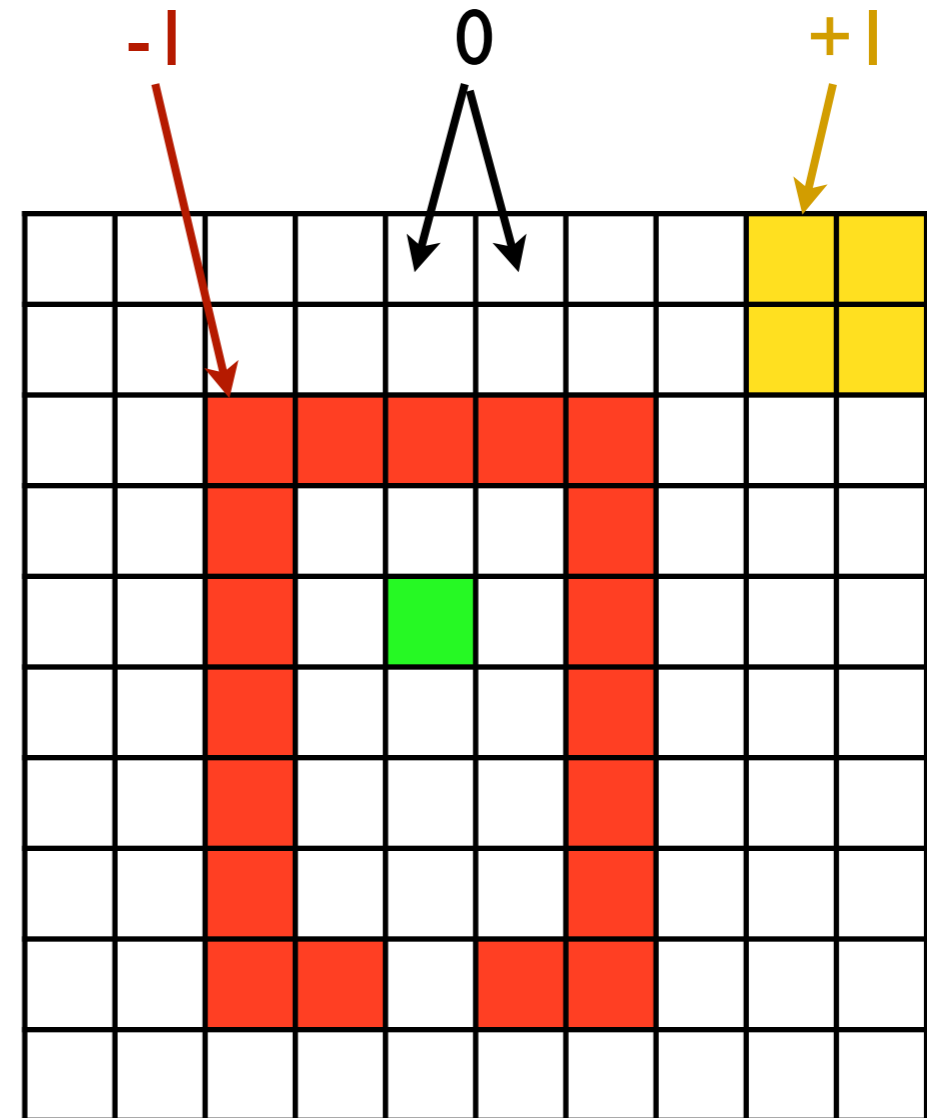
$$\mathcal{T}_{ss'}^a = p(s_{t+1} | s_t, a_t)$$

$$r_t \sim \mathcal{R}(s_{t+1}, a_t, s_t)$$

$$\pi(a|s) = p(a|s)$$

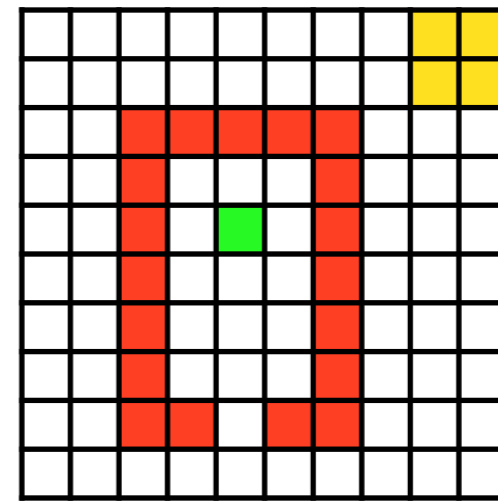


$$\begin{aligned} s_t &\in \mathcal{S} \\ a_t &\in \mathcal{A} \\ \mathcal{T}_{ss'}^a &= p(s_{t+1} | s_t, a_t) \\ r_t &\sim \mathcal{R}(s_{t+1}, a_t, s_t) \\ \pi(a|s) &= p(a|s) \end{aligned}$$



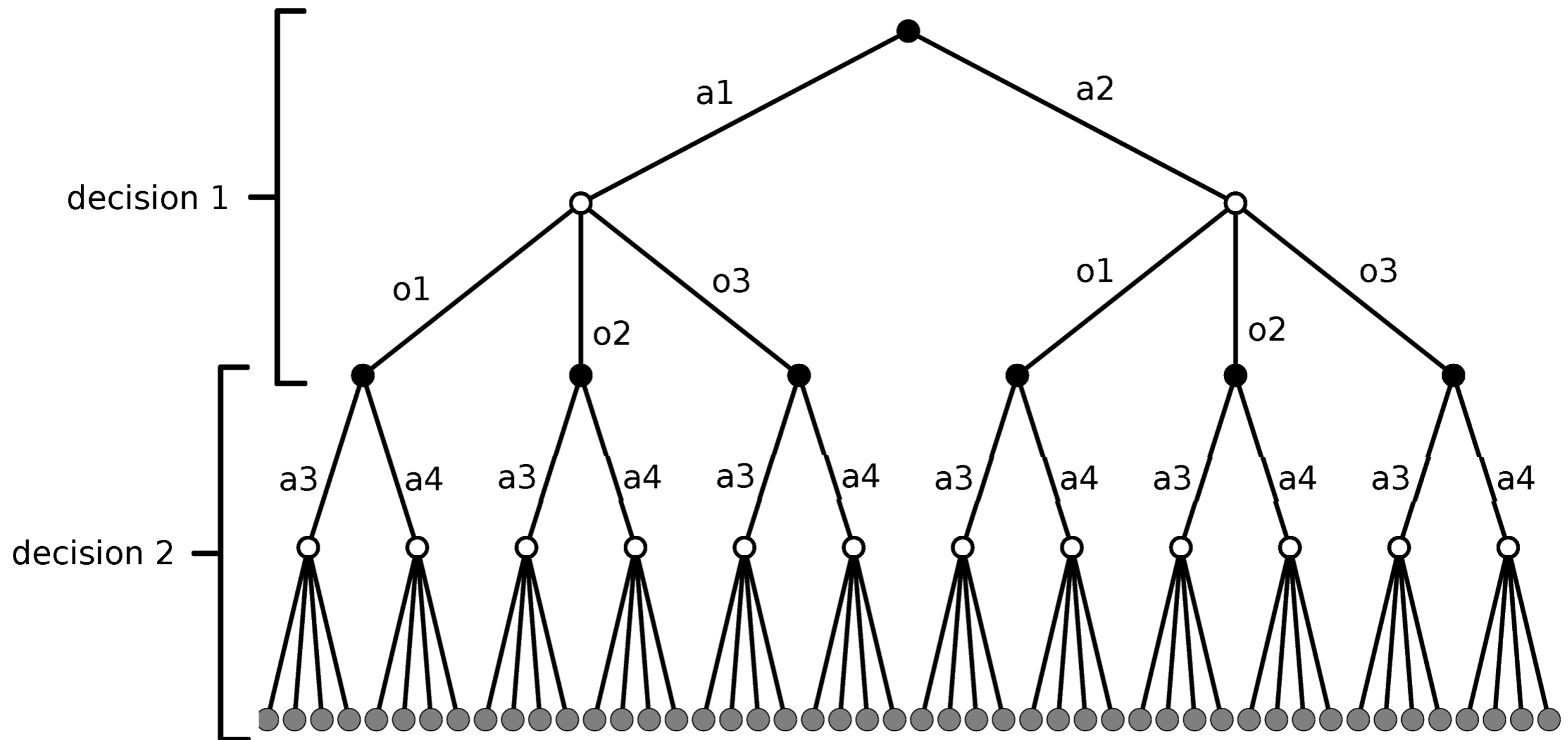
- ▶ Aim: maximise total future reward

$$\sum_{t=1}^{\infty} r_t$$

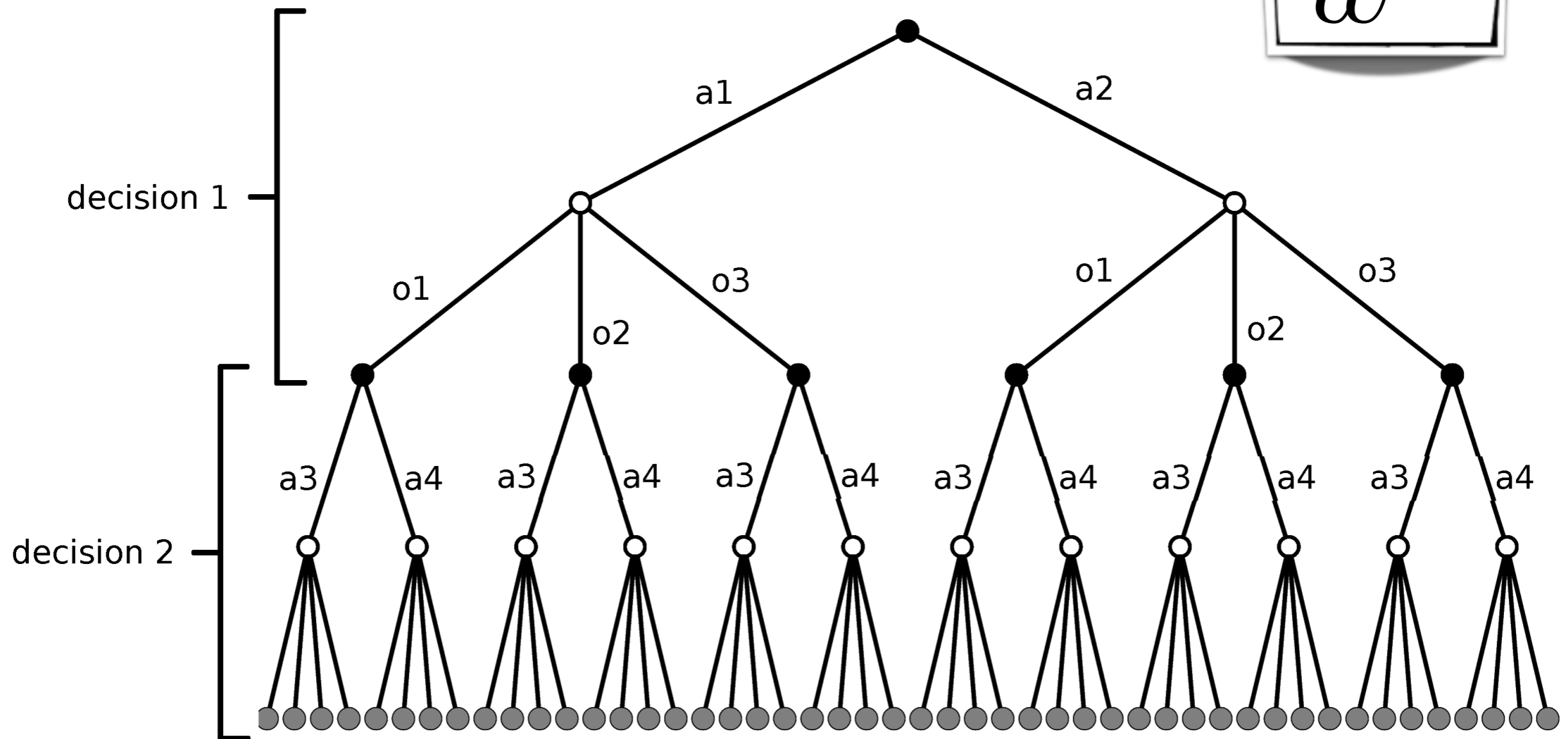
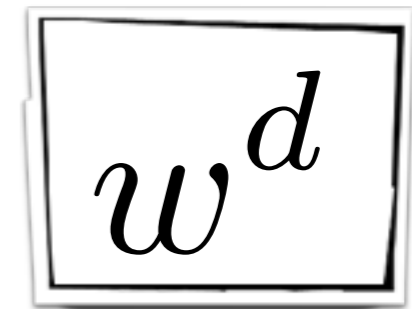


- ▶ i.e. we have to sum over paths through the future and weigh each by its probability
- ▶ Best policy achieves best long-term reward

# Exhaustive tree search

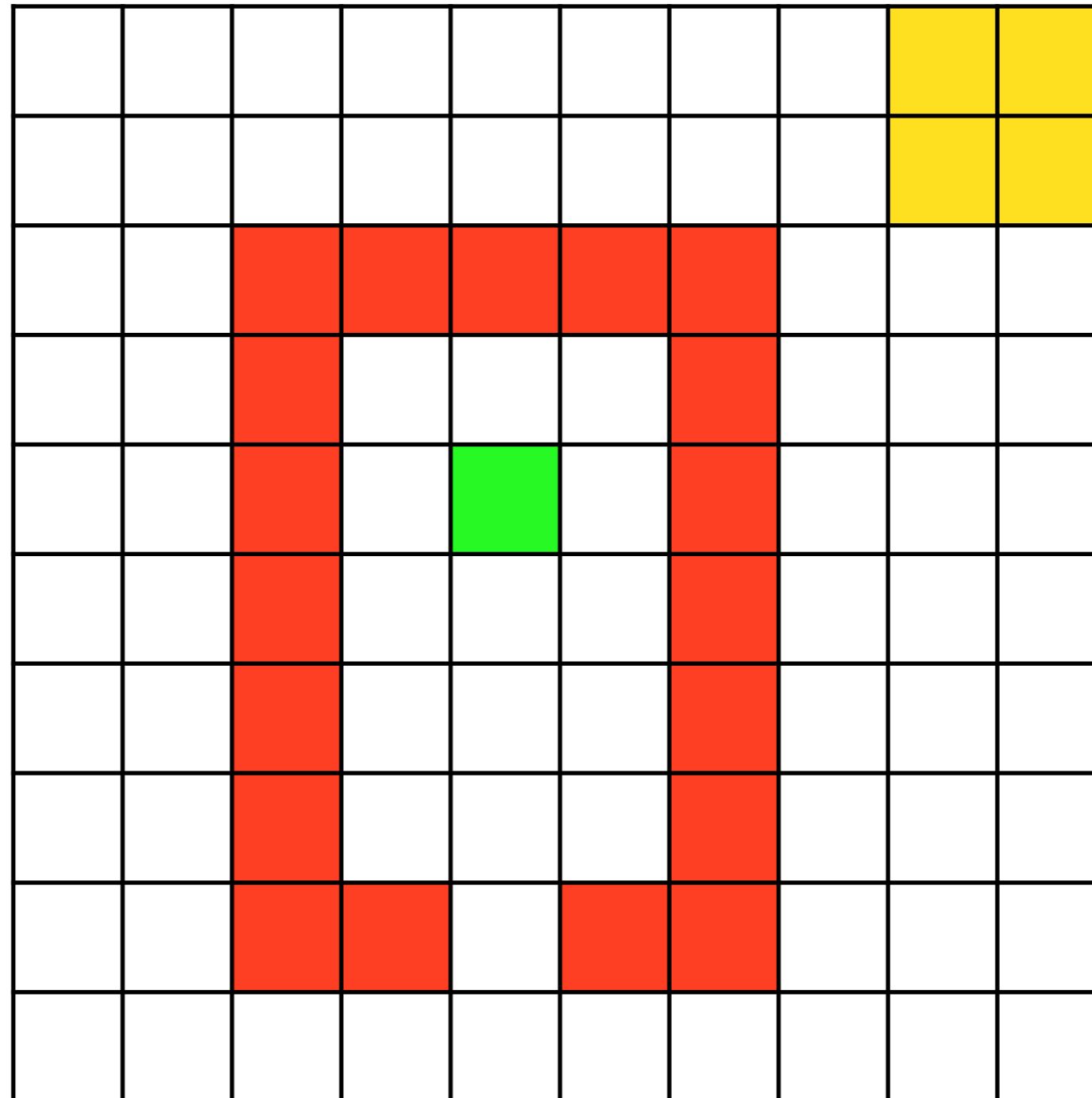


# Exhaustive tree search



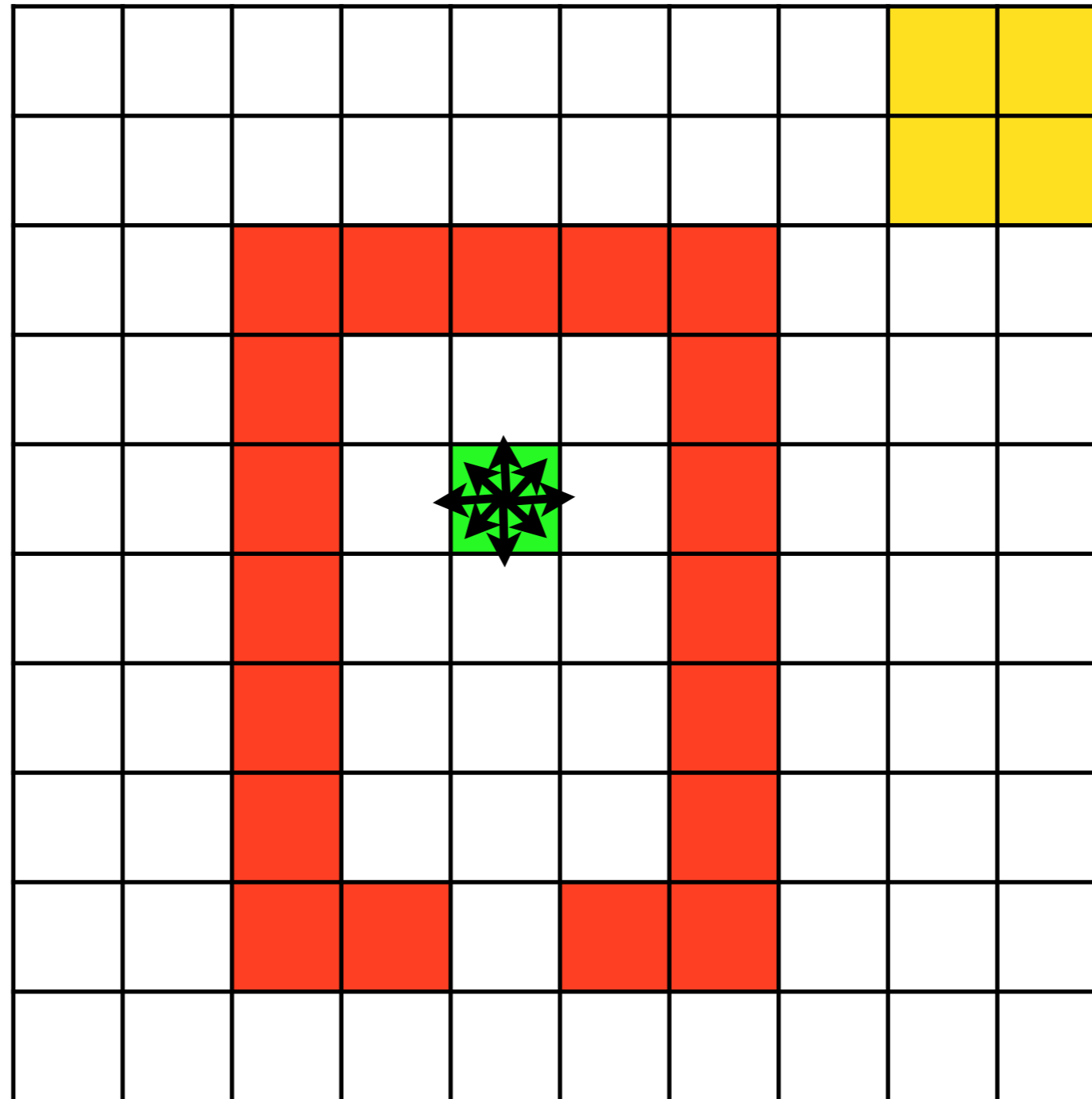


$$\sum_{t=1}^{\infty} r_t$$



$$\sum_{t=1}^{\infty} r_t$$

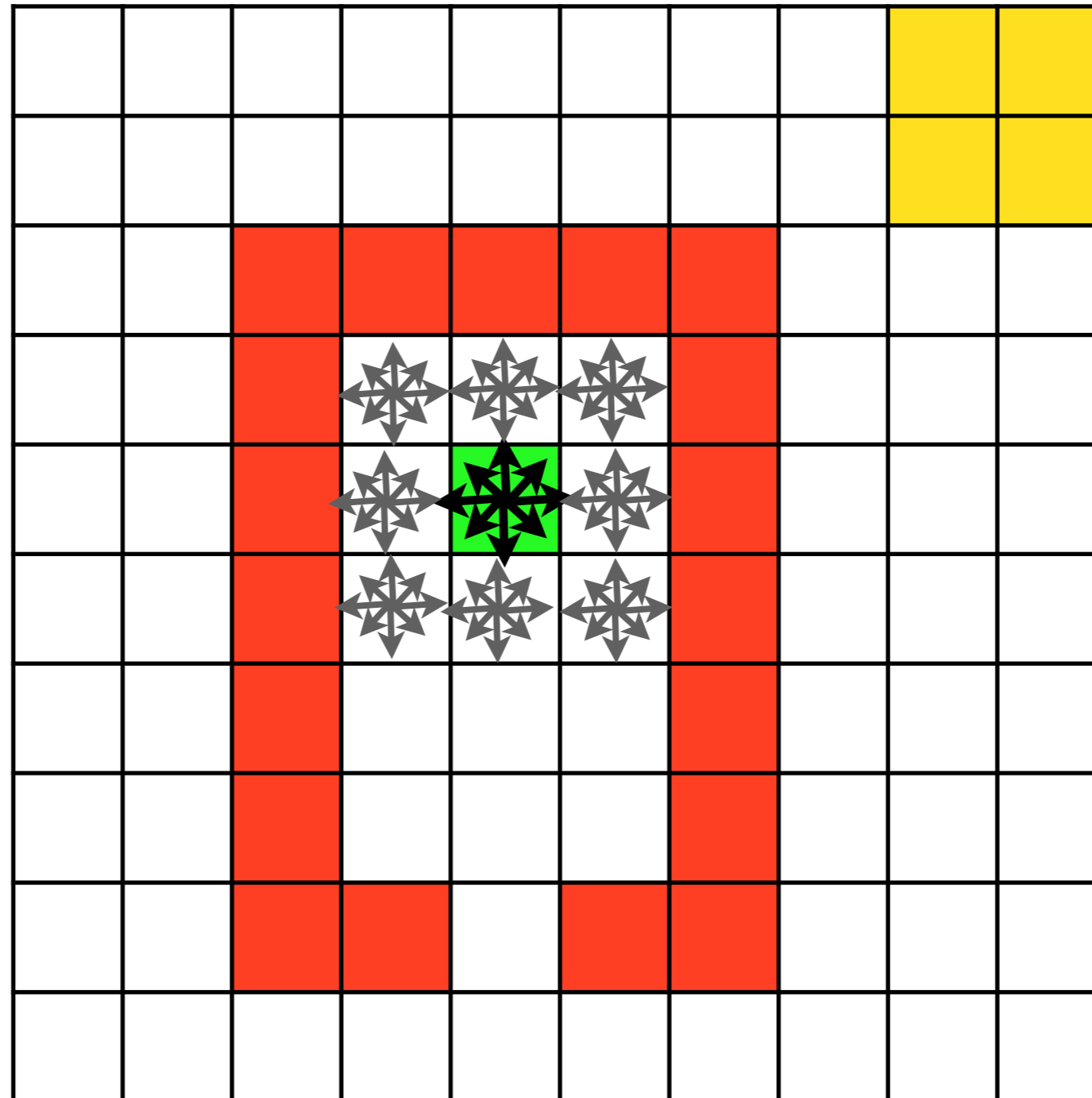
8



$$\sum_{t=1}^{\infty} r_t$$

8

64



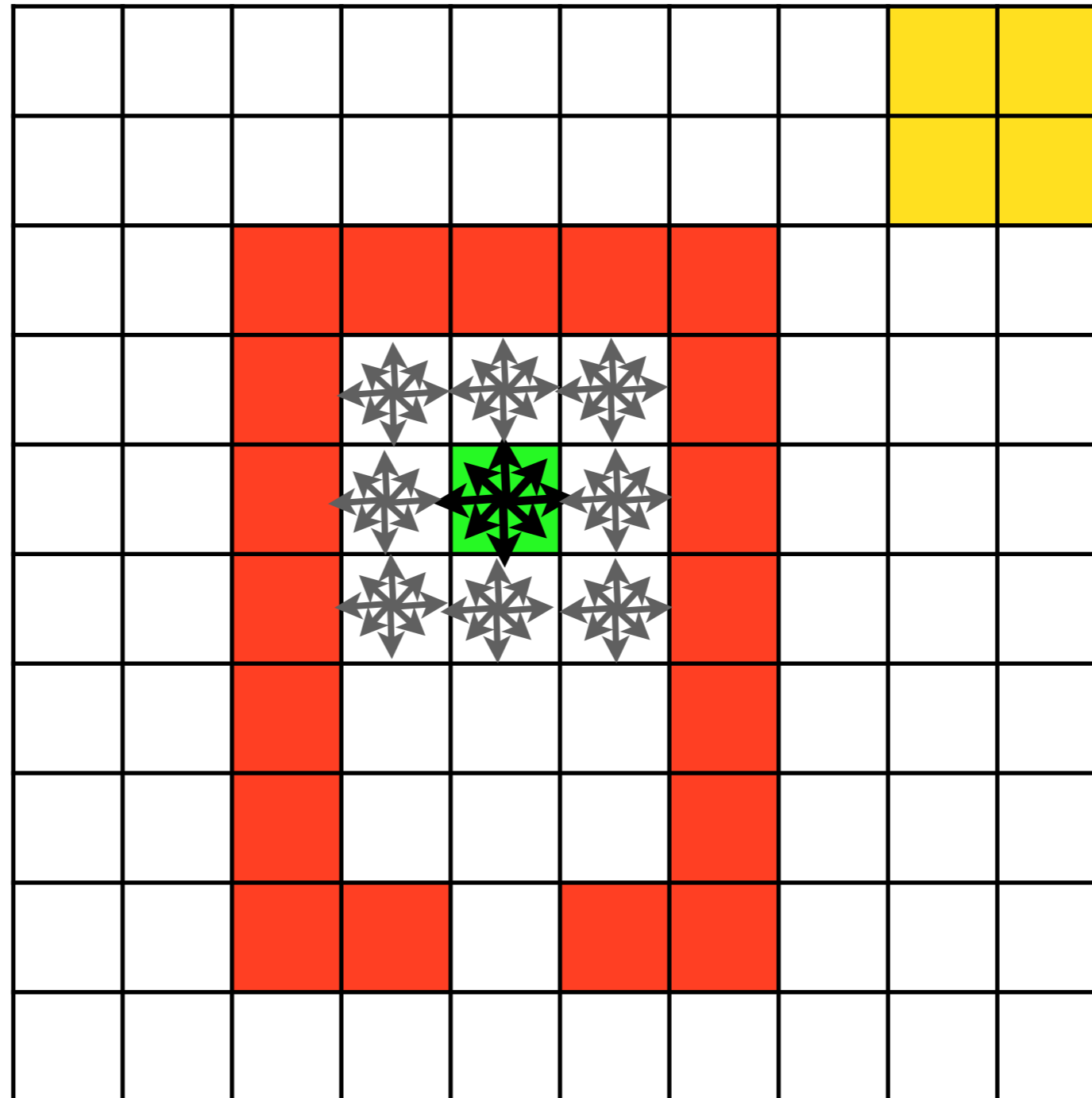
$$\sum_{t=1}^{\infty} r_t$$

8

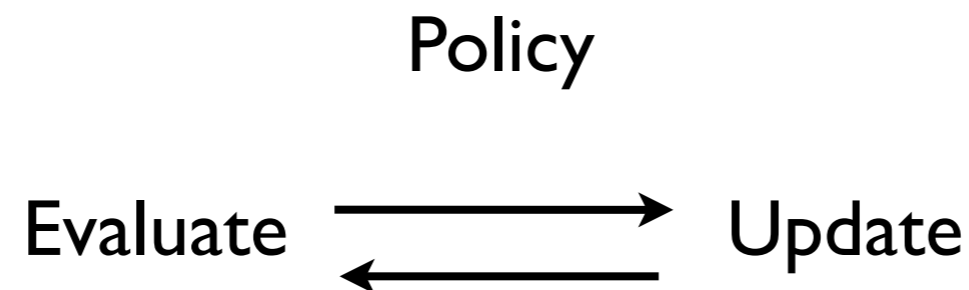
64

512

...



- ▶ Pose the problem mathematically
- ▶ Policy evaluation
- ▶ Policy iteration
- ▶ Monte Carlo techniques: experience samples
- ▶ TD learning

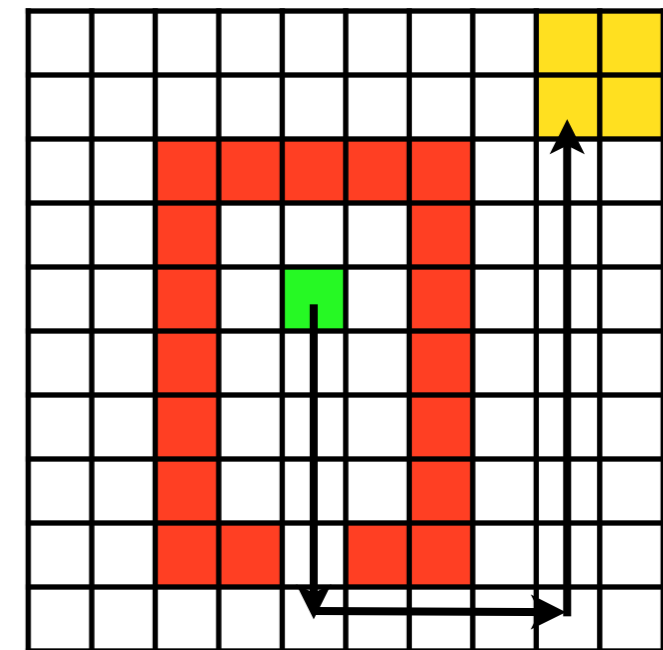


- ▶ Aim: maximise total future reward

$$\sum_{t=1}^{\infty} r_t$$

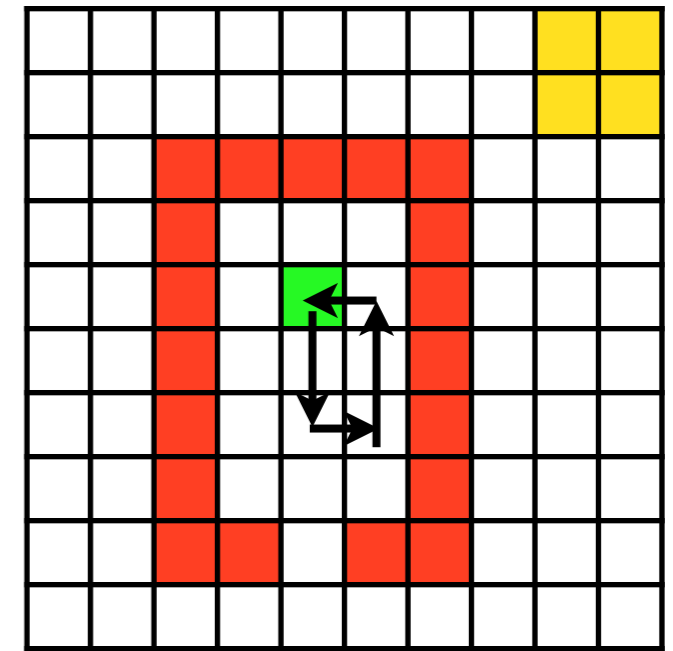
- ▶ To know which is best, evaluate it first
- ▶ The policy determines the expected reward from each state

$$V^{\pi}(s_1) = \mathbb{E} \left[ \sum_{t=1}^{\infty} r_t \mid s_1 = 1, a_t \sim \pi \right]$$



- ▶ Given a policy, each state has an expected value

$$V^\pi(s_1) = \mathbb{E} \left[ \sum_{t=1}^{\infty} r_t \mid s_1 = 1, a_t \sim \pi \right]$$



- ▶ But:  $\sum_{t=0}^{\infty} r_t = \infty$

- ▶ Episodic  $\sum_{t=0}^T r_t < \infty$

- ▶ Discounted  $\sum_{t=0}^{\infty} \gamma^t r_t < \infty$

- infinite horizons

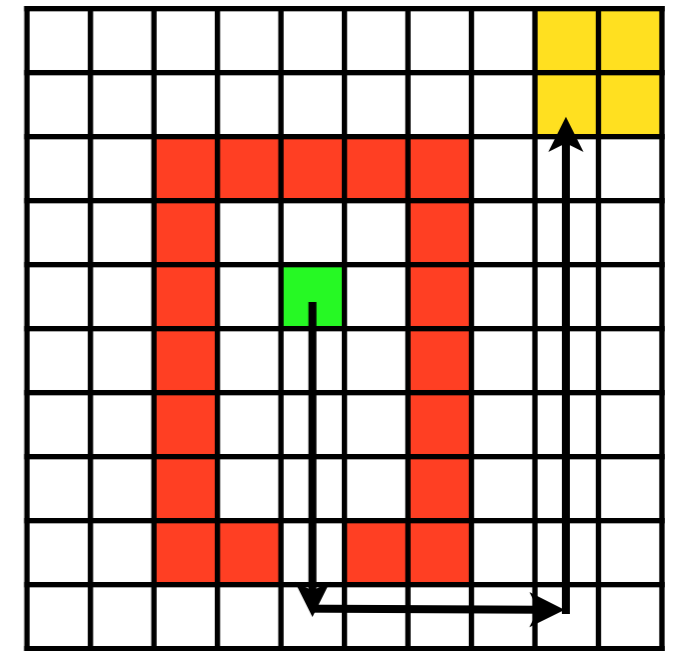
$$\sum_{t=0}^T \gamma^t r_t$$

- finite, exponentially distributed horizons

$$T \sim \frac{1}{\tau} e^{t/\tau}$$

- ▶ Given a policy, each state has an expected value

$$v^\pi(s_1) = \mathbb{E} \left[ \sum_{t=1}^{\infty} r_t \mid s_1 = 1, a_t \sim \pi \right]$$



- ▶ But:  $\sum_{t=0}^{\infty} r_t = \infty$

- ▶ Episodic  $\sum_{t=0}^T r_t < \infty$

- ▶ Discounted  $\sum_{t=0}^{\infty} \gamma^t r_t < \infty$

- infinite horizons

$$\sum_{t=0}^T \gamma^t r_t$$

- finite, exponentially distributed horizons

$$T \sim \frac{1}{\tau} e^{t/\tau}$$







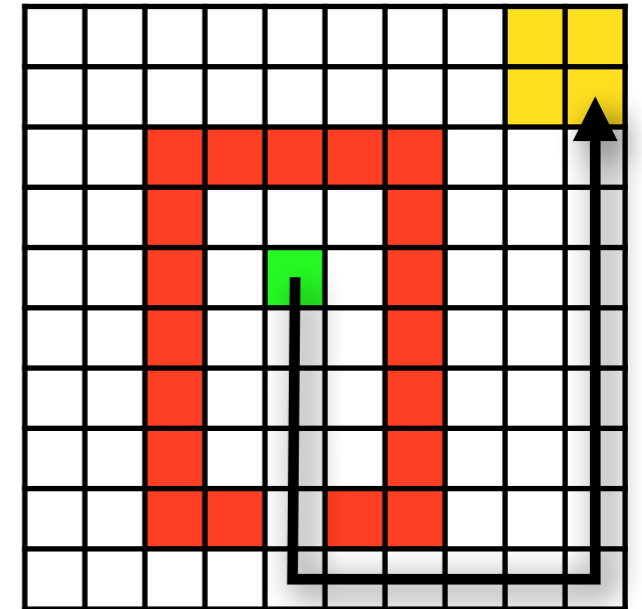
$$V^\pi(s_t) = \mathbb{E}[r_1 | s_t = s, \pi] + \mathbb{E}[V(s_{t+1}), \pi]$$

$$r_1 \sim \mathcal{R}(s_2, a_1, s_1)$$

$$\mathbb{E}[r_1 | s_t = s, \pi] = \mathbb{E} \left[ \sum_{s_{t+1}} p(s_{t+1} | s_t, a_t) \mathcal{R}(s_{t+1}, a_t, s_t) \right]$$

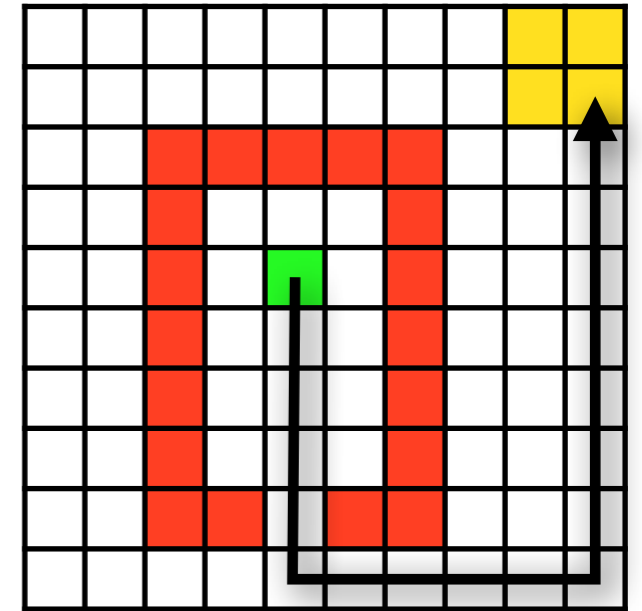
$$= \sum_{a_t} p(a_t | s_t) \left[ \sum_{s_{t+1}} p(s_{t+1} | s_t, a_t) \mathcal{R}(s_{t+1}, a_t, s_t) \right]$$

$$= \sum_{a_t} \pi(a_t, s_t) \left[ \sum_{s_{t+1}} \mathcal{T}_{s_t s_{t+1}}^{a_t} \mathcal{R}(s_{t+1}, a_t, s_t) \right]$$



$$V^\pi(s_t) = \mathbb{E}[r_1 | s_t = s, \pi] + \mathbb{E}[V(s_{t+1}), \pi]$$

$$\mathbb{E}[r_1 | s_t, \pi] = \sum_a \pi(a, s_t) \left[ \sum_{s_{t+1}} \mathcal{T}_{s_t s_{t+1}}^a \mathcal{R}(s_{t+1}, a, s_t) \right]$$



$$\mathbb{E}[V^\pi(s_{t+1}), \pi, s_t] = \sum_a \pi(a, s_t) \left[ \sum_{s_{t+1}} \mathcal{T}_{s_t s_{t+1}}^a V^\pi(s_{t+1}) \right]$$

$$V^\pi(s) = \sum_a \pi(a|s) \left[ \sum_{s'} \mathcal{T}_{ss'}^a [\mathcal{R}(s', a, s) + V^\pi(s')] \right]$$

$$V^\pi(s) = \sum_a \pi(a|s) \left[ \sum_{s'} \mathcal{T}_{ss'}^a [\mathcal{R}(s', a, s) + V^\pi(s')] \right]$$

All future reward from state  $s$  =  $E$  [ Immediate reward + All future reward from next state  $s'$  ]

$$V^\pi(s) = \sum_a \pi(a|s) \left[ \sum_{s'} \mathcal{T}_{ss'}^a [\mathcal{R}(s', a, s) + V^\pi(s')] \right]$$

$$V^\pi(s) = \sum_a \pi(a|s) \left[ \sum_{s'} \mathcal{T}_{ss'}^a [\mathcal{R}(s', a, s) + V^\pi(s')] \right]$$

All future reward from state  $s$  = E [ Immediate reward + All future reward from next state  $s'$  ]

$$V^\pi(s) = \sum_a \pi(a|s) \underbrace{\left[ \sum_{s'} \mathcal{T}_{ss'}^a [\mathcal{R}(s', a, s) + V^\pi(s')] \right]}_{Q^\pi(s, a)}$$

- ▶ so we can define state-action values as:

$$\begin{aligned} Q(s, a) &= \sum_{s'} \mathcal{T}_{ss'}^a [\mathcal{R}(s', a, s) + V(s')] \\ &= \mathbb{E} \left[ \sum_{t=1}^{\infty} r_t | s, a \right] \end{aligned}$$

- ▶ and state values are average state-action values:

$$V(s) = \sum_a \pi(a|s) Q(s, a)$$



$$V^\pi(s) = \sum_a \pi(a|s) \left[ \sum_{s'} \mathcal{T}_{ss'}^a [\mathcal{R}(s', a, s) + V^\pi(s')] \right]$$

- ▶ to evaluate a policy, we need to solve the above equation, i.e. find the self-consistent state values
- ▶ options for policy evaluation
  - exhaustive tree search - outwards, inwards, depth-first
  - value iteration: iterative updates
  - linear solution in 1 step
  - sampling

Option 1: turn it into update equation

Option 2: linear solution

(w/ absorbing states)

$$\begin{aligned} V(s) &= \sum_a \pi(a, s_t) \left[ \sum_{s'} \mathcal{T}_{ss'}^a [\mathcal{R}(s', a, s) + V(s')] \right] \\ \Rightarrow \mathbf{v} &= \mathbf{R}^\pi + \mathbf{T}^\pi \mathbf{v} \\ \Rightarrow \mathbf{v}^\pi &= (\mathbf{I} - \mathbf{T}^\pi)^{-1} \mathbf{R}^\pi \quad \mathcal{O}(|\mathcal{S}|^3) \end{aligned}$$

Option 1: turn it into update equation

$$V^{k+1}(s) = \sum_a \pi(a, s_t) \left[ \sum_{s'} \mathcal{T}_{ss'}^a [\mathcal{R}(s', a, s) + V^k(s')] \right]$$

Option 2: linear solution

(w/ absorbing states)

$$V(s) = \sum_a \pi(a, s_t) \left[ \sum_{s'} \mathcal{T}_{ss'}^a [\mathcal{R}(s', a, s) + V(s')] \right]$$

$$\Rightarrow \mathbf{v} = \mathbf{R}^\pi + \mathbf{T}^\pi \mathbf{v}$$

$$\Rightarrow \mathbf{v}^\pi = (\mathbf{I} - \mathbf{T}^\pi)^{-1} \mathbf{R}^\pi \quad \mathcal{O}(|\mathcal{S}|^3)$$

Given the value function for a policy, say via linear solution

$$V^\pi(s) = \sum_a \pi(a|s) \underbrace{\left[ \sum_{s'} \mathcal{T}_{ss'}^a [\mathcal{R}(s', a, s) + V^\pi(s')] \right]}_{Q^\pi(s, a)}$$

Given the values  $V$  for the policy, we can improve the policy by always choosing the best action:

$$\pi'(a|s) = \begin{cases} 1 & \text{if } a = \operatorname{argmax}_a Q^\pi(s, a) \\ 0 & \text{else} \end{cases}$$

It is guaranteed to improve:

$$Q^\pi(s, \pi'(s)) = \max_a Q^\pi(s, a) \geq Q^\pi(s, \pi(s)) = \mathcal{V}^\pi(s)$$

← for deterministic policy

Policy evaluation

$$\mathbf{v}^\pi = (\mathbf{I} - \mathbf{T}^\pi)^{-1} \mathbf{R}^\pi$$

$$\pi(a|s) = \begin{cases} 1 & \text{if } a = \operatorname{argmax}_a \sum_{s'} \mathcal{T}_{ss'}^a [\mathcal{R}_{ss}^a + V^{\pi}(s')] \\ 0 & \text{else} \end{cases}$$

Policy evaluation

$$\mathbf{v}^\pi = (\mathbf{I} - \mathbf{T}^\pi)^{-1} \mathbf{R}^\pi$$

greedy policy improvement

$$\pi(a|s) = \begin{cases} 1 & \text{if } a = \operatorname{argmax}_a \sum_{s'} \mathcal{T}_{ss'}^a [\mathcal{R}_{ss}^a + V^{\pi}(s')] \\ 0 & \text{else} \end{cases}$$

Policy evaluation

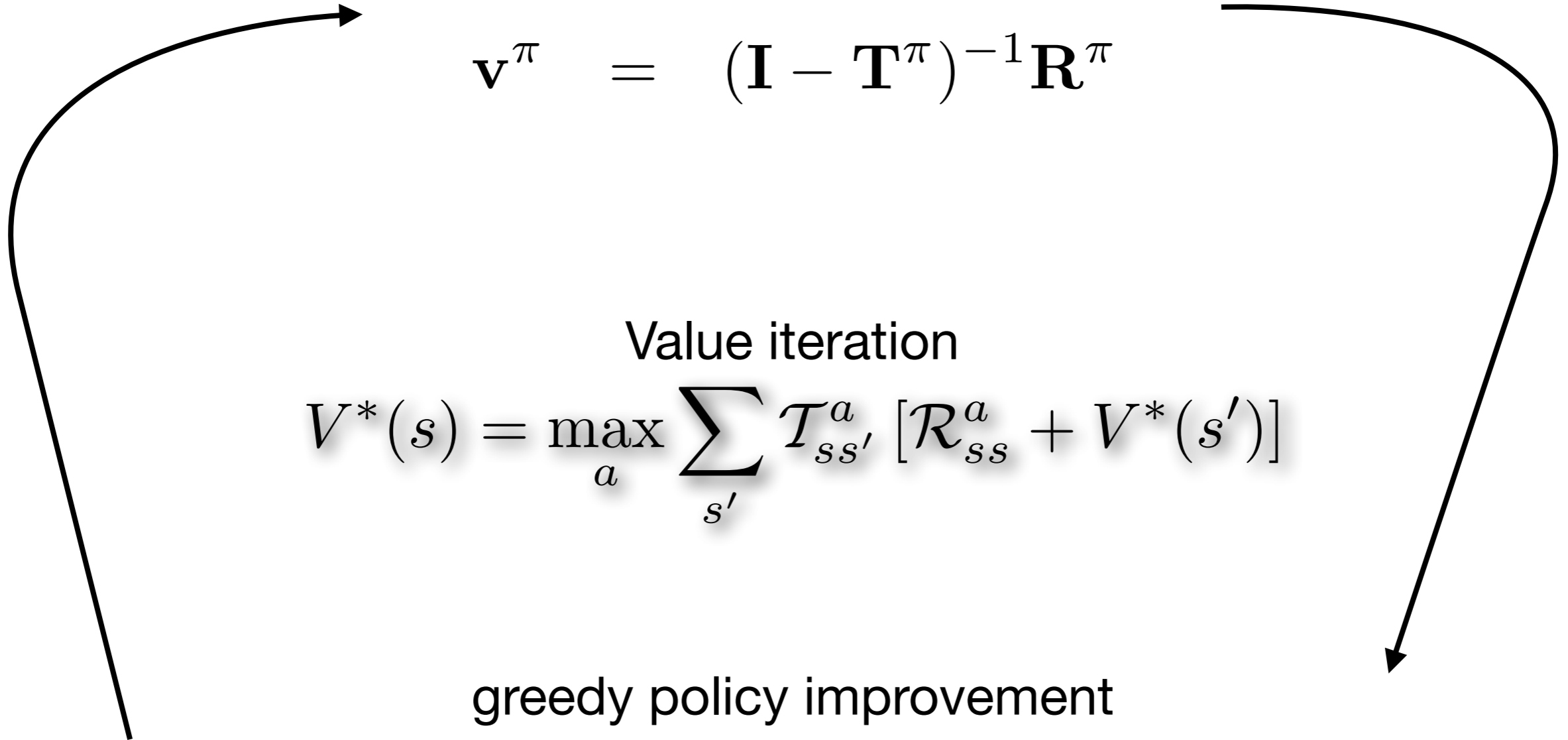
$$\mathbf{v}^\pi = (\mathbf{I} - \mathbf{T}^\pi)^{-1} \mathbf{R}^\pi$$

Value iteration

$$V^*(s) = \max_a \sum_{s'} \mathcal{T}_{ss'}^a [\mathcal{R}_{ss}^a + V^*(s')]$$

greedy policy improvement

$$\pi(a|s) = \begin{cases} 1 & \text{if } a = \operatorname{argmax}_a \sum_{s'} \mathcal{T}_{ss'}^a [\mathcal{R}_{ss}^a + V^{pi}(s')] \\ 0 & \text{else} \end{cases}$$



- ▶ So far we have assumed knowledge of  $R$  and  $T$ 
  - $R$  and  $T$  are the ‘model’ of the world, so we assume full knowledge of the dynamics and rewards in the environment
- ▶ What if we don’t know them?
- ▶ We can still learn from state-action-reward samples
  - we can learn  $R$  and  $T$  from them, and use our estimates to solve as above
  - alternatively, we can directly estimate  $V$  or  $Q$



Option 3: sampling

$$V(s) = \sum_a \pi(a, s_t) \left[ \sum_{s'} \mathcal{T}_{ss'}^a [\mathcal{R}(s', a, s) + V(s')] \right]$$

So we can just draw some samples from the policy and the transitions and average over them:

$$a = \sum_k f(x_k) p(x_k)$$
$$x^{(i)} \sim p(x) \rightarrow \hat{a} = \frac{1}{N} \sum_i f(x^{(i)})$$

## Option 3: sampling

So we can just draw some samples from the policy and the transitions and average over them:

$$a = \sum_k f(x_k)p(x_k)$$
$$x^{(i)} \sim p(x) \rightarrow \hat{a} = \frac{1}{N} \sum_i f(x^{(i)})$$

## Option 3: sampling

this is an expectation over policy and transition samples.

So we can just draw some samples from the policy and the transitions and average over them:

$$a = \sum_k f(x_k) p(x_k)$$
$$x^{(i)} \sim p(x) \rightarrow \hat{a} = \frac{1}{N} \sum_i f(x^{(i)})$$

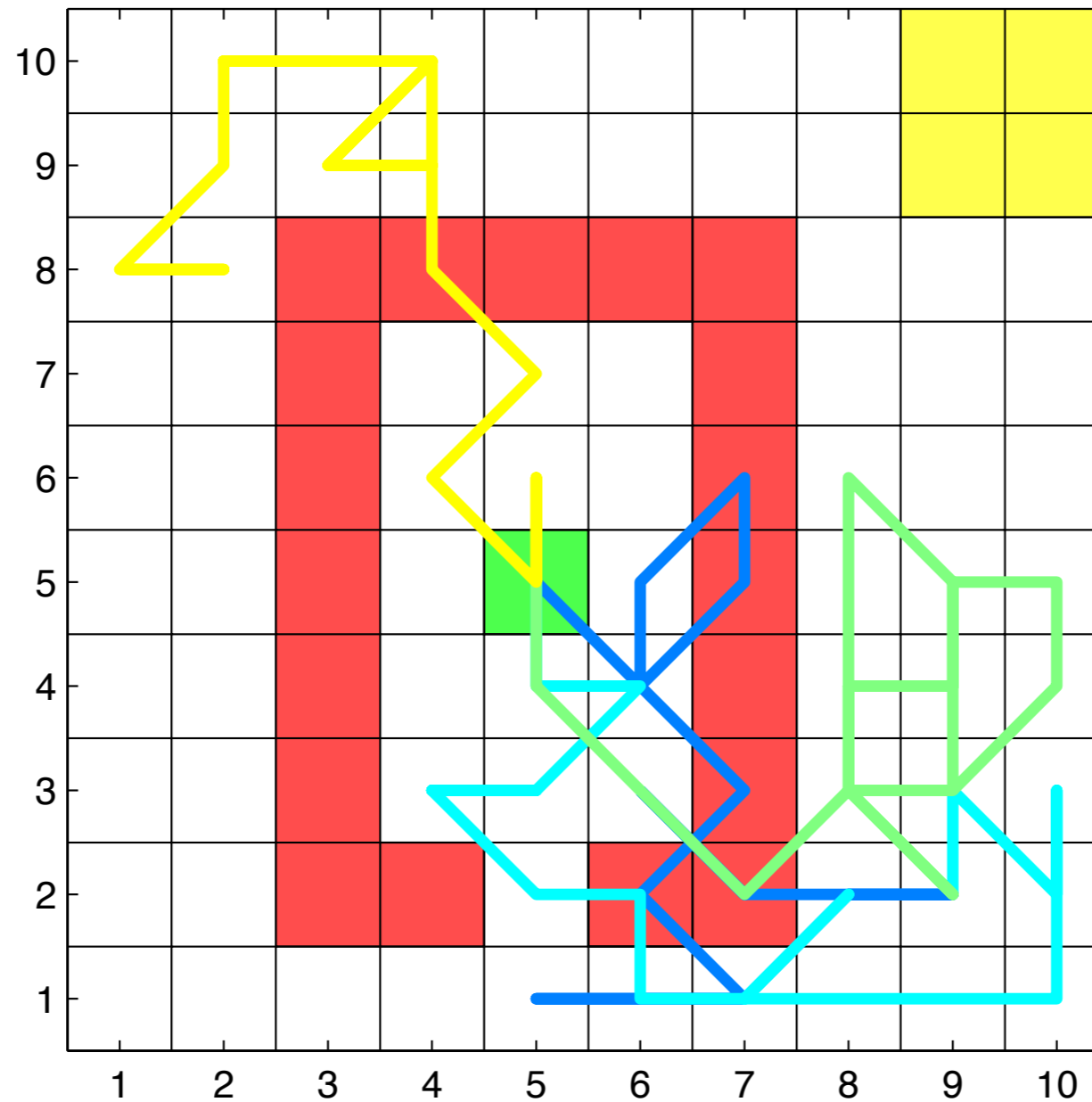
## Option 3: sampling

this is an expectation over policy and transition samples.

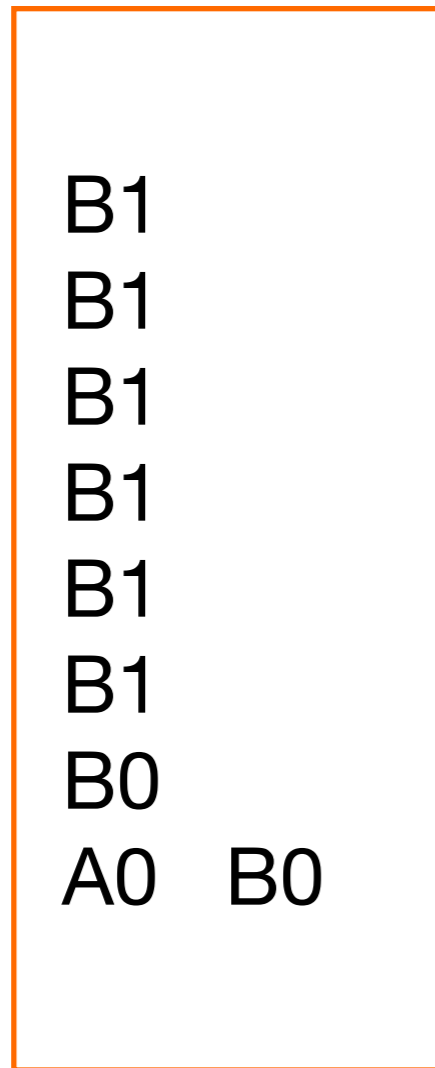
So we can just draw some samples from the policy and the transitions and average over them:

$$a = \sum_k f(x_k) p(x_k)$$
$$x^{(i)} \sim p(x) \rightarrow \hat{a} = \frac{1}{N} \sum_i f(x^{(i)})$$

more about this later...



A new problem: exploration versus exploitation



Markov (every visit)

$$V(B)=3/4$$

$$V(A)=0$$

TD

$$V(B)=3/4$$

$$V(A)=\sim 3/4$$

- ▶ Average over various bootstrappings:  $TD(\lambda)$

## ► First visit MC

- randomly start in all states, generate paths, average for starting state only

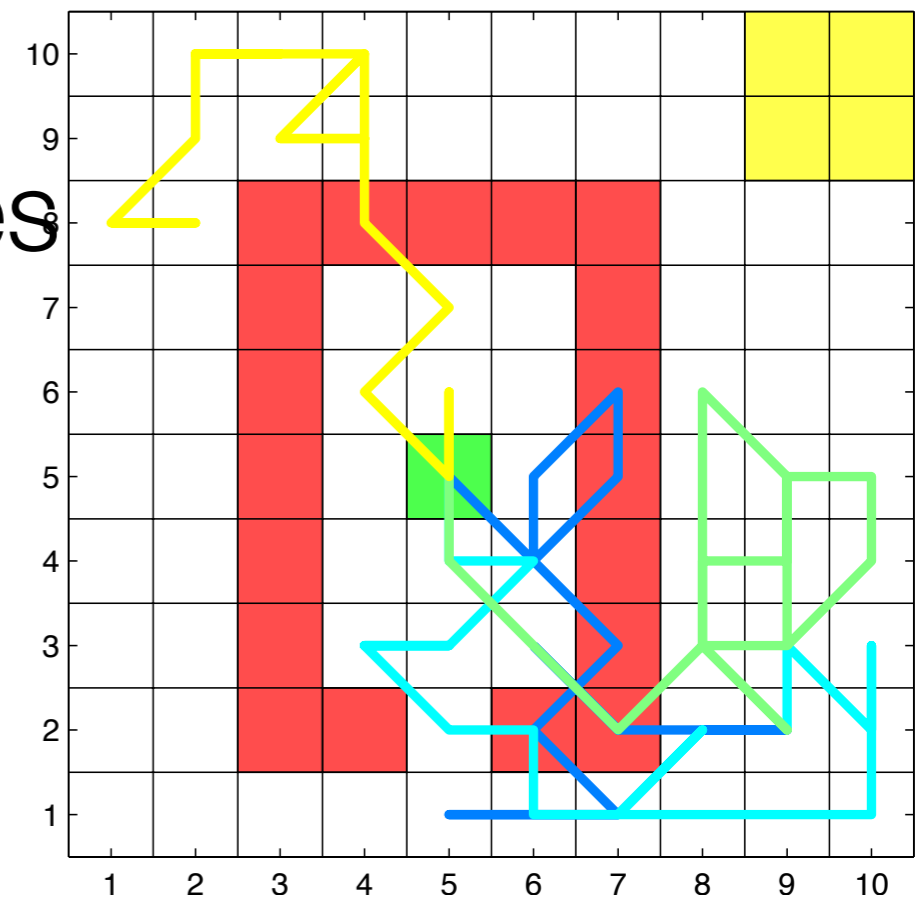
$$\mathcal{V}(s) = \frac{1}{N} \sum_i \left\{ \sum_{t'=1}^T r_{t'}^i \mid s_0 = s \right\}$$

## ► More efficient use of samples

- Every visit MC
- Bootstrap: TD
- Dyna

## ► Better samples

- on policy versus off policy
- Stochastic search, UCT...



Bellman equation

$$V(s) = \sum_a \pi(a, s) \left[ \sum_{s'} \mathcal{T}_{ss'}^a [\mathcal{R}(s', a, s) + V(s')] \right]$$

Not yet converged, so it doesn't hold:

$$dV(s) = -V(s) + \sum_a \pi(a, s) \left[ \sum_{s'} \mathcal{T}_{ss'}^a [\mathcal{R}(s', a, s) + V(s')] \right]$$

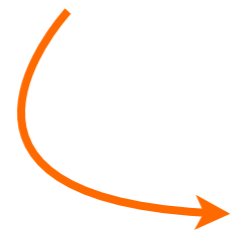
And then use this to update

$$V^{i+1}(s) = V^i(s) + dV(s)$$



$$dV(s) = -V(s) + \sum_a \pi(a, s) \left[ \sum_{s'} \mathcal{T}_{ss'}^a [\mathcal{R}(s', a, s) + V(s')] \right]$$

$$dV(s) = -V(s) + \sum_a \pi(a, s) \left[ \sum_{s'} \mathcal{T}_{ss'}^a [\mathcal{R}(s', a, s) + V(s')] \right]$$



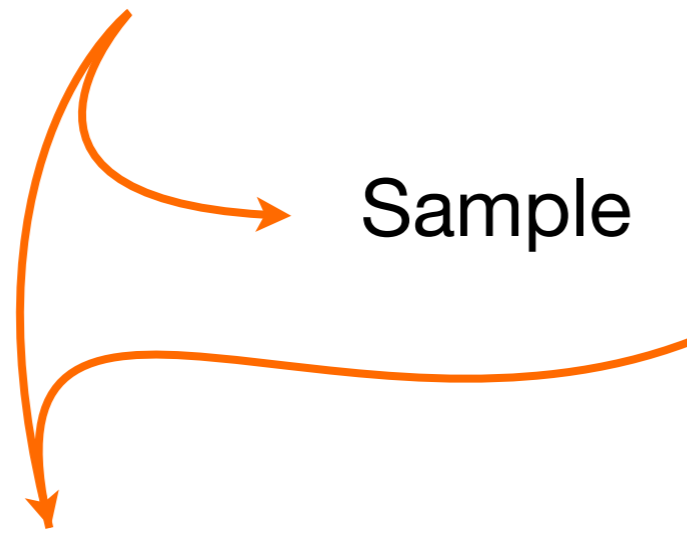
Sample

$$a_t \sim \pi(a|s_t)$$

$$s_{t+1} \sim \mathcal{T}_{s_t, s_{t+1}}^{a_t}$$

$$r_t = \mathcal{R}(s_{t+1}, a_t, s_t)$$

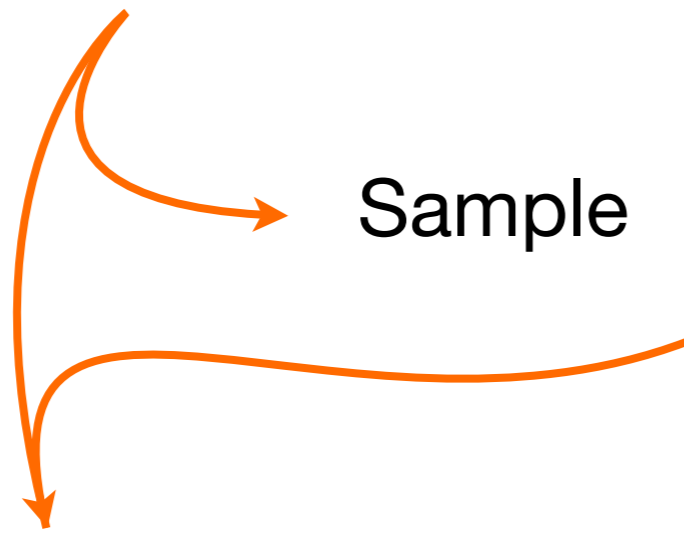
$$dV(s) = -V(s) + \sum_a \pi(a, s) \left[ \sum_{s'} \mathcal{T}_{ss'}^a [\mathcal{R}(s', a, s) + V(s')] \right]$$



$$\begin{aligned} a_t &\sim \pi(a|s_t) \\ s_{t+1} &\sim \mathcal{T}_{s_t, s_{t+1}}^{a_t} \\ r_t &= \mathcal{R}(s_{t+1}, a_t, s_t) \end{aligned}$$

$$\delta_t = -V_{t-1}(s_t) + r_t + V_{t-1}(s_{t+1})$$

$$dV(s) = -V(s) + \sum_a \pi(a, s) \left[ \sum_{s'} \mathcal{T}_{ss'}^a [\mathcal{R}(s', a, s) + V(s')] \right]$$



$$\begin{aligned} a_t &\sim \pi(a|s_t) \\ s_{t+1} &\sim \mathcal{T}_{s_t, s_{t+1}}^{a_t} \\ r_t &= \mathcal{R}(s_{t+1}, a_t, s_t) \end{aligned}$$

$$\delta_t = -V_{t-1}(s_t) + r_t + V_{t-1}(s_{t+1})$$

$$V^{i+1}(s) = V^i(s) + dV(s) \quad \longrightarrow \quad V_t(s_t) = V_{t-1}(s_t) + \alpha \delta_t$$

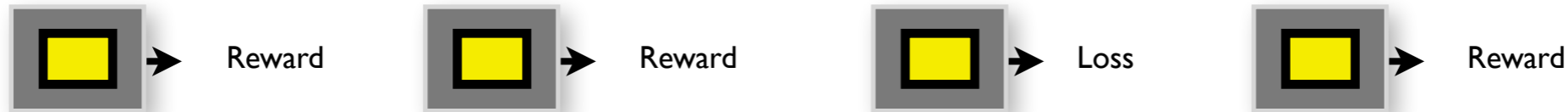
$$\begin{aligned}a_t &\sim \pi(a|s_t) \\s_{t+1} &\sim \mathcal{T}_{s_t, s_{t+1}}^{a_t} \\r_t &= \mathcal{R}(s_{t+1}, a_t, s_t) \\ \delta_t &= -V_t(s_t) + r_t + V_t(s_{t+1}) \\ V_{t+1}(s_t) &= V_t(s_t) + \alpha\delta_t\end{aligned}$$

- ▶ Pavlovian conditioning

## ▶ Pavlovian conditioning



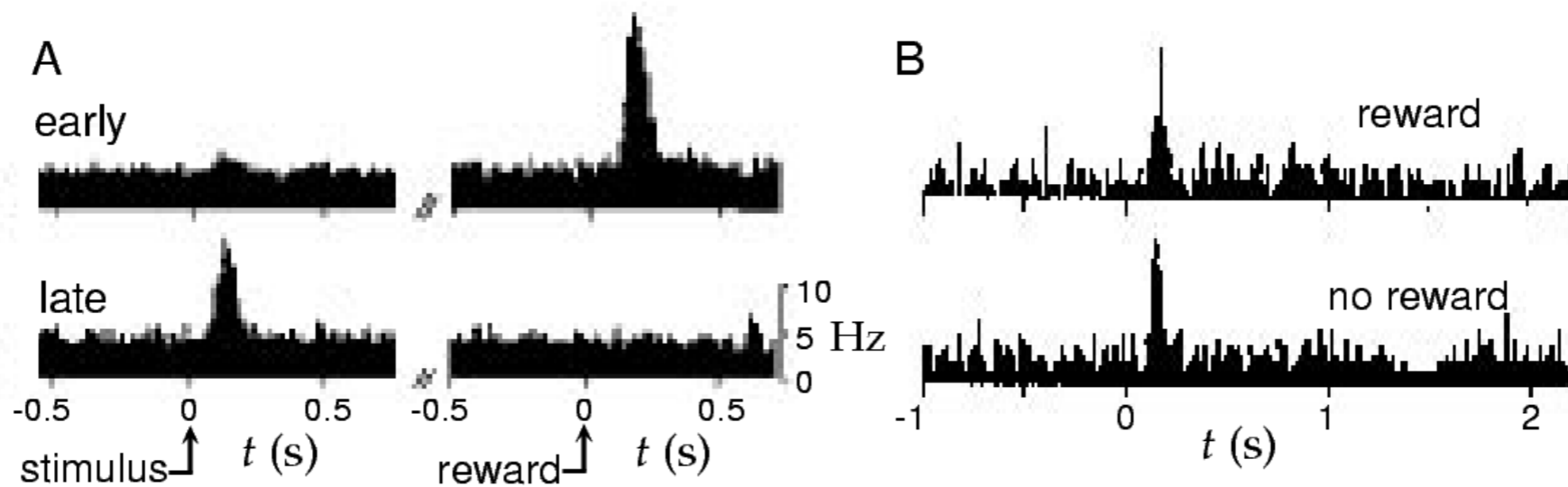
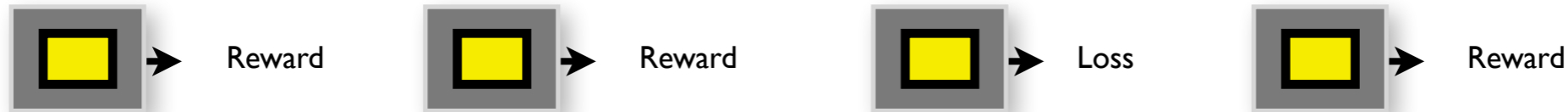
## ► Pavlovian conditioning



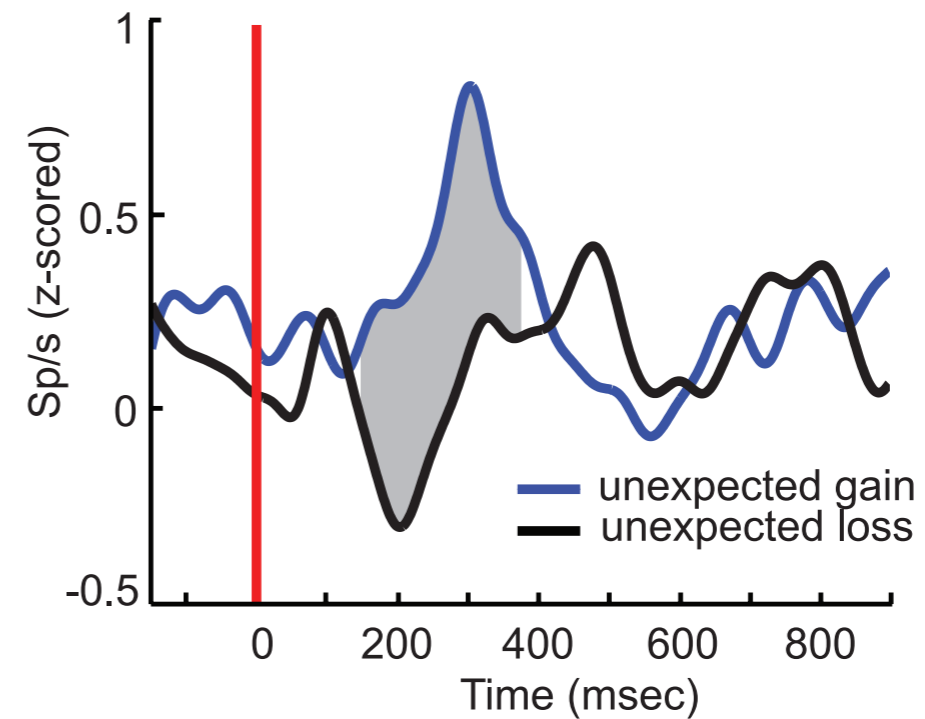
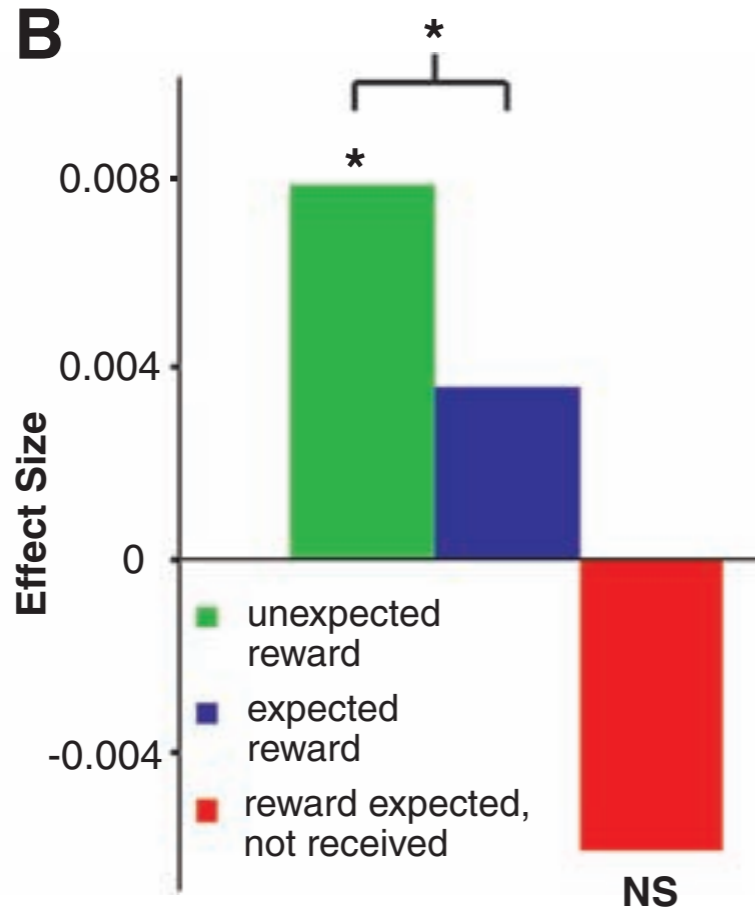
$$\mathcal{V}_{t+1}(s) = \mathcal{V}_t(s) + \epsilon \underbrace{(\mathcal{R}_t - \mathcal{V}_t(s))}_{= \text{Prediction error}}$$



## ► Pavlovian conditioning

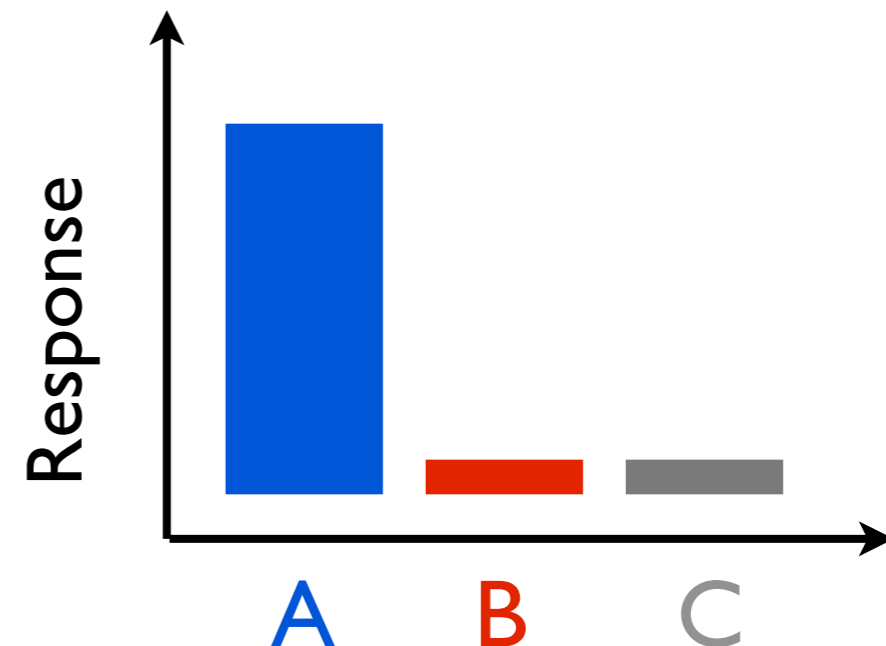


$$V_{t+1}(s) = V_t(s) + \epsilon \underbrace{(\mathcal{R}_t - V_t(s))}_{= \text{Prediction error}}$$

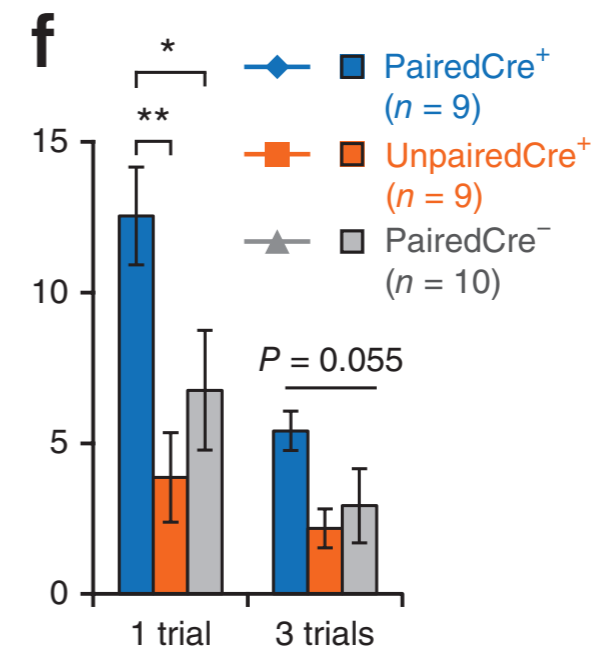
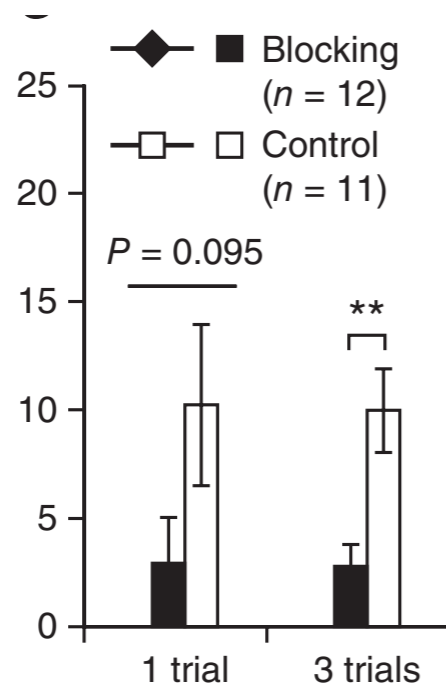
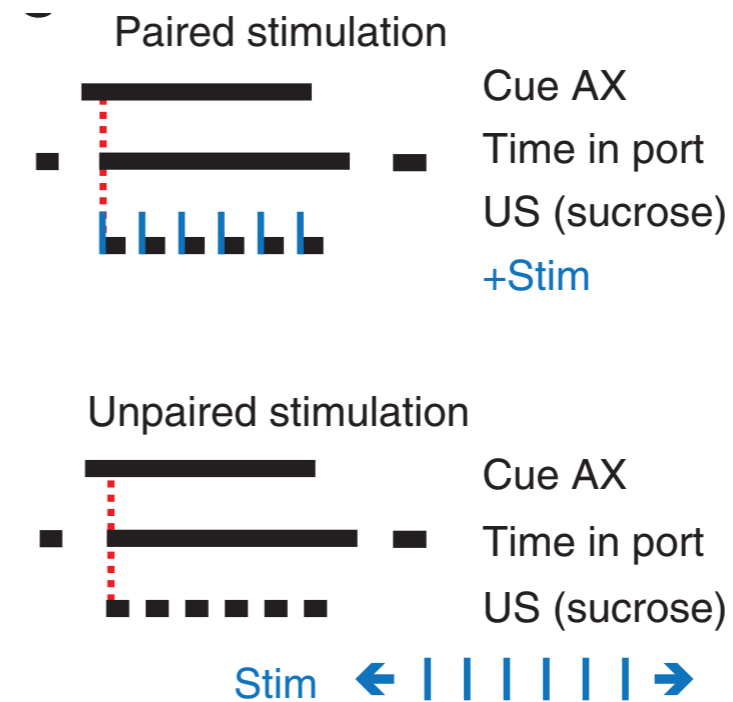
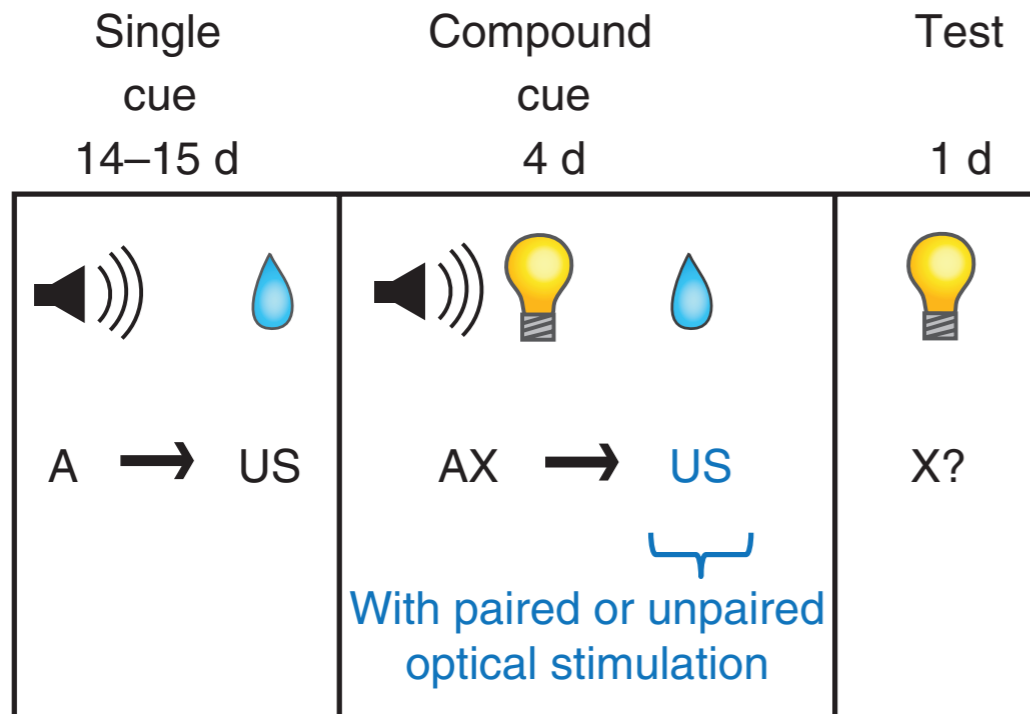


- ▶ Are predictions and prediction errors really causally important in learning?

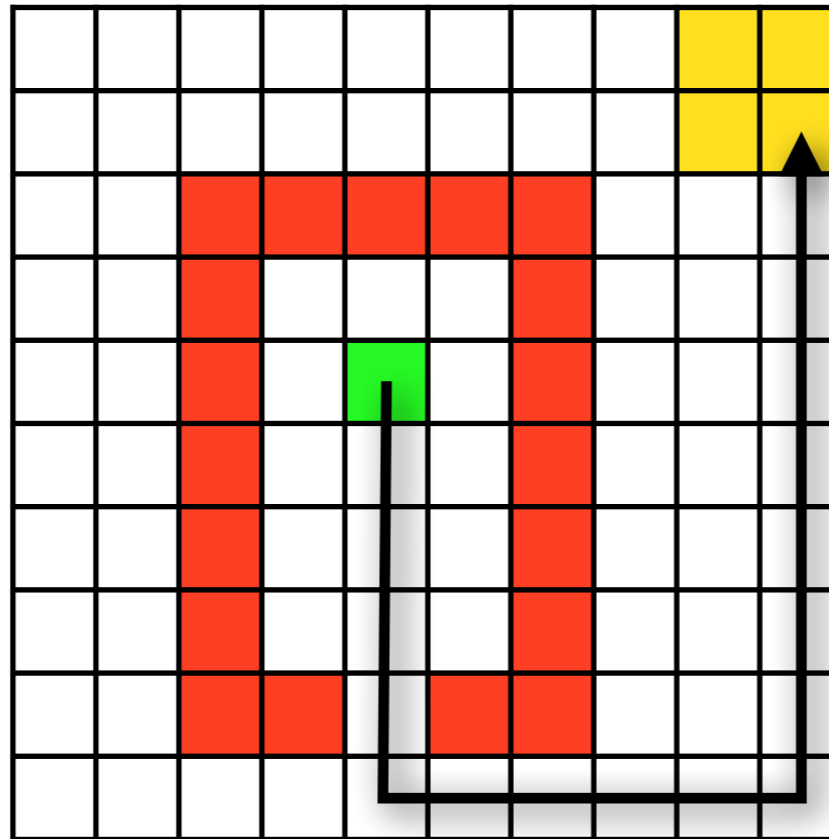
- 1: A -> Reward
- 2: A+B -> Reward
- 3: A -> ? approach
- B -> ? approach



# Causal role of phasic DA in learning



Steinberg et al., 2013 Nat. Neurosci.



$$V(s_t) = \mathbb{E}[r_t + r_{t+1} + r_{t+2} + \dots]$$

$$= \mathbb{E}[r_t] + \mathbb{E}[r_{t+1} + r_{t+2} + r_{t+3} \dots]$$

$$\Rightarrow V(s_t) = \mathbb{E}[r_t] + V(s_{t+1})$$

# “Cached” solutions to MDPs

- ▶ Learn from experience
- ▶ If we have true values  $V$ , then this is true every trial:

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

- ▶ If it is not true (we don't know true  $V$ ), then we get an error:

$$\delta = (\mathbb{E}[r_t] + V(s_{t+1})) - V(s_t) \neq 0$$

- ▶ So now we can update with our **experience**

$$V(s_t) \leftarrow V(s_t) + \epsilon \delta$$

- ▶ This is an **average over past experience**

- ▶ Do TD for state-action values instead:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha[r_t + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)]$$

$$s_t, a_t, r_t, s_{t+1}, a_{t+1}$$

- ▶ convergence guarantees - will estimate  $Q^\pi(s, a)$



- ▶ Learn off-policy
  - draw from some policy
  - “only” require extensive sampling

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \left[ \underbrace{r_t + \gamma \max_a Q(s_{t+1}, a)}_{\text{update towards optimum}} - Q(s_t, a_t) \right]$$

- ▶ will estimate  $Q^*(s, a)$

	<b>Model-free</b>
Pavlovian (state) values	$V^{\text{MF}}(s)$
Instrumental (state-action) values	$Q^{\text{MF}}(s, a)$

*There are both Pavlovian state and instrumental state-action values, and both of these can be either model-free (cached) or model-based.*

## ▶ “Cached” learning

- average experience
- do again what worked in the past
- averages are cheap to compute - no computational curse
- averages move slowly

If you have an average over large number of subjects, it won't move much if you add one more.

## ▶ “Cached” learning

- average experience
- do again what worked in the past
- averages are cheap to compute - no computational curse
- averages move slowly

If you have an average over large number of subjects, it won't move much if you add one more.

## ▶ “Goal-directed” or “Model-based” decisions

- Think through possible options and choose the best
- Requires detailed model of the world
- Requires huge computational resources
- Learning = building the model, extracting structure

# MF and MB learning of V and Q values

	<b>Model-free</b>	<b>Model-based</b>
Pavlovian (state) values	$V^{\text{MF}}(s)$	$V^{\text{MB}}(s)$
Instrumental (state-action) values	$Q^{\text{MF}}(s, a)$	$Q^{\text{MB}}(s, a)$

*There are both Pavlovian state and instrumental state-action values, and both of these can be either model-free (cached) or model-based.*

- ▶ Pavlovian model-free learning:

$$V_t(s) = V_{t-1}(s) + \epsilon(r_t - V_{t-1}(s))$$

$$p(a|s, V) \propto f(a, V(s)) p(a|s)$$



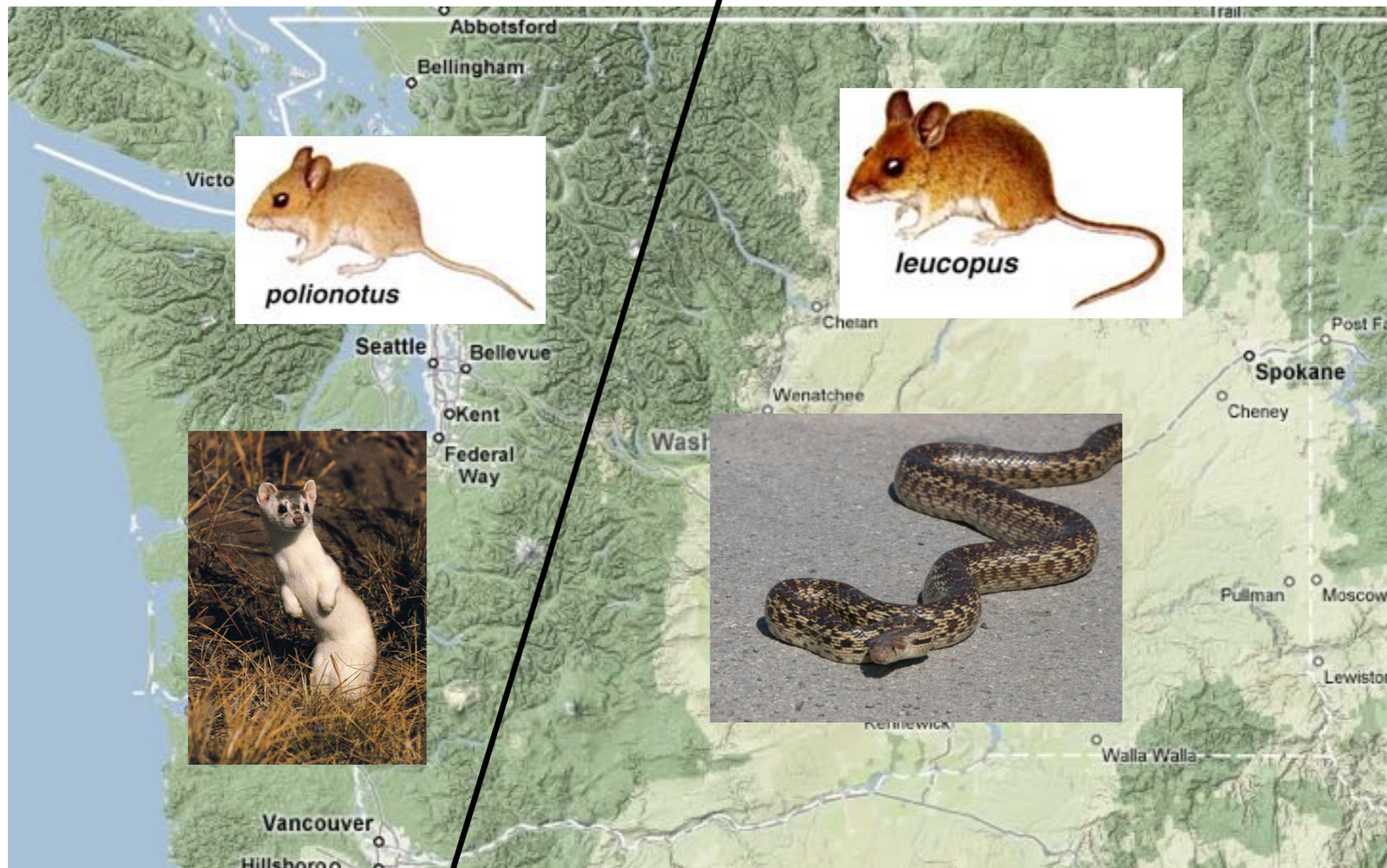
- ▶ Instrumental model-free learning:



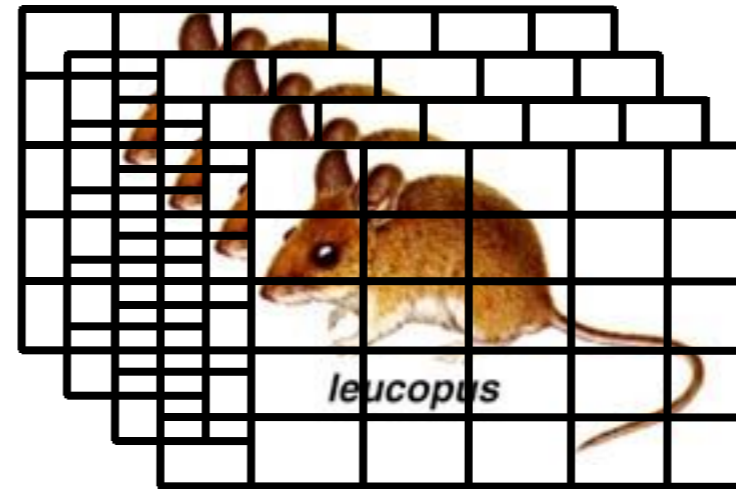
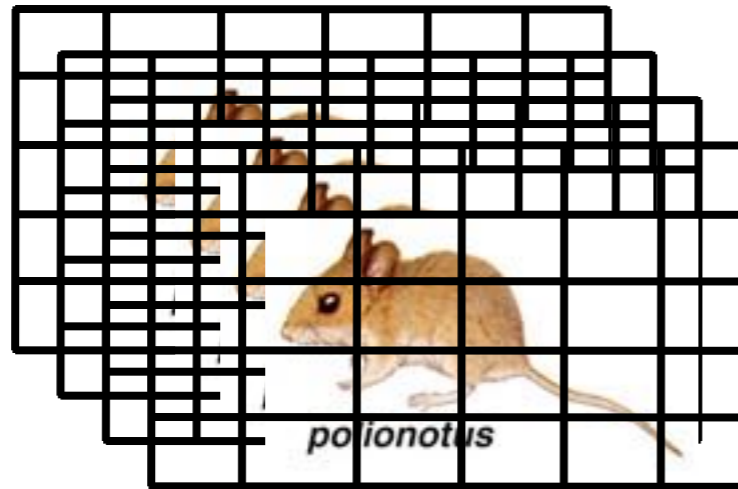
$$Q_t(a, s) = Q_{t-1}(a, s) + \epsilon(r_t - Q_{t-1}(a, s))$$



# Innate evolutionary strategies



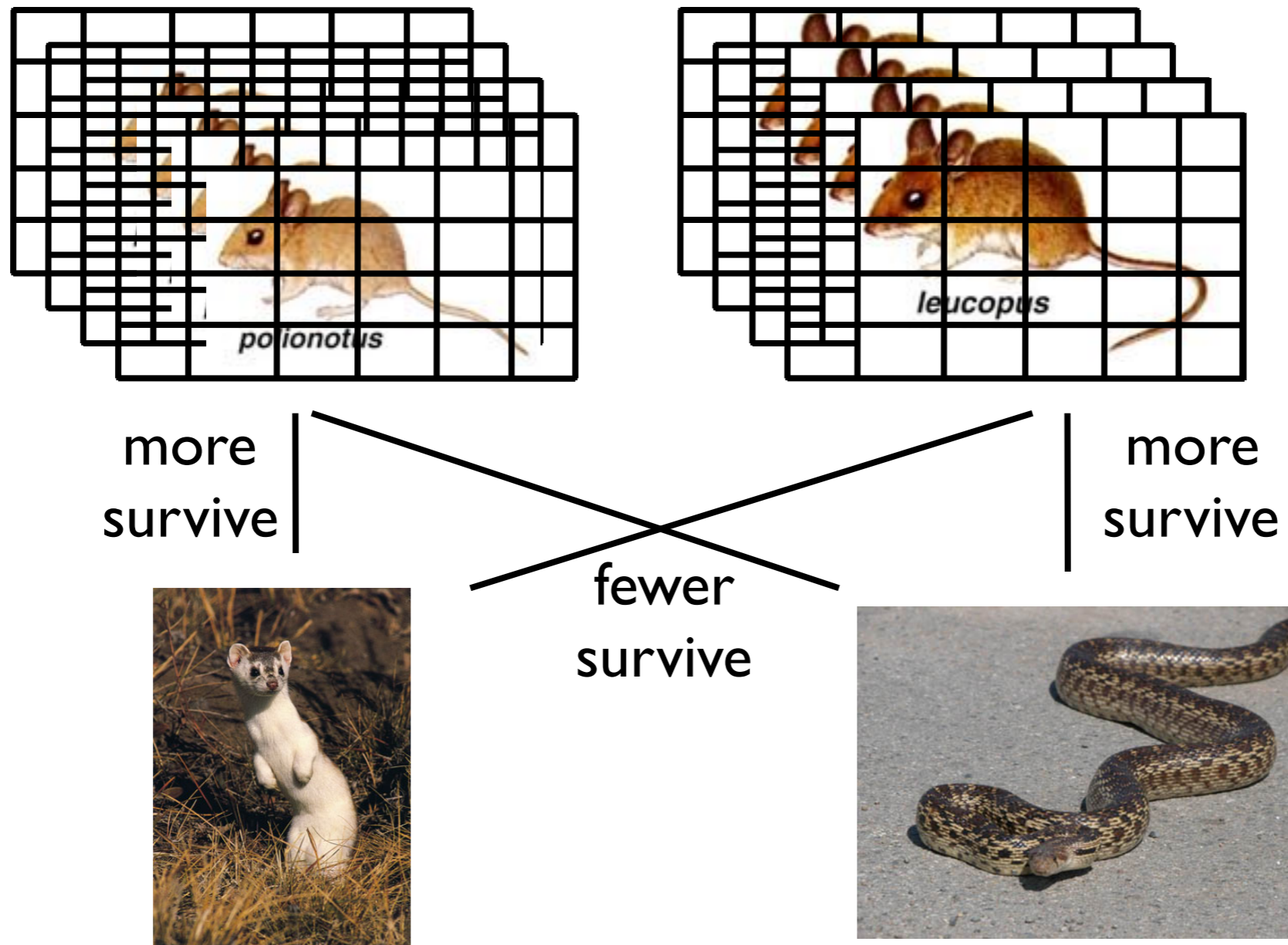
Hirsch and Bolles 1980





# Innate evolutionary strategies

are quite sophisticated...



Hirsch and Bolles 1980



- powerful
- inflexible over short timescale
- adaptive on evolutionary scale

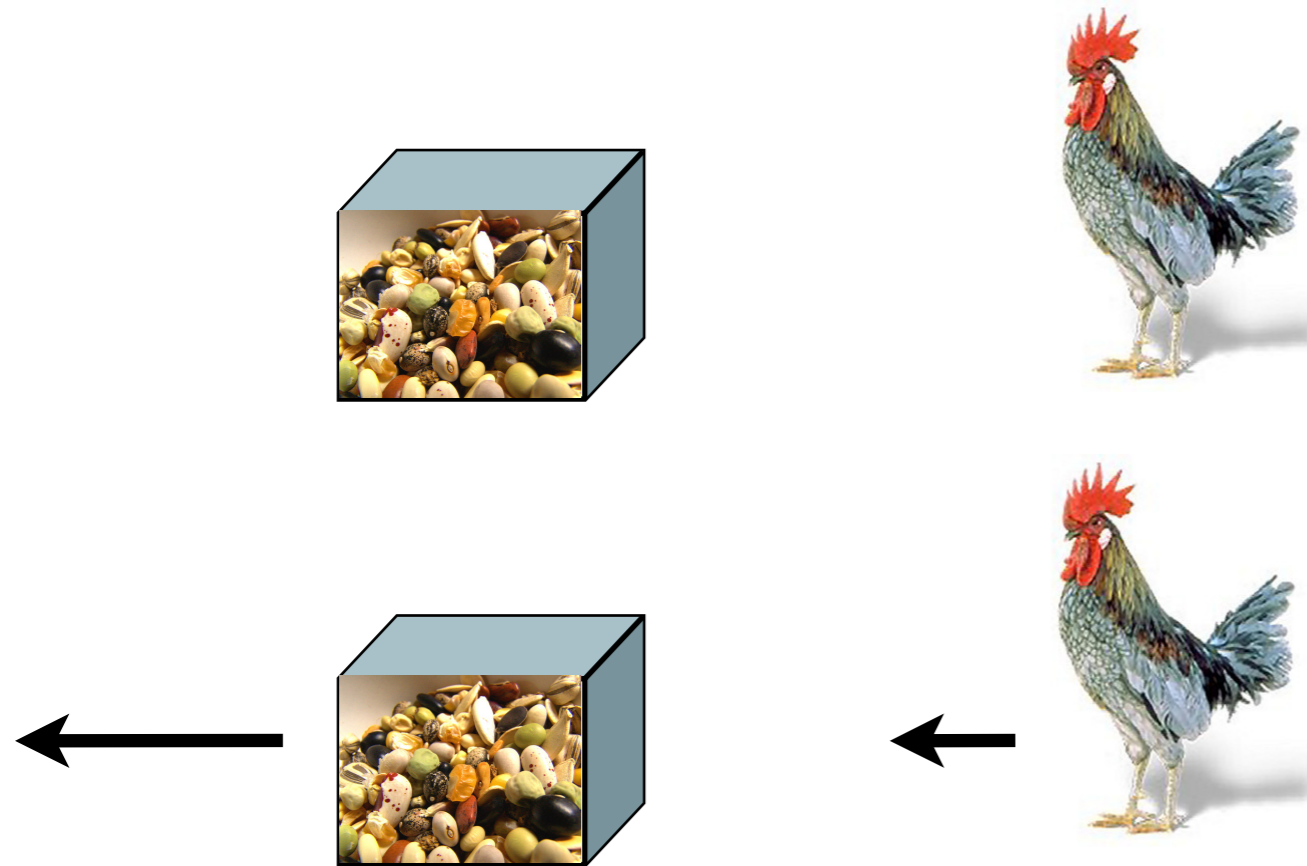
Hershberger 1986



- powerful
- inflexible over short timescale
- adaptive on evolutionary scale

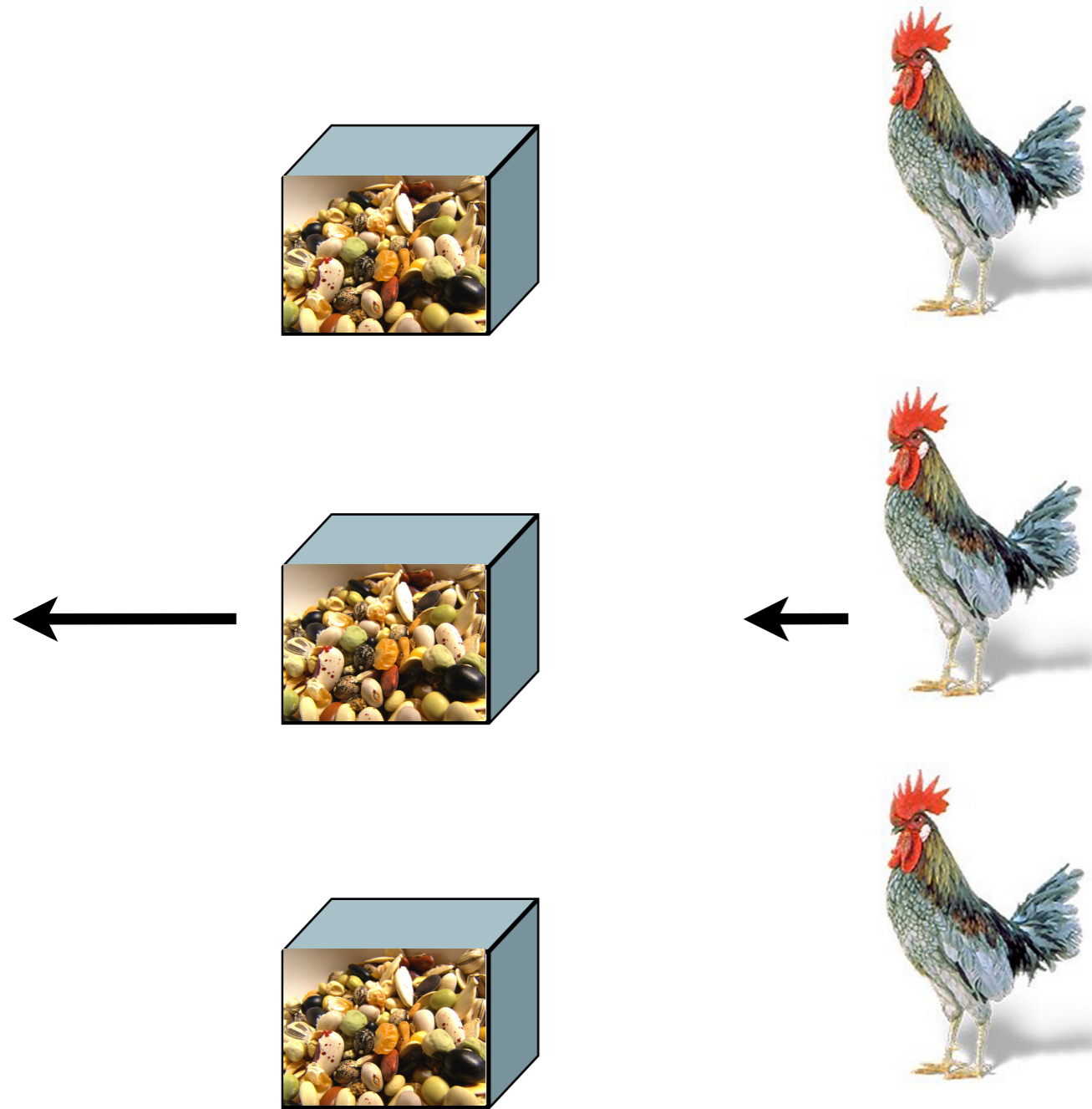


Hershberger 1986



- powerful
- inflexible over short timescale
- adaptive on evolutionary scale

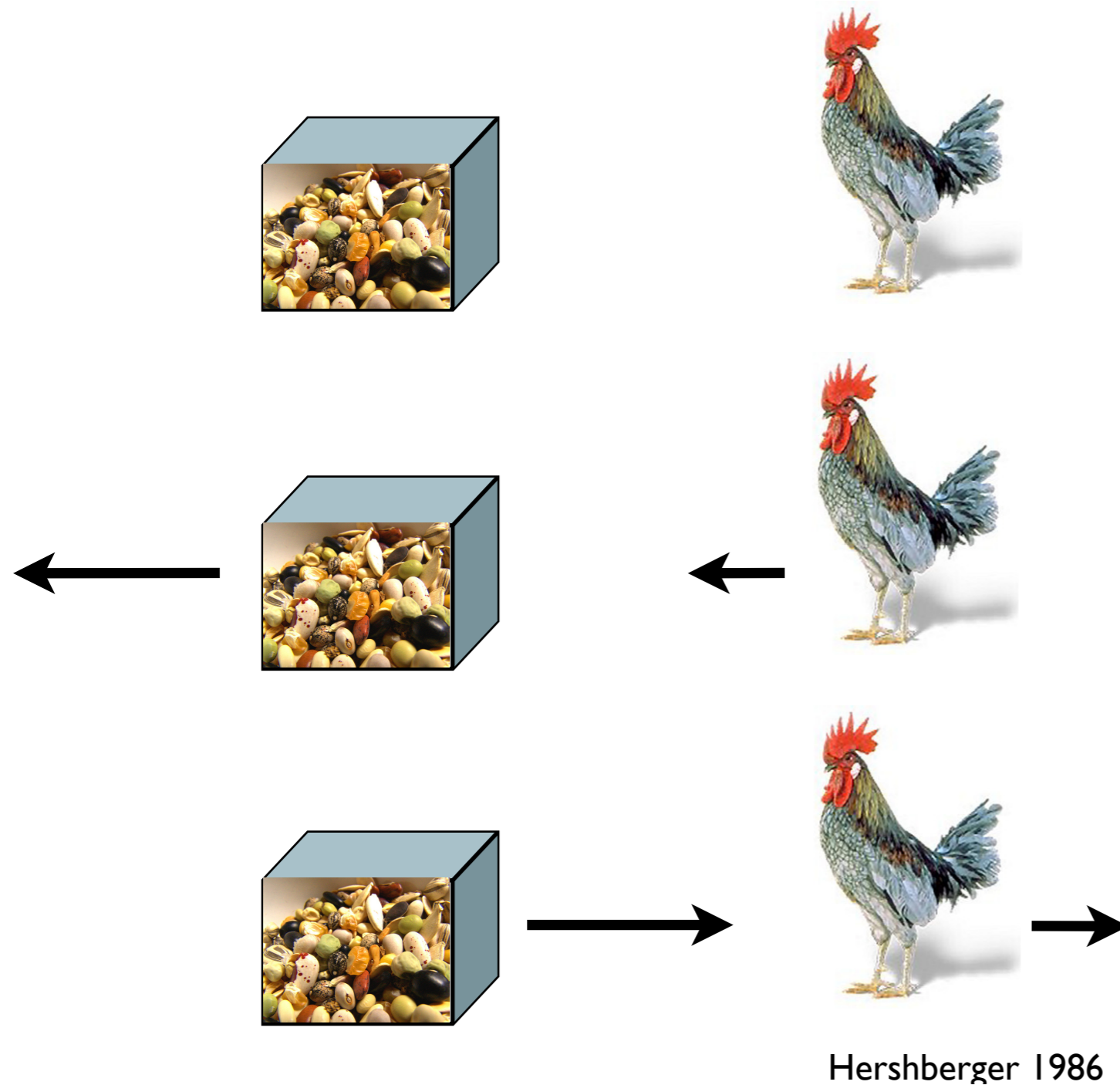
Hershberger 1986



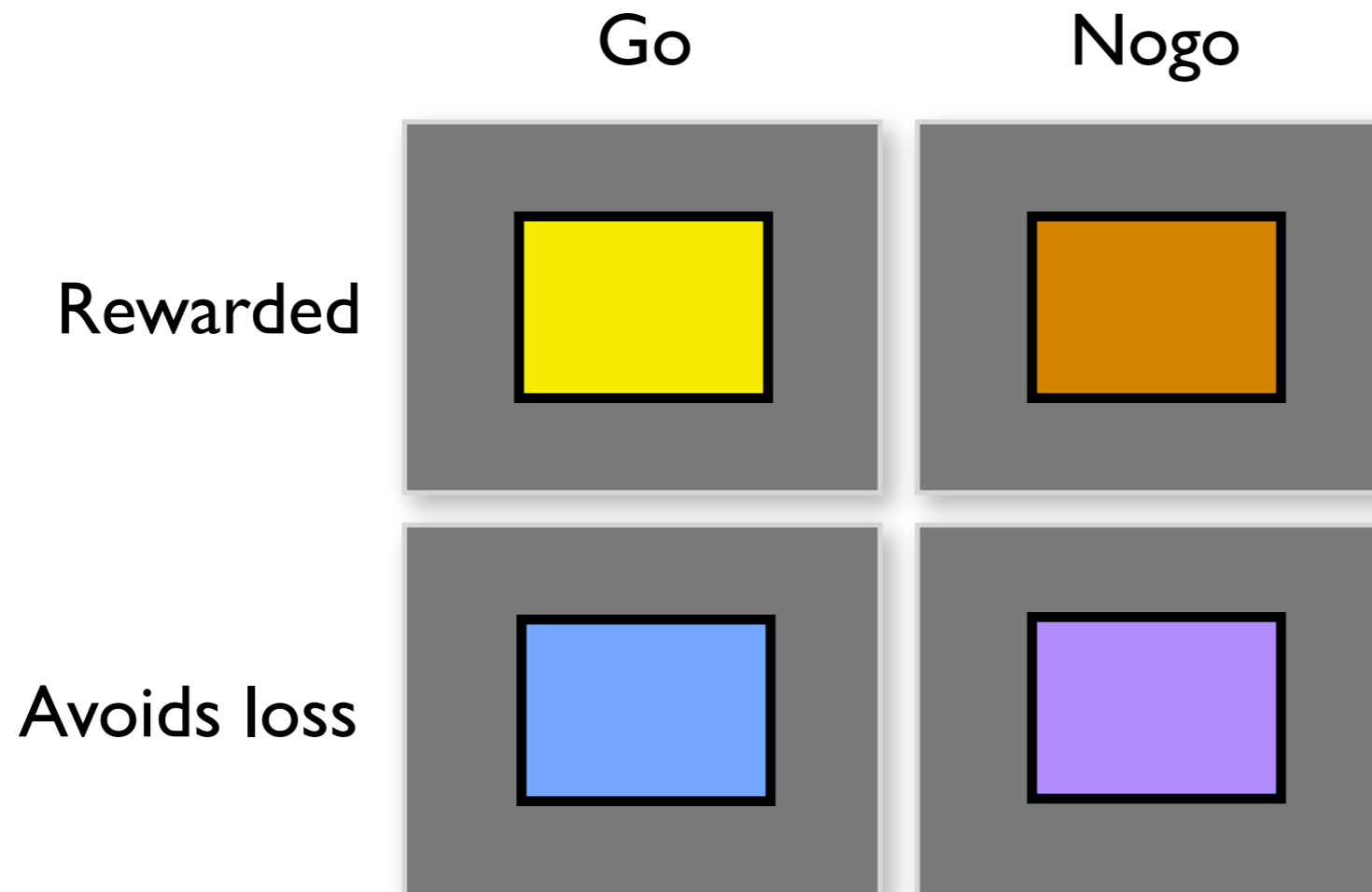
- powerful
- inflexible over short timescale
- adaptive on evolutionary scale

Hershberger 1986

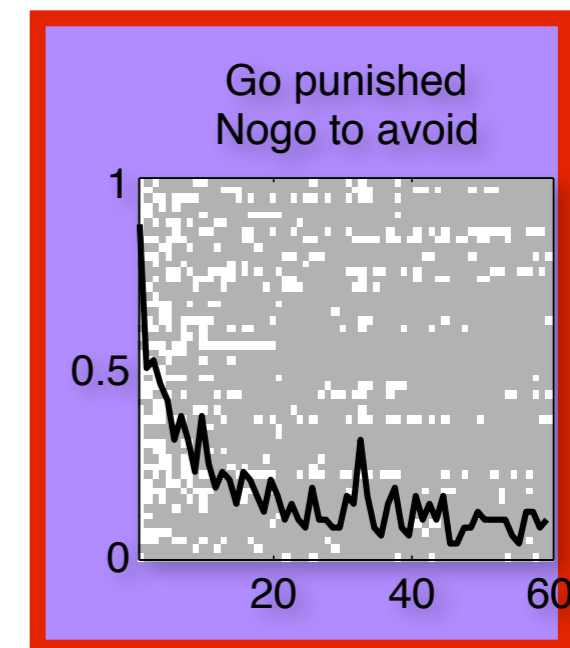
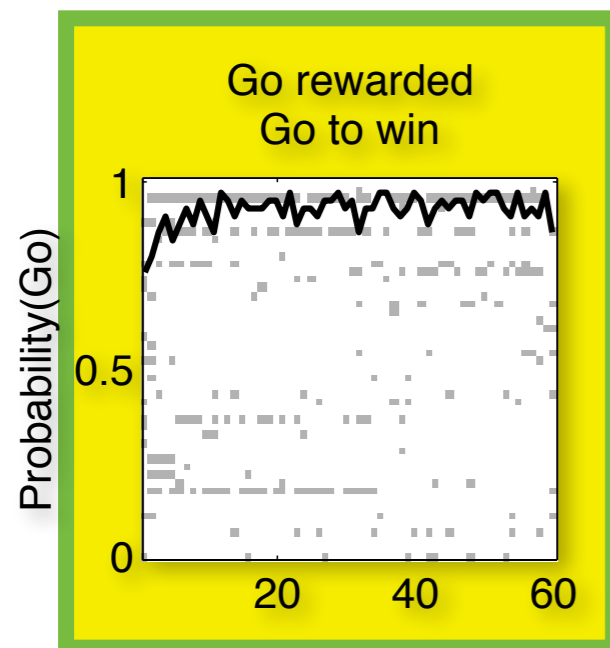
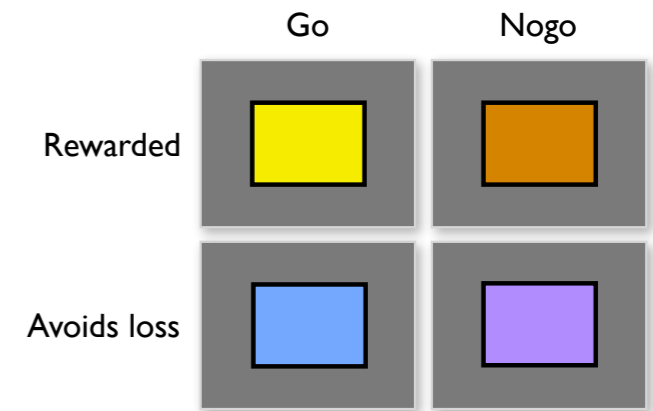
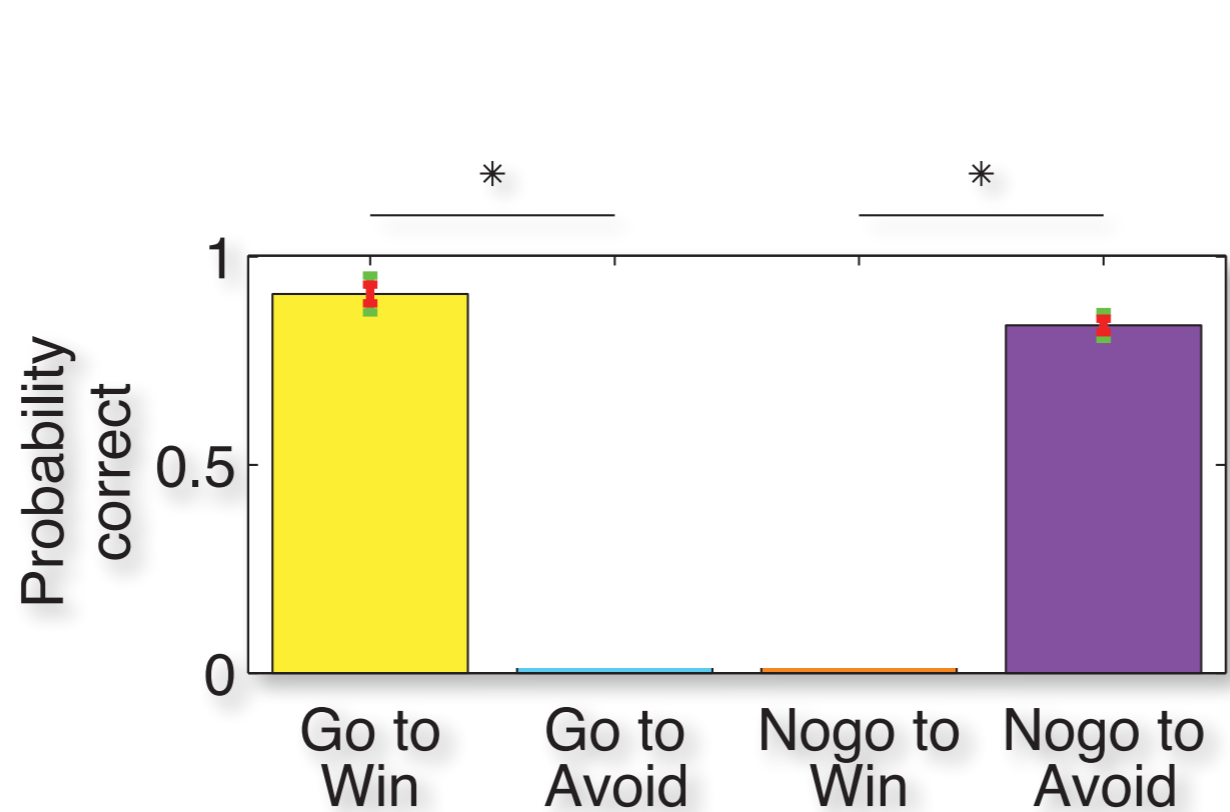




- powerful
- inflexible over short timescale
- adaptive on evolutionary scale

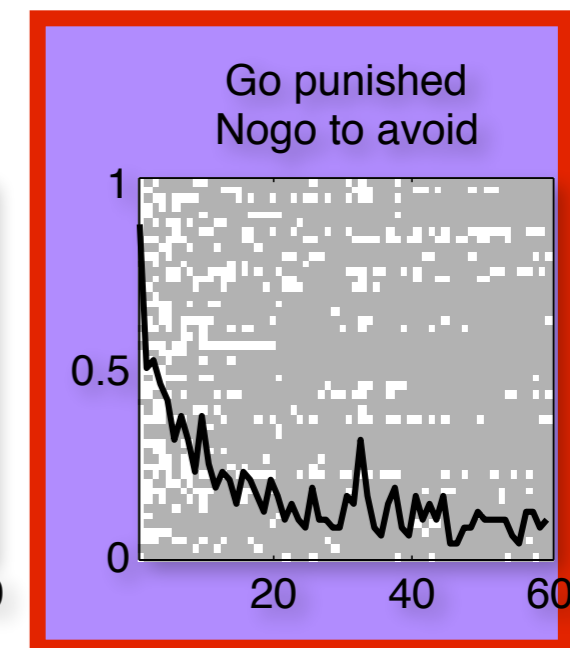
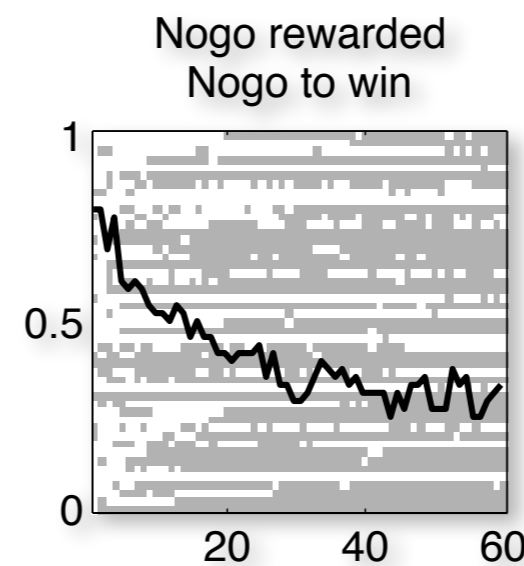
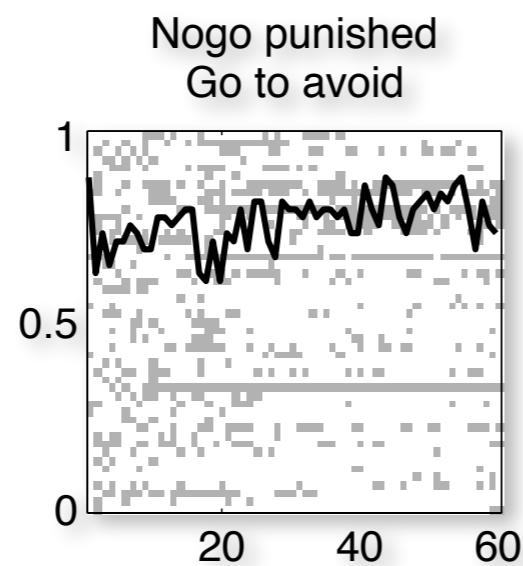
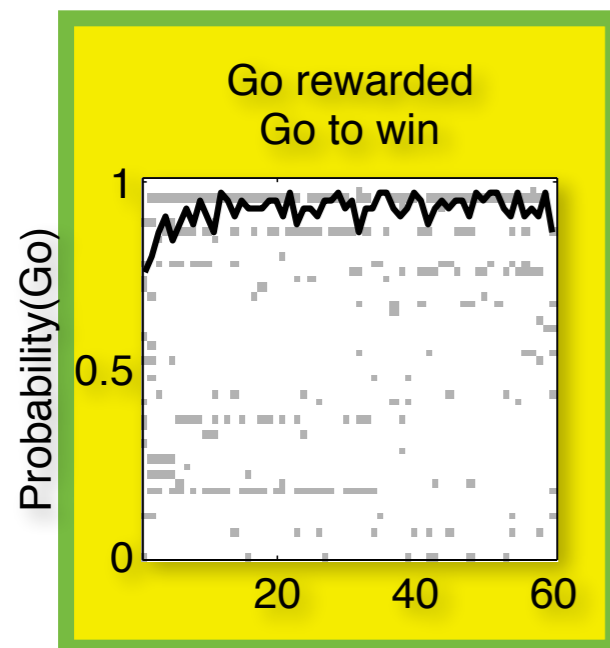
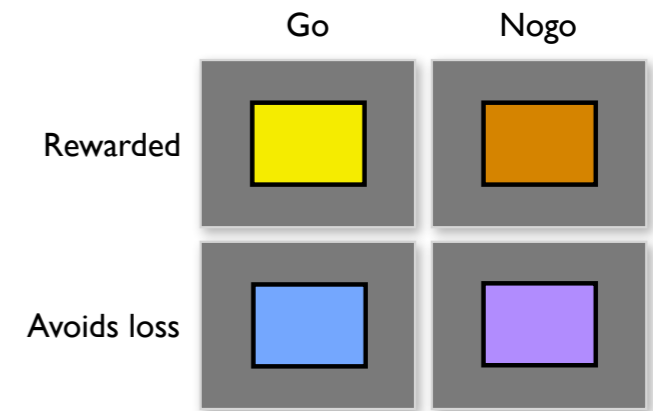
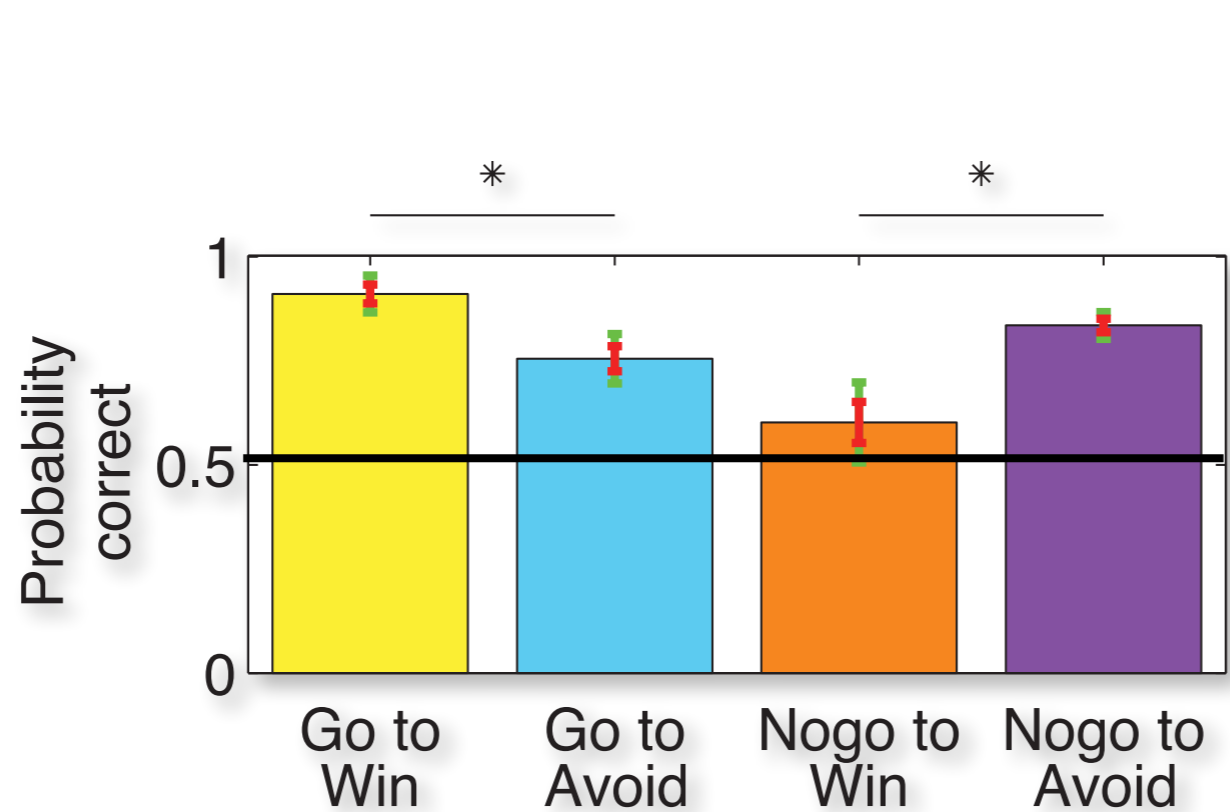


# Affective go / nogo task





# Affective go / nogo task



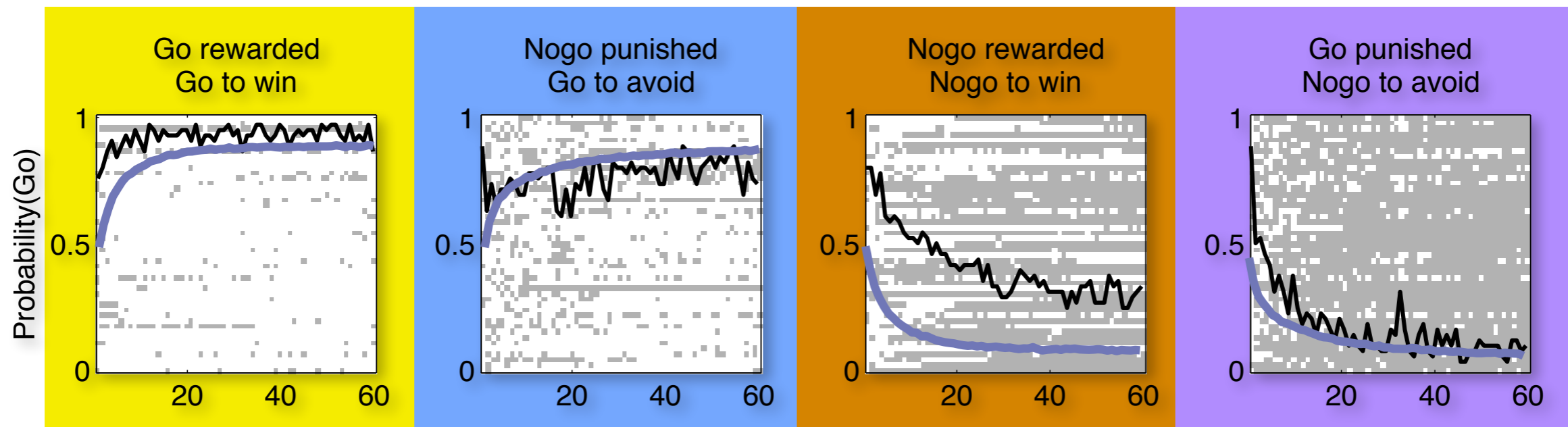
# Models



## ► Instrumental

$$p_t(a|s) \propto Q_t(s, a)$$

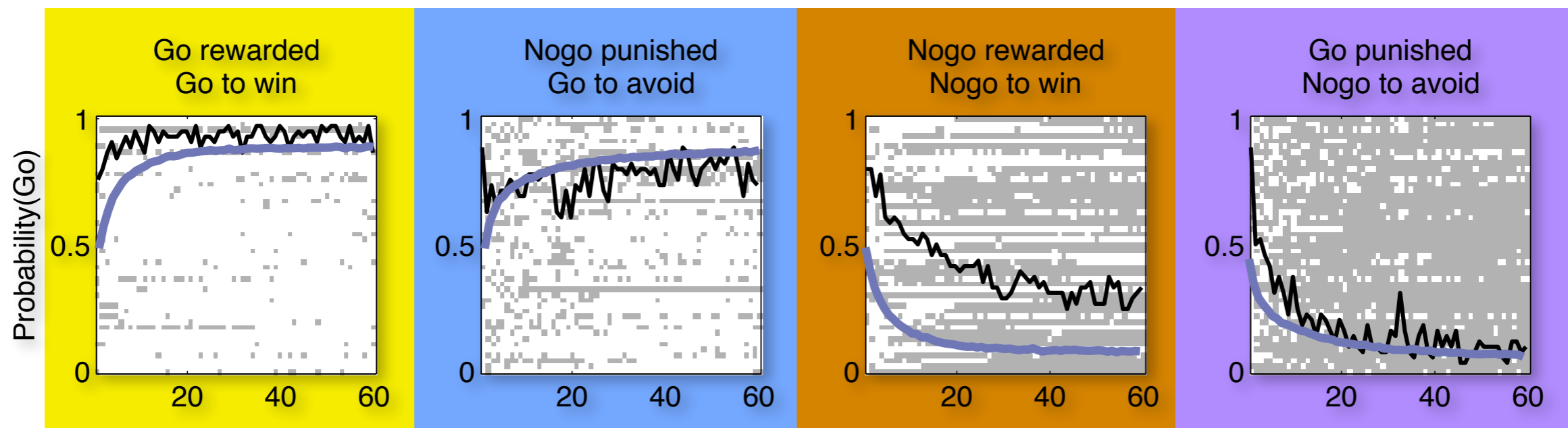
$$Q_{t+1}(s, a) = Q_t(s, a) + \alpha(r_t - Q_t(s, a))$$



# Models



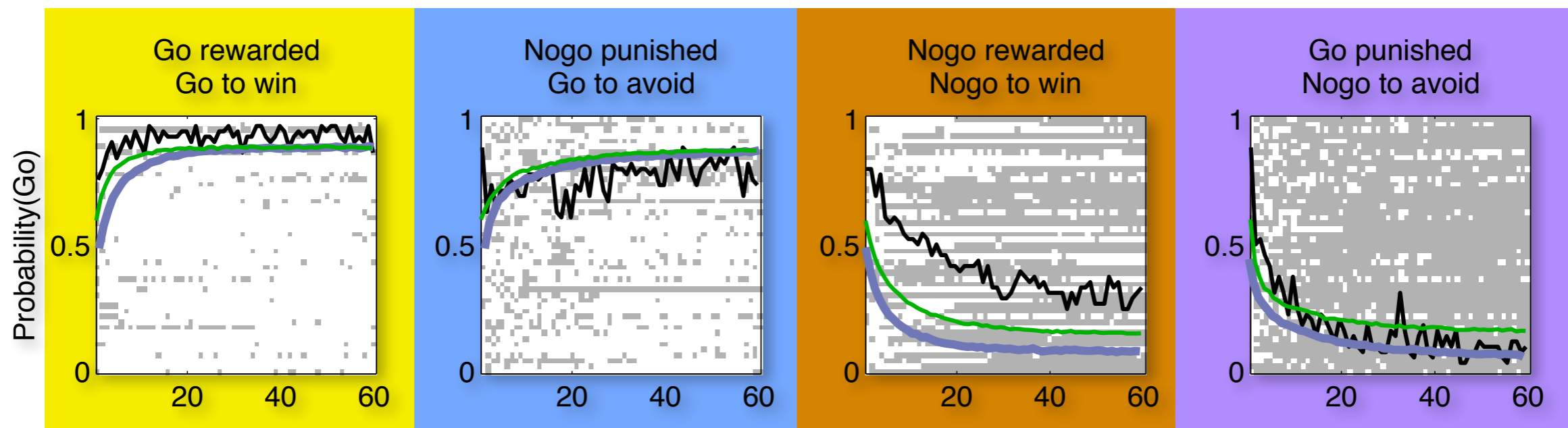
## ► Instrumental + bias



# Models



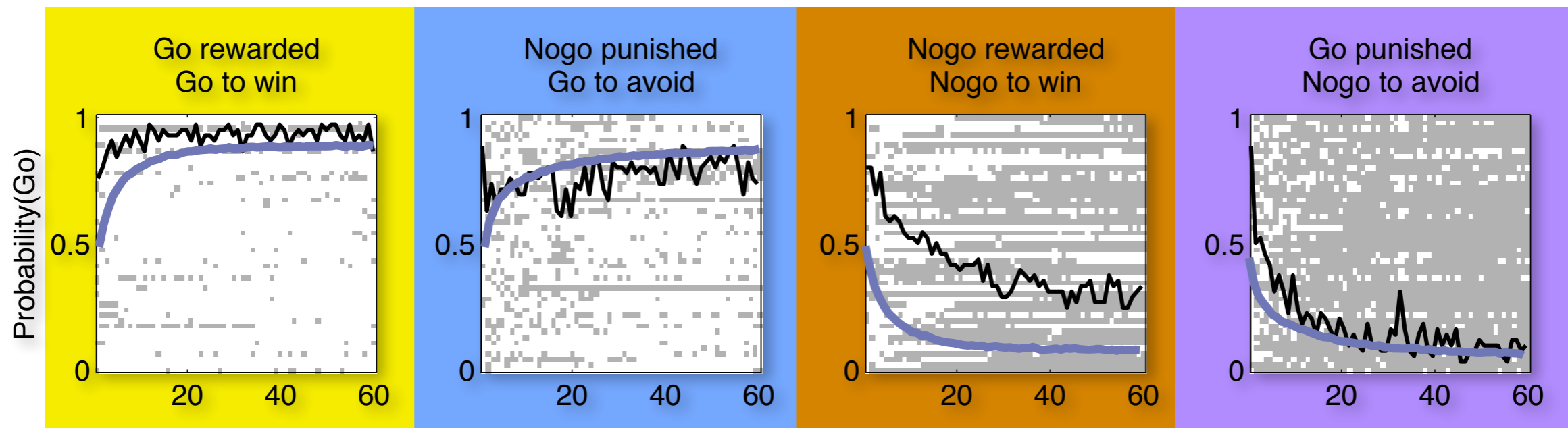
## ► Instrumental + bias



# Models



## ► Instrumental + bias + Pavlovian

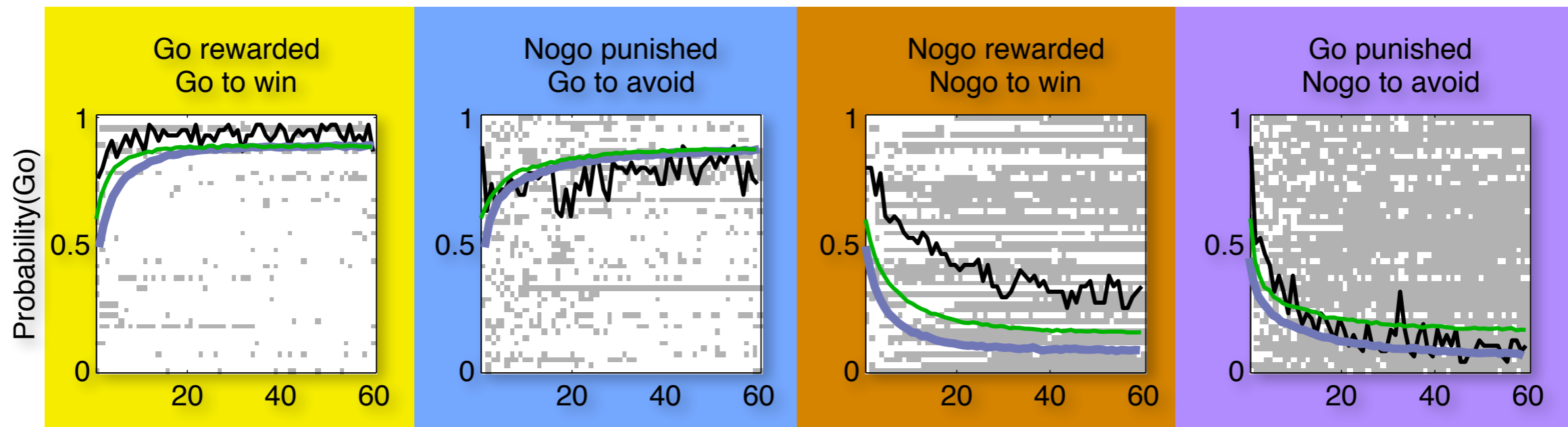


Guitart-Masip et al., 2012 J Neurosci

# Models



## ► Instrumental + bias + Pavlovian

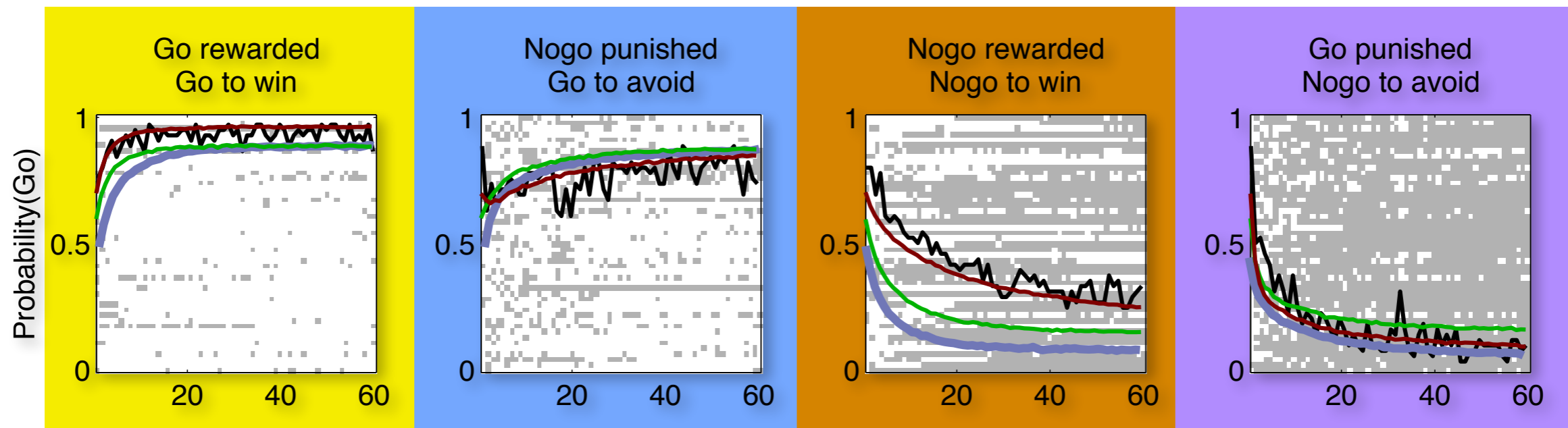


Guitart-Masip et al., 2012 J Neurosci

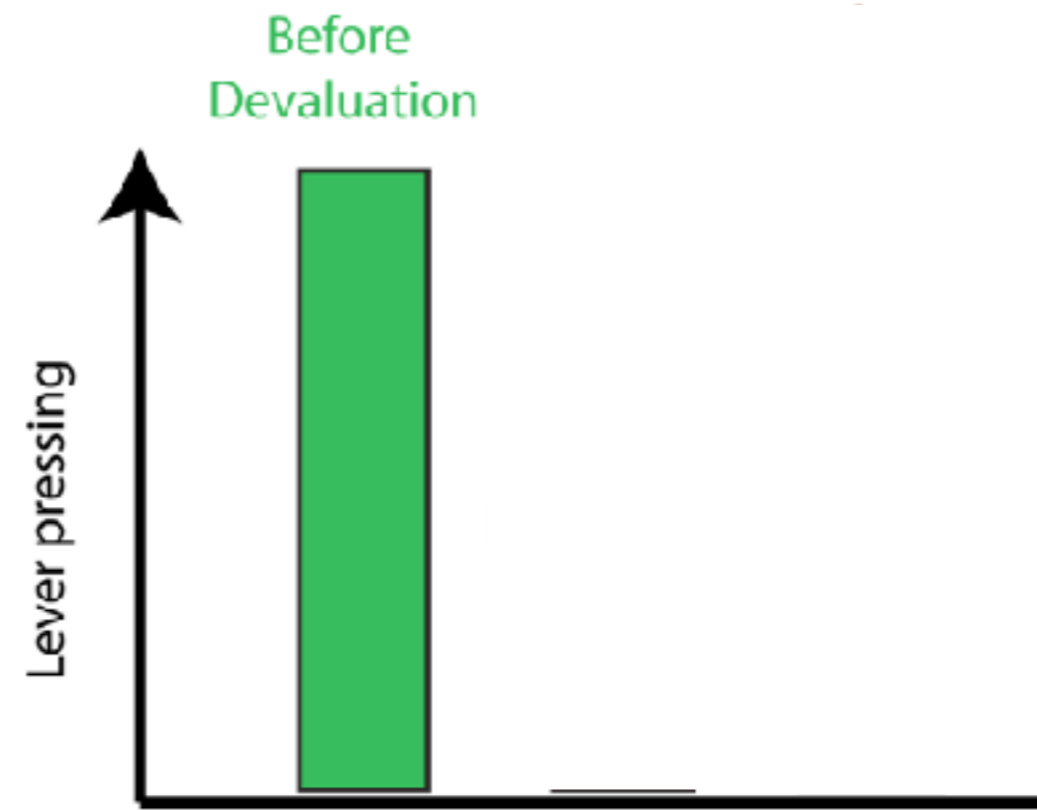
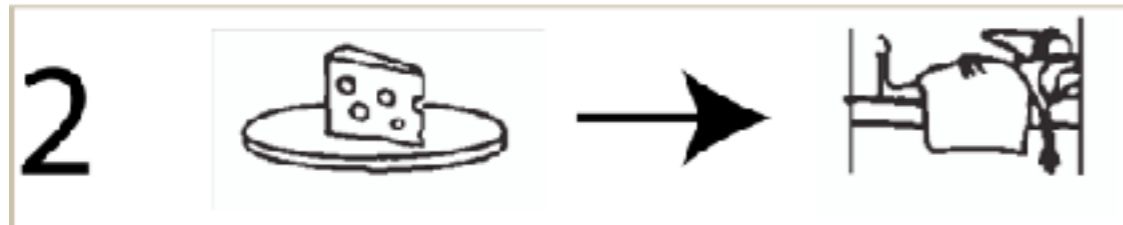
# Models



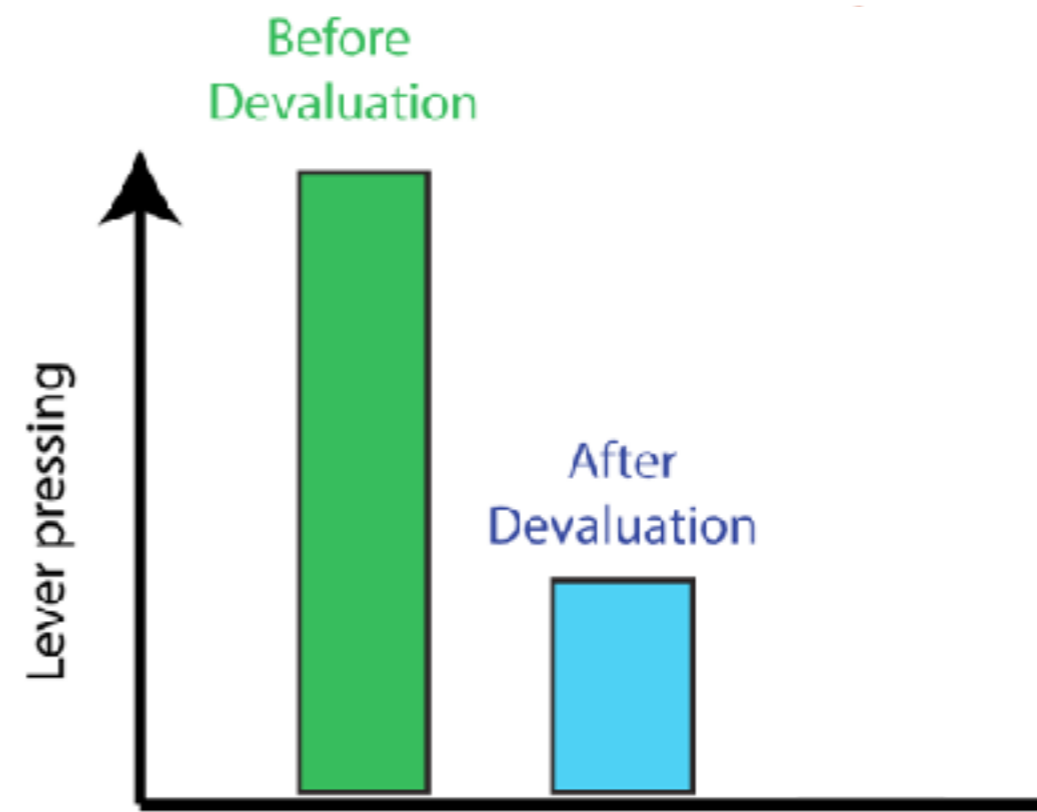
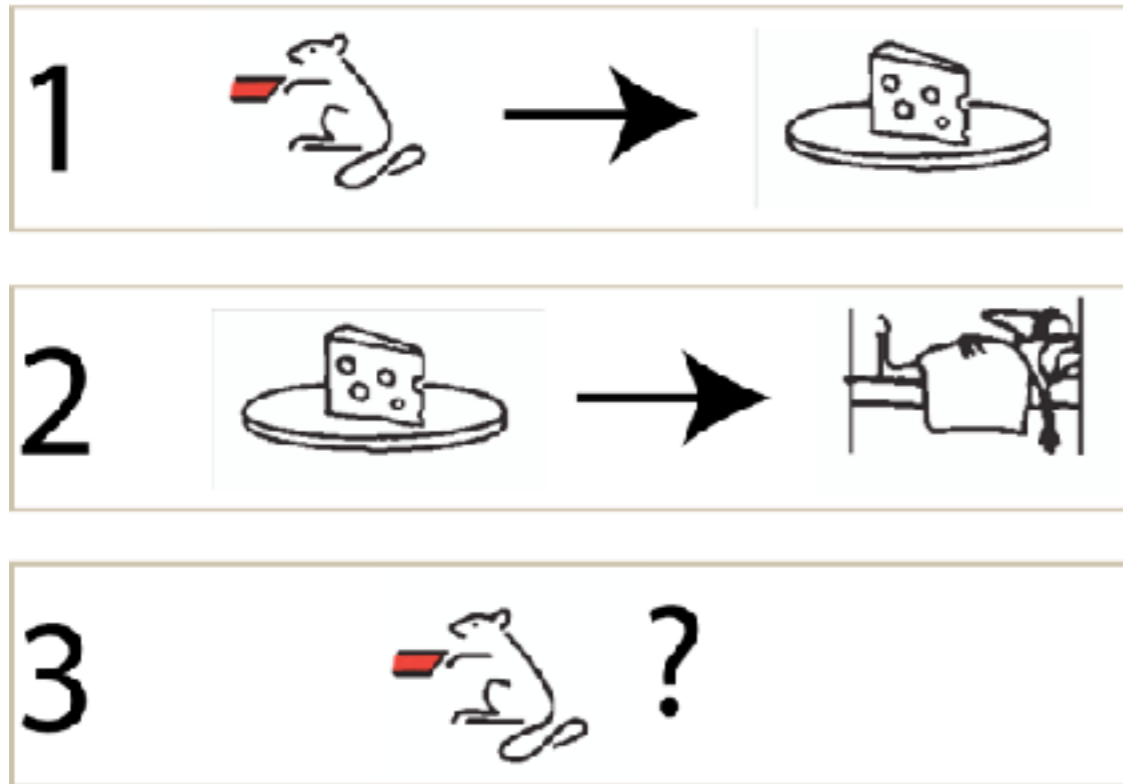
- ▶ Instrumental + bias + Pavlovian

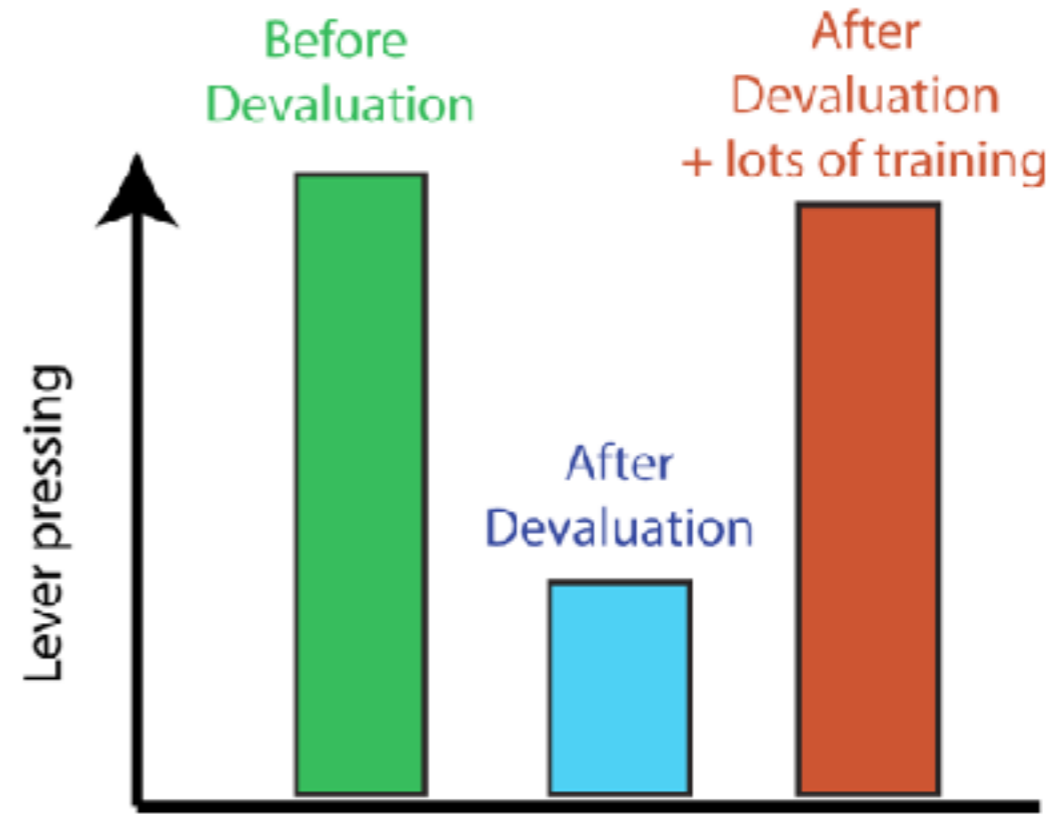
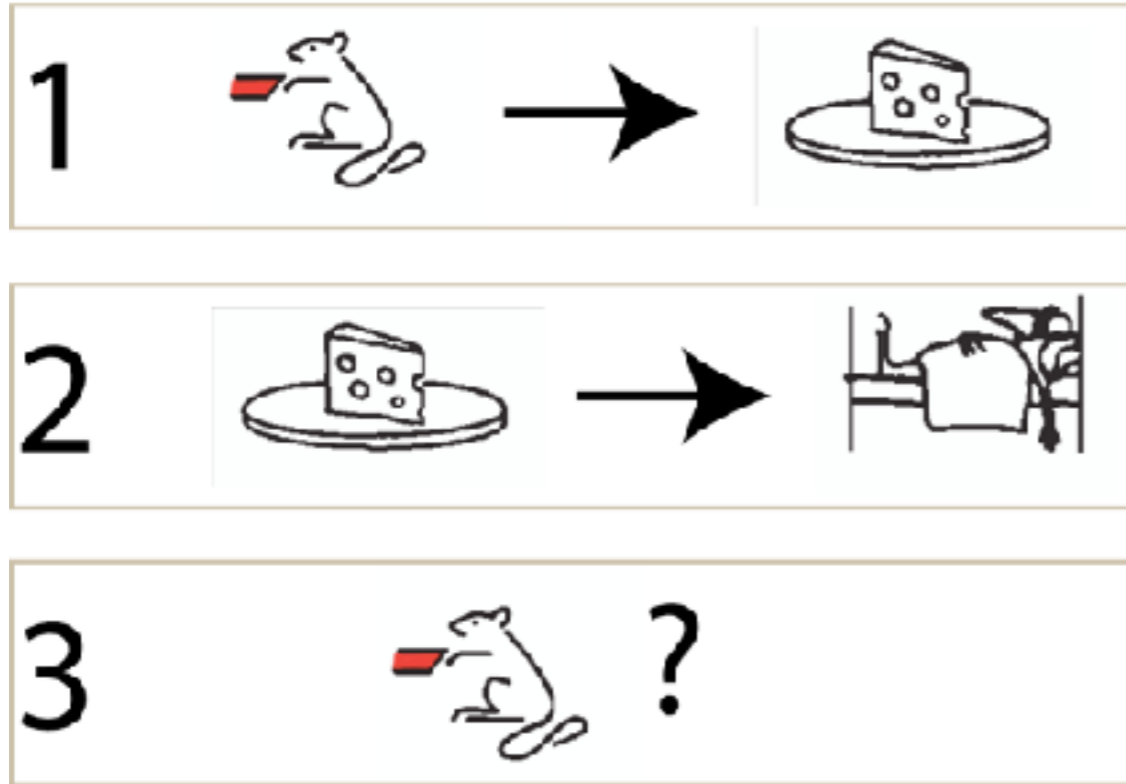


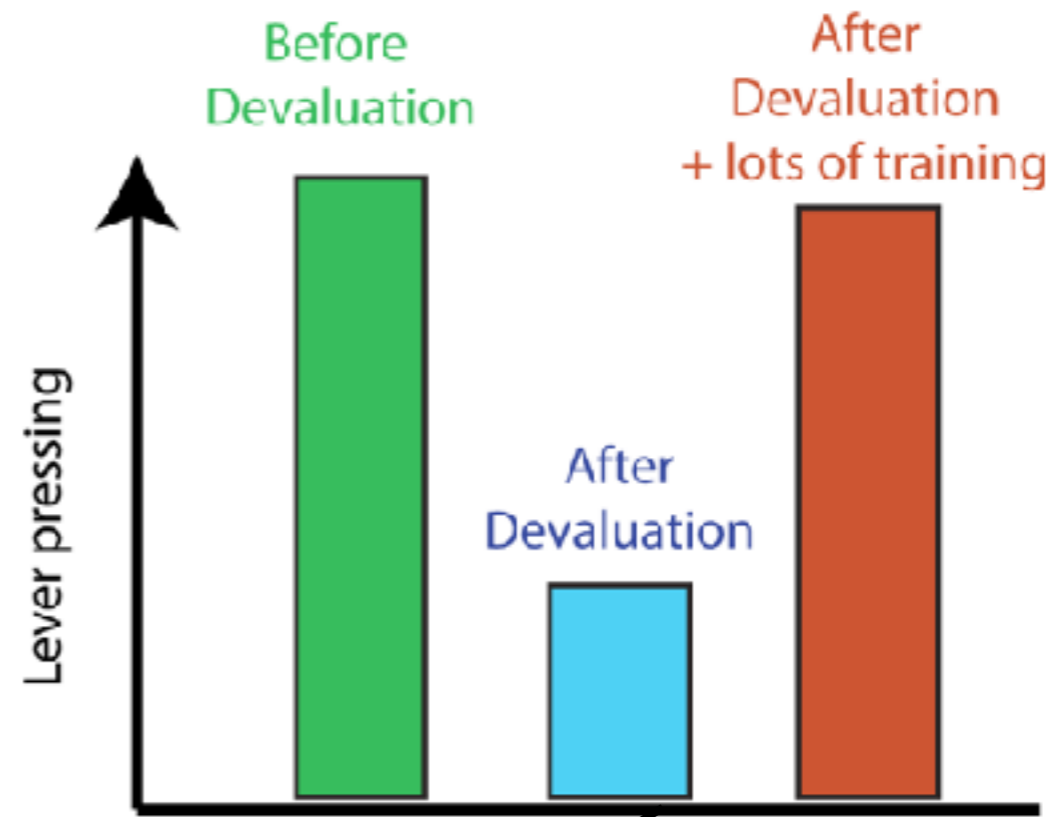
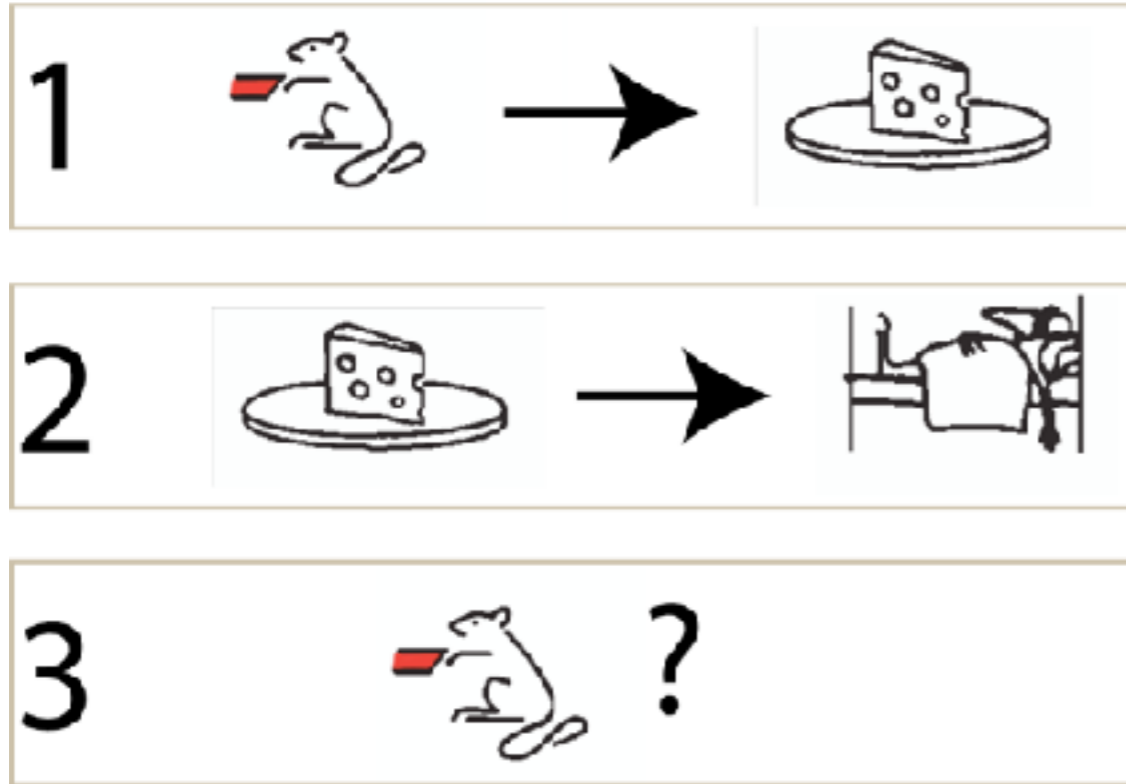
Guitart-Masip et al., 2012 J Neurosci



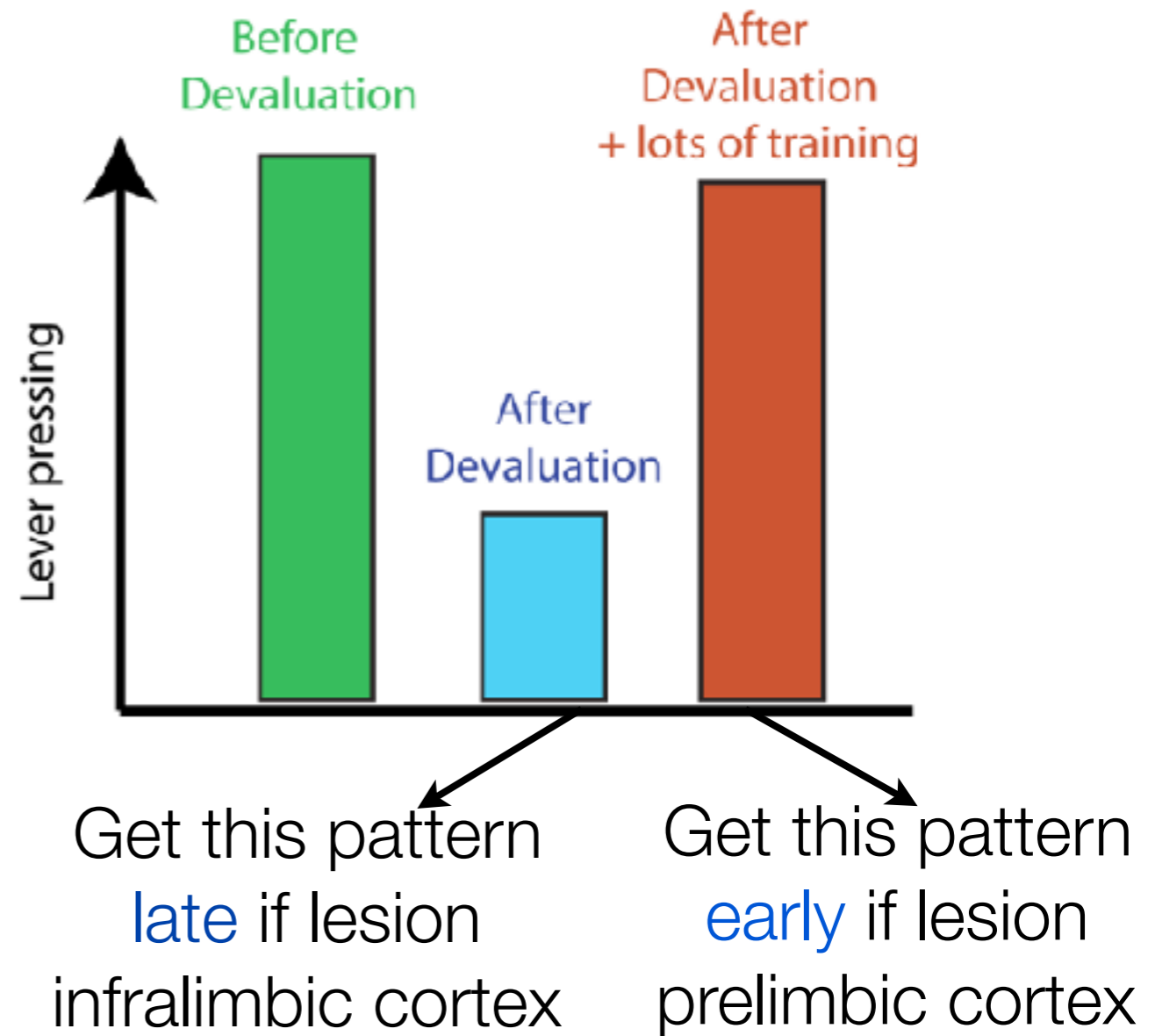
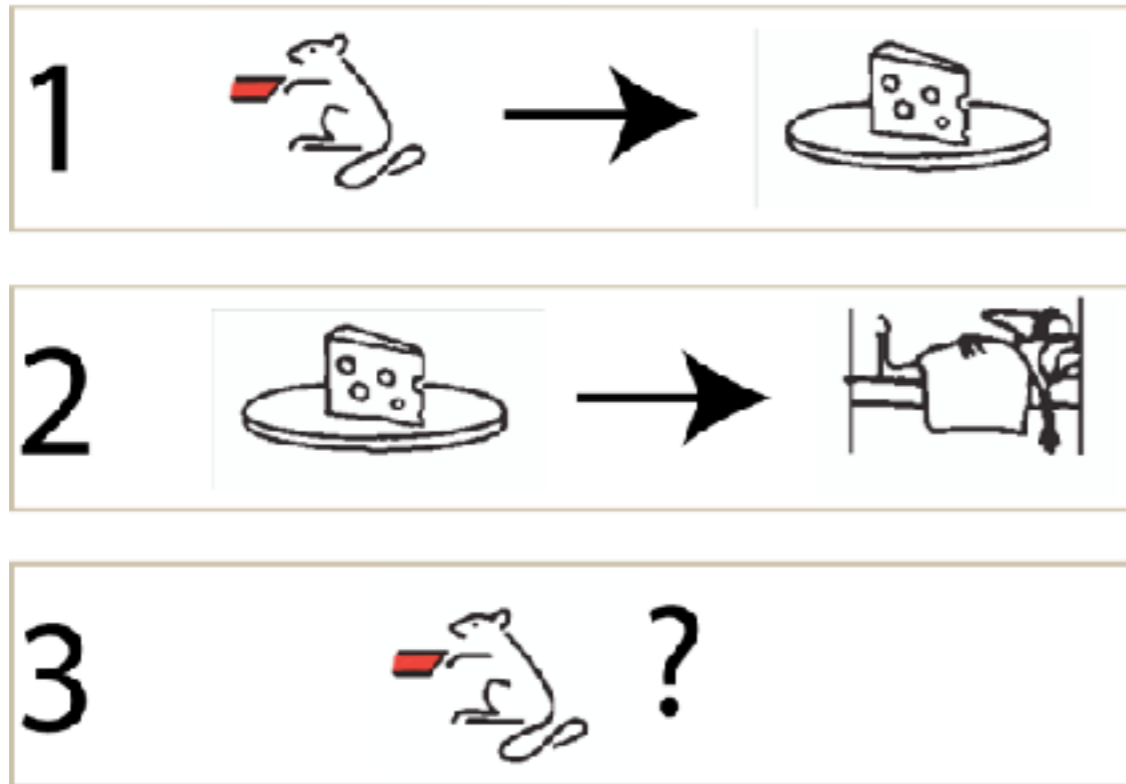


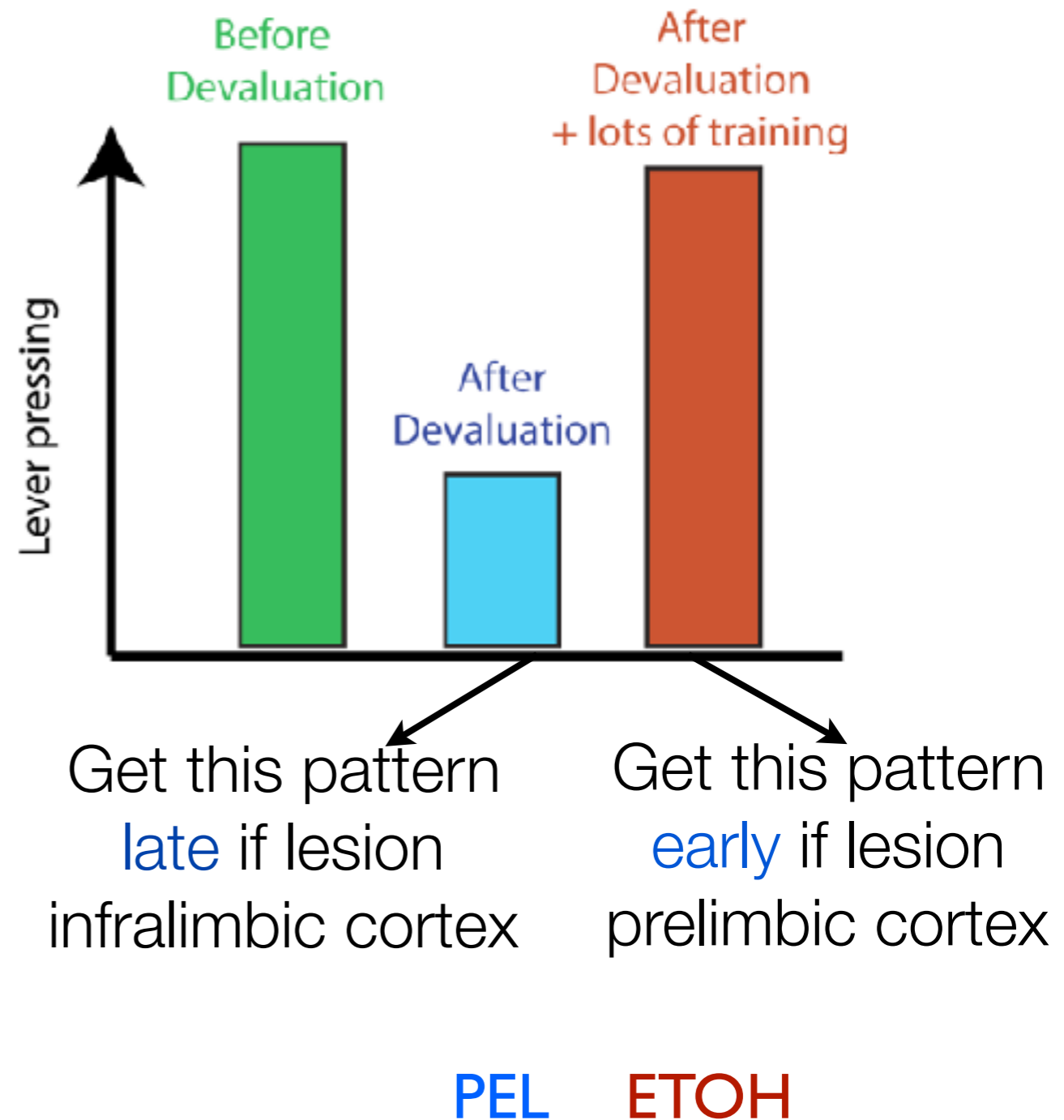
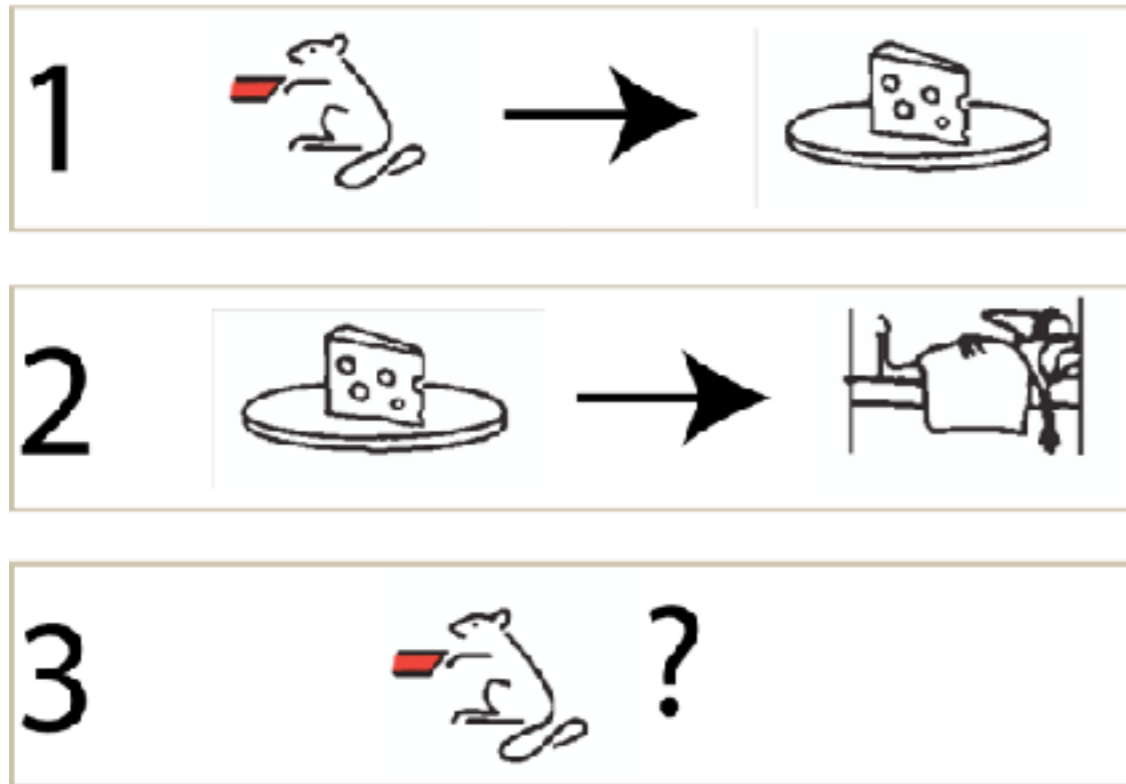




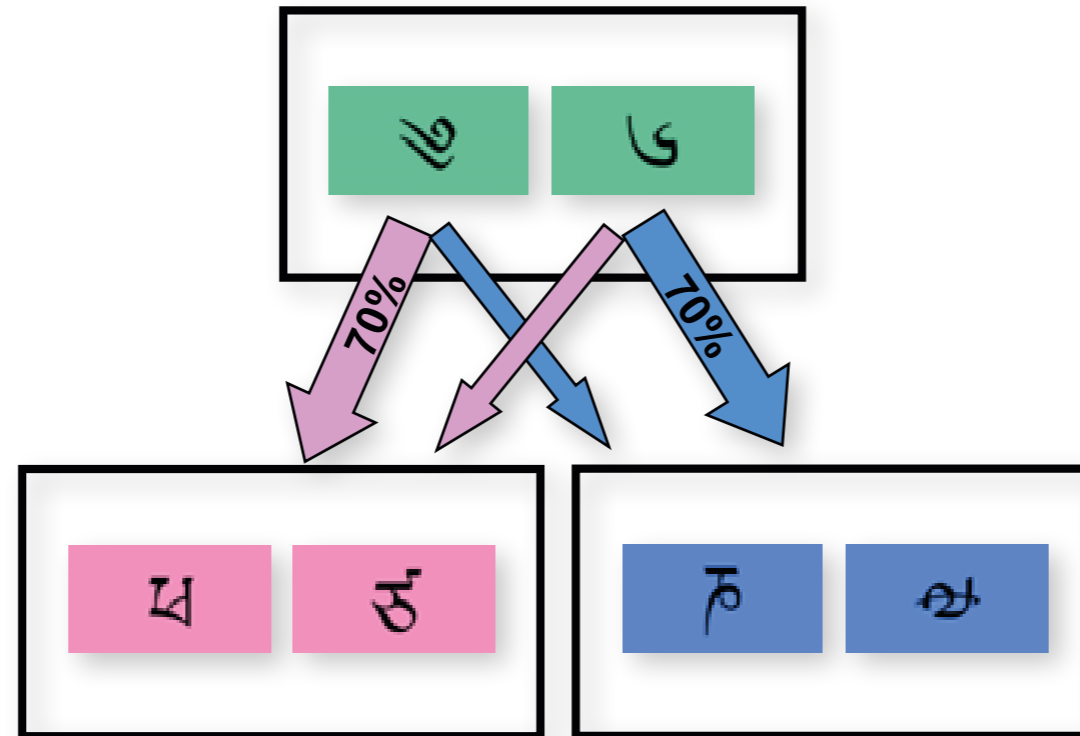


Get this pattern  
late if lesion  
infralimbic cortex

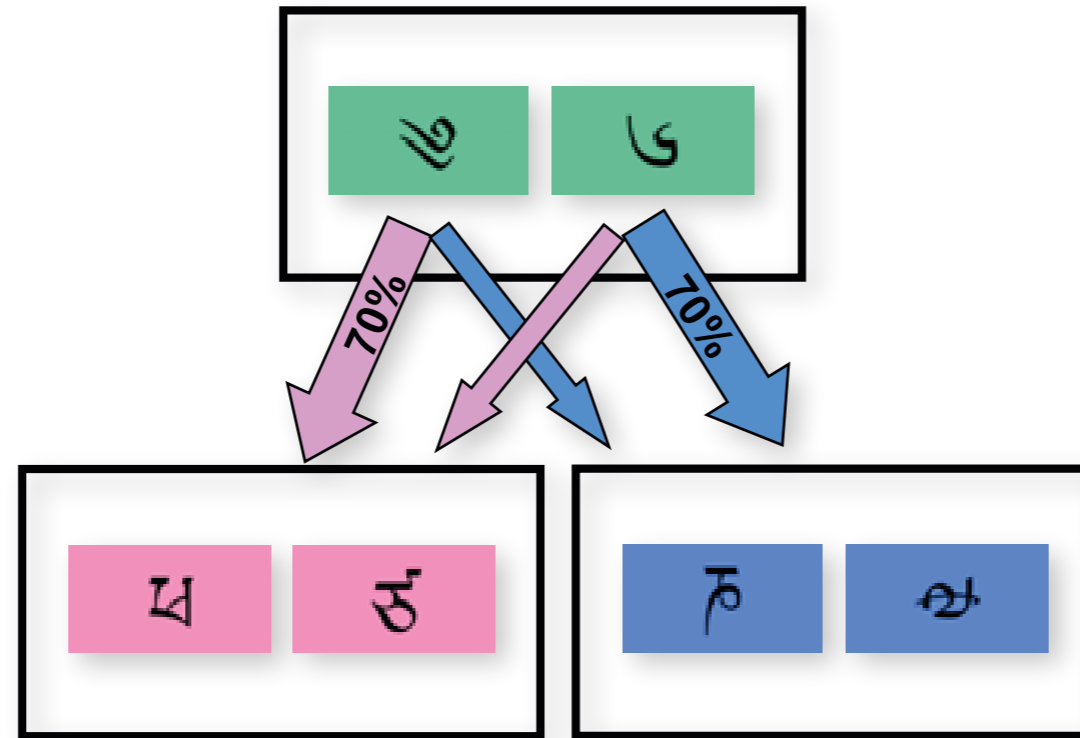




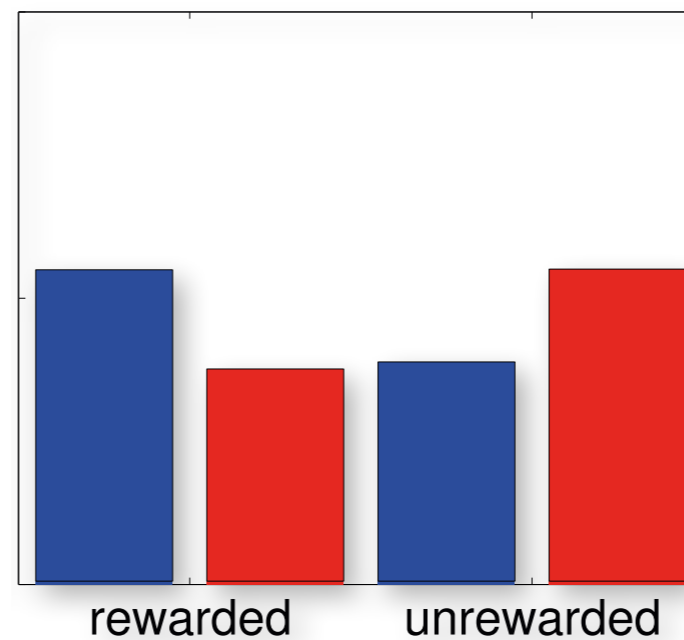
# Two-step task



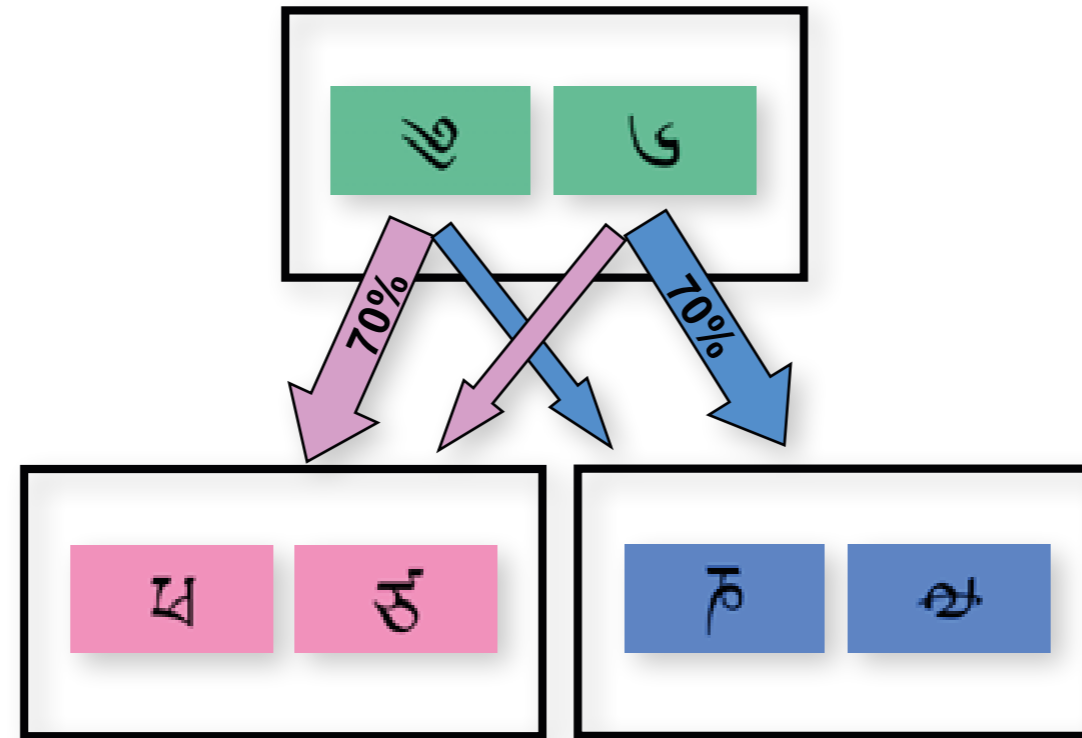
# Two-step task



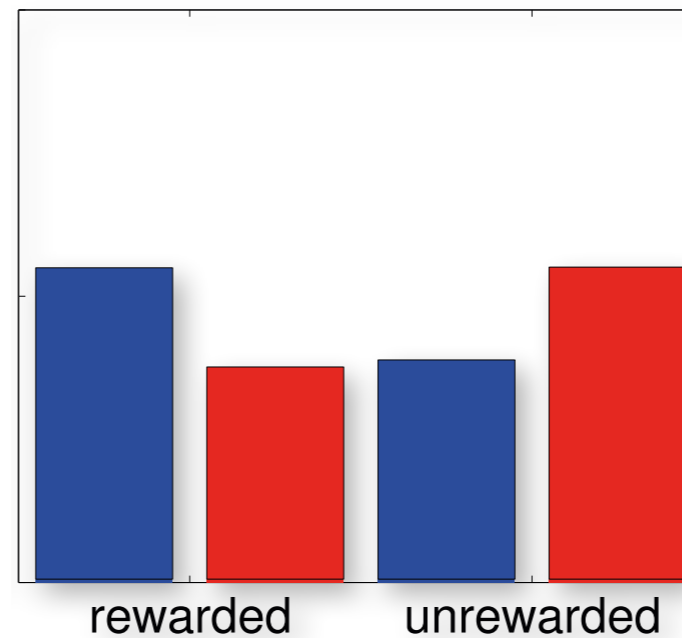
**B** model-based



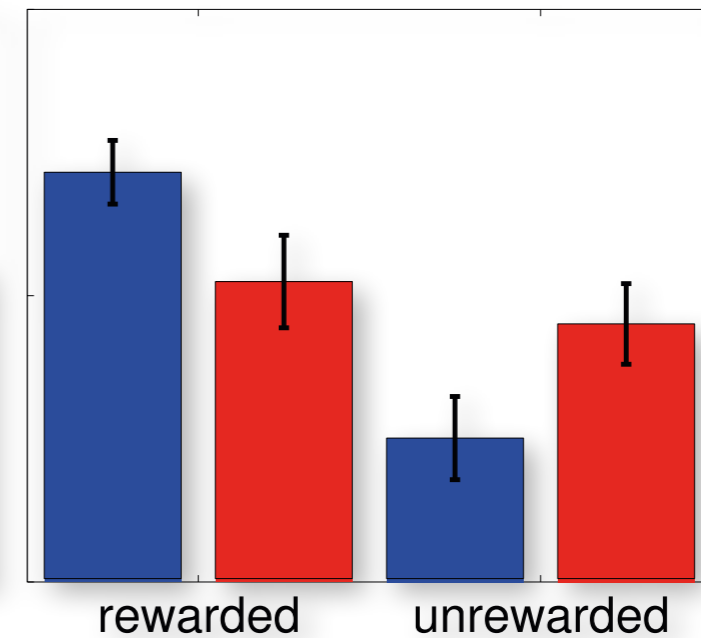
# Two-step task



**B** model-based

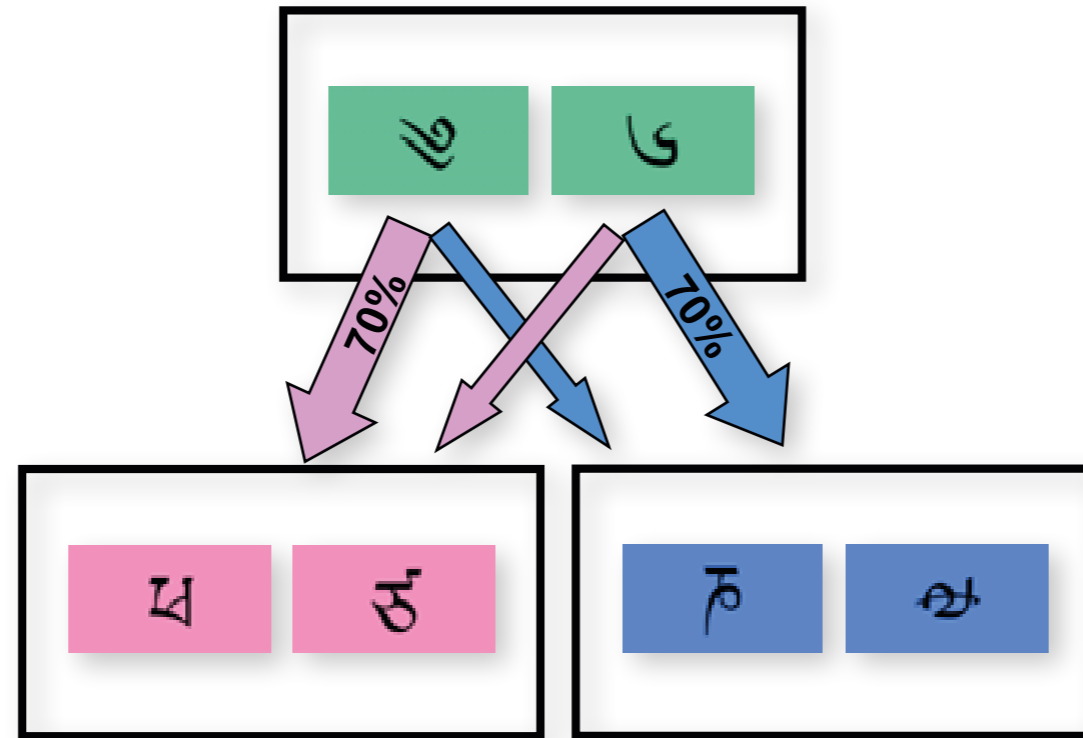


**C** data





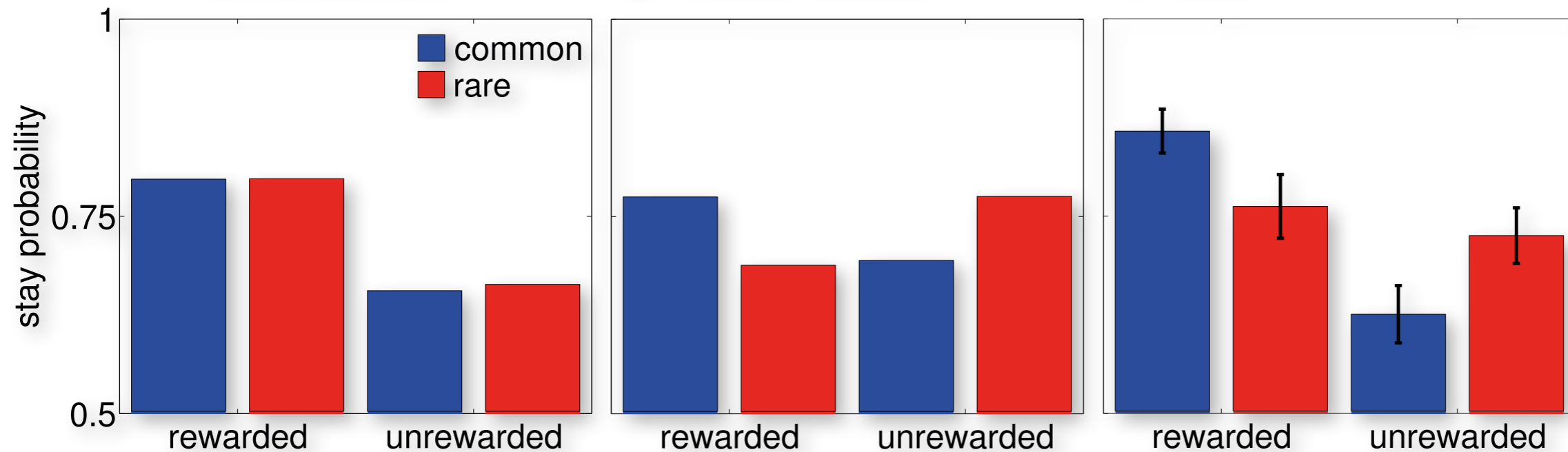
# Two-step task



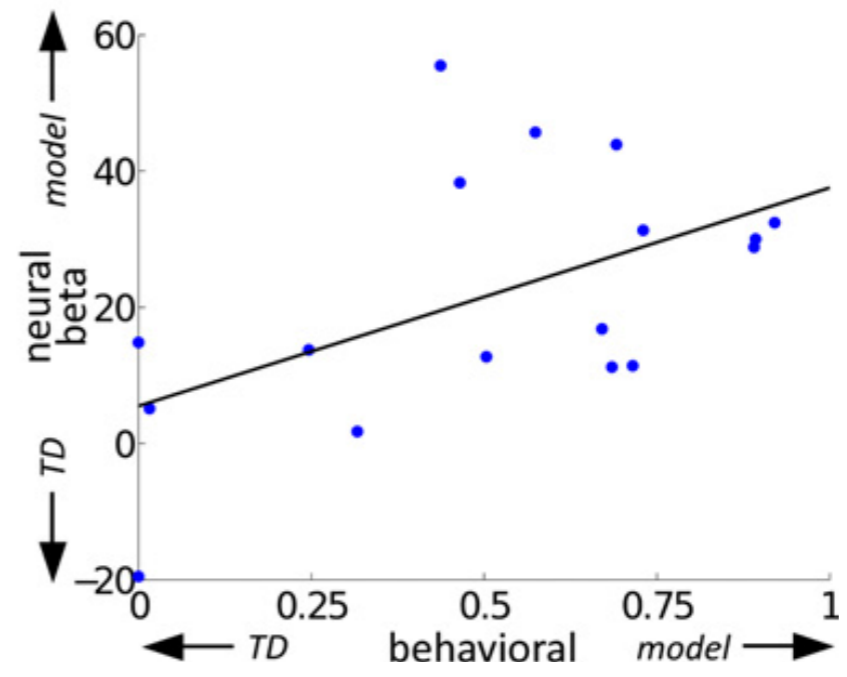
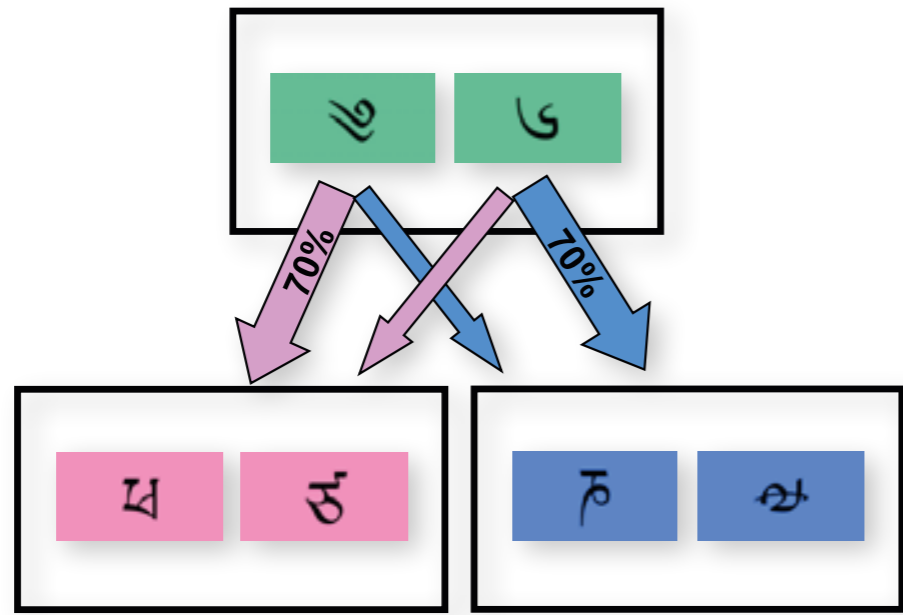
**A** reinforcement

**B** model-based

**C** data



# Fault line 1: Balance of cached and goal



**A** reinforcement

**B** model-based

**C** data

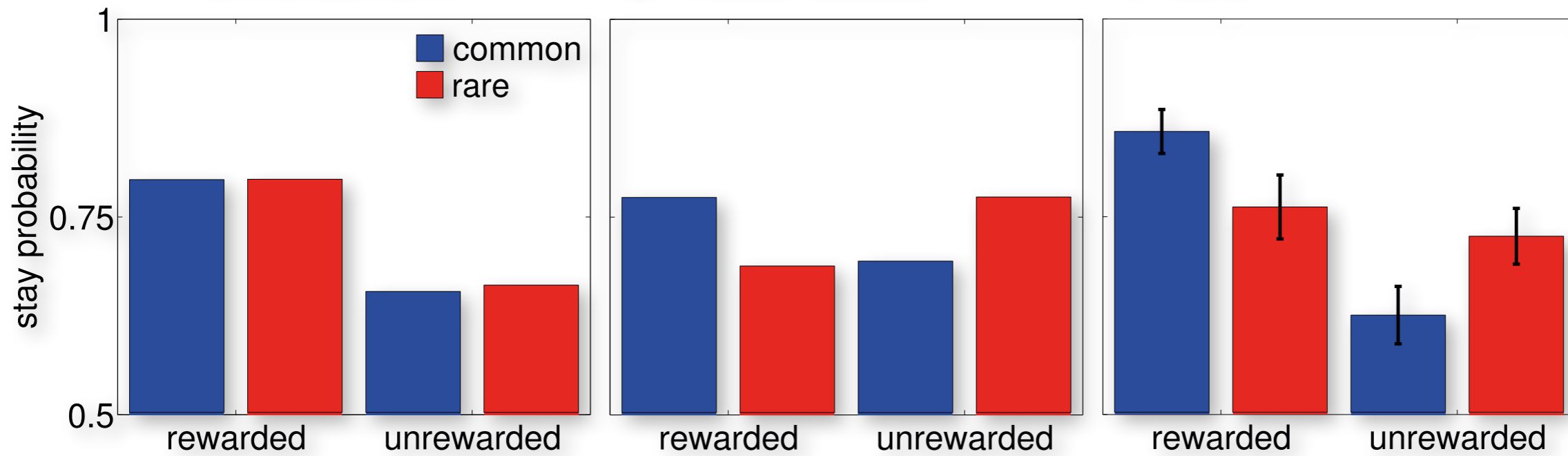
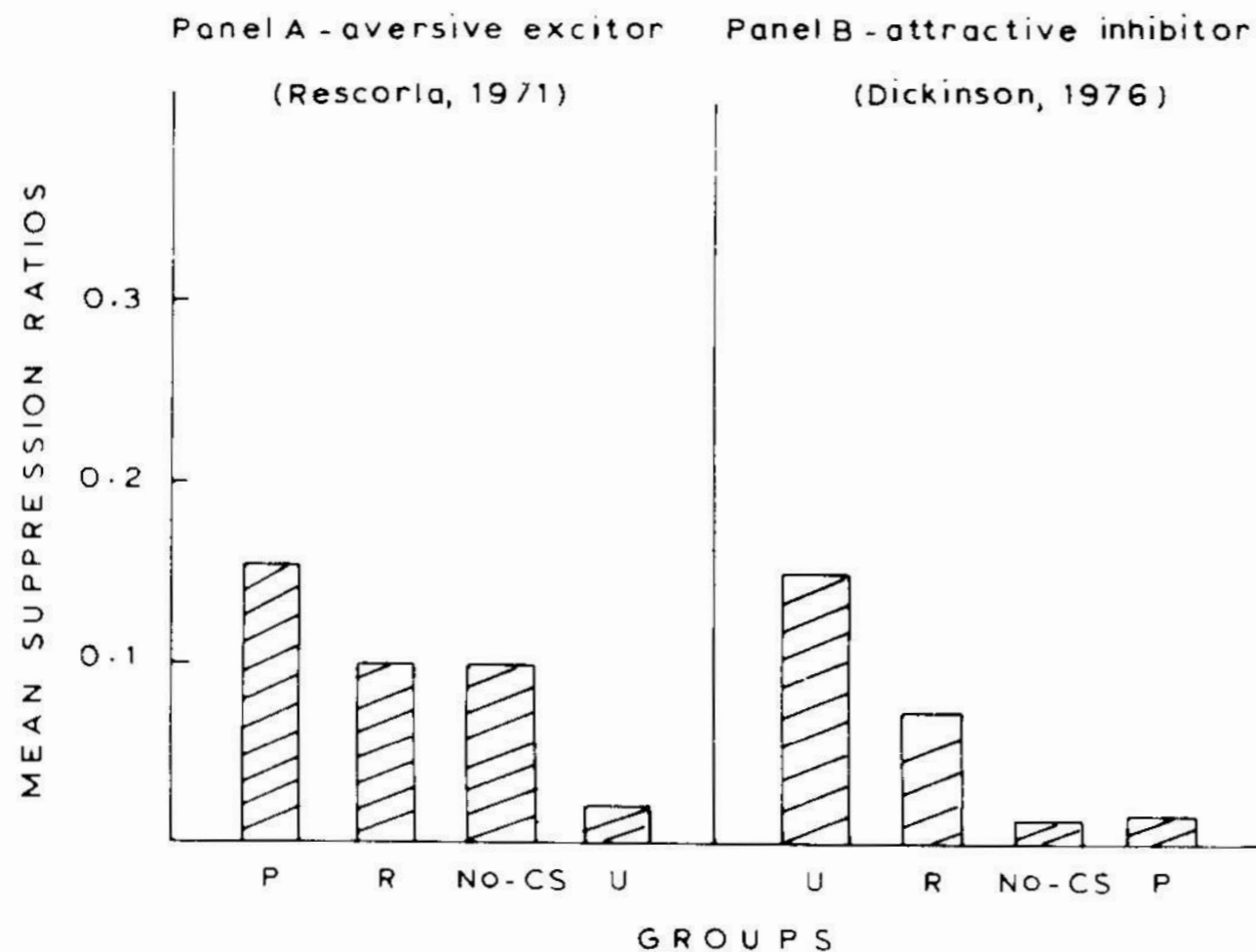


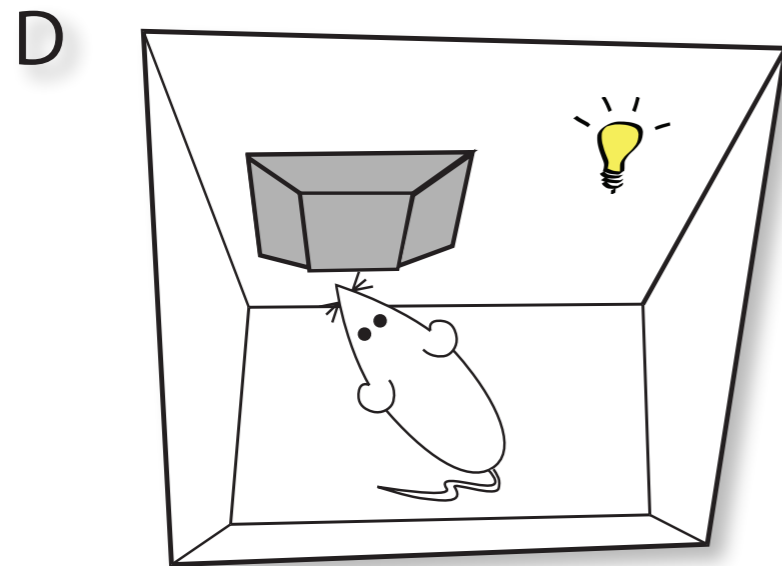
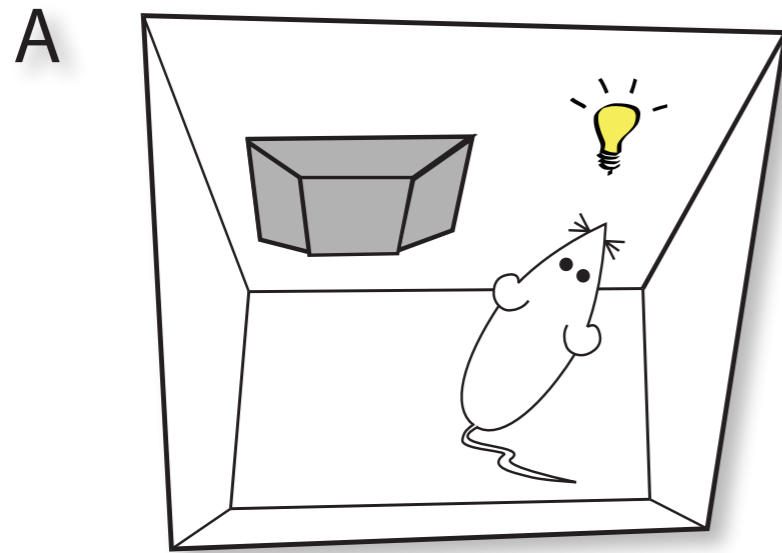
TABLE 8.4  
Blocking of aversive conditioning

Conditions	Stage 1	Stage 2	test
1	A → shock	AX → shock	X
2	control treatments	AX → shock	X
3	A → food omission	AX → shock	X

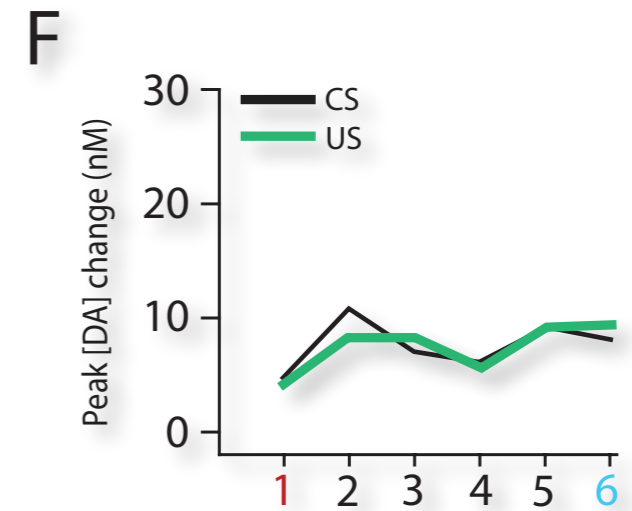
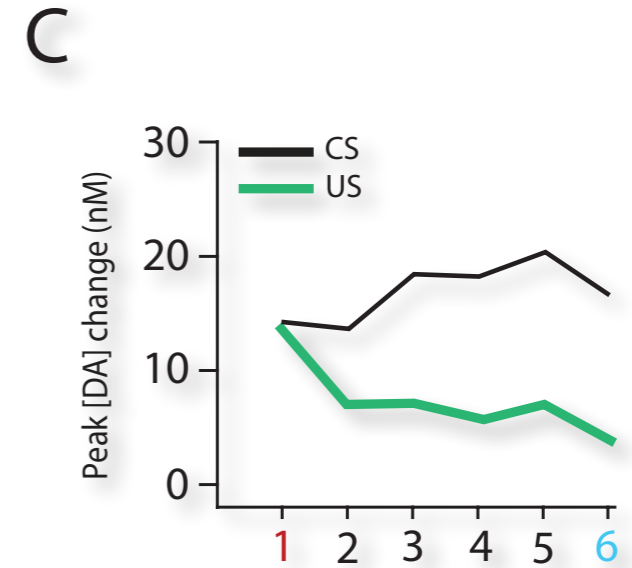
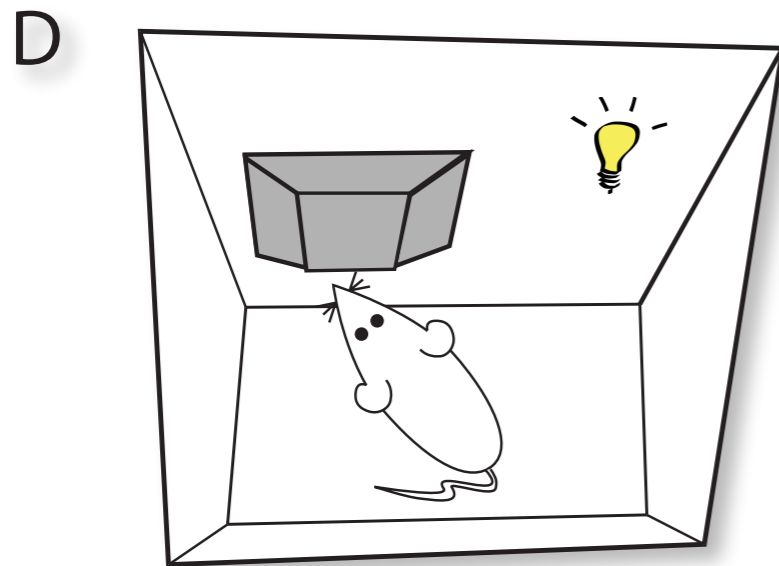
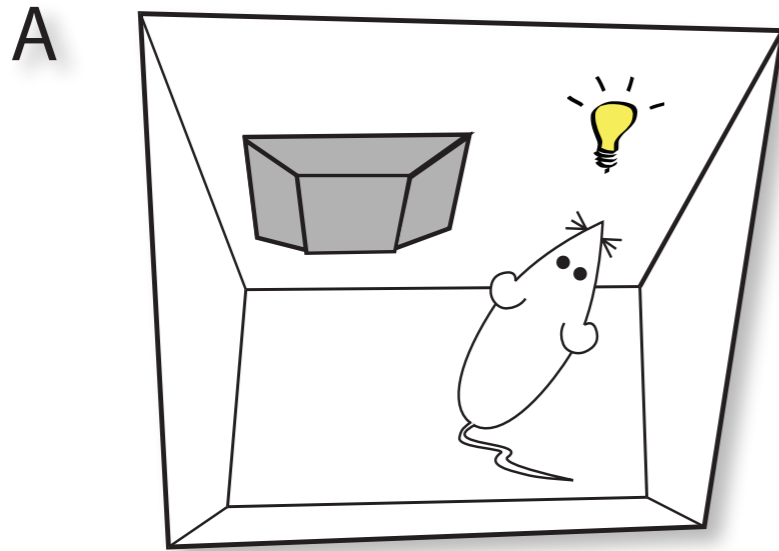


“bad”  
vs  
“good”

Dickinson and Dearing 1979

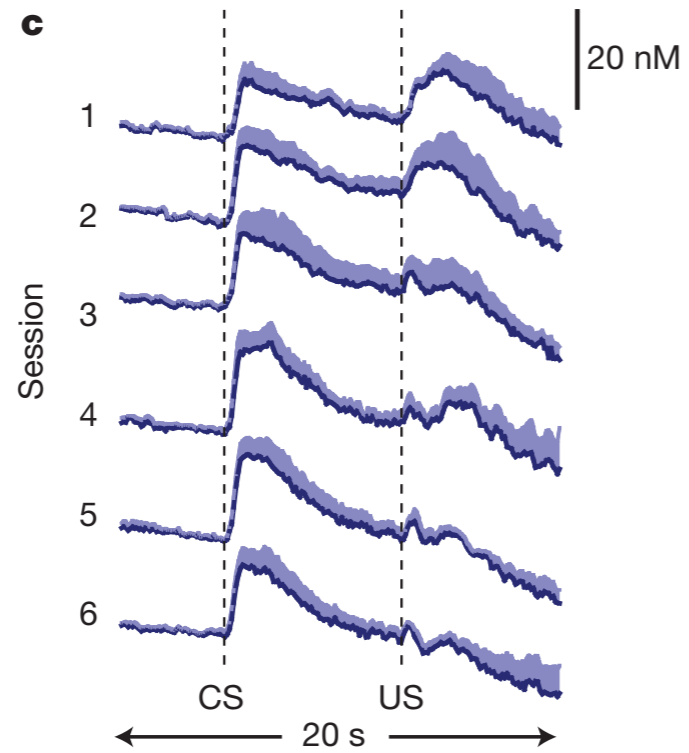


Flagel et al., 2011 Nature, Huys et al., 2014 Prog. Neurobiol.

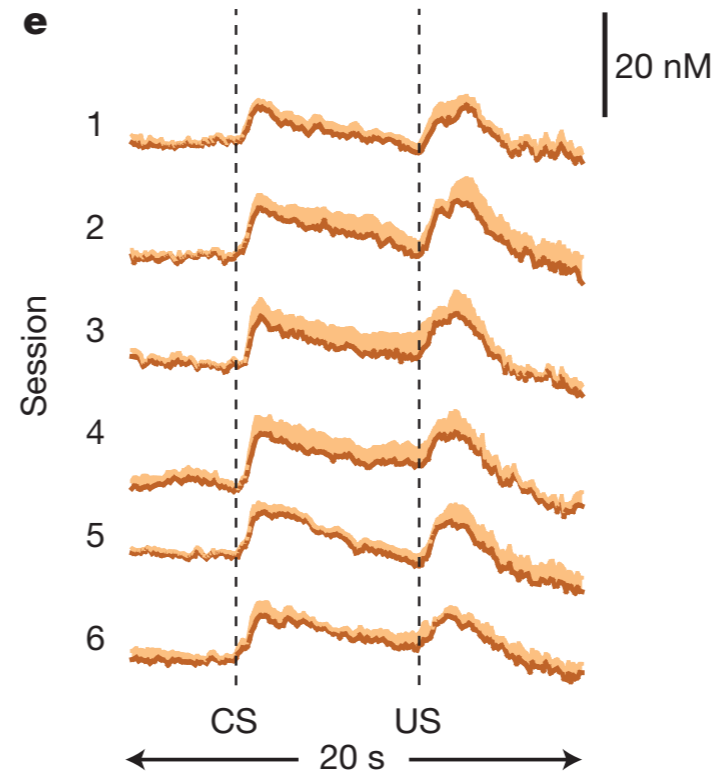


Flagel et al., 2011 Nature, Huys et al., 2014 Prog. Neurobiol.

## Sign trackers



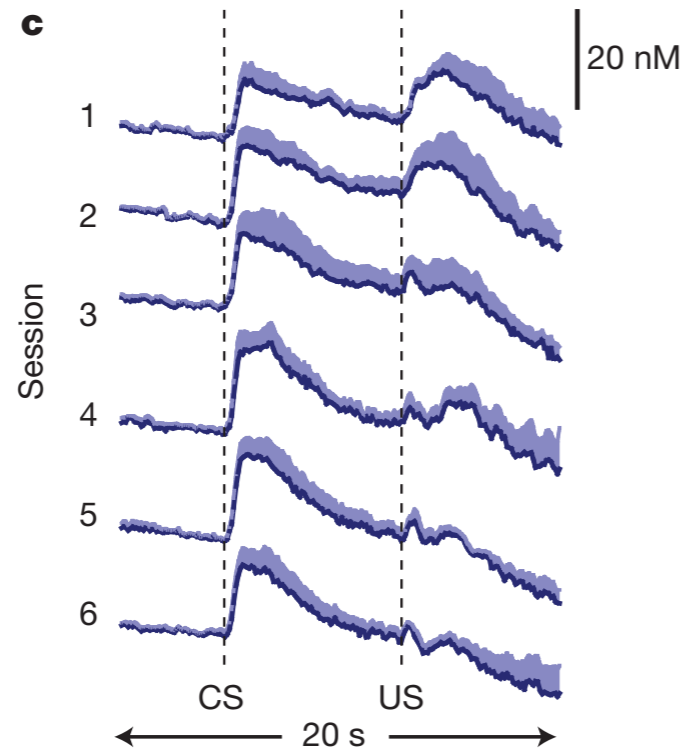
## Goal trackers



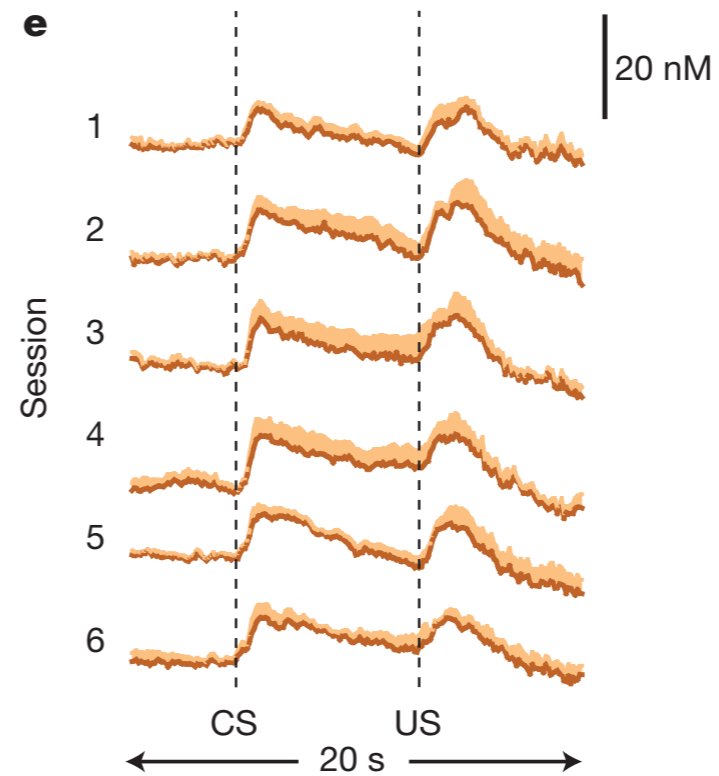
$$\delta = r - Q$$

# Absent model?

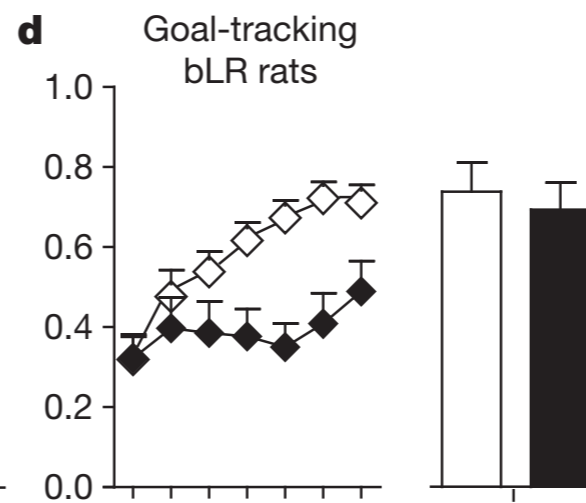
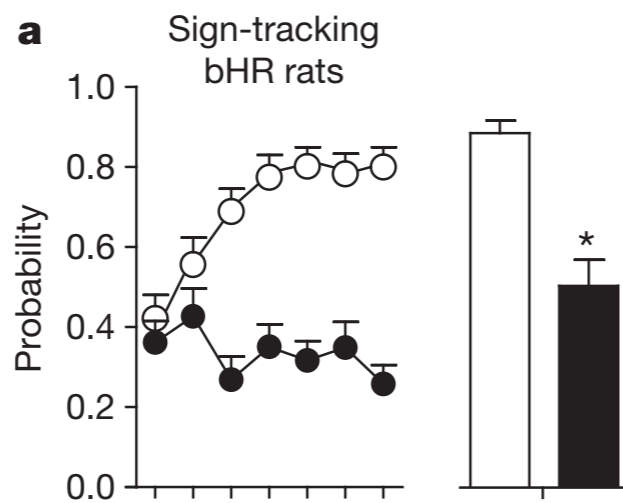
## Sign trackers



## Goal trackers

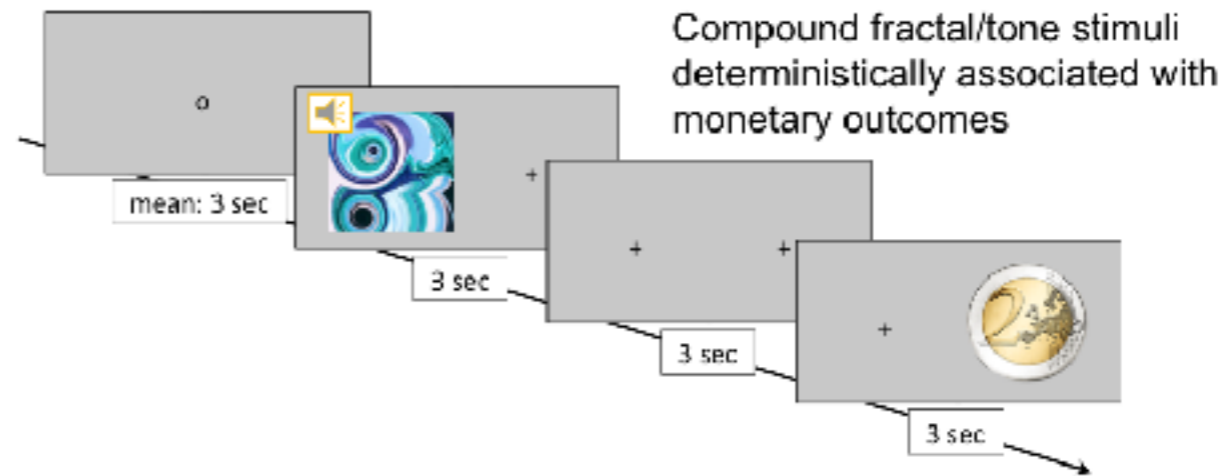


$$\delta = r - Q$$



# Sign-tracking in humans?

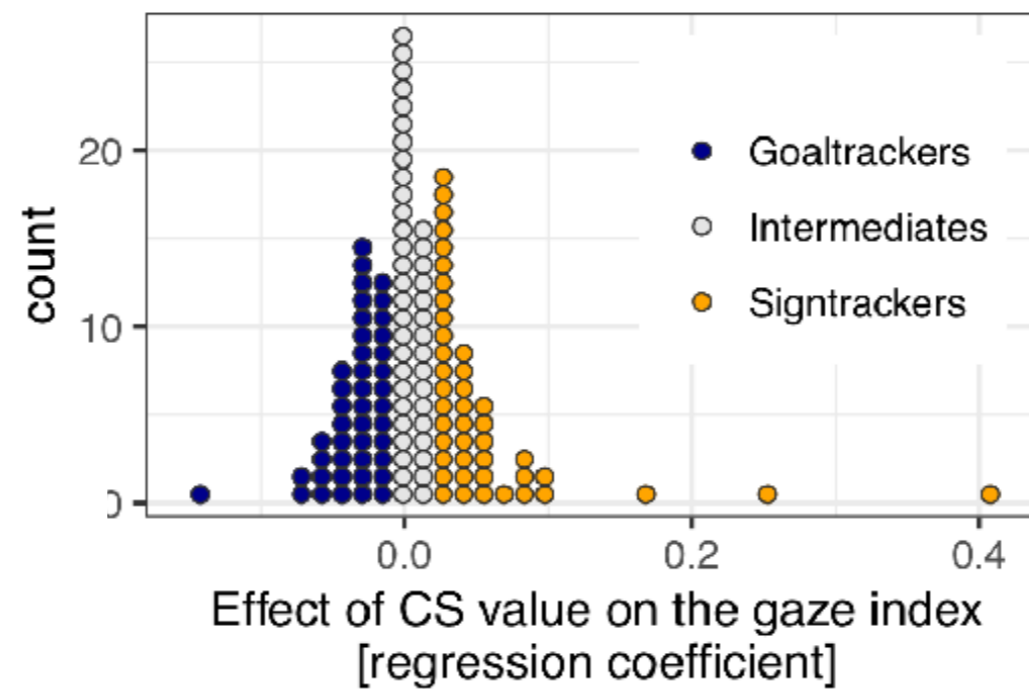
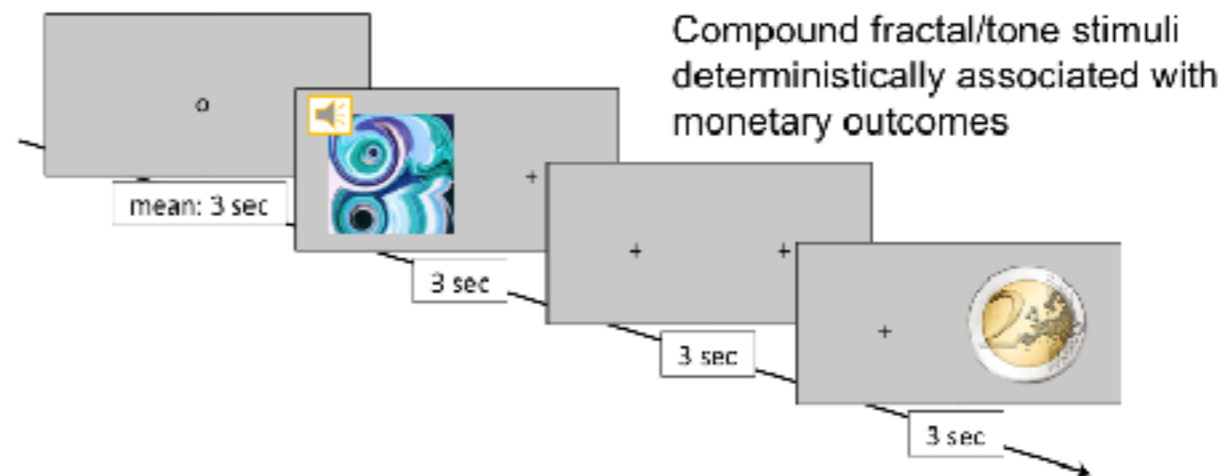
Experimental Paradigm  
Pavlovian Conditioning



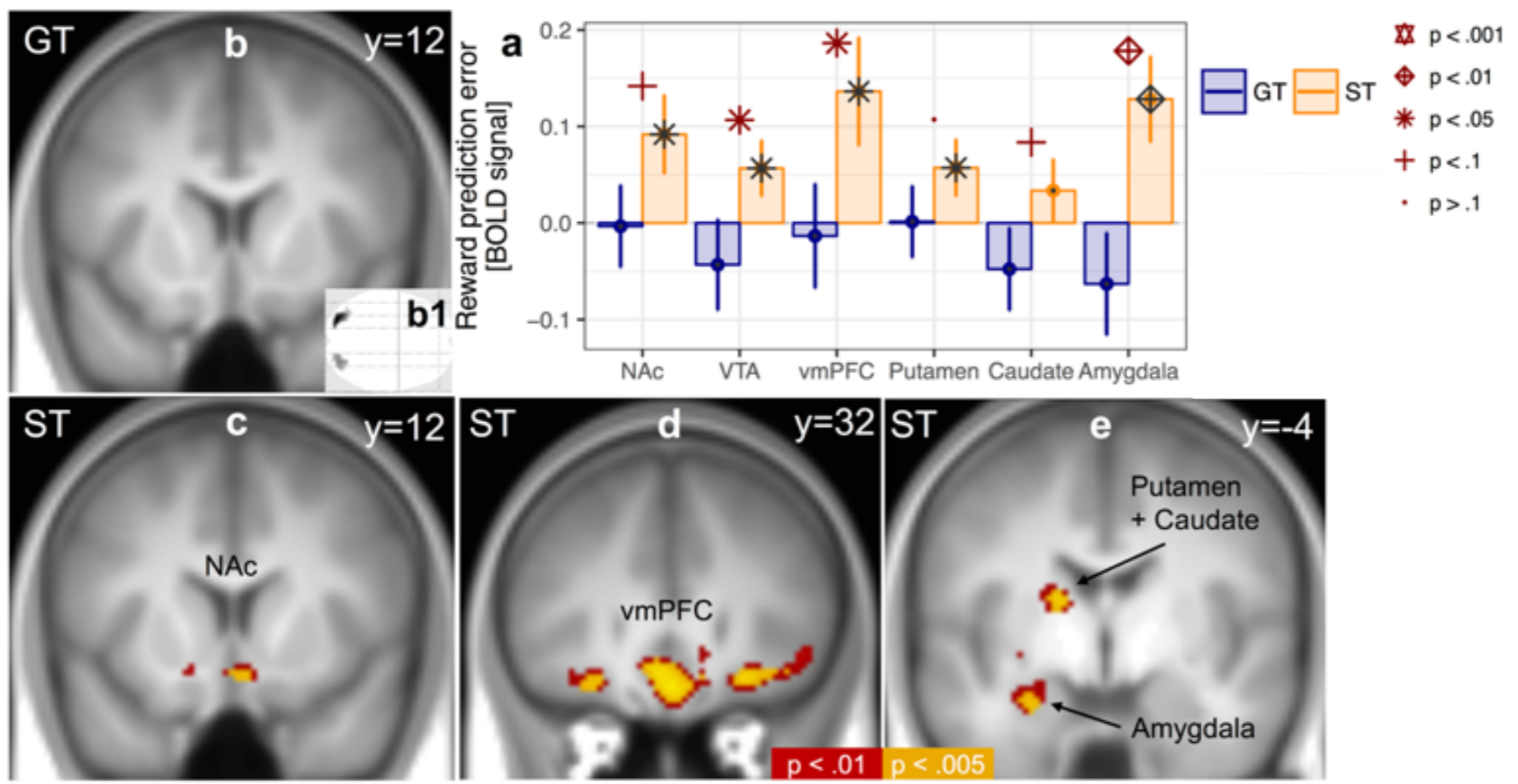


# Sign-tracking in humans?

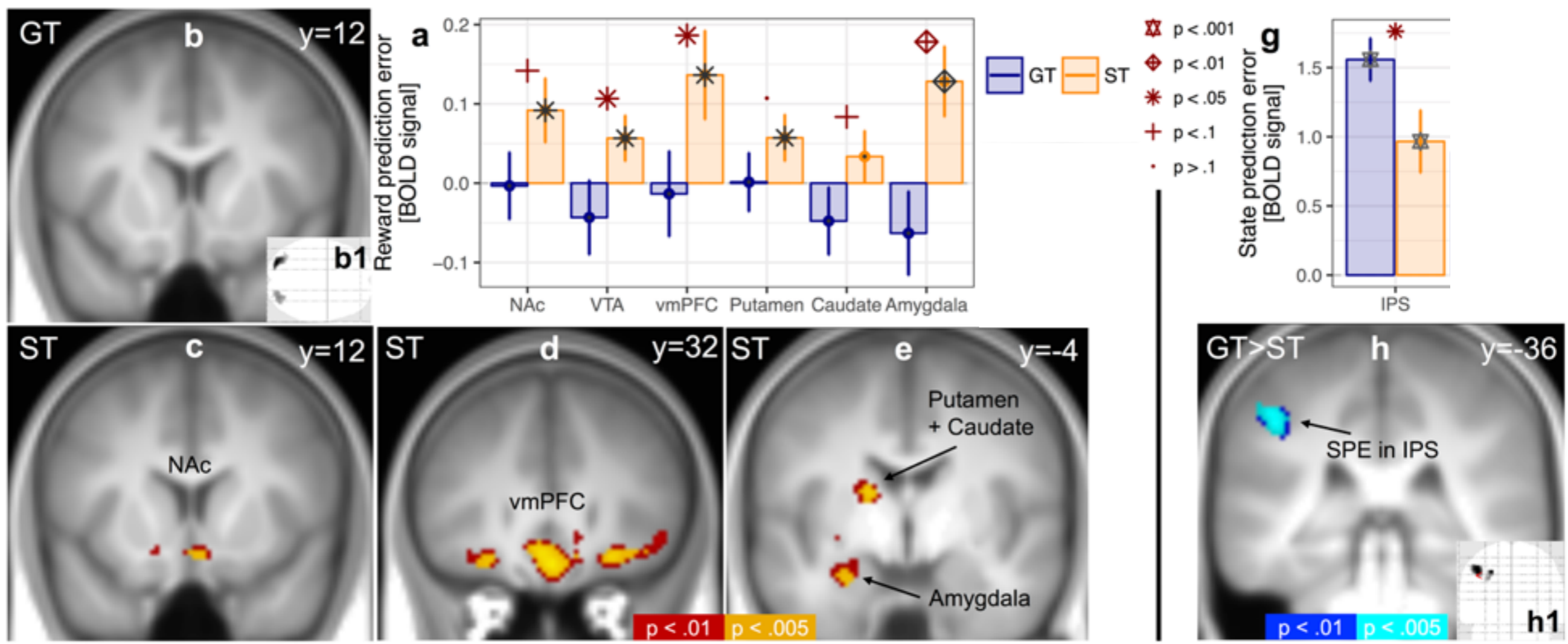
Experimental Paradigm  
Pavlovian Conditioning

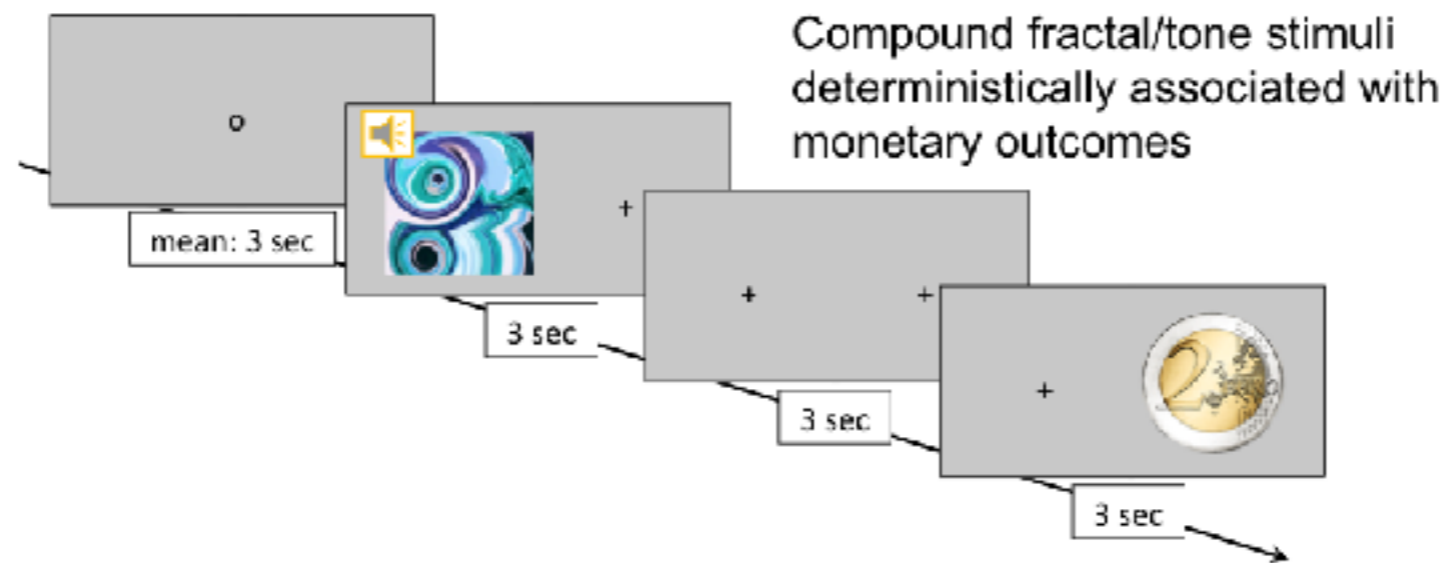


# Double dissociation between ST and GT



# Double dissociation between ST and GT



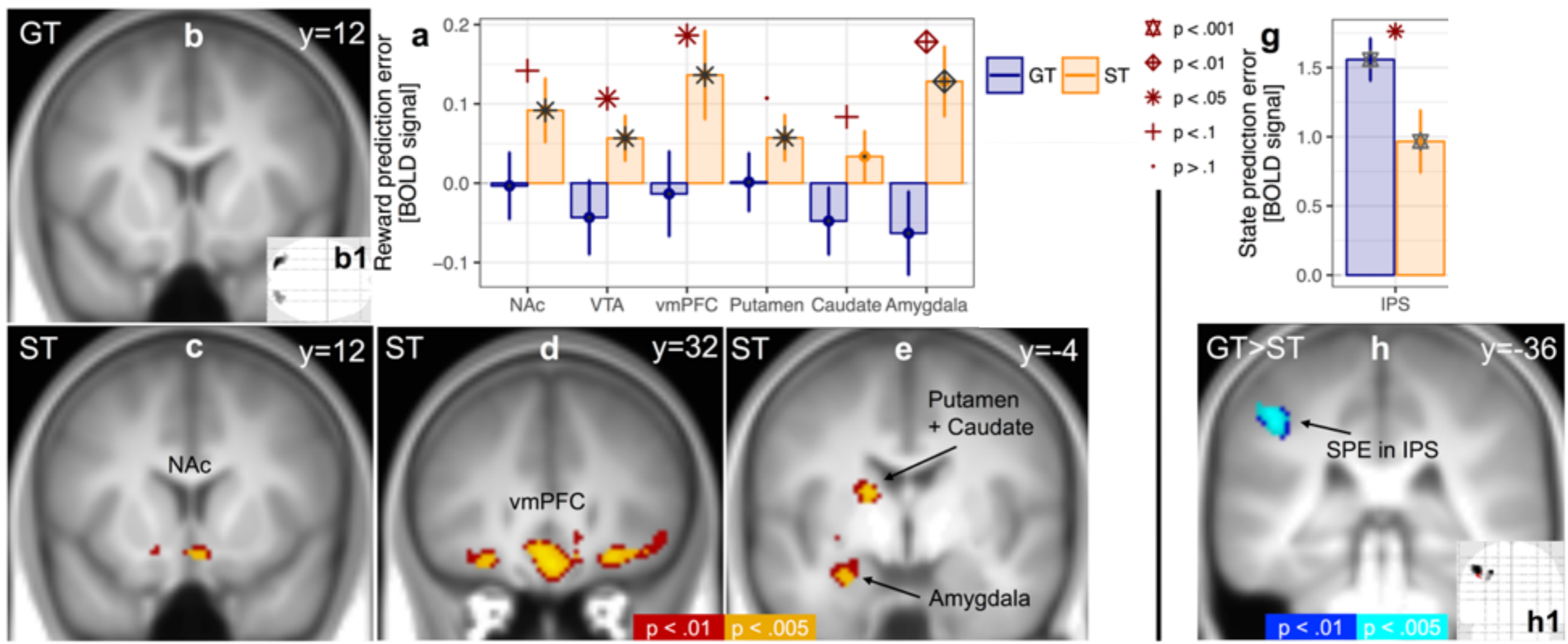


ST: learn expected value  $V$

GT: learn mappings  $T$  from CS to US identity

$$\mathcal{V}(s) = \sum_a \pi(a; s) \sum_{s'} \mathcal{T}(s' | s, a) [\mathcal{R}(s', a, s) + \mathcal{V}(s')]$$

# Double dissociation between ST and GT



$$V^\pi(s) = \sum_a \pi(a|s) \sum_{s'} \mathcal{T}_{ss'}^a [\mathcal{R}_{ss'}^a + V(s')]$$

$$\mathbf{v}^\pi = \mathbf{R}^\pi + \mathbf{T}^\pi \mathbf{v}^\pi$$

$$\mathbf{v}^\pi = (\mathbf{I} - \mathbf{T}^\pi)^{-1} \mathbf{R}^\pi$$

$$\hat{\mathbf{v}} = \mathbf{M}\mathbf{w}$$

$$M^\pi(s, :) = \mathbf{1}_s + \gamma \sum_{s'} T^\pi(s, s') M^\pi(s', :),$$

$$M^\pi(s, :) \leftarrow M^\pi(s, :) + \alpha_{SR} [\mathbf{1}_s + \gamma M^\pi(s', :) - M^\pi(s, :)],$$

$$M^\pi = (I - \gamma T^\pi)^{-1}$$



## ► “Model-free learning”

$$M^\pi(s, :) = \mathbf{1}_s + \gamma \sum_{s'} T^\pi(s, s') M^\pi(s', :),$$

$$M^\pi(s, :) \leftarrow M^\pi(s, :) + \alpha_{SR} [\mathbf{1}_s + \gamma M^\pi(s', :) - M^\pi(s, :)],$$

$$M^\pi = (I - \gamma T^\pi)^{-1}$$

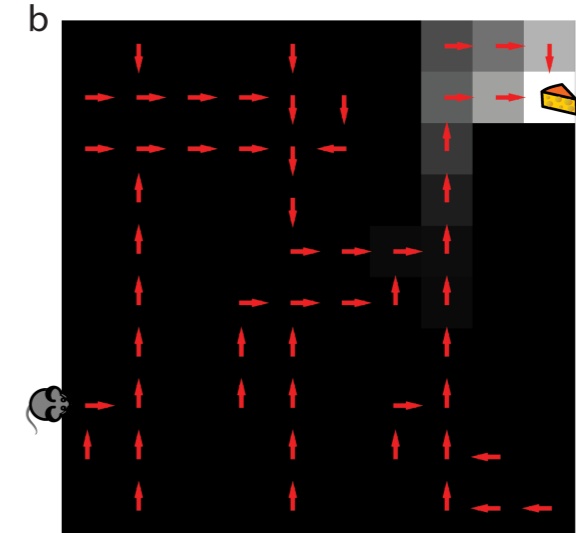


## ► “Model-free learning”

$$M^\pi(s, :) = \mathbf{1}_s + \gamma \sum_{s'} T^\pi(s, s') M^\pi(s', :),$$

$$M^\pi(s, :) \leftarrow M^\pi(s, :) + \alpha_{SR} [\mathbf{1}_s + \gamma M^\pi(s', :) - M^\pi(s, :)],$$

$$M^\pi = (I - \gamma T^\pi)^{-1}$$

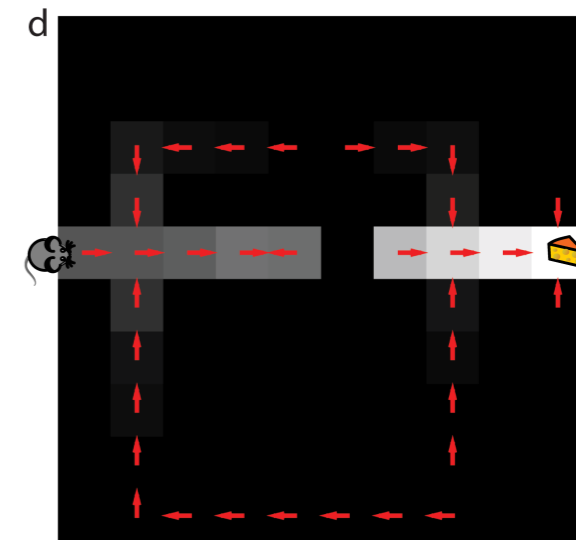
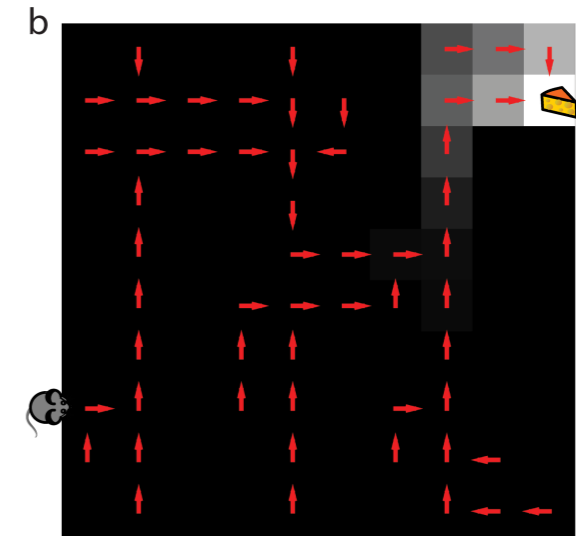


► “Model-free learning”

$$M^\pi(s, :) = \mathbf{1}_s + \gamma \sum_{s'} T^\pi(s, s') M^\pi(s', :),$$

$$M^\pi(s, :) \leftarrow M^\pi(s, :) + \alpha_{SR} [\mathbf{1}_s + \gamma M^\pi(s', :) - M^\pi(s, :)],$$

$$M^\pi = (I - \gamma T^\pi)^{-1}$$



## ► “Model-free learning”

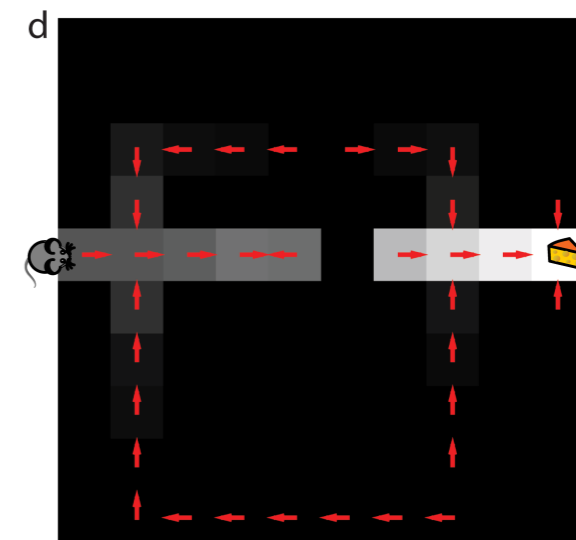
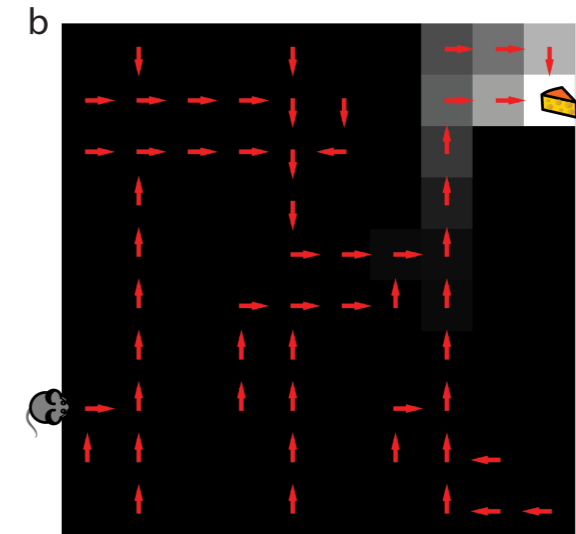
$$M^\pi(s, :) = \mathbf{1}_s + \gamma \sum_{s'} T^\pi(s, s') M^\pi(s', :),$$

$$M^\pi(s, :) \leftarrow M^\pi(s, :) + \alpha_{SR} [\mathbf{1}_s + \gamma M^\pi(s', :) - M^\pi(s, :)],$$

## ► “Model-based learning”

- Estimate transition and compute

$$M^\pi = (I - \gamma T^\pi)^{-1}$$



► “Model-free learning”

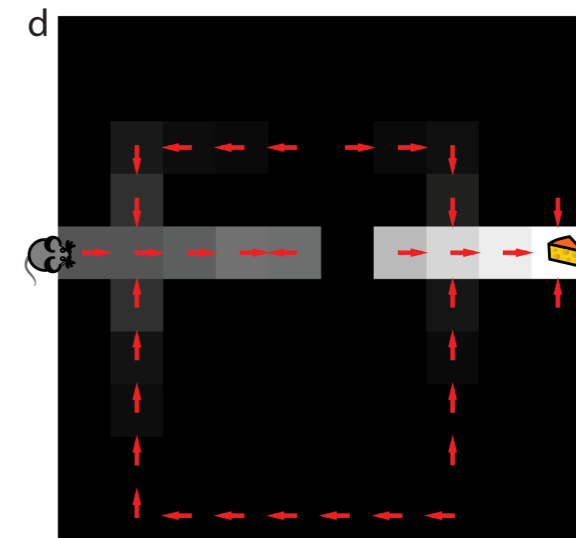
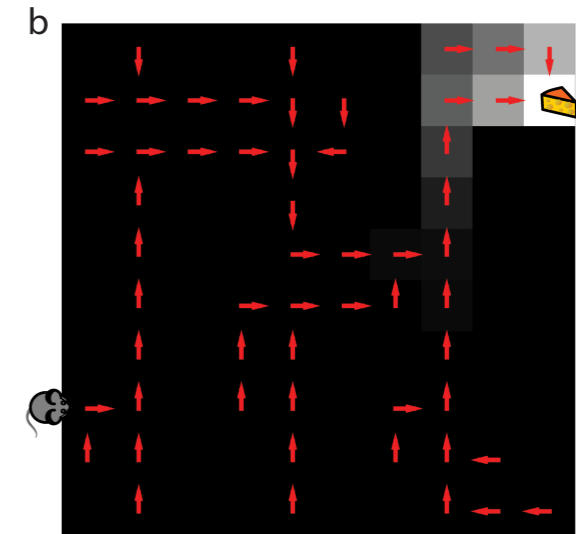
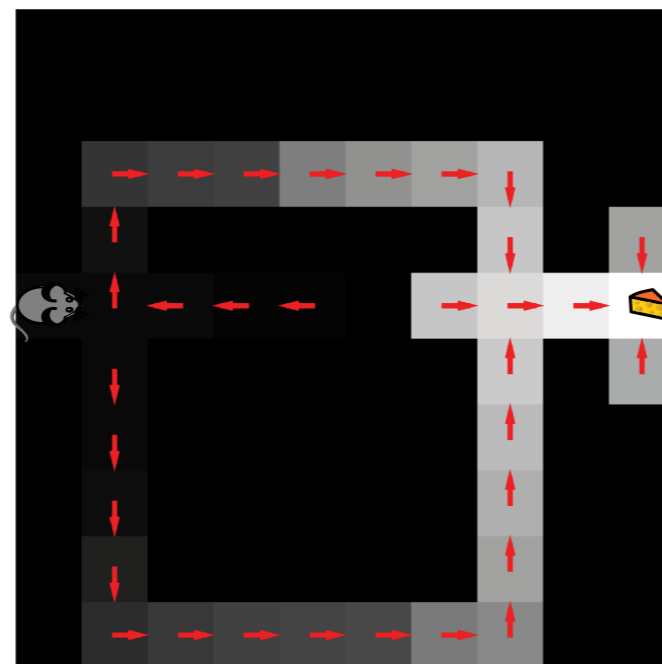
$$M^\pi(s, :) = \mathbf{1}_s + \gamma \sum_{s'} T^\pi(s, s') M^\pi(s', :),$$

$$M^\pi(s, :) \leftarrow M^\pi(s, :) + \alpha_{SR} [\mathbf{1}_s + \gamma M^\pi(s', :) - M^\pi(s, :)],$$

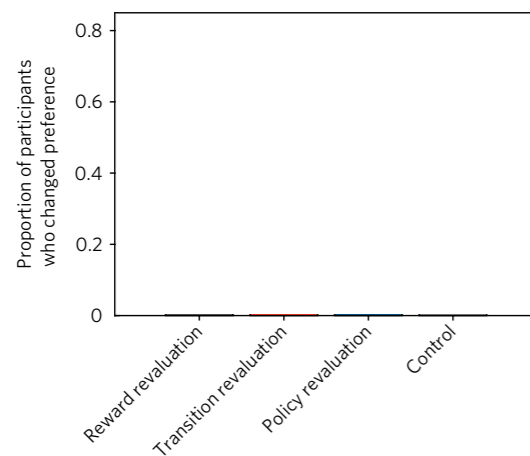
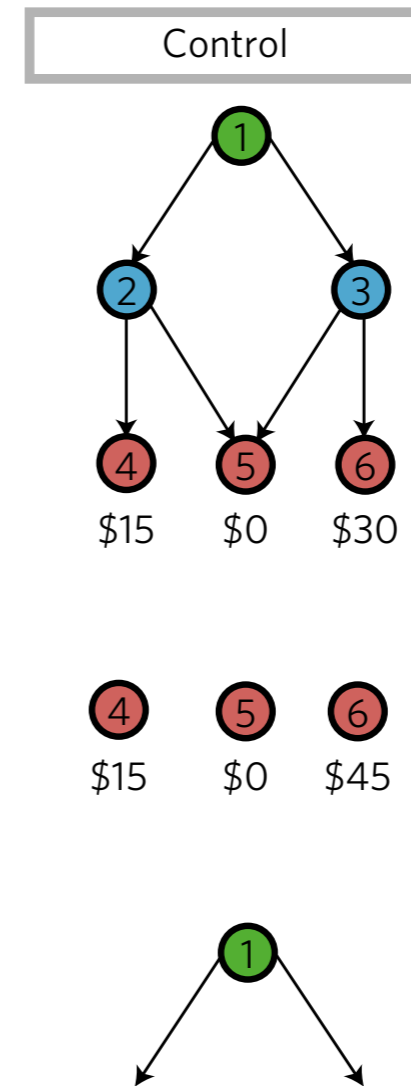
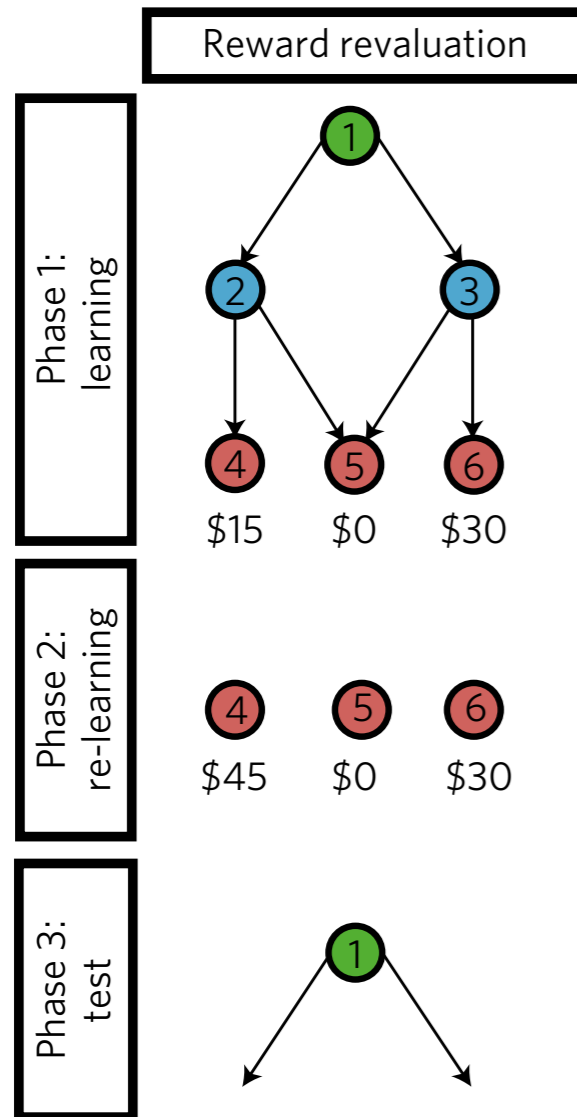
► “Model-based learning”

- Estimate transition and compute

$$M^\pi = (I - \gamma T^\pi)^{-1}$$

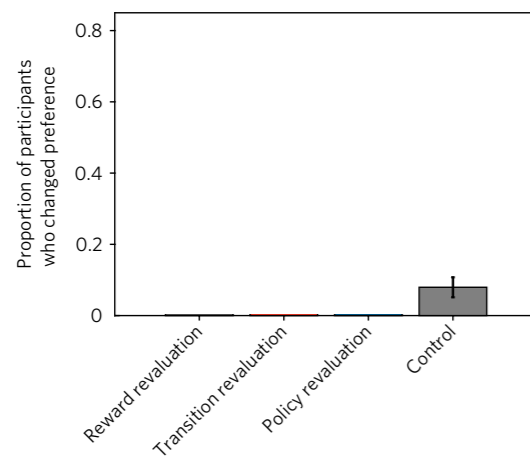
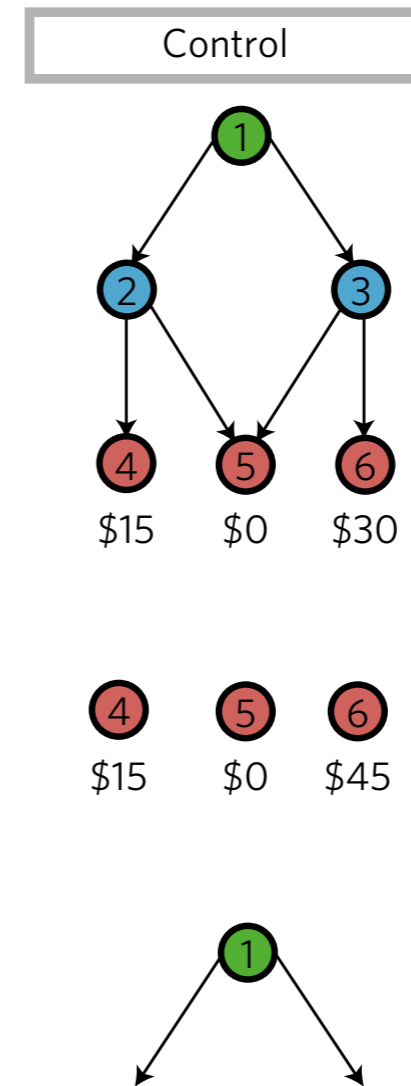
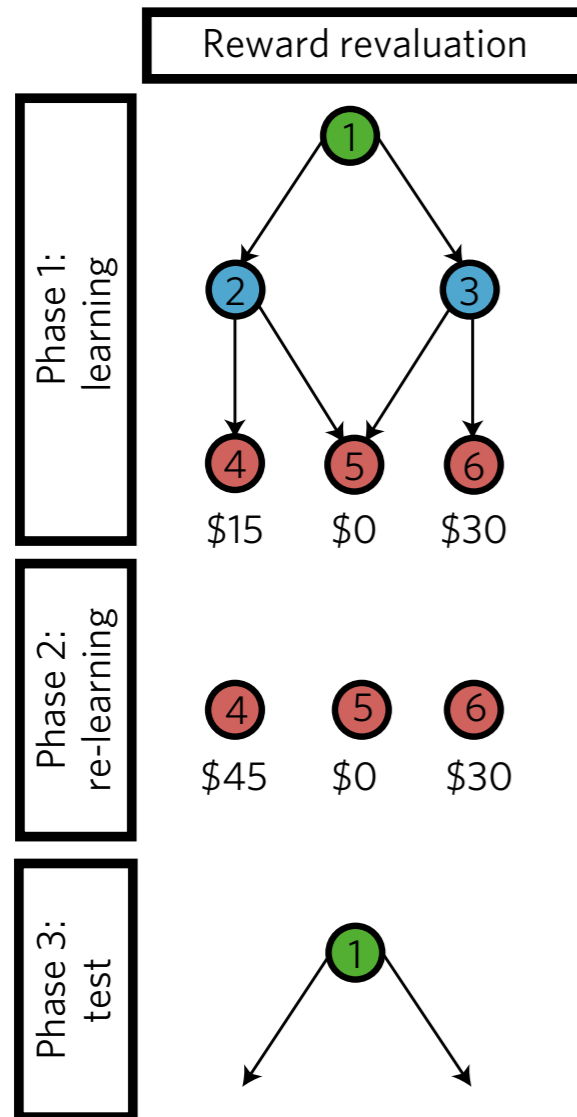


# Human successor learning



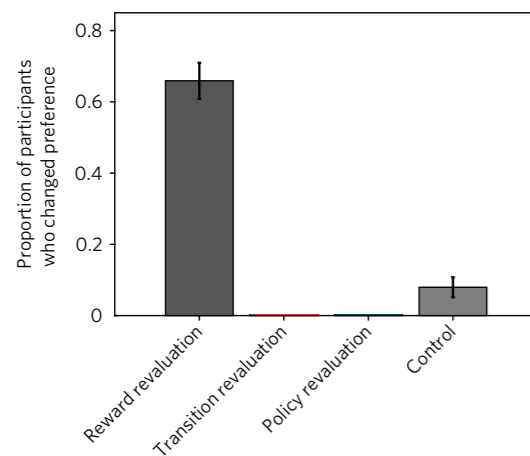
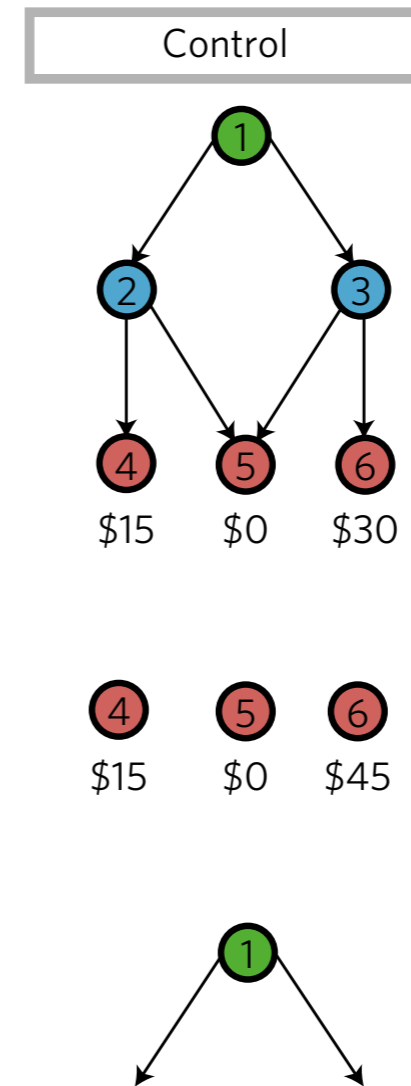
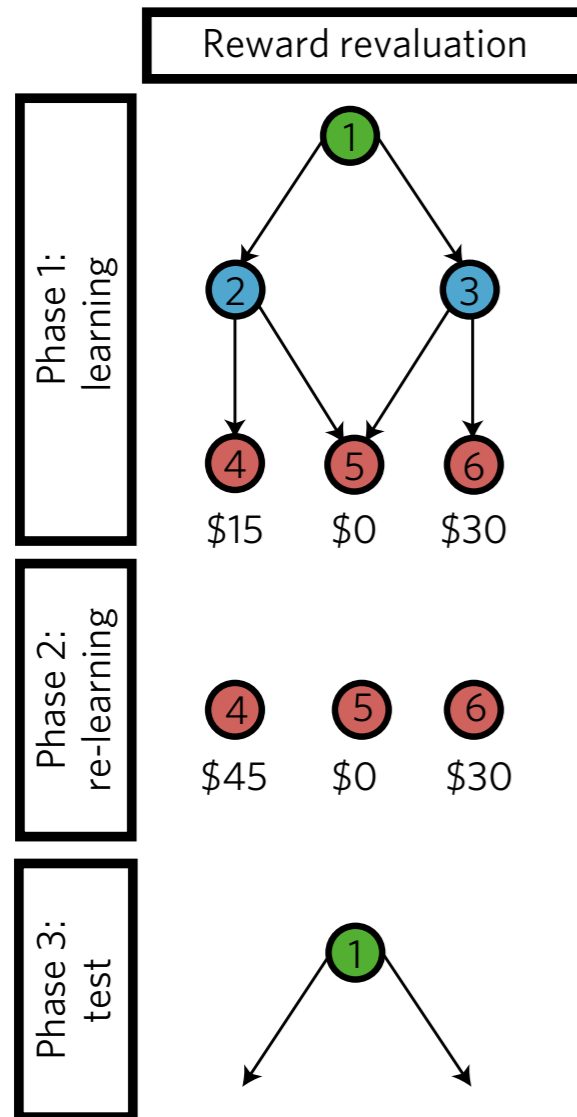
Momennejad et al., 2017 Nat. Hum. Beh.

# Human successor learning



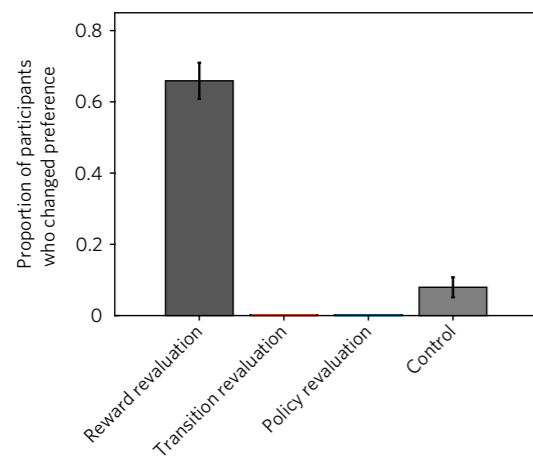
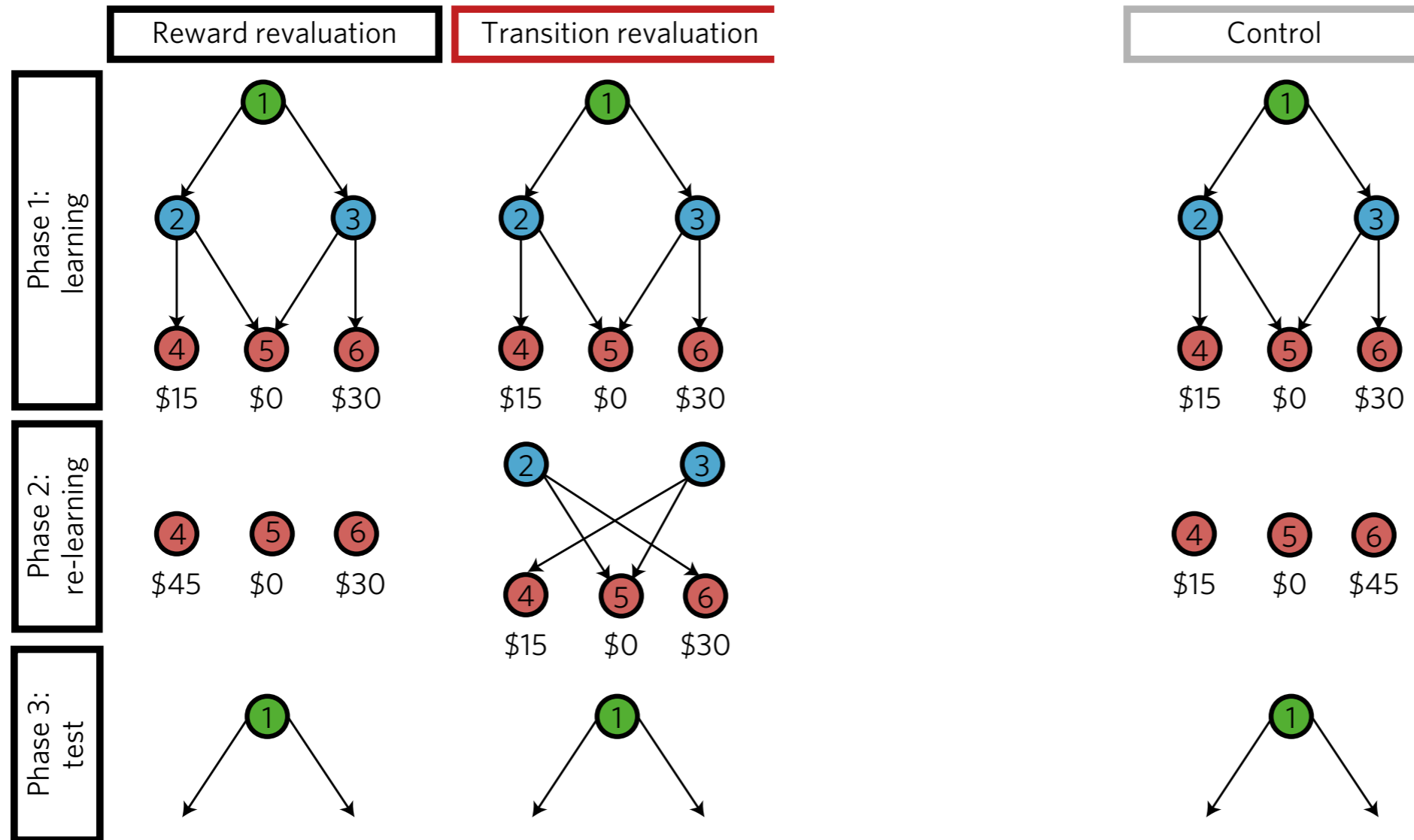
Momennejad et al., 2017 Nat. Hum. Beh.

# Human successor learning



Momennejad et al., 2017 Nat. Hum. Beh.

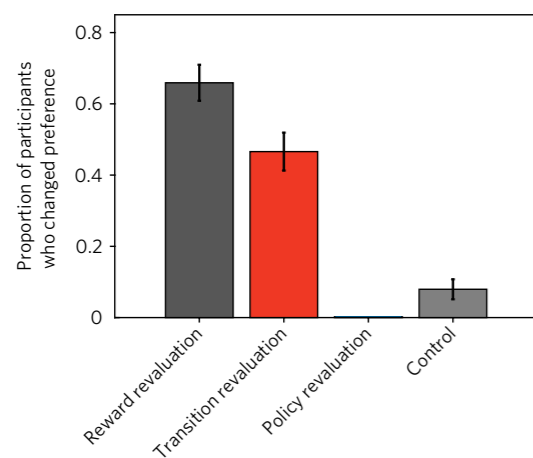
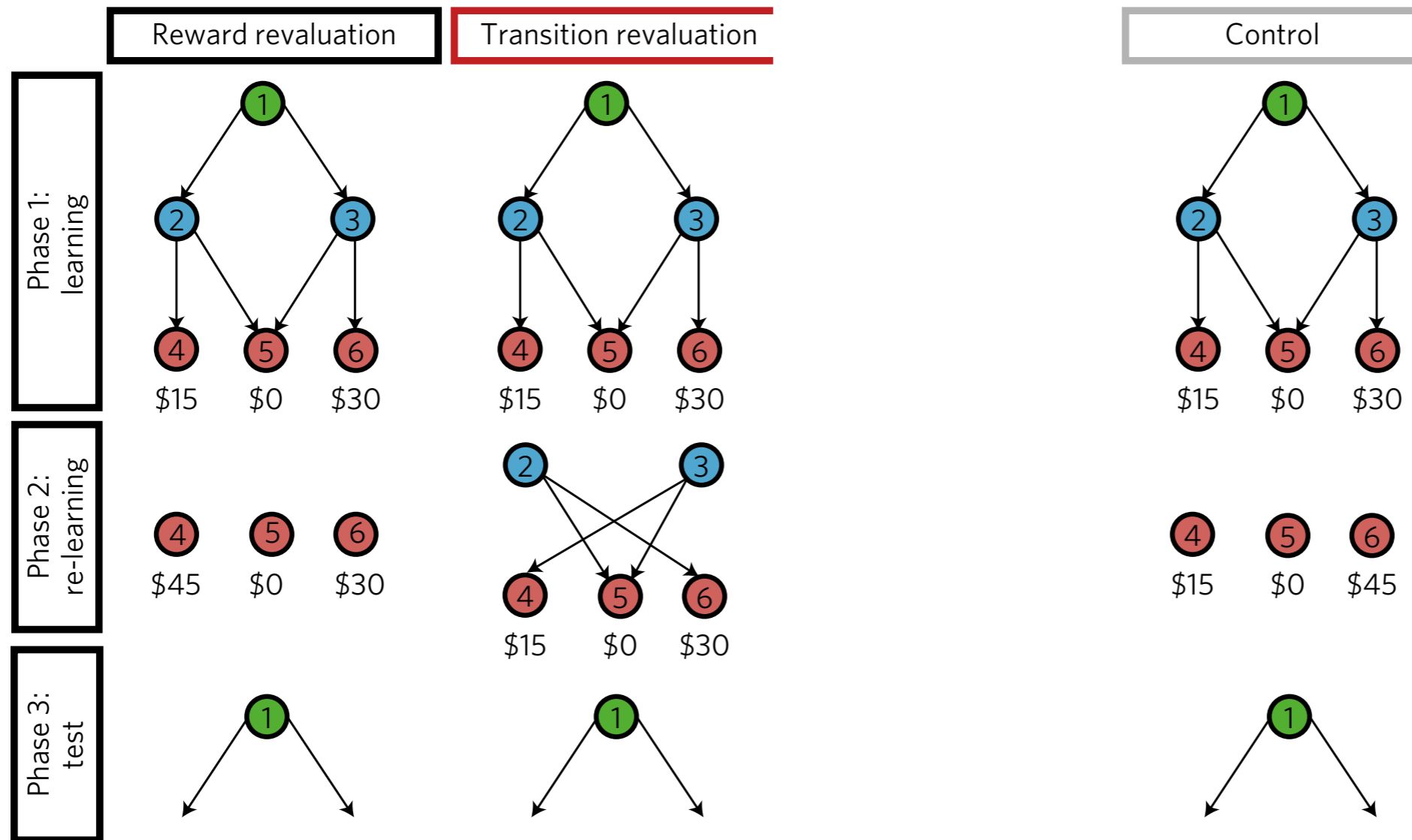
# Human successor learning



Momennejad et al., 2017 Nat. Hum. Beh.

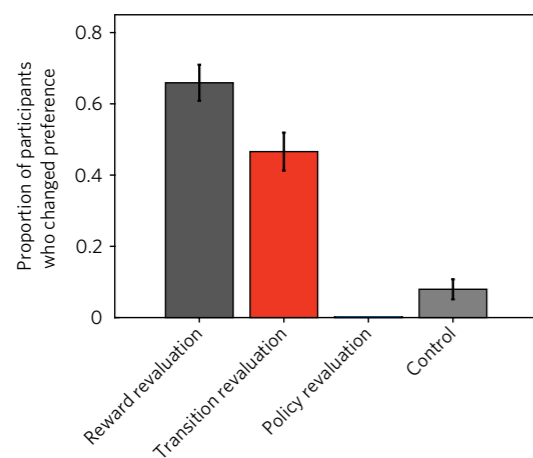
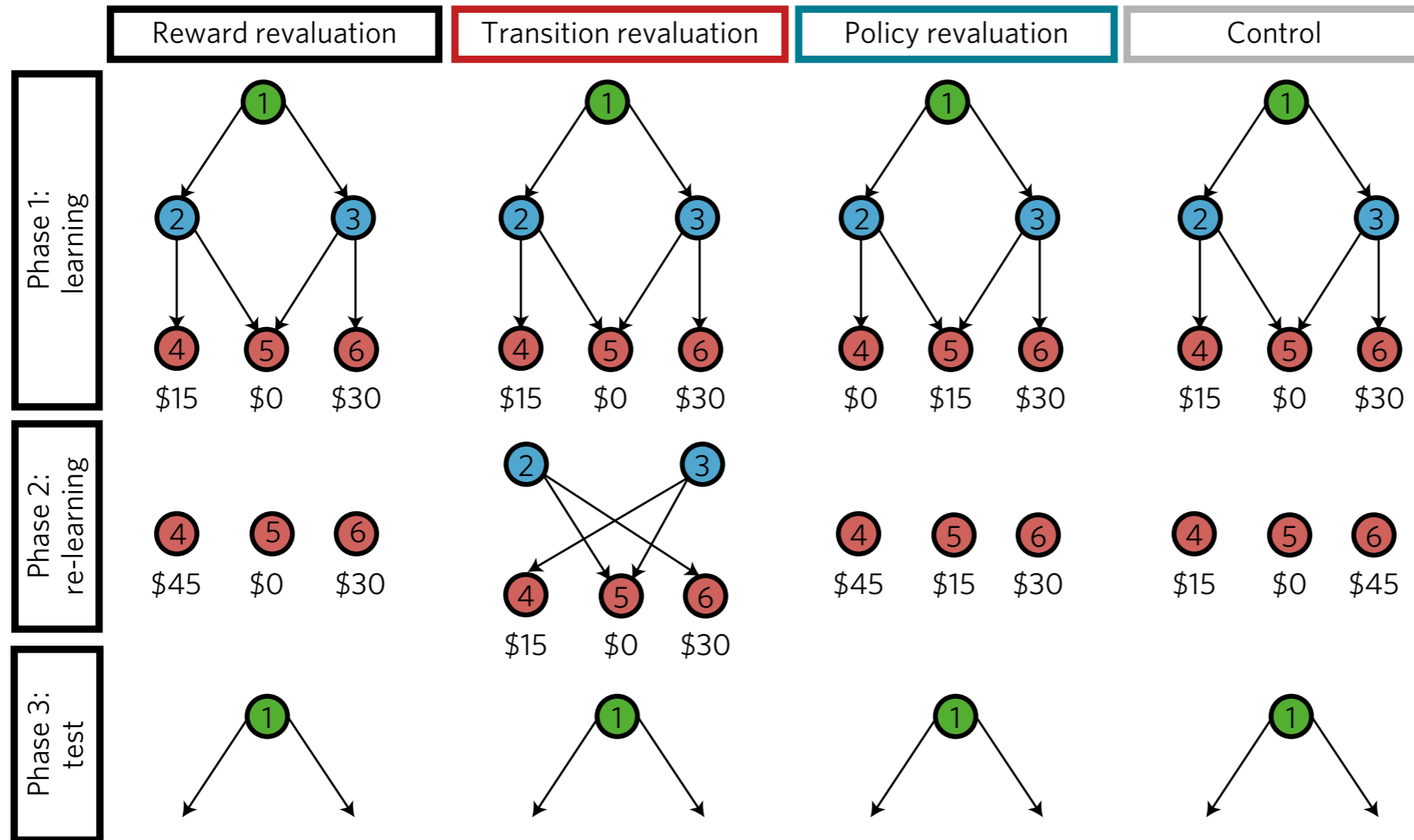


# Human successor learning



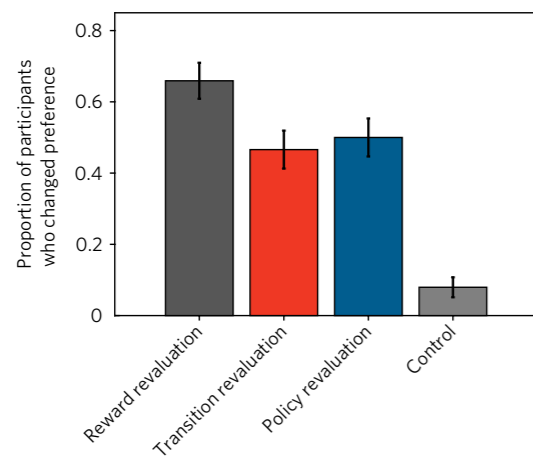
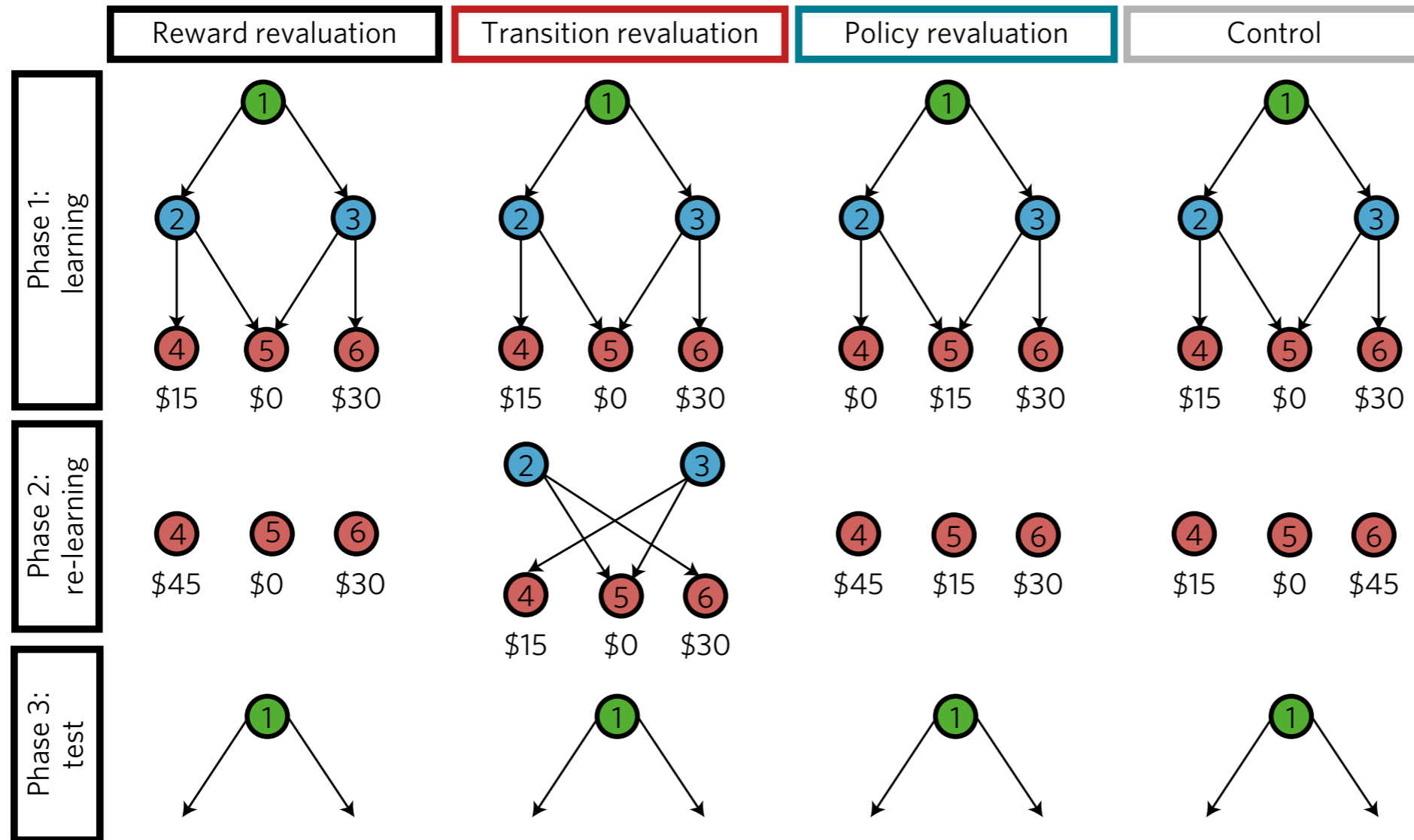
Momennejad et al., 2017 Nat. Hum. Beh.

# Human successor learning



Momennejad et al., 2017 Nat. Hum. Beh.

# Human successor learning



Momennejad et al., 2017 Nat. Hum. Beh.