
Lecture 6: Computational Cognitive Modeling

Reinforcement Learning (pt. 3)

course website:

[**https://brendenlake.github.io/CCM-site/**](https://brendenlake.github.io/CCM-site/)

Three levels of description (*David Marr, 1982*)

Computational

Why do things work the way they do?
What is the goal of the computation?
What are the unifying principles?



Algorithmic

What representations can implement such computations?
How does the choice of representations determine the algorithm?

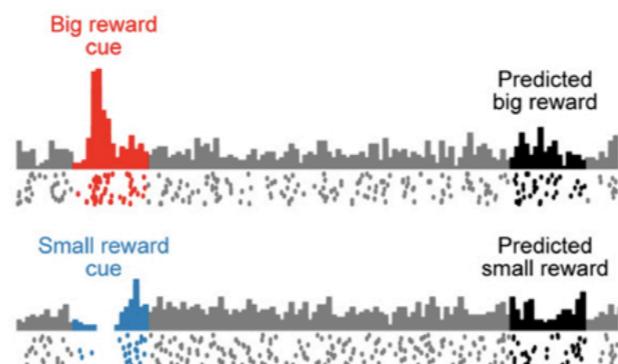
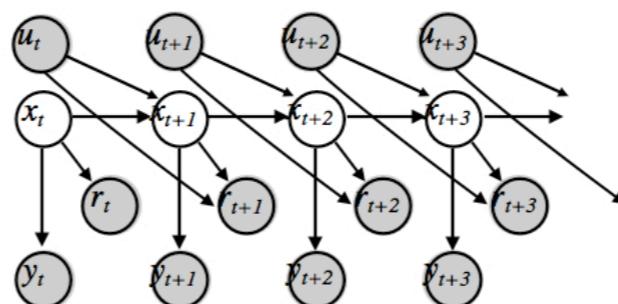
Implementational

How can such a system be built in hardware?
How can neurons carry out the computations?

maximize:

$$R_t = r_{t+1} + r_{t+2} + \dots + r_T$$

Bellman

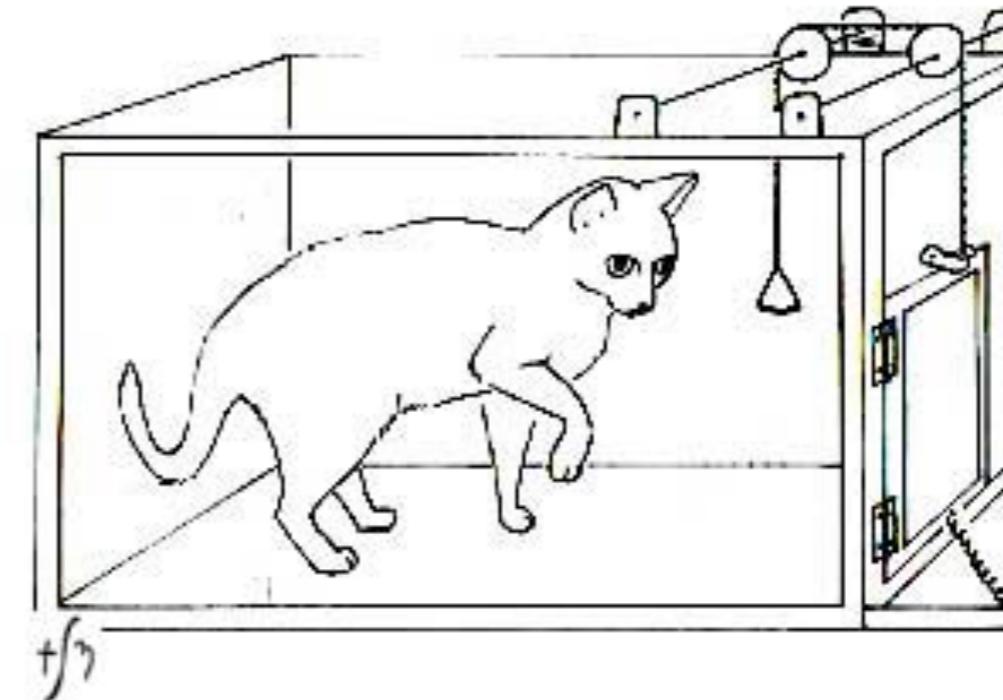
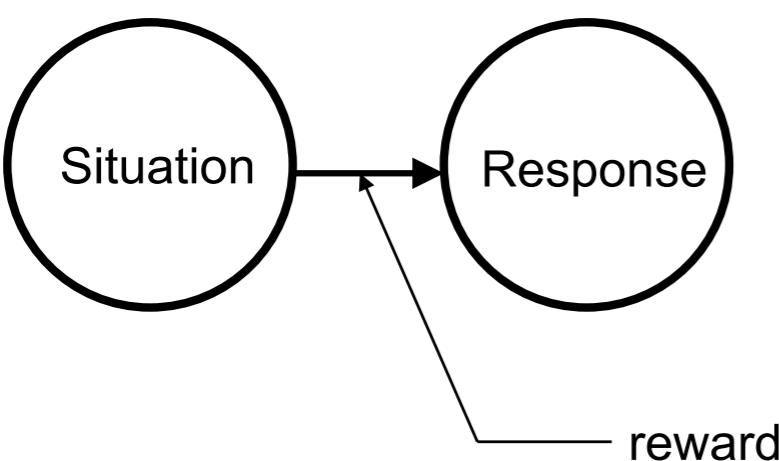


Dynamic programming,
TD methods, Monte
Carlo

Neural firing patterns,
prediction errors,
system level
neuroscience

Edward Thorndike

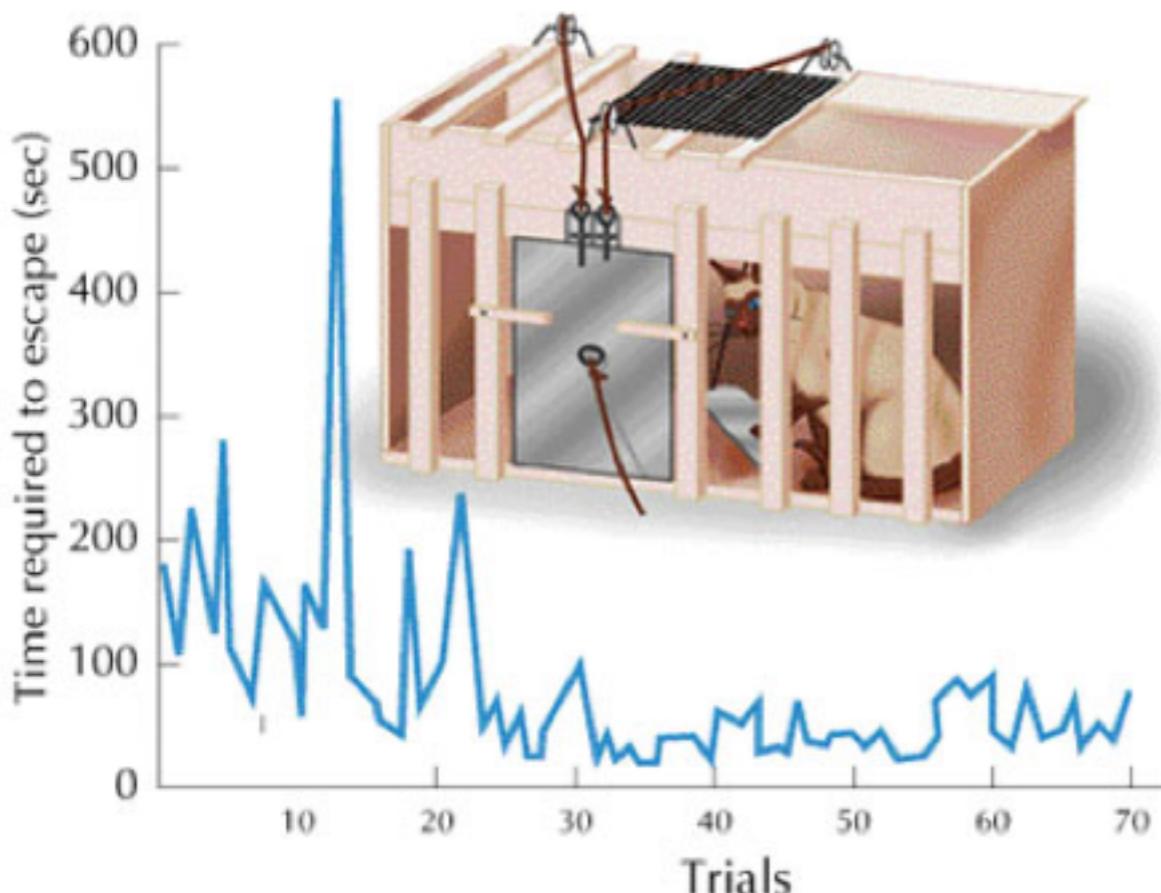
Law of Effect



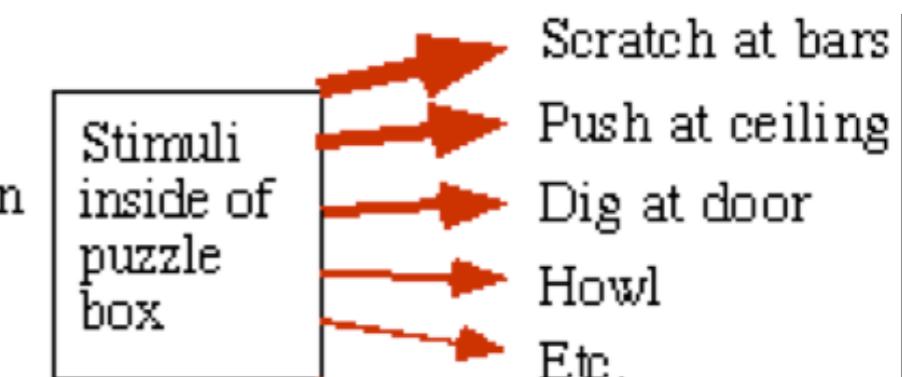
“Of several responses made to the same situation, those which are accompanied or closely **followed by satisfaction** to the animal **will**, other things being equal, **be more firmly connected with the situation**, so that, when it recurs, they will be more likely to recur; those which are accompanied or closely followed by discomfort to the animal will, other things being equal, have their connections with that situation weakened, so that, when it recurs, they will be less likely to occur” (Thorndike, 1911)

Edward Thorndike

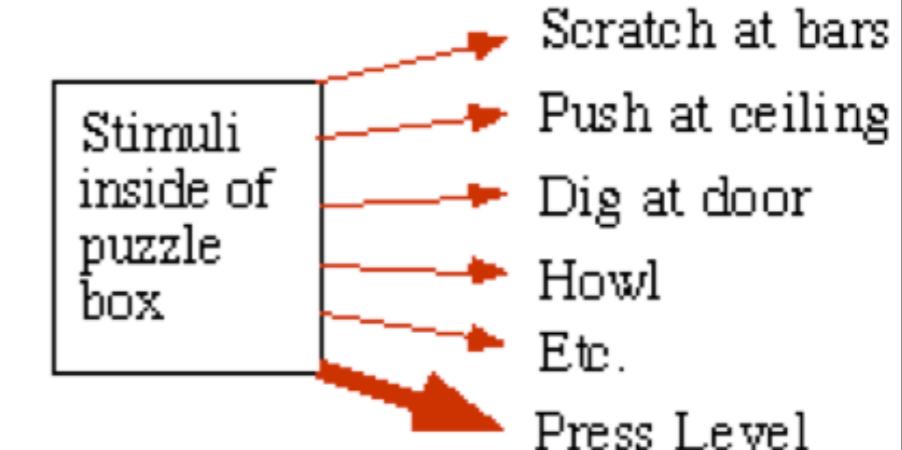
Law of Effect



First trial in
puzzle box



After many
trials in
puzzle box



classical conditioning



- prediction
- ... revealed by behavior
- ... shaped by learning

Maximizing Reward

The learning algorithm which approximates the full Bellman Solution (temporal difference):

$$V_t = E[r_{t+1}] + V_{t+1}$$

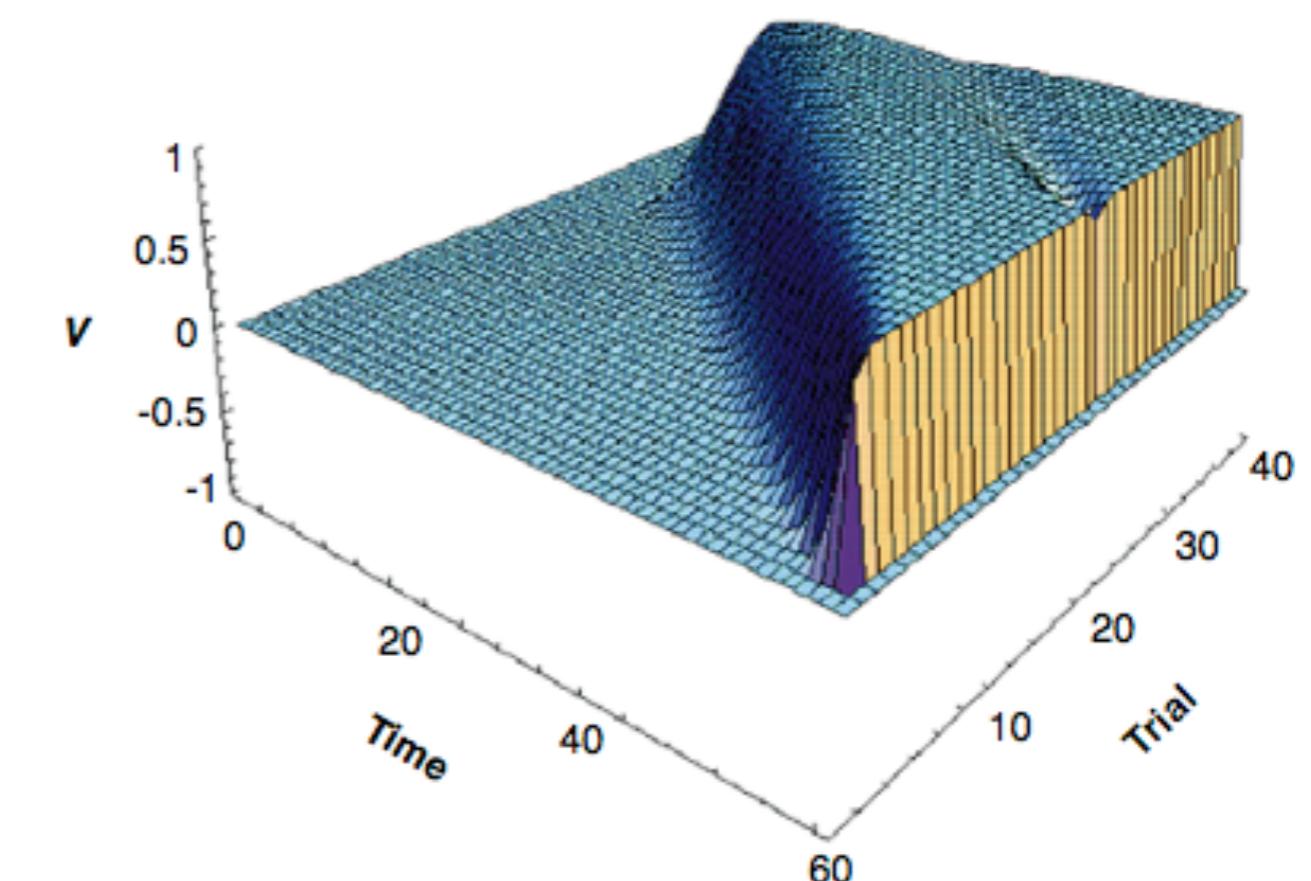
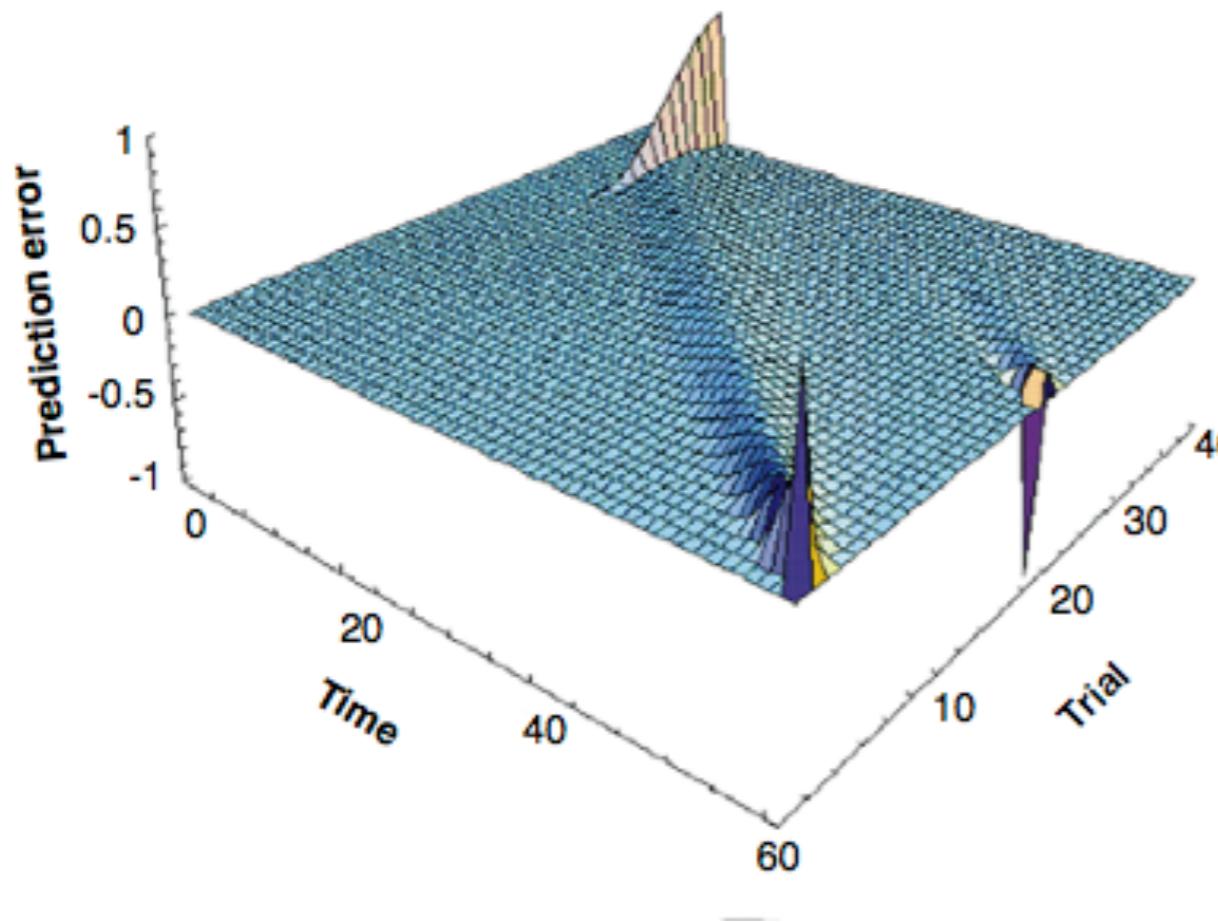
$$V_t \leftarrow V_t + \eta(r_{t+1} + V_{t+1} - V_t)$$

$$V_t \leftarrow (1 - \eta)V_t + \eta(r_{t+1} + V_{t+1})$$

*Compare with **Rescorla-Wagner***

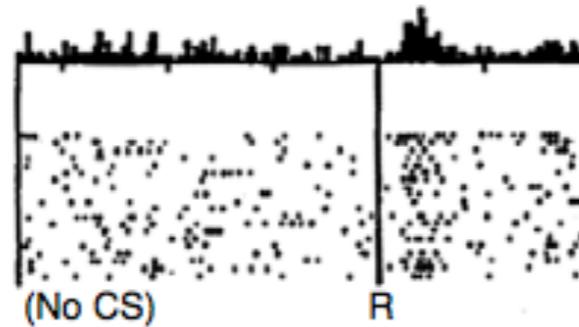
$$V_t \leftarrow V_t + \eta(r_{t+1} - V_t)$$

Temporal-difference learning and “backing up”

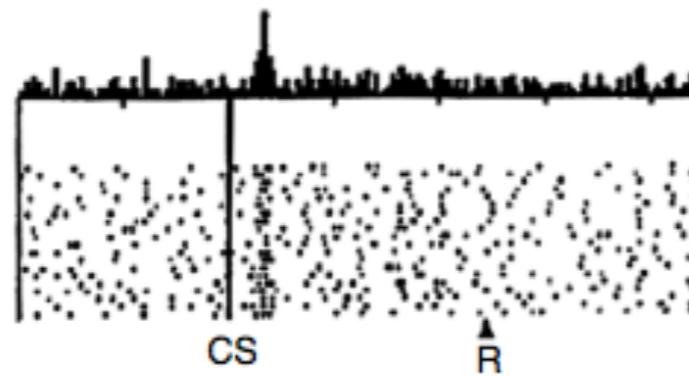


**Do dopamine neurons report an error
in the prediction of reward?**

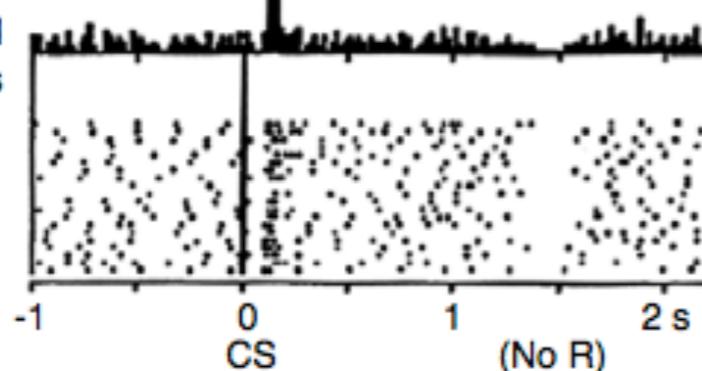
No prediction
Reward occurs



Reward predicted
Reward occurs



Reward predicted
No reward occurs

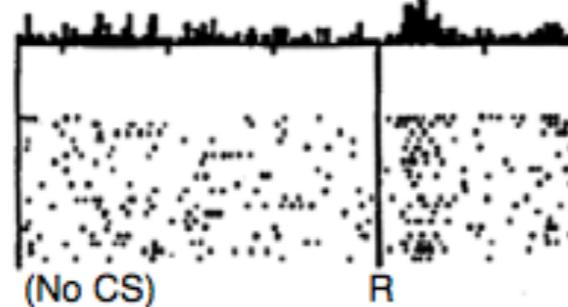


**Dopamine
neurons do it
too...**

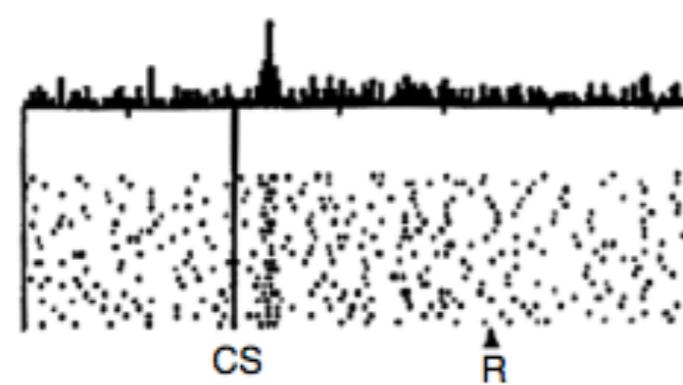
Dopamine neurons do it too...

Do dopamine neurons report an error
in the prediction of reward?

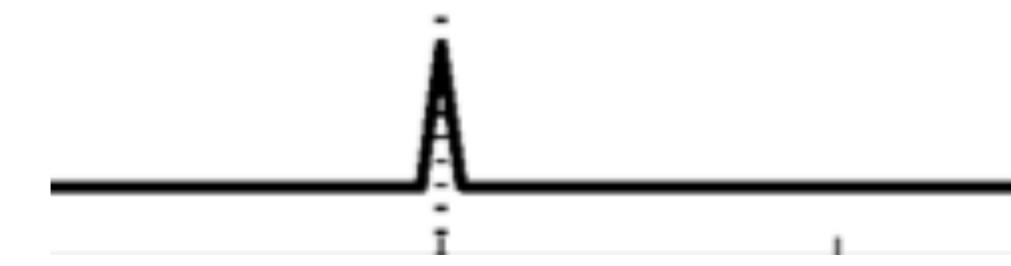
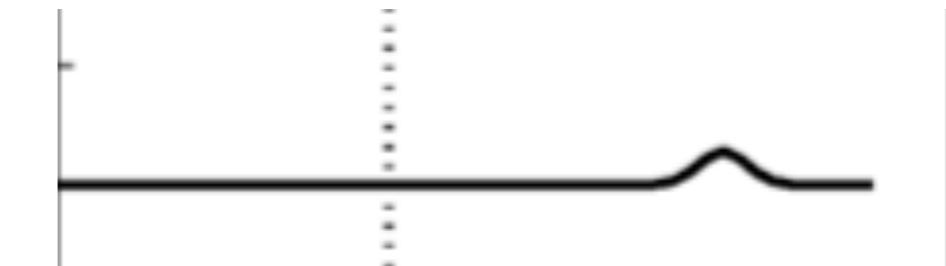
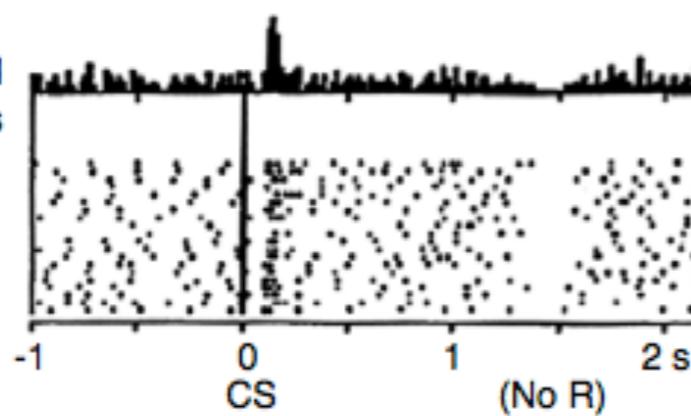
No prediction
Reward occurs



Reward predicted
Reward occurs



Reward predicted
No reward occurs



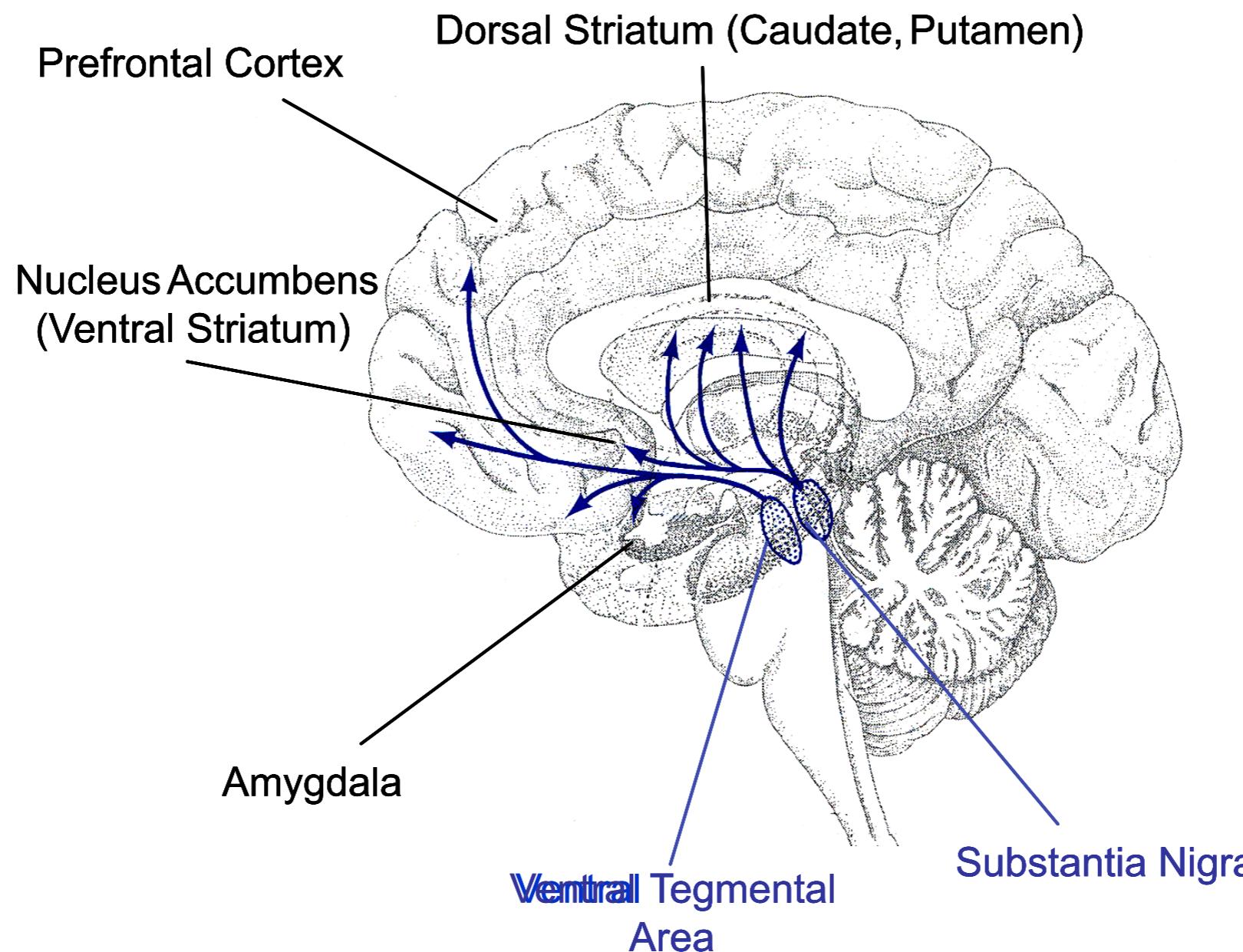
δ_t

δ_t

δ_t

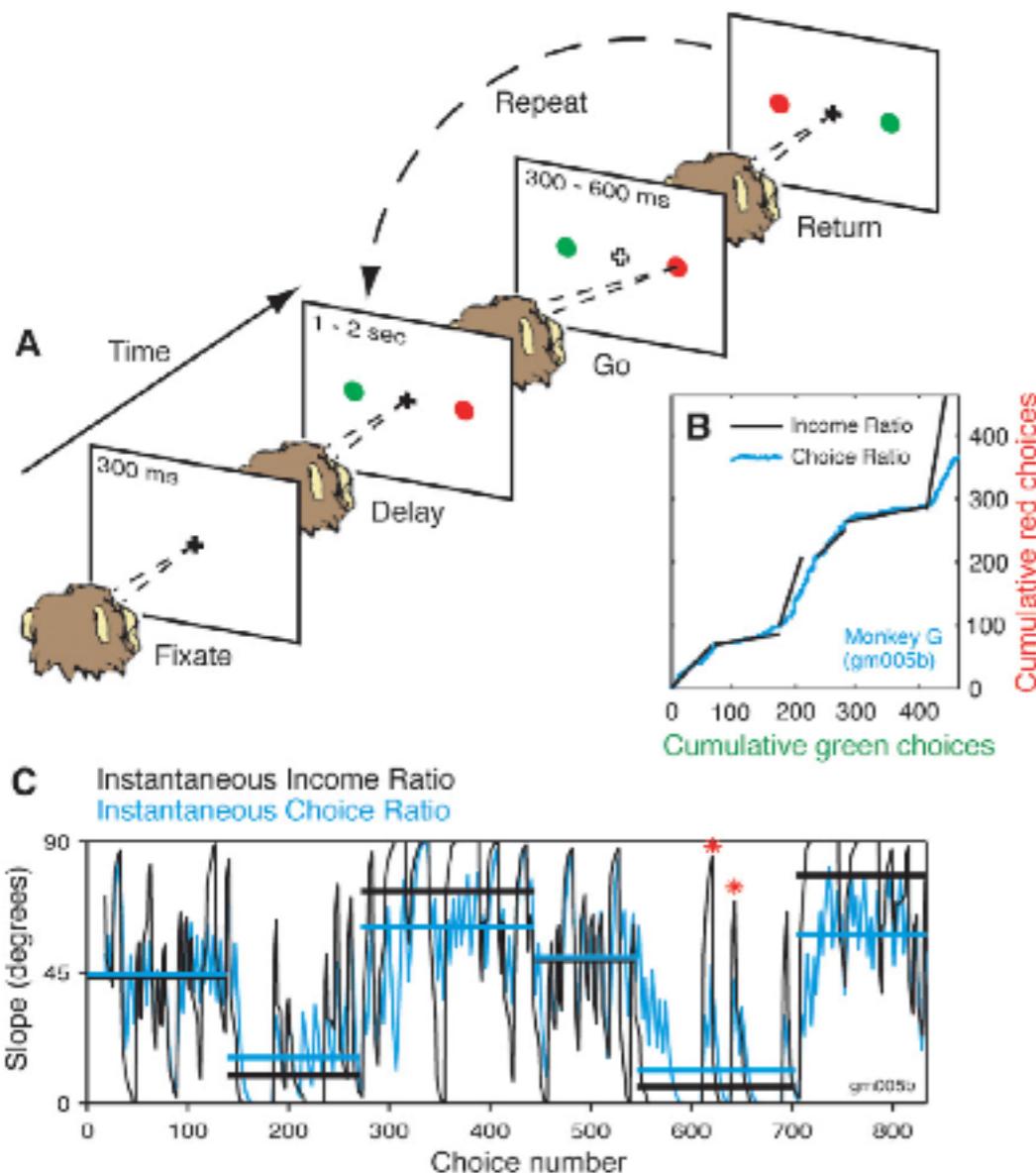
*Also correlates with magnitude of reward, degree
of intermittent reinforcement, etc...*

The role of dopamine in reward learning



Matching Behavior and the Representation of Value in the Parietal Cortex

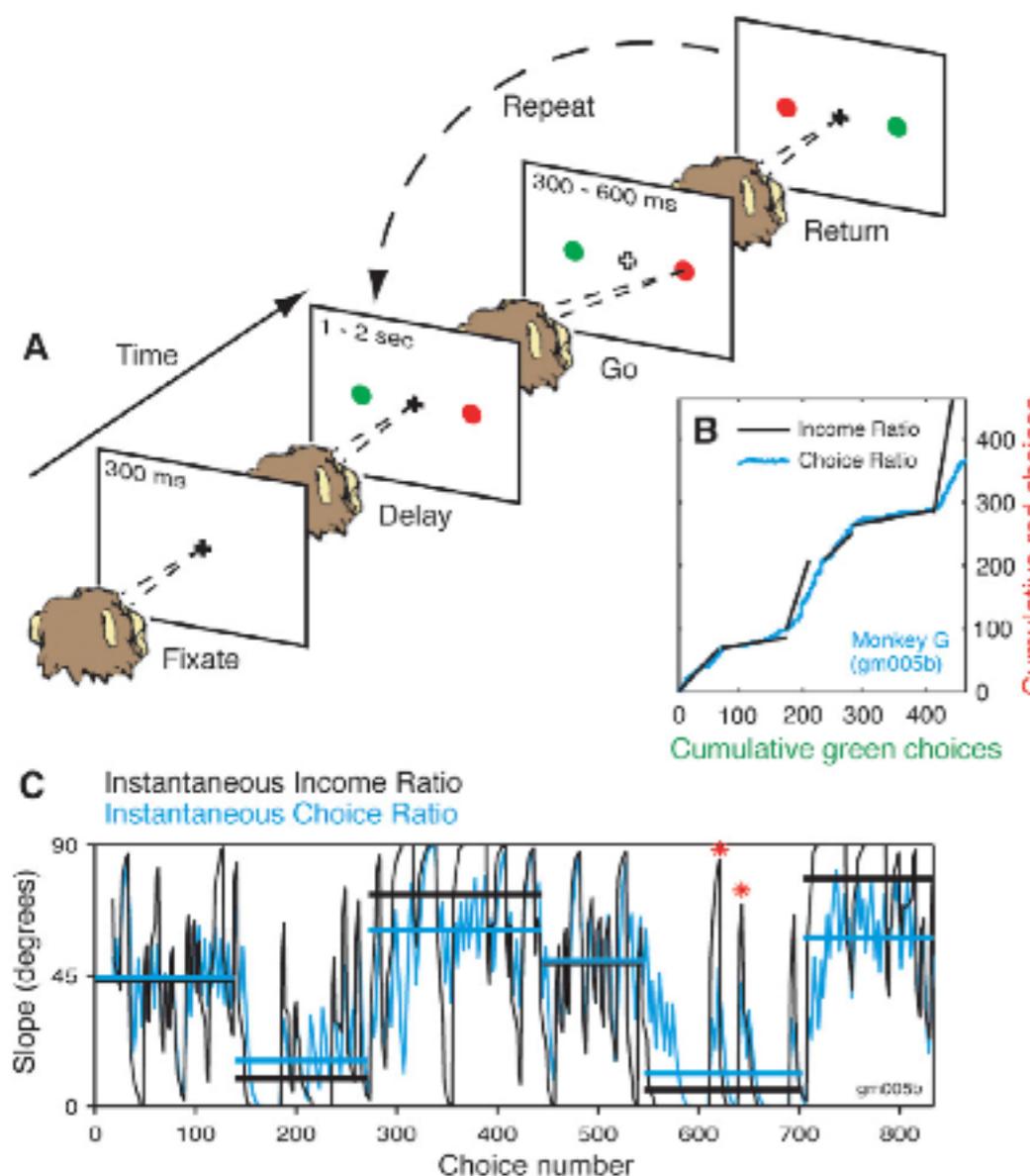
Leo P. Sugrue,* Greg S. Corrado, William T. Newsome



- Rhesus monkeys make eye fixations to one of two locations which get rewarded at different rates
- Dynamic foraging task since the reward from each option “drifts” in time (Poisson arrival times, but one are baited are always available... like VI schedule, makes matching the optimal strategy)
- Ever 200 or so trials, the reward magnitudes or rates change suddenly causing discontinuous shifts in the reward rates

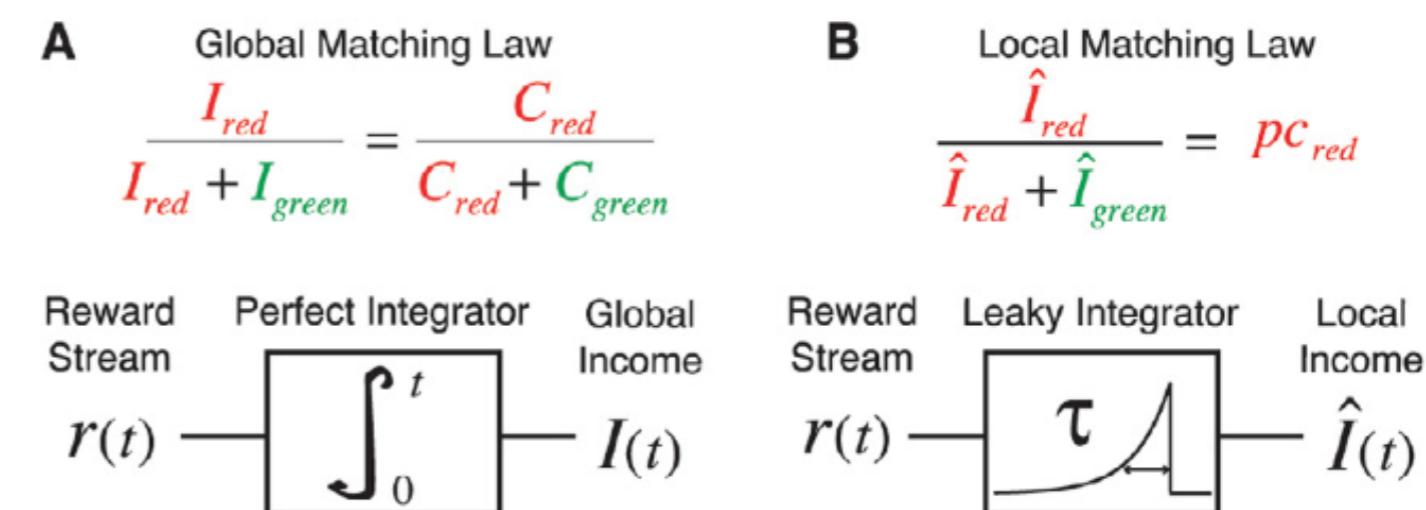
Matching Behavior and the Representation of Value in the Parietal Cortex

Leo P. Sugrue,* Greg S. Corrado, William T. Newsome



Learning

- Optimal integration of the reward history would be incapable of this quick adaptation
- Instead, suggest some “leaky integrator” which is pooling information in time (recency biased)

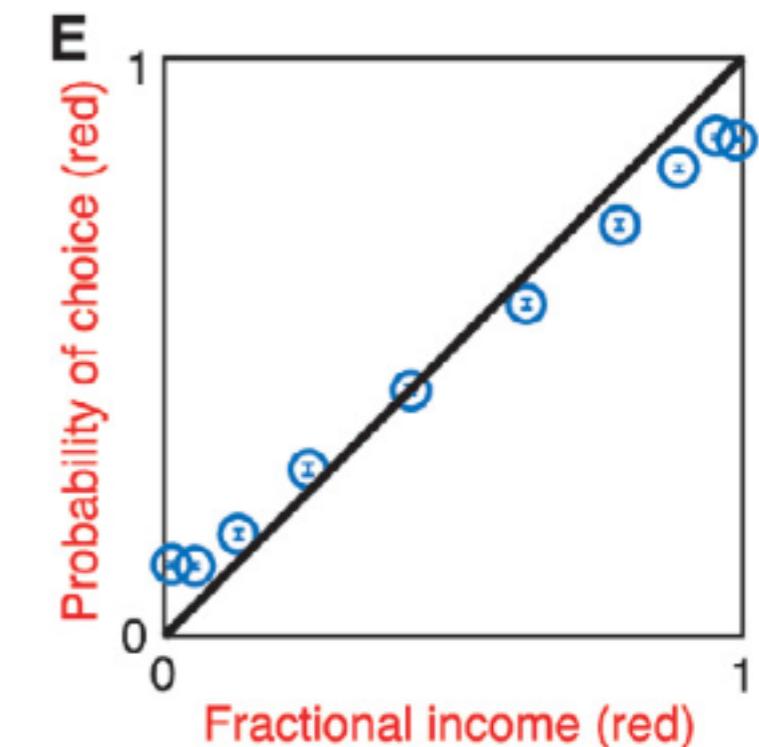
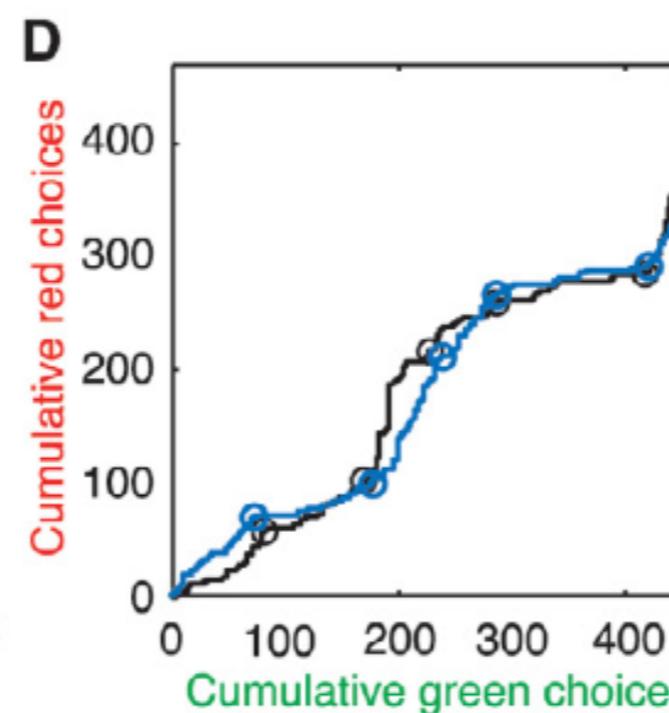
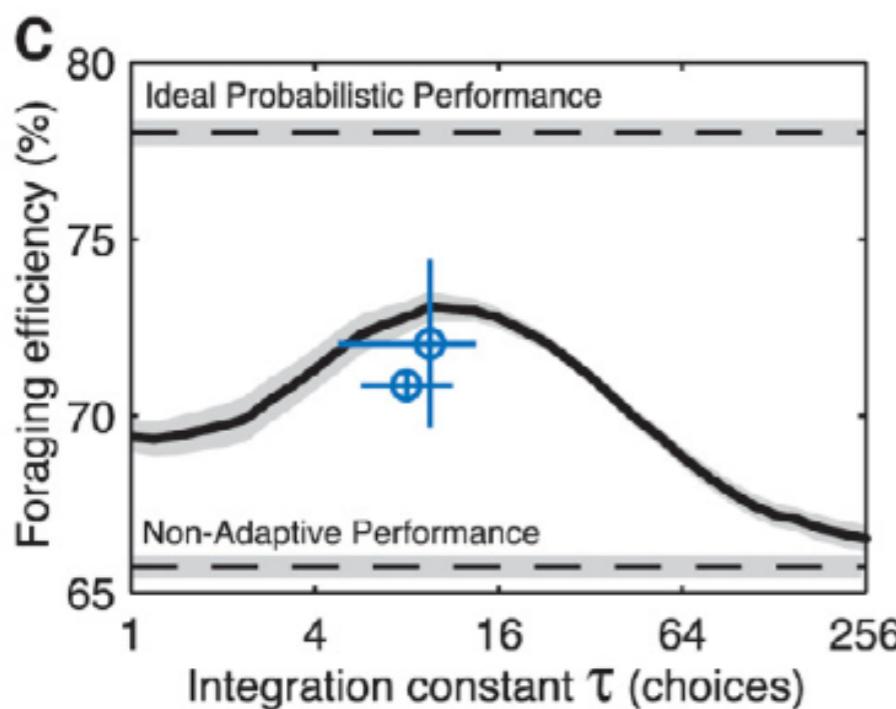


- Distinct from traditional matching law which integrates “globally”

Matching Behavior and the Representation of Value in the Parietal Cortex

Leo P. Sugrue,* Greg S. Corrado, William T. Newsome

How do the monkeys do?



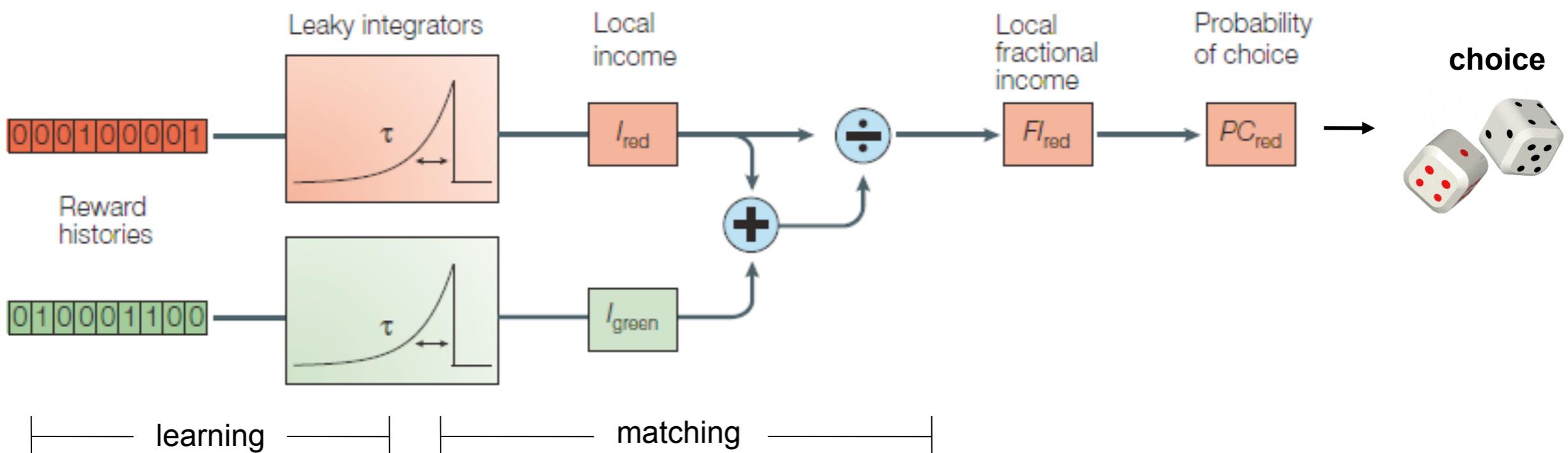
Near optimal setting of tau

The monkey's match!

Matching Behavior and the Representation of Value in the Parietal Cortex

Leo P. Sugrue,* Greg S. Corrado, William T. Newsome

The model

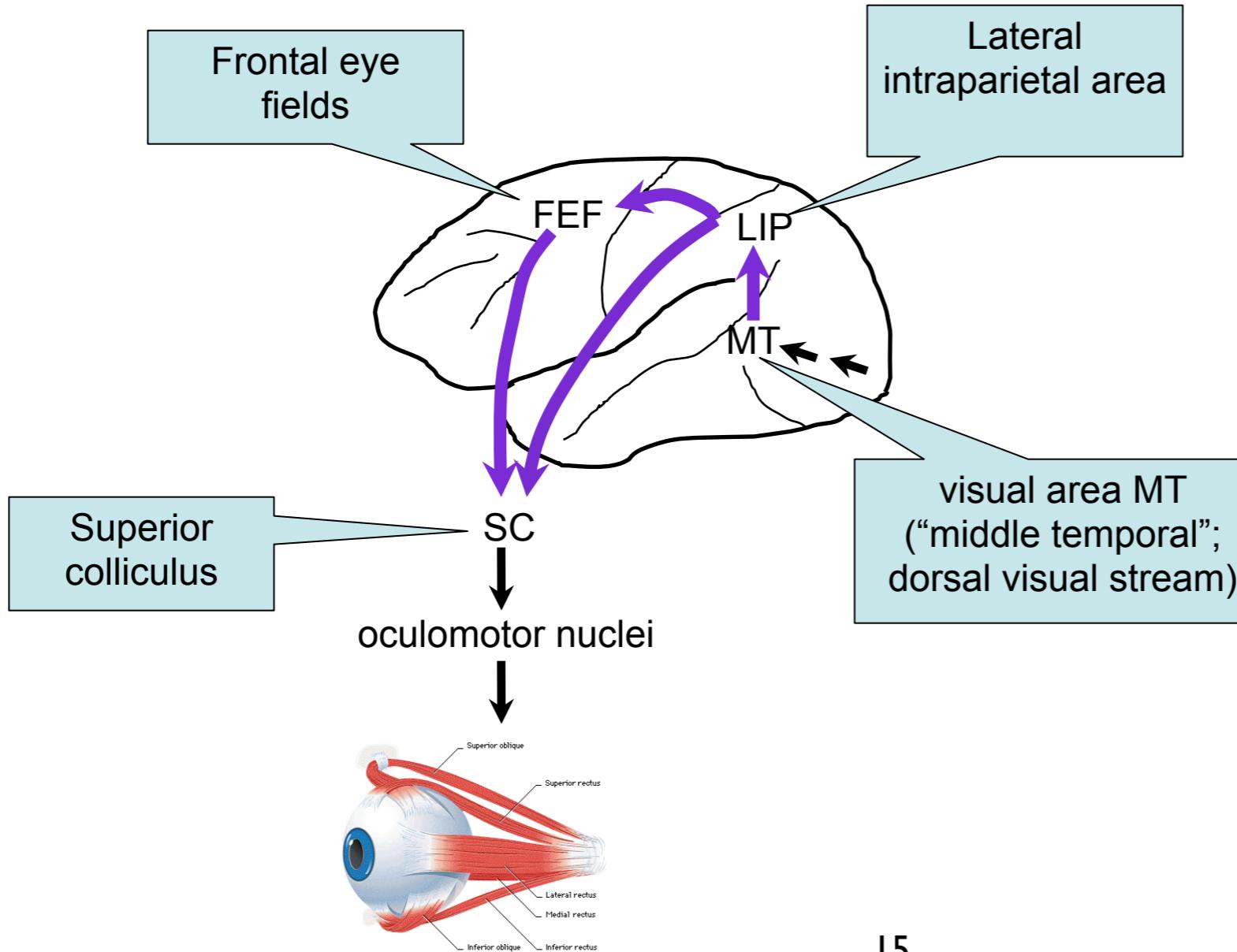


Matching Behavior and the Representation of Value in the Parietal Cortex

Leo P. Sugrue,* Greg S. Corrado, William T. Newsome

Recording from target-selective neuron is LIP

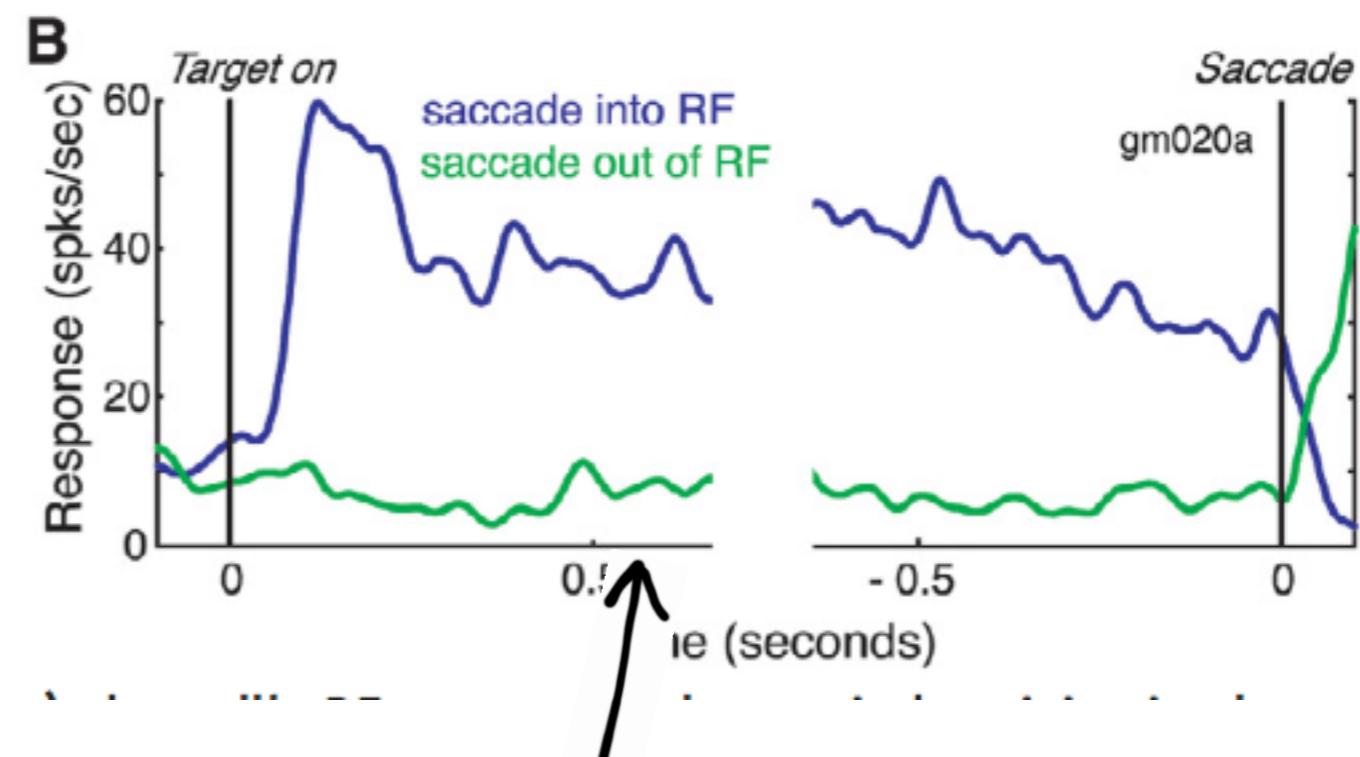
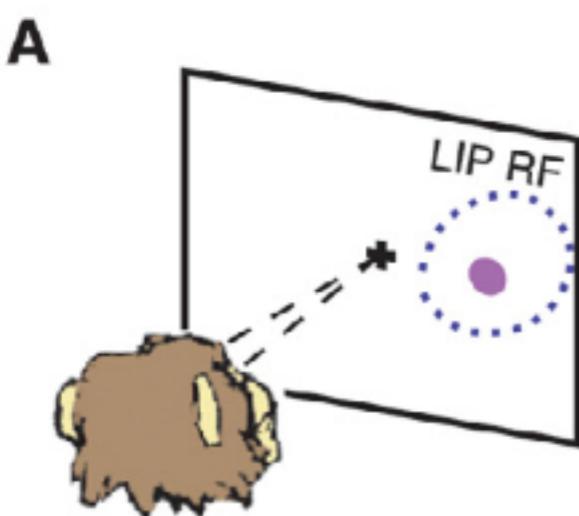
- Region involved in planning and executing saccades



Matching Behavior and the Representation of Value in the Parietal Cortex

Leo P. Sugrue,* Greg S. Corrado, William T. Newsome

Recording from target-selective neuron is LIP



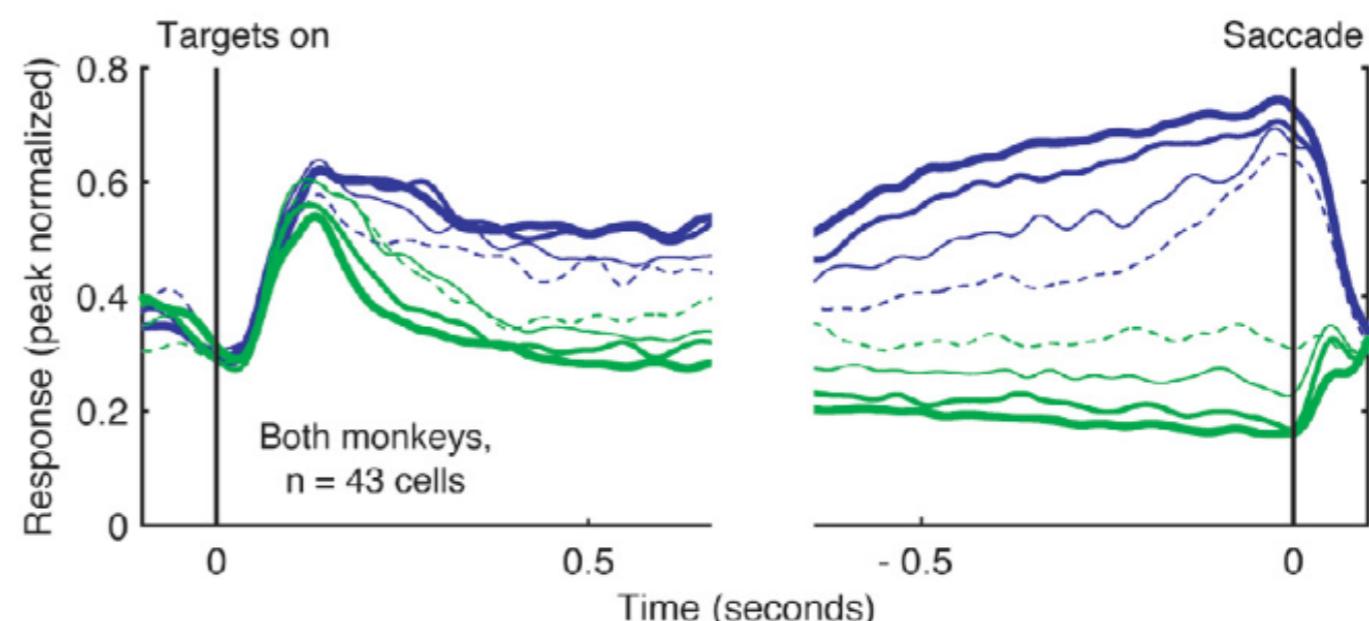
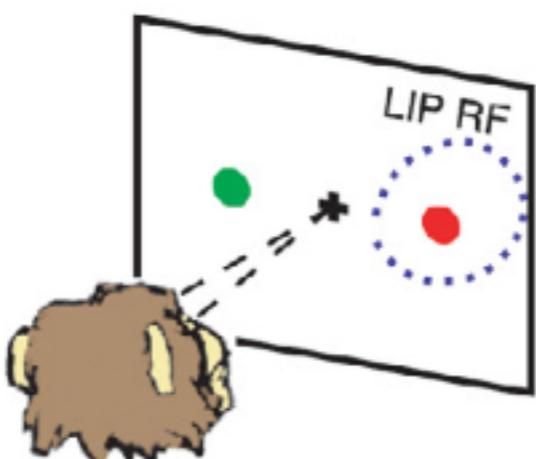
- Selected on the basis of delayed saccade task to find RF selective cells (like a forced choice version of matching game)

Spatially selective response

Matching Behavior and the Representation of Value in the Parietal Cortex

Leo P. Sugrue,* Greg S. Corrado, William T. Newsome

Recording from target-selective neuron in LIP

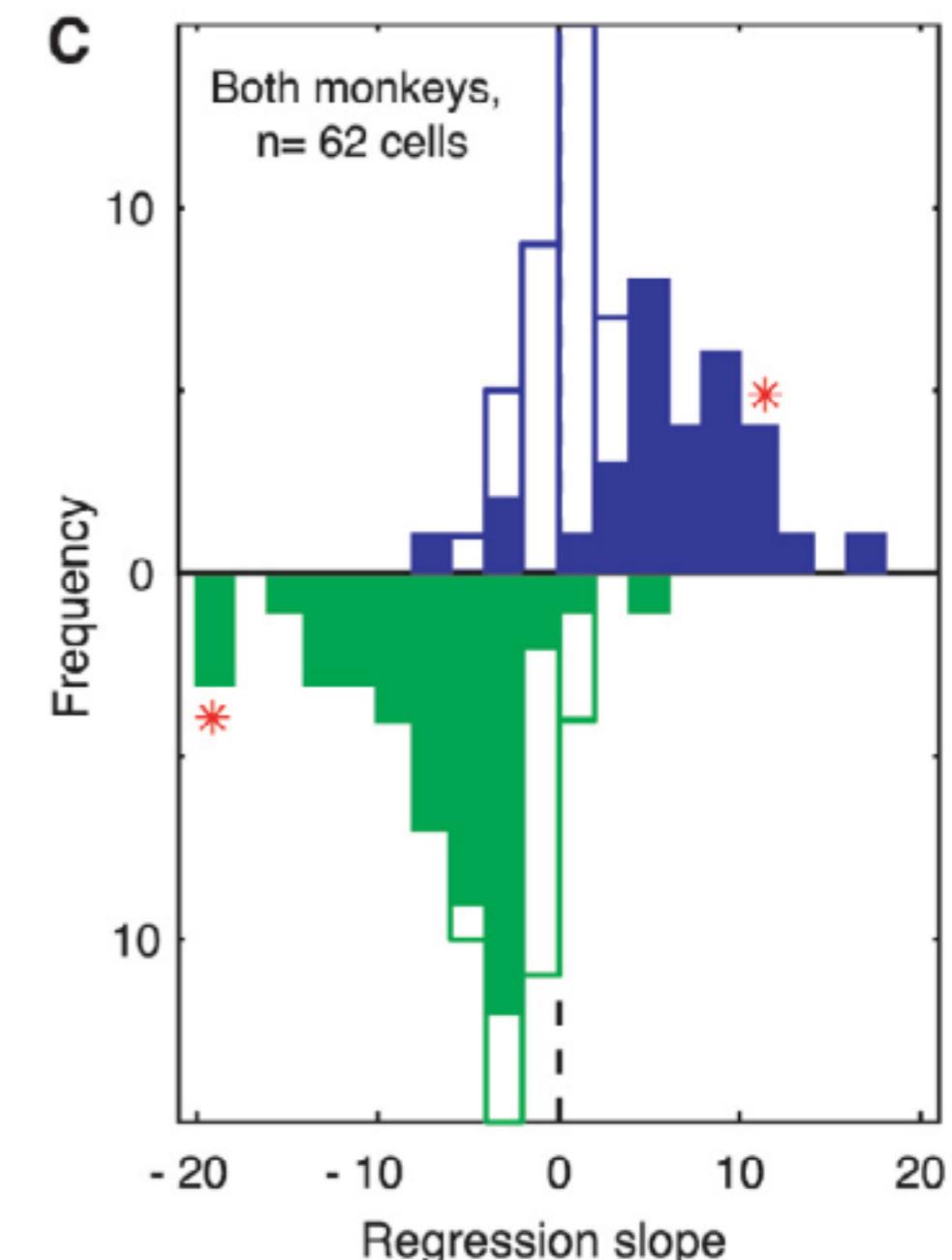
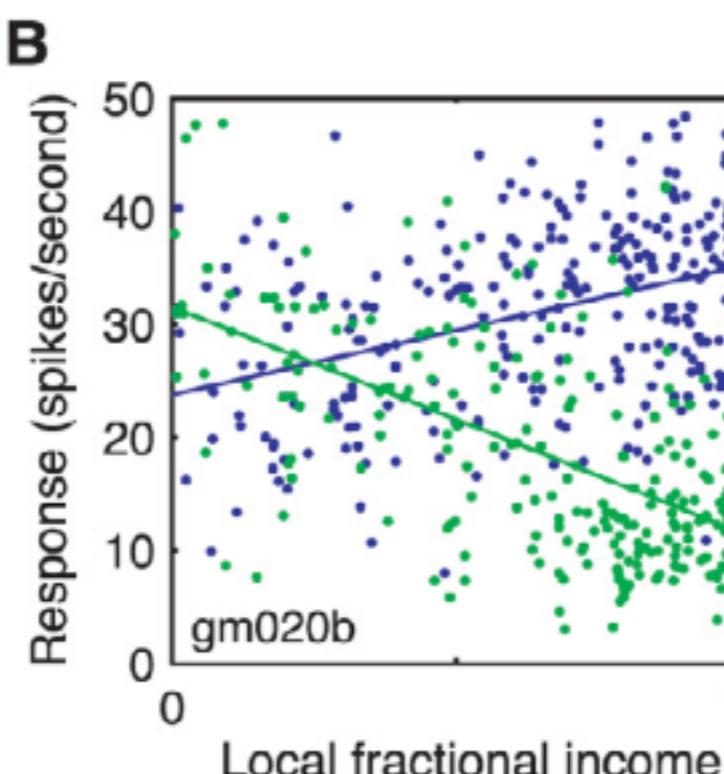
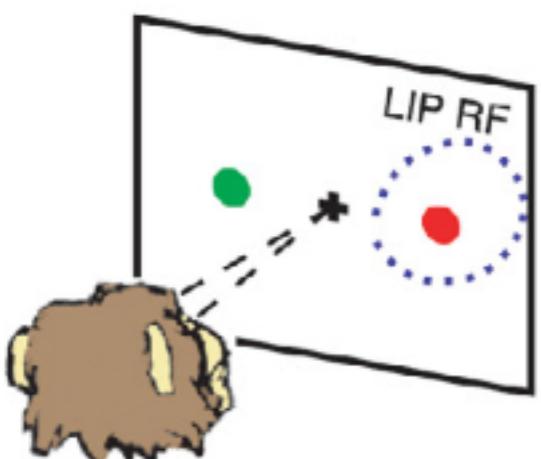


- Cells modulated by trial-to-trial utility of target in the receptive field

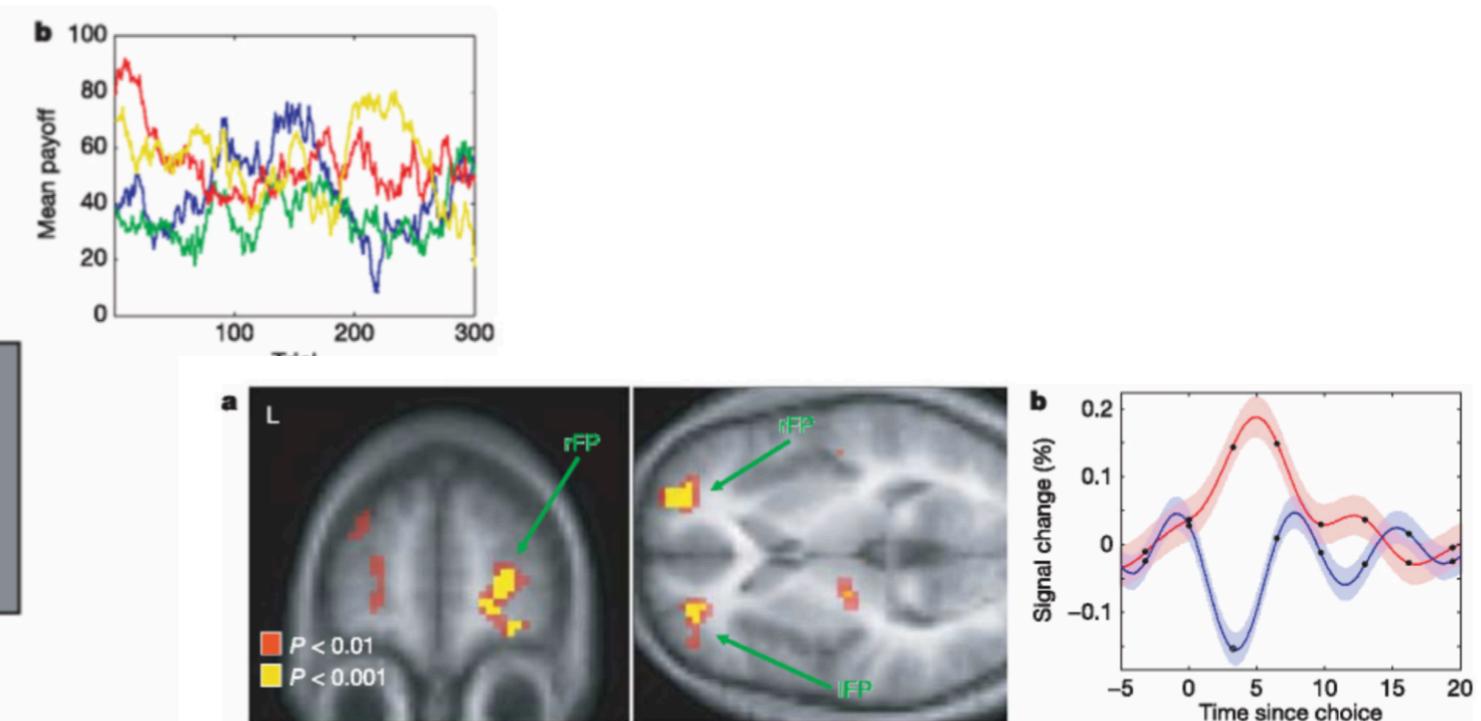
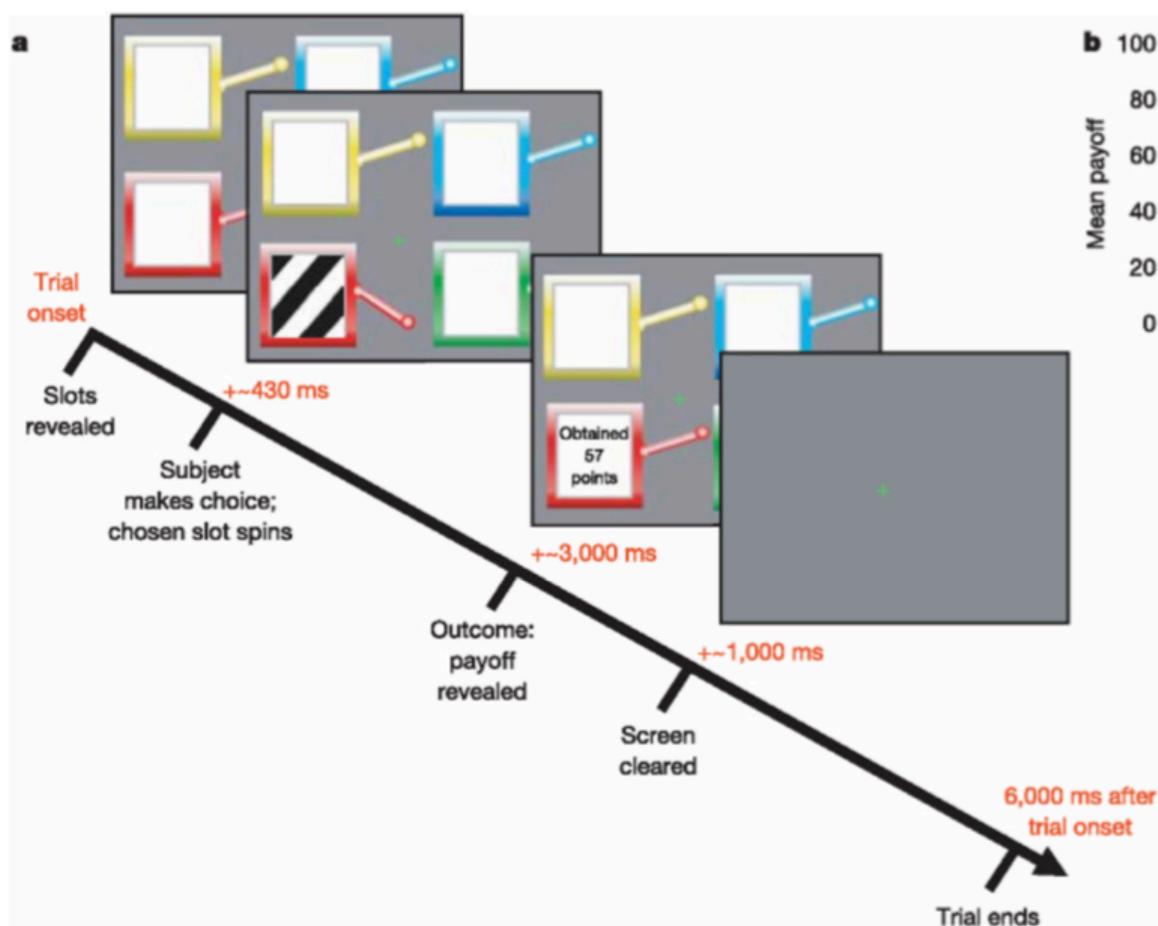
Matching Behavior and the Representation of Value in the Parietal Cortex

Leo P. Sugrue,* Greg S. Corrado, William T. Newsome

Recording from target-selective neuron is LIP



- Cells modulated by trial-to-trial utility of target in the receptive field

**Figure 3.**

Exploration-related activity in frontopolar cortex. a, Regions of left and right frontopolar cortex (IFP, rFP) showing significantly increased activation on exploratory compared with exploitative trials. Activation maps (yellow, $P < 0.001$; red, $P < 0.01$) are superimposed on a subject-averaged structural scan. The coordinates of activated areas are $[-27,48,4, \text{peak } z = 3.49]$ for IFP and $[27,57,6, \text{peak } z = 4.13]$ for rFP. b, rFP BOLD time courses averaged over 1,515 exploratory (red line) and 2,646 exploitative (blue line) decisions. Black dots indicate the sampling frequency (although, because sample alignment varied from trial to trial, time courses were upscaled). Coloured fringes show error bars (representing s.e.m.).

Published in final edited form as:

Nature. 2006 June 15; 441(7095): 876–879. doi:10.1038/nature04766.

Cortical substrates for exploratory decisions in humans

Nathaniel D. Daw, John P. O'Doherty, Peter Dayan, Ben Seymour, and Raymond J. Dolan

Generalization and Function Approximation

- Tabular value functions

state	action \longrightarrow			
	0.1	0.3	0.7	0.2
0.4
0.2
	$Q(s,a)$			

What happens if the size of the state/action space is large?

- Large numbers of states/actions?
- Continuously-valued states/actions?
- Most states never experienced exactly before

Memory

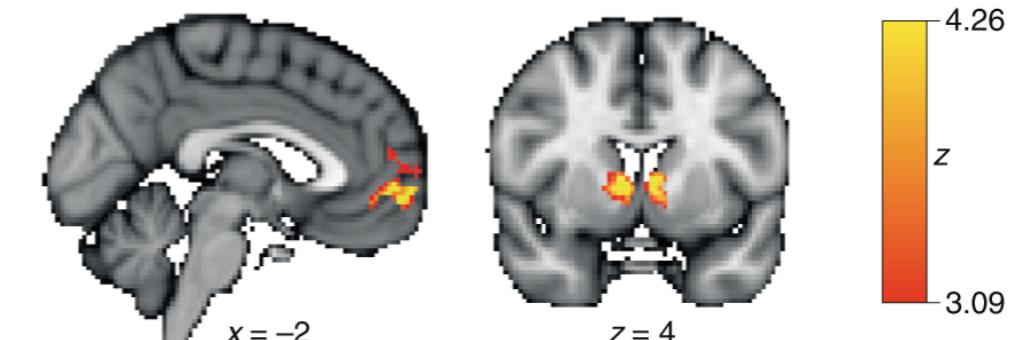
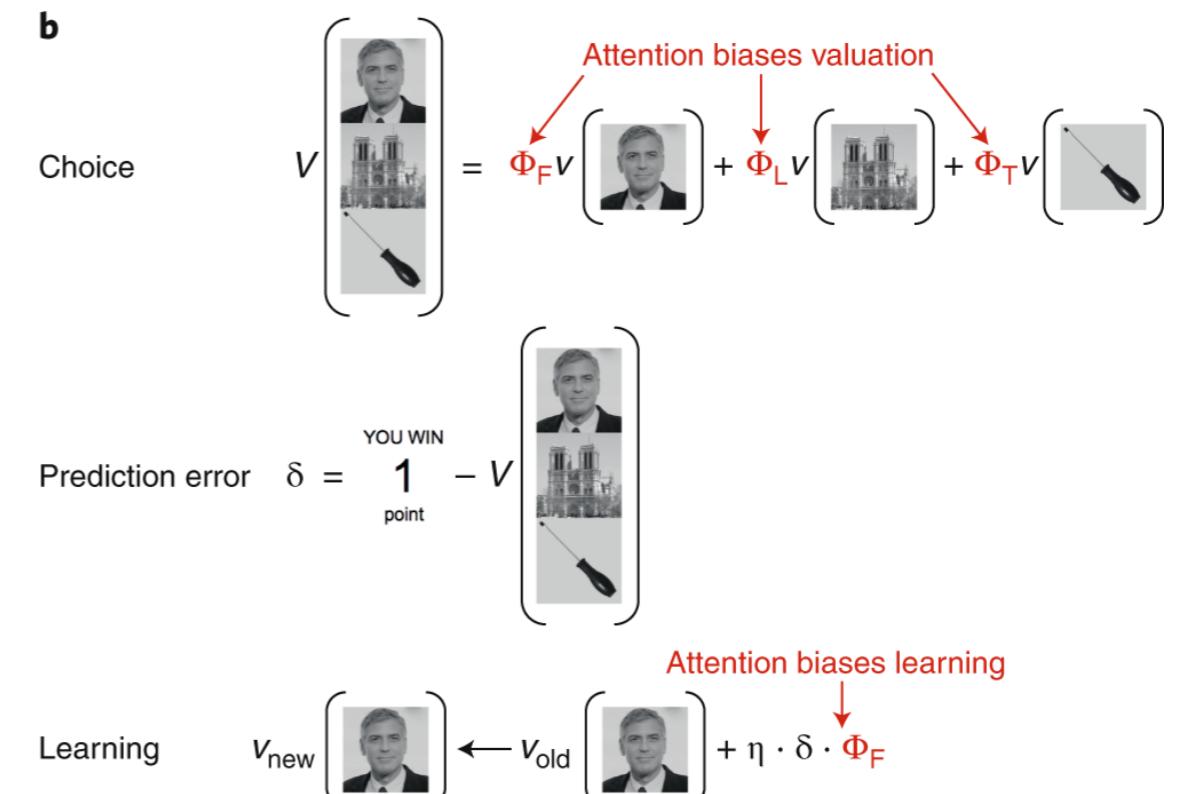
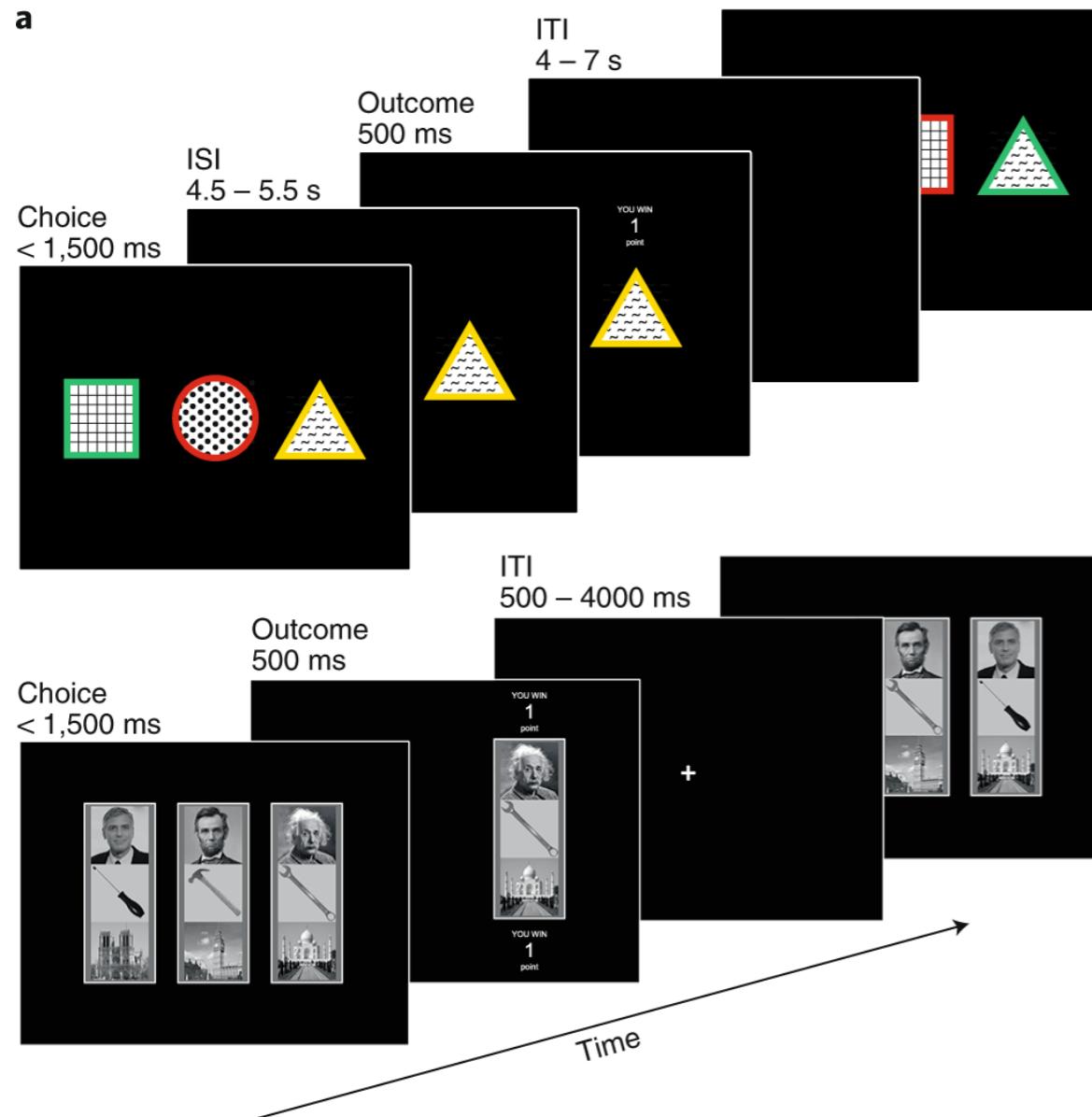
Time

Data

GENERALISATION: how experience with small part of state space is used to produce good behaviour over large part of state space

Learning task-state representations

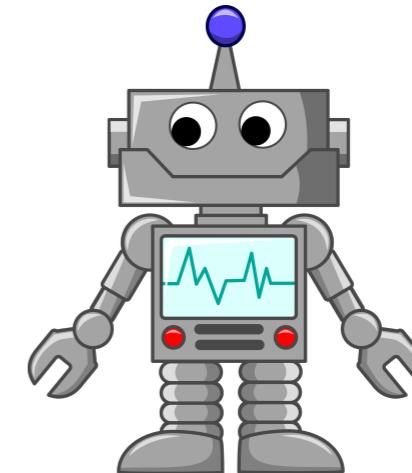
Yael Niv 



Generalization and Function Approximation

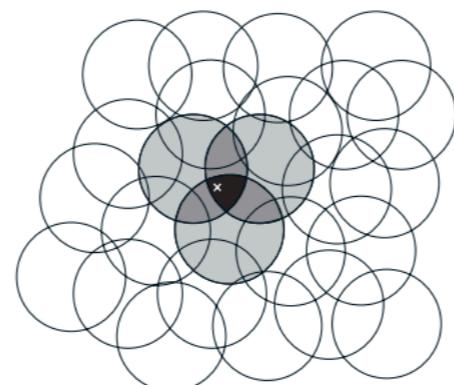
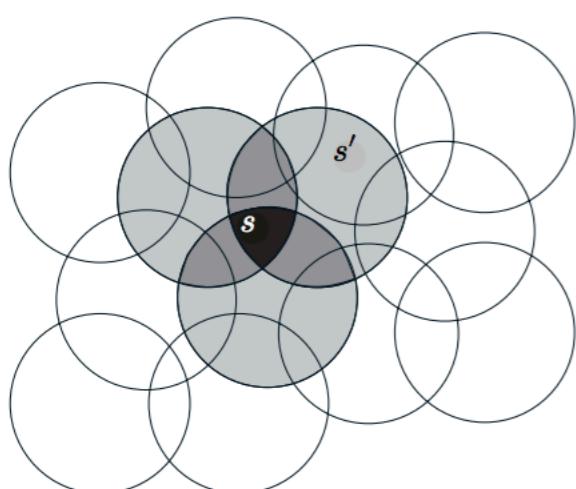
- Instead of directly looking up the states in a kind of addressable-memory, we represent different features of the states

$$\vec{\phi}_s = \begin{pmatrix} x \\ y \\ \text{heading} \\ \text{batterypower} \end{pmatrix}$$

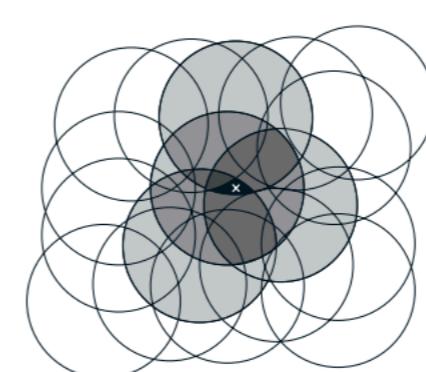


- Position in x, y coordinates (real numbers)
- Heading in degrees w.r.t. north (real number or quantised)
- Battery power = some real number

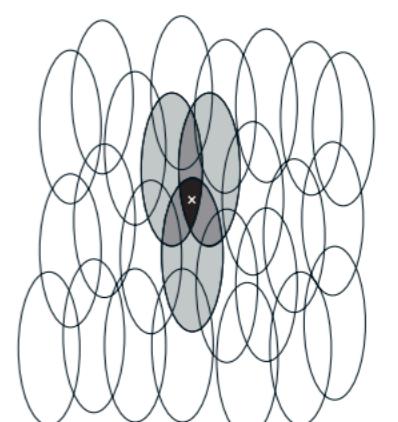
Coarse coding/tiling



Narrow generalization



Broad generalization



Asymmetric generalization

Figure 9.7: Generalization in linear function approximation methods is determined by the sizes and shapes of the features' receptive fields. All three of these cases have roughly the same number and density of features.

Generalization and Function Approximation

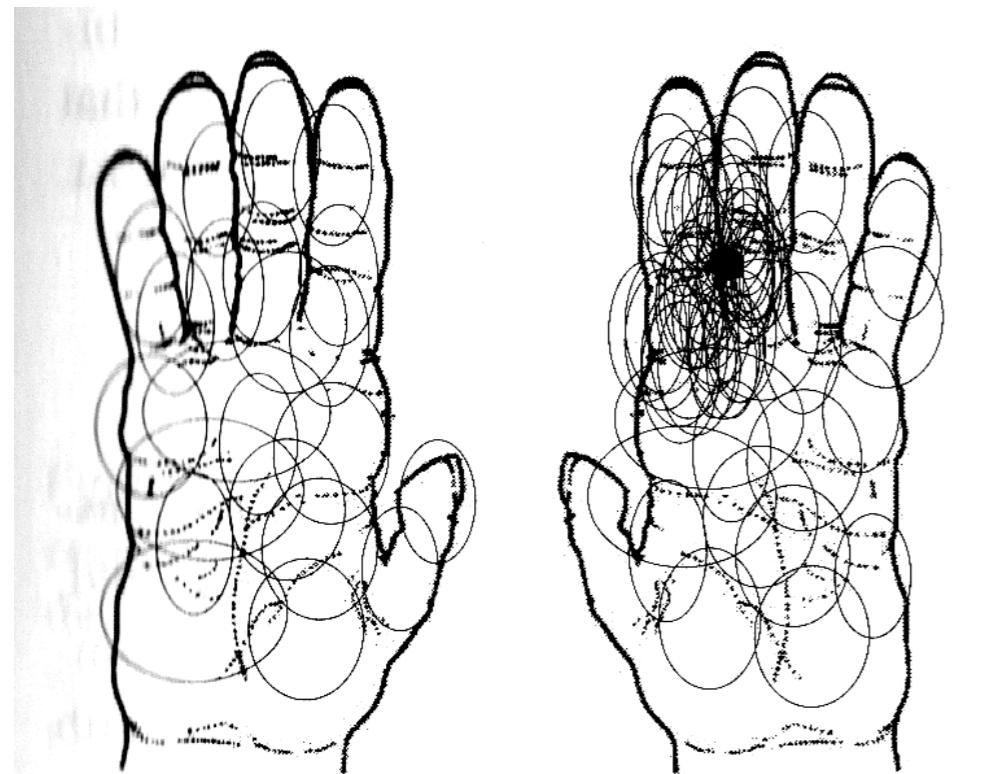
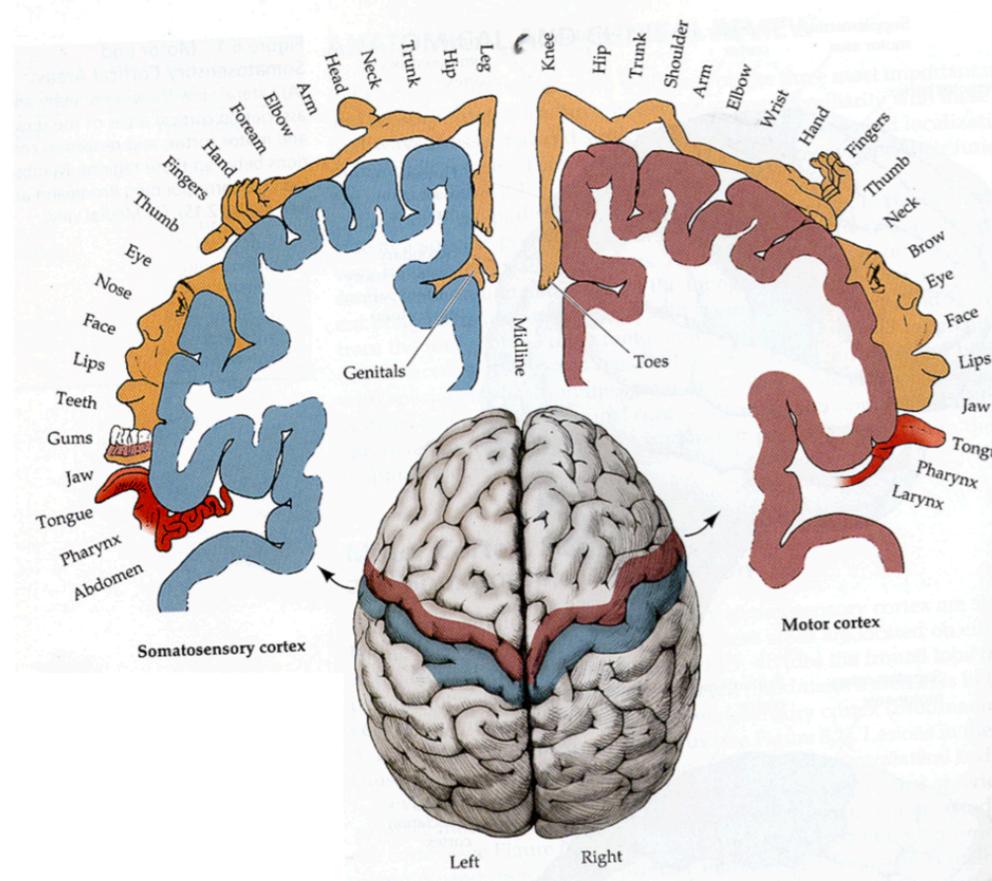


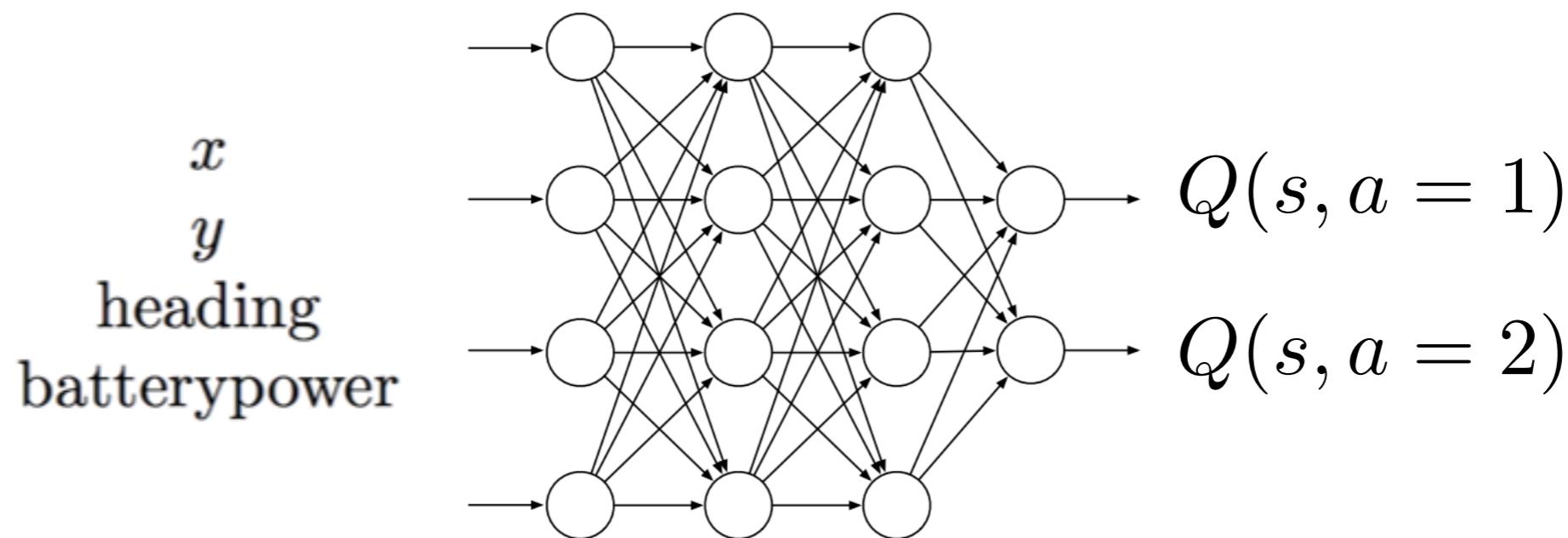
Figure 3.6 — CORTICAL PLASTICITY. The density of receptive fields of somatosensory cortical neurons, mapped onto the monkey palm surface, after prolonged exposure to local stimulation (*right*; the arrow indicates the stimulated spot). The map for the unstimulated hand is shown as a control (*left*). Merely touching the skin repeatedly at the designated location for the period of several days caused that location to become over-represented in the brain.



*More neuron (finer tiling) =
higher fidelity, better
discrimination*

Generalization and Function Approximation

- Neural networks make a useful approximation scheme due to the universal approximation properties



- Train with backprop as in past stuff but the loss function is the Q-learning update rule

$$L = \frac{1}{2} [r + \max_{a'} Q(s', a') - Q(s, a)]^2$$

Reminder: Backpropagation algorithm for computing gradient

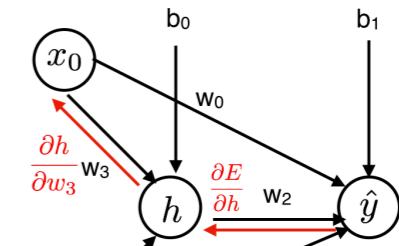
Multi-step strategy:

$$\frac{\partial E}{\partial w_3} = \frac{\partial E}{\partial h} \frac{\partial h}{\partial w_3}$$

Step 1) Compute how error changes as a function of hidden unit activation

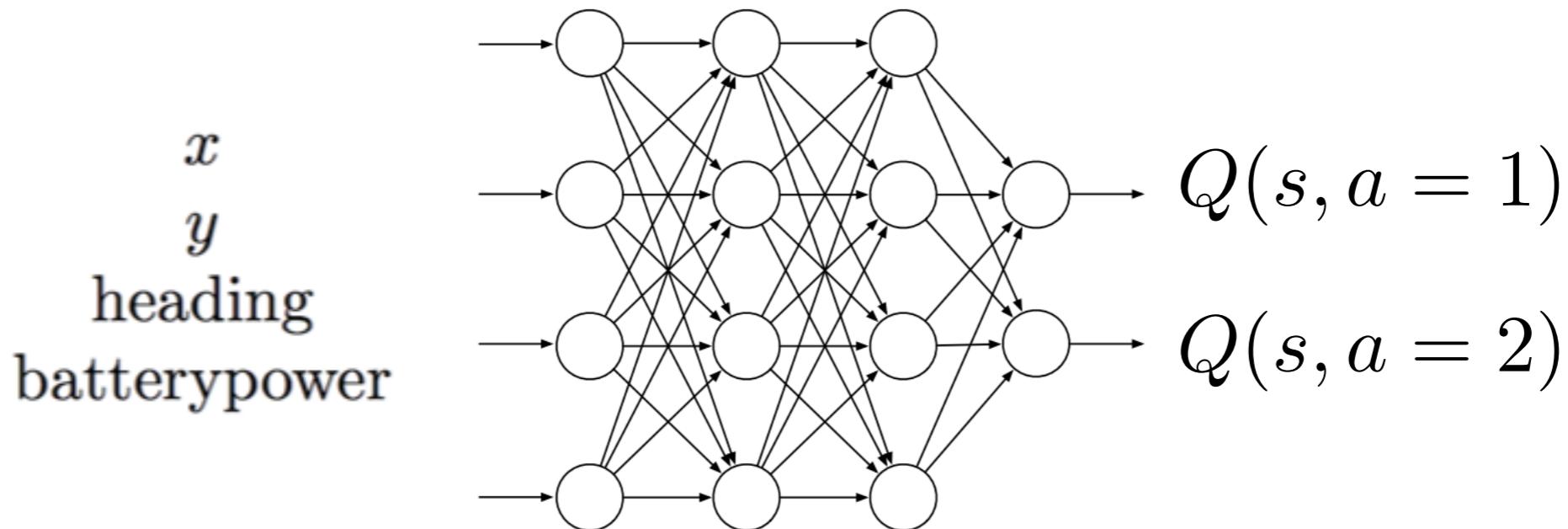
Step 2) Compute how hidden unit activation changes as a function of weight

$$E(w, b) = (\hat{y} - y)^2 = (g(\text{net}) - y)^2$$

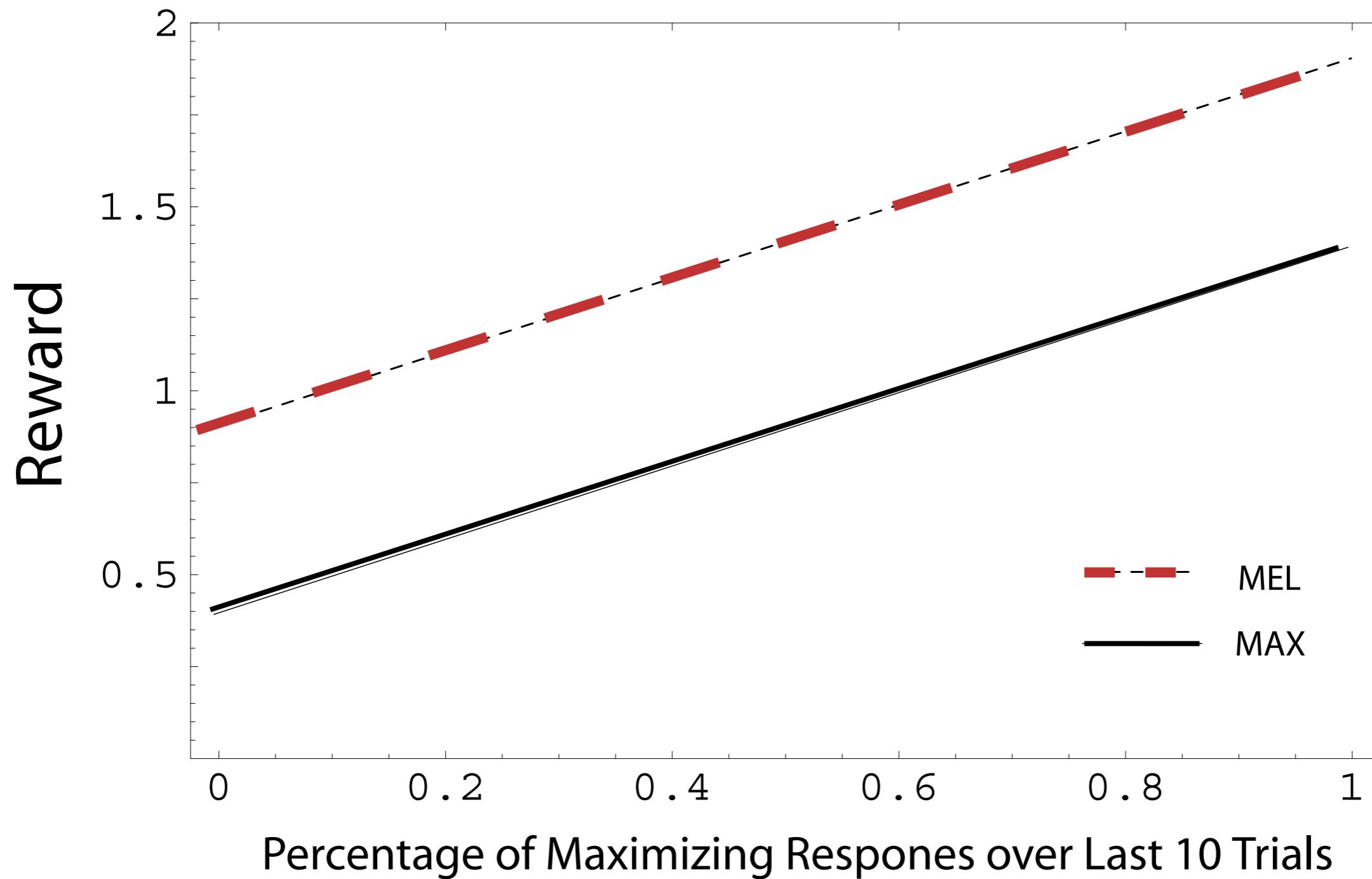


Q-learning gradient updates

- Do a feedforward pass of the network for the current state s and get predicted Q values for all actions (output nodes)
- Choose an action using e.g., softmax or epsilon greedy to encourage exploration
- Do a feedforward pass for the network for the next state (s') and get the maximum possible Q-value
- Set Q-value target for the chosen action to $r + \max_{a'} Q(s', a')$ for the chosen action and zero for all others
- Backpropagate that error as for normal network

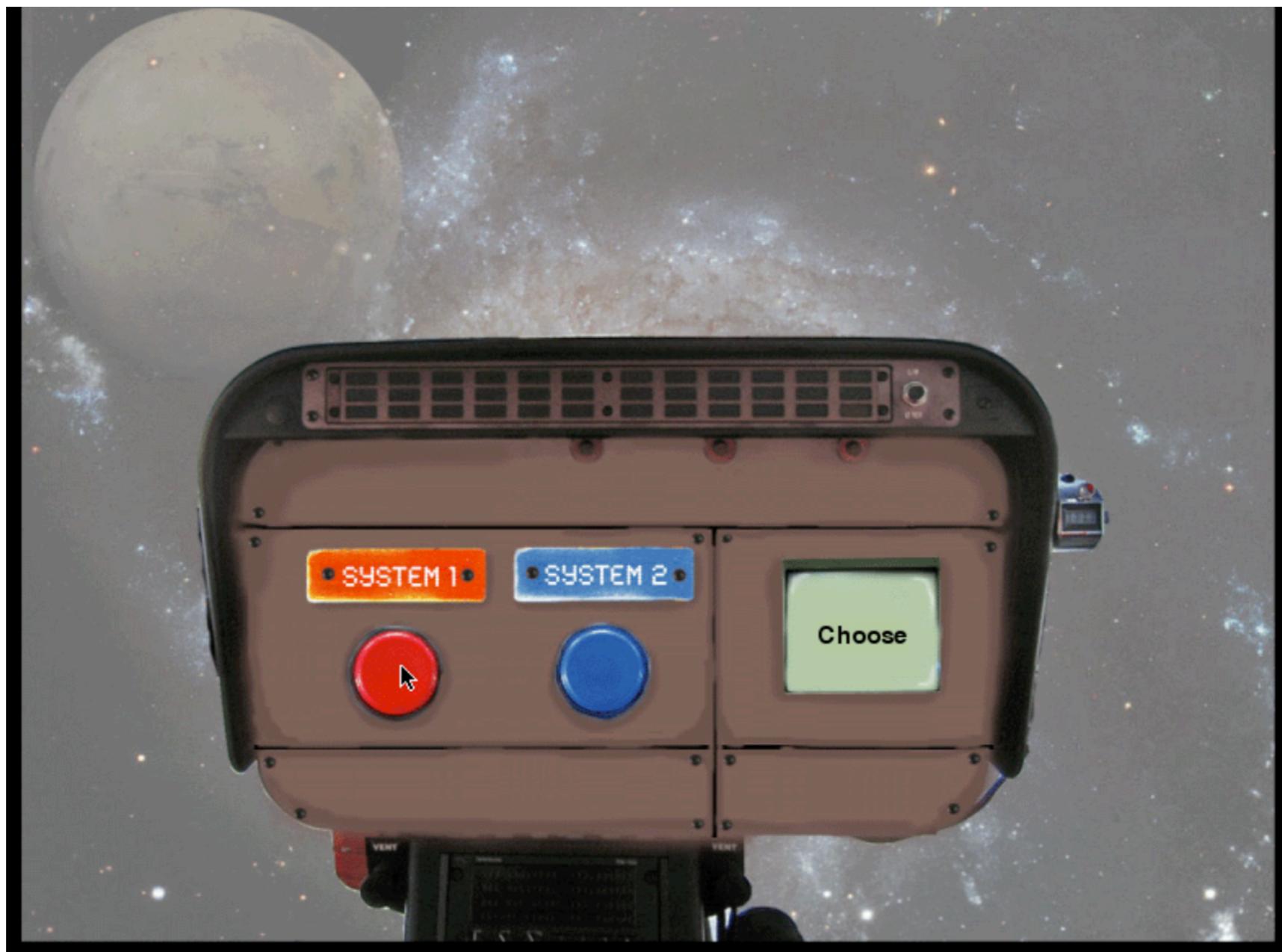


$$L = \frac{1}{2} [r + \max_{a'} Q(s', a') - Q(s, a)]^2$$

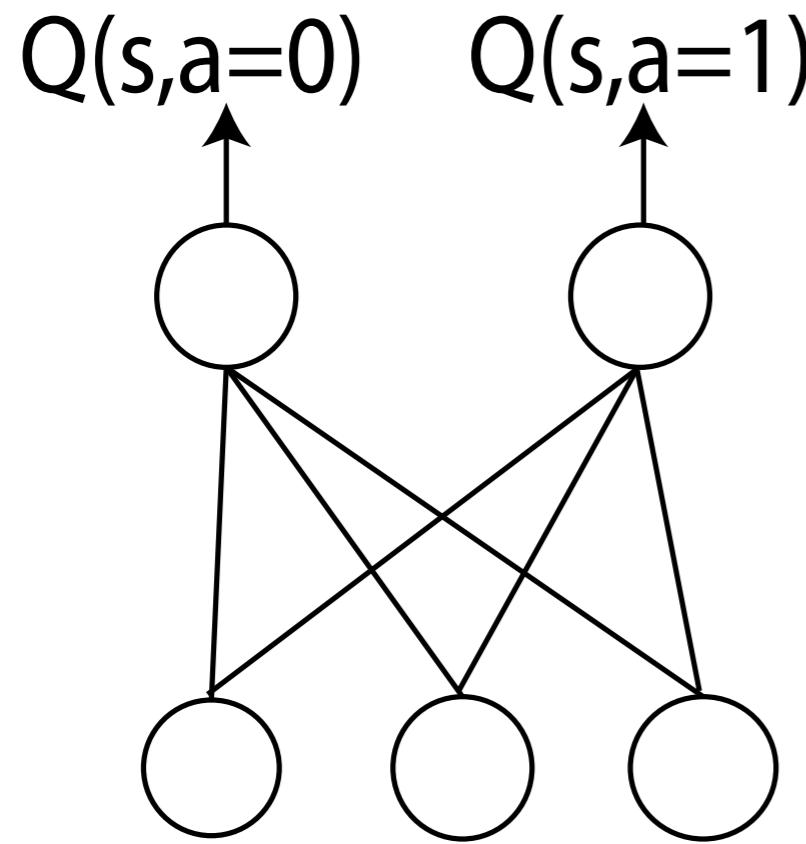




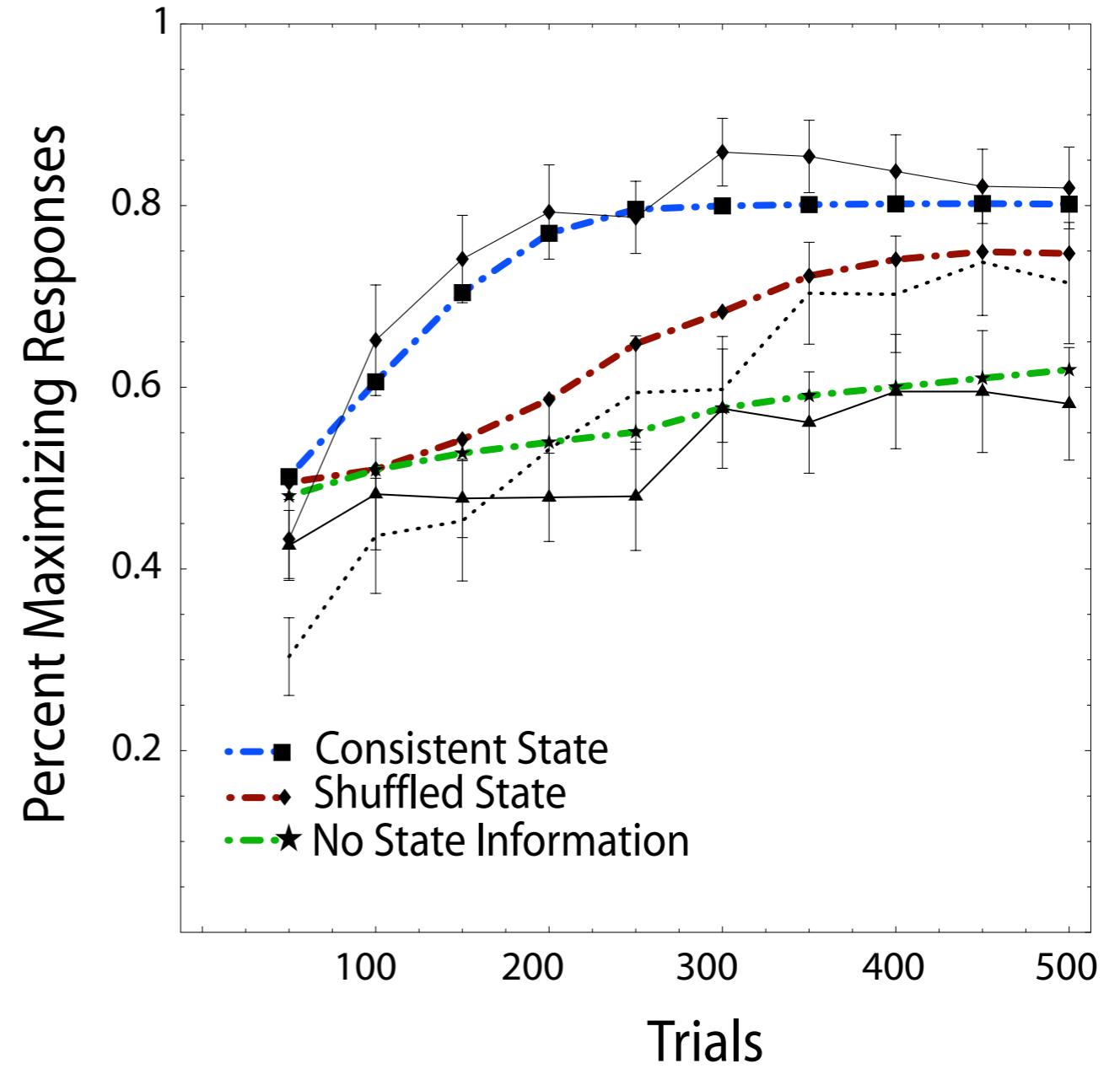
WELCOME TO THE N.A.S.A. MARS FARMING PROJECT







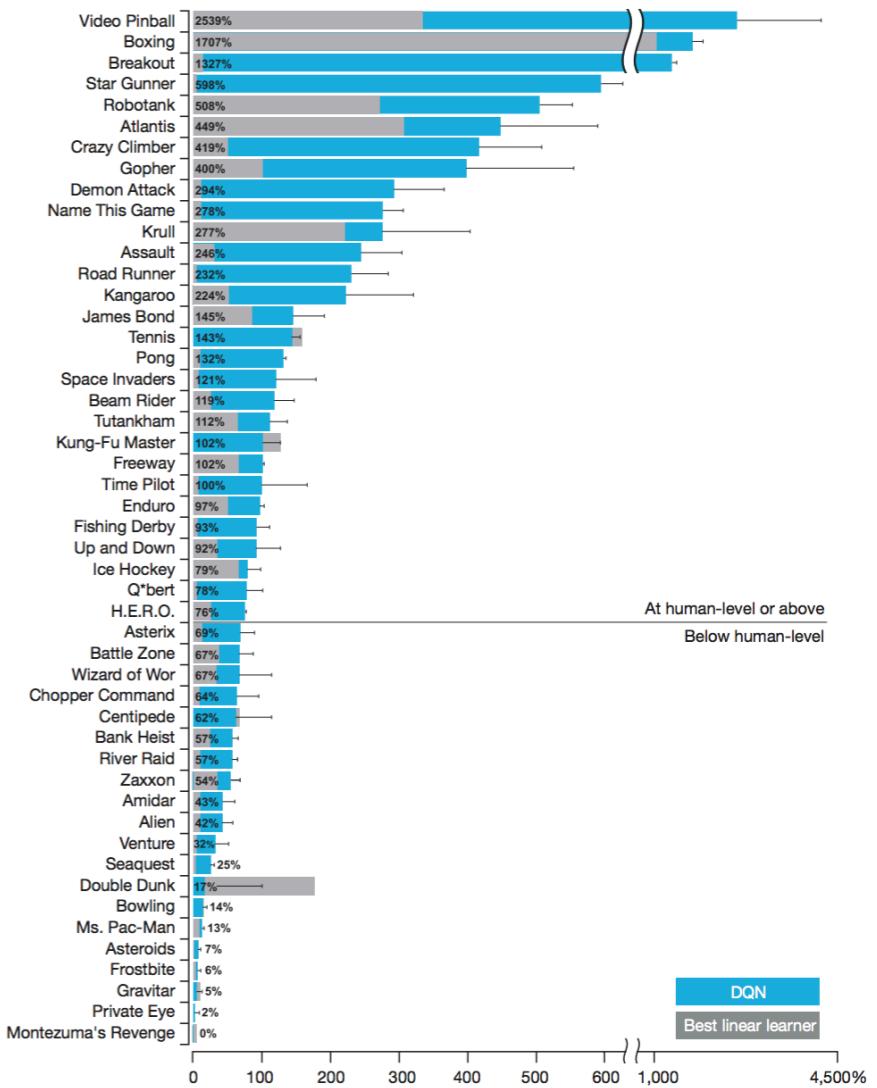
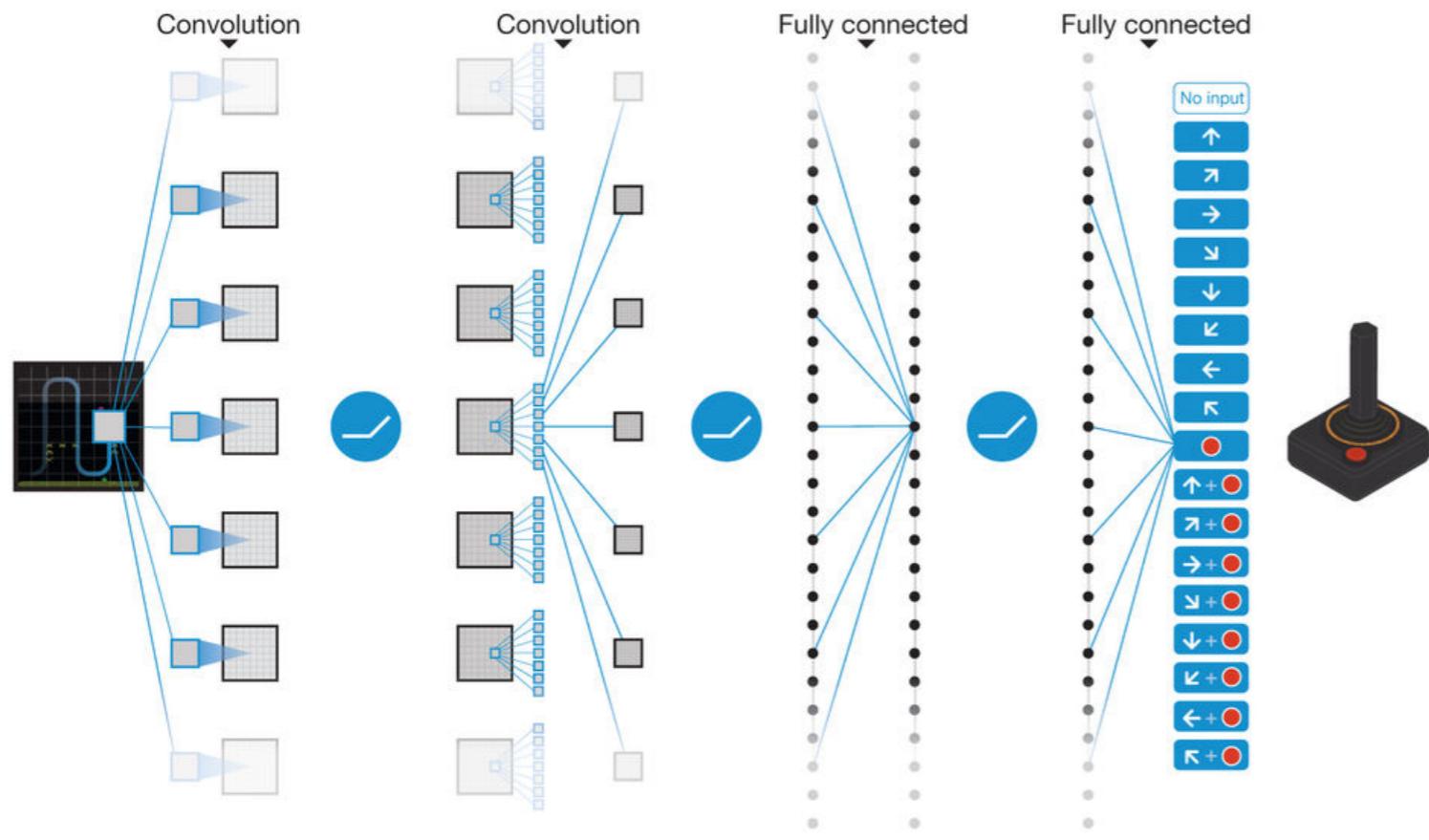
Position of Light	Last Choice [0/1]	Bias Unit [0-10]
-------------------	-------------------	------------------



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

$$L = \frac{1}{2} [r + \max_{a'} Q(s', a') - Q(s, a)]^2$$

The way the problem is represented to the network makes it easier to learn the state representation matters!



LETTER

doi:10.1038/nature14236

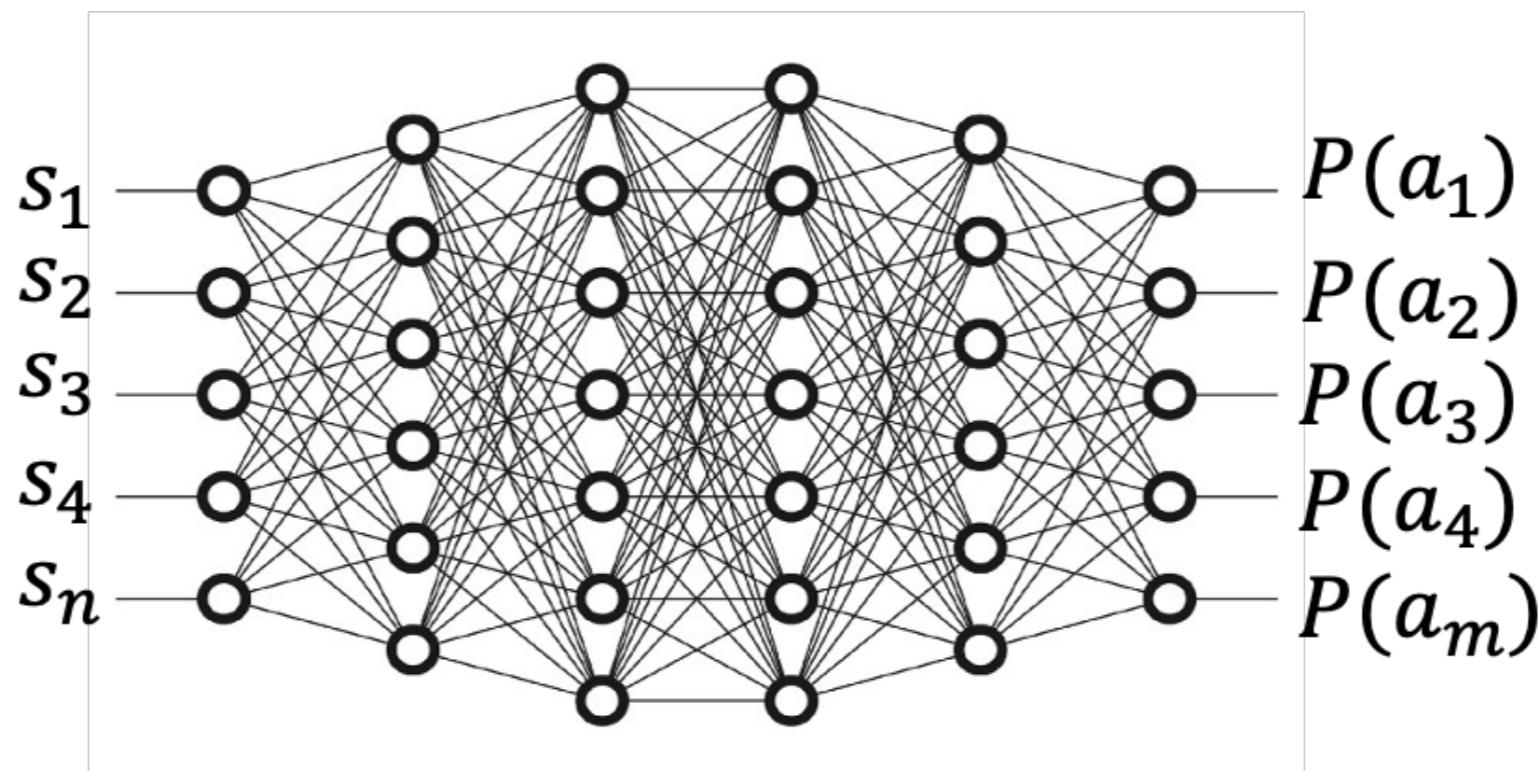
Human-level control through deep reinforcement learning

Volodymyr Mnih^{1*}, Koray Kavukcuoglu^{1*}, David Silver^{1*}, Andrei A. Rusu¹, Joel Veness¹, Marc G. Bellemare¹, Alex Graves¹, Martin Riedmiller¹, Andreas K. Fidjeland¹, Georg Ostrovski¹, Stig Petersen¹, Charles Beattie¹, Amir Sadik¹, Ioannis Antonoglou¹, Helen King¹, Dharshan Kumaran¹, Daan Wierstra¹, Shane Legg¹ & Demis Hassabis¹

DQN!

Learning Policies Directly Using Policy Gradient

- So far we have focused on learning values (v-values or Q-values) with the policy treated separately on top (e.g., epsilon-greedy or softmax)
- Instead what if we **learn policies directly?**
- Represent the policy as a continuous, differentiable function (e.g., an artificial neural network)



Learning Policies Directly Using Policy Gradient

- So far we have focused on learning values (v-values or Q-values) with the policy treated separately on top (e.g., epsilon-greedy or softmax)
- Instead what if we **learn policies directly?**
- Represent the policy as a continuous, differentiable function (e.g., an artificial neural network)

$$R_t(\tau) = \sum_{t'=t}^T \gamma^{t'-t} r_{t'}$$

The return

$$J(\pi_\theta) = \mathbb{E}_{\tau \sim \pi_\theta} [R(\tau)] = \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=0}^T \gamma^t r_t \right]$$

The objective

$$\max_{\theta} J(\pi_\theta) = \mathbb{E}_{\tau \sim \pi_\theta} [R(\tau)]$$

Aka, the problem to optimize

Learning Policies Directly Using Policy Gradient

- So far we have focused on learning values (v-values or Q-values) with the policy treated separately on top (e.g., epsilon-greedy or softmax)
- Instead what if we **learn policies directly?**
- Represent the policy as a continuous, differentiable function (e.g., an artificial neural network)

$$\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\pi_{\theta})$$

Gradient ascent in parameter Space to maximize the return

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^T R_t(\tau) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right]$$

The policy gradient

** see several sources including Sutton and Barto book for proof of policy gradient theorem, Williams (1992) REINFORCE algorithm

Learning Policies Directly Using Policy Gradient

- So far we have focused on learning values (v-values or Q-values) with the policy treated separately on top (e.g., epsilon-greedy or softmax)
 - Instead what if we **learn policies directly?**
 - Represent the policy as a continuous, differentiable function (e.g., an artificial neural network)
-
- If the Return is greater than zero than probability of action is increased, if less than zero probability of action is decreased (Thorndike law of effect!)
 - Can alter discrete action probabilities but also apply to continuous action parameters like the mean and standard deviation of a Gaussian!
 - Combined neural networks with the problem of dynamic control!

Algorithm 2.1 REINFORCE algorithm

```
1: Initialize learning rate  $\alpha$ 
2: Initialize weights  $\theta$  of a policy network  $\pi_\theta$ 
3: for  $episode = 0, \dots, MAX\_EPISODE$  do
4:   Sample a trajectory  $\tau = s_0, a_0, r_0, \dots, s_T, a_T, r_T$ 
5:   Set  $\nabla_\theta J(\pi_\theta) = 0$ 
6:   for  $t = 0, \dots, T$  do
7:      $R_t(\tau) = \sum_{t'=t}^T \gamma^{t'-t} r'_t$ 
8:      $\nabla_\theta J(\pi_\theta) = \nabla_\theta J(\pi_\theta) + R_t(\tau) \nabla_\theta \log \pi_\theta(a_t | s_t)$ 
9:   end for
10:   $\theta = \theta + \alpha \nabla_\theta J(\pi_\theta)$ 
11: end for
```

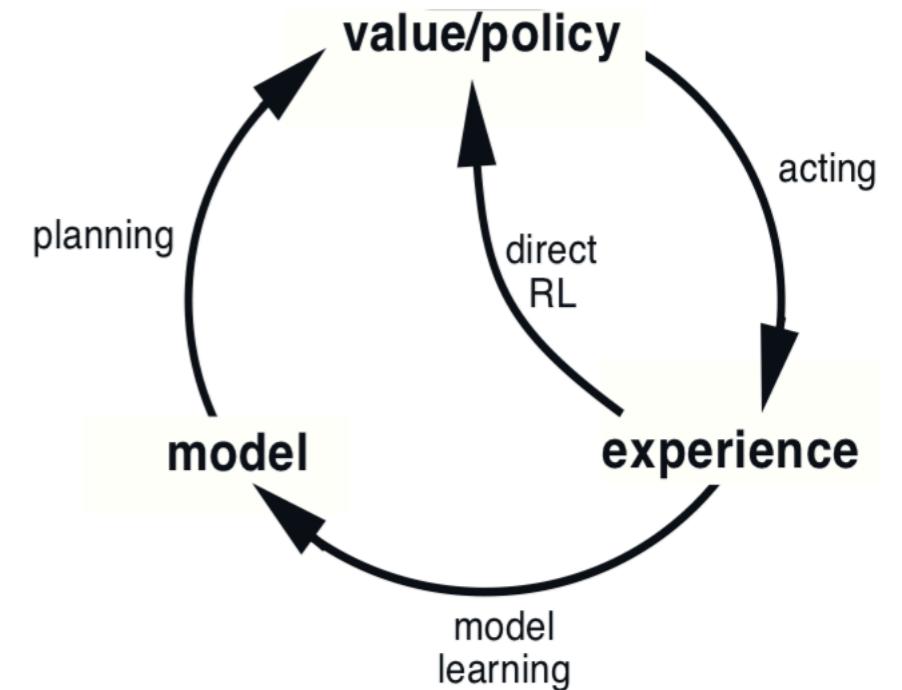
Planning or Model-based Learning

Temporal-Difference/Monte-carlo methods we have considered so far are called “model-free” because they don’t maintain or learn information about the state transitions and rewards explicitly.

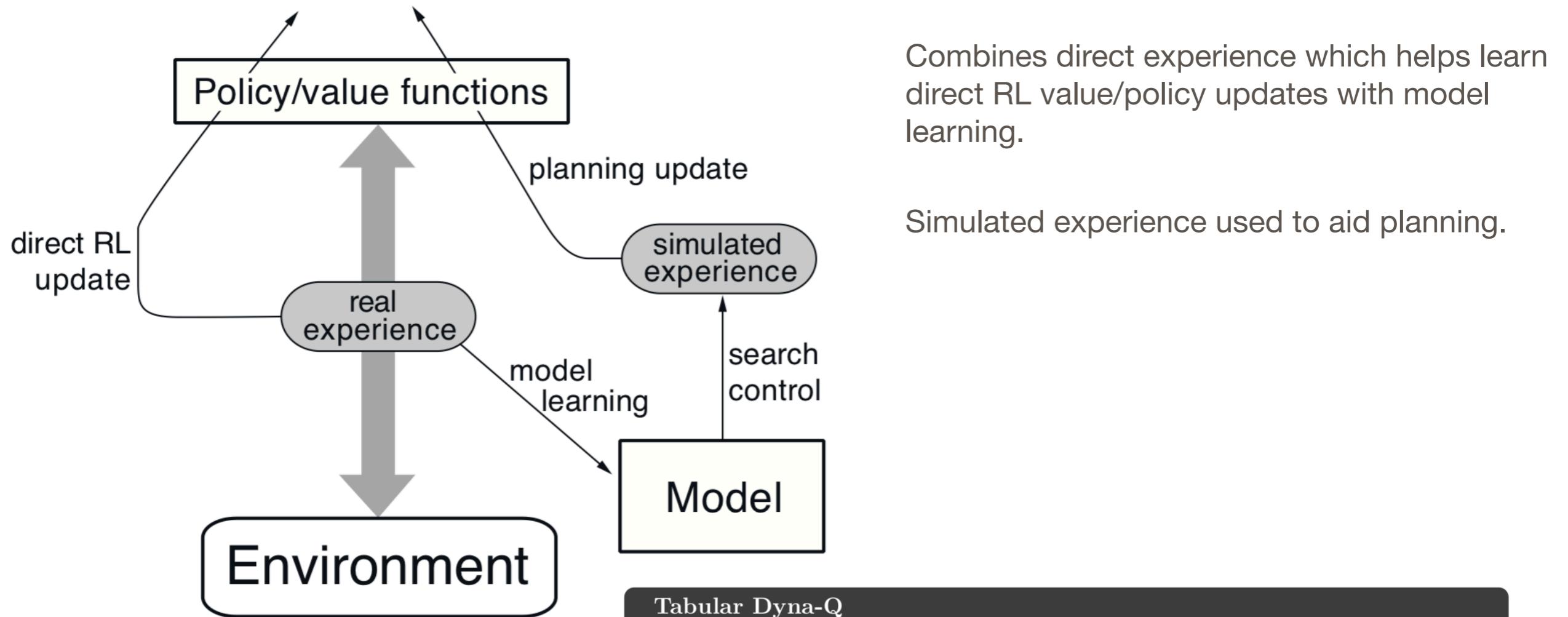
In contrast, Dynamic Programming needs a model of the environment in terms of the transition probabilities and rewards.

If you have a model you can distinguish **planning** from **learning**.

Planning - predict what will happen in the future... what state will I end up in, what reward will I get?



The DYNA algorithm



Tabular Dyna-Q

Initialize $Q(s, a)$ and $Model(s, a)$ for all $s \in \mathcal{S}$ and $a \in \mathcal{A}(s)$

Loop forever:

- (a) $S \leftarrow$ current (nonterminal) state
- (b) $A \leftarrow \varepsilon\text{-greedy}(S, Q)$
- (c) Take action A ; observe resultant reward, R , and state, S'
- (d) $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$
- (e) $Model(S, A) \leftarrow R, S'$ (assuming deterministic environment)
- (f) Loop repeat n times:

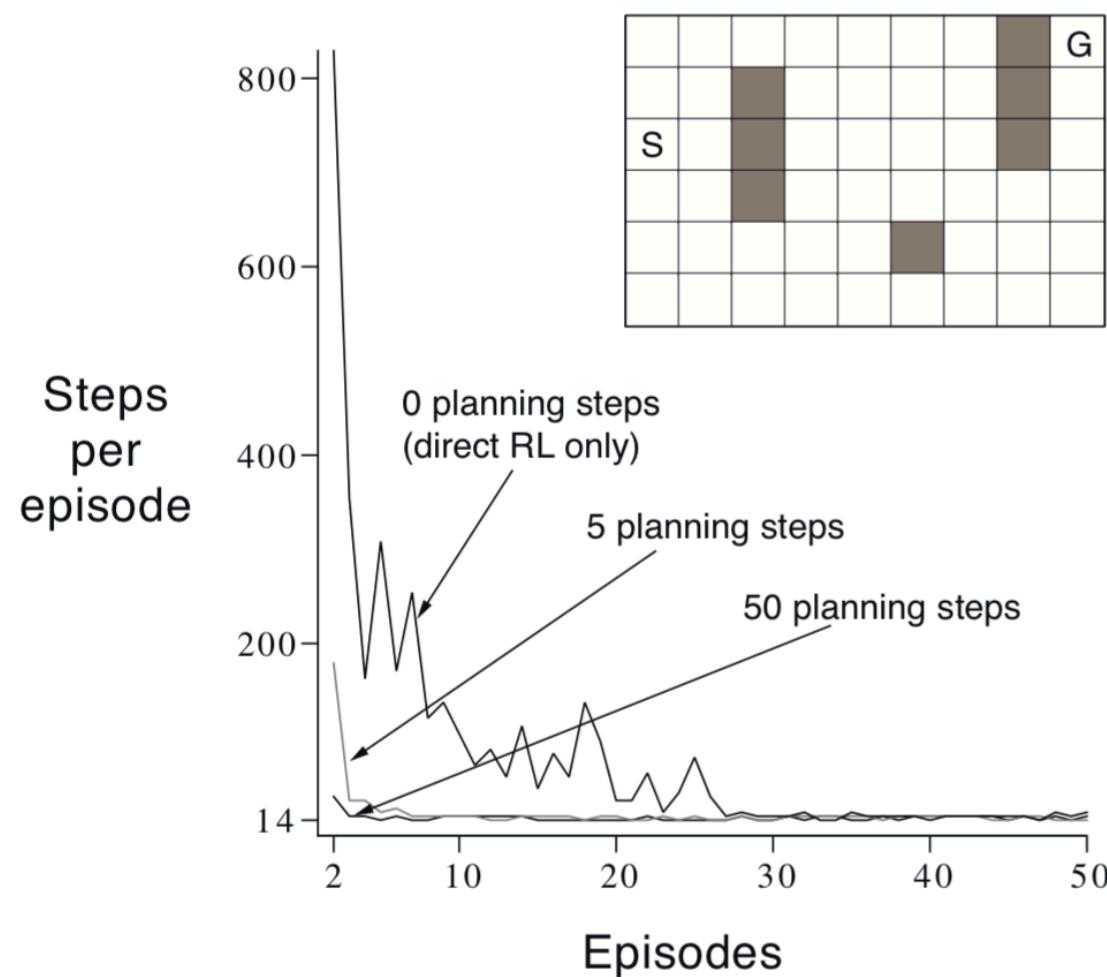
$S \leftarrow$ random previously observed state

$A \leftarrow$ random action previously taken in S

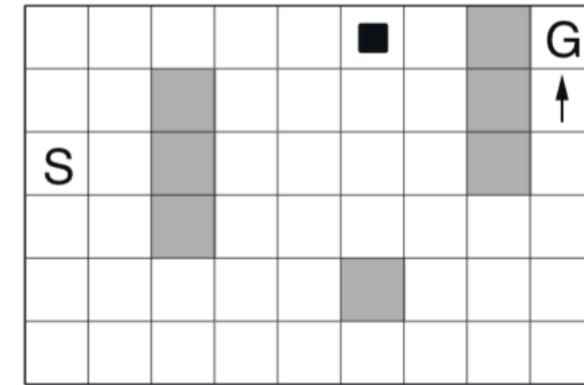
$R, S' \leftarrow Model(S, A)$

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$

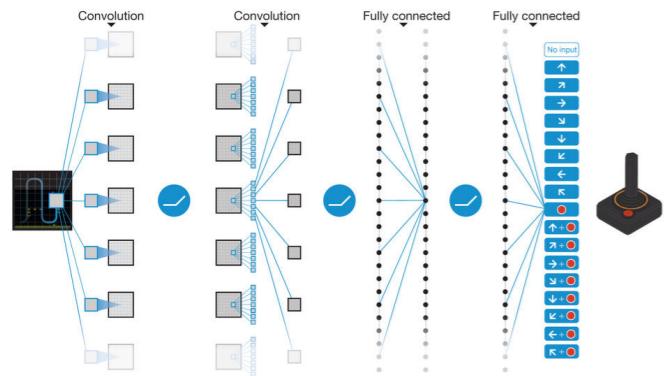
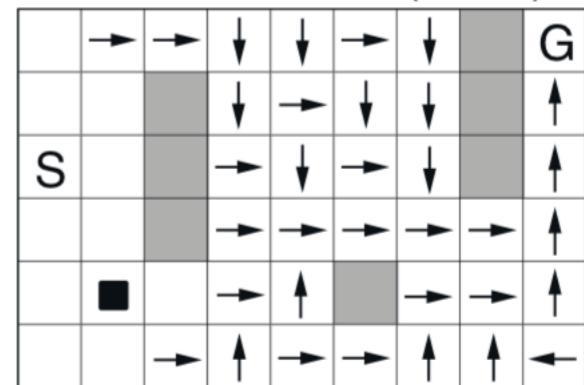
The DYNA algorithm



WITHOUT PLANNING ($n=0$)

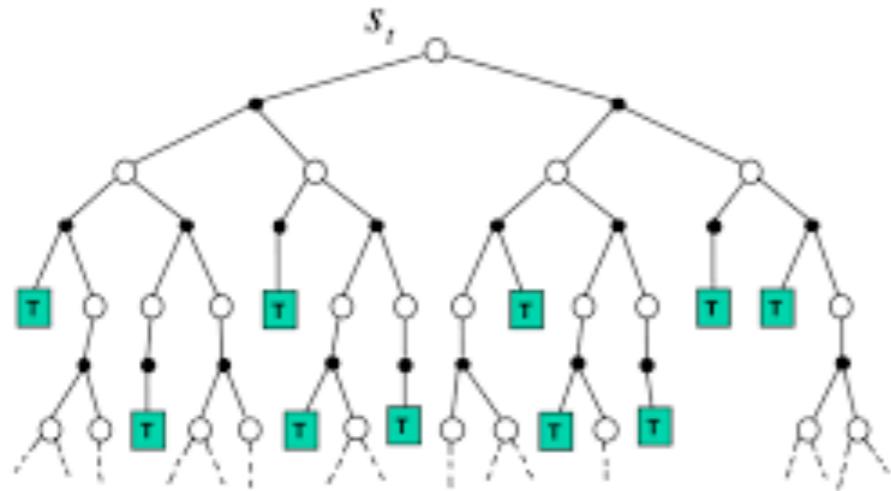


WITH PLANNING ($n=50$)



“experience replay” also used in Atari Solution

Online planning

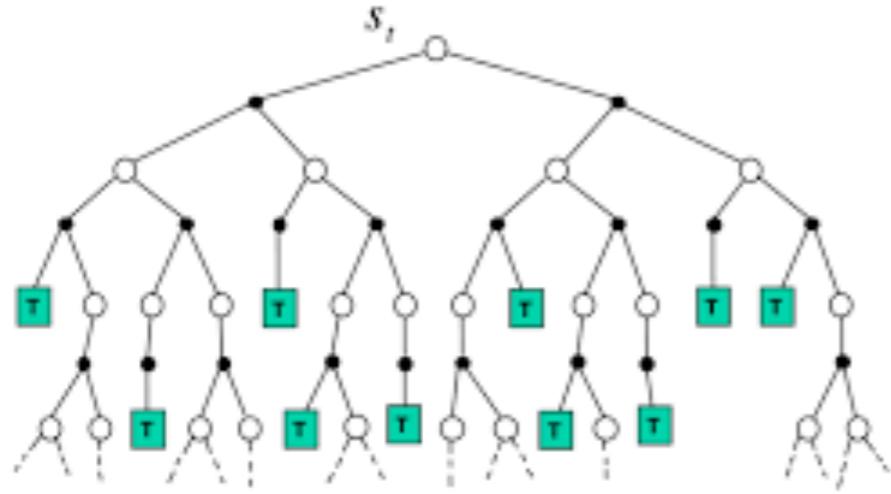


From each state, if we have a good model of environment we can simulate forward many times to determine the value of the current state or action.

However, the branching factor of this tree can be a major problem.

One recent innovation is **Monte Carlo Tree Search (MCTS)**

Monte Carlo Tree Search



$$\bar{x}_i \pm \sqrt{\frac{2 \ln n}{n_i}}$$

- \bar{x}_i : the mean payout for machine i
- n_i : the number of plays of machine i
- n : the total number of plays
- Pick option with highest confidence bound

UCB1

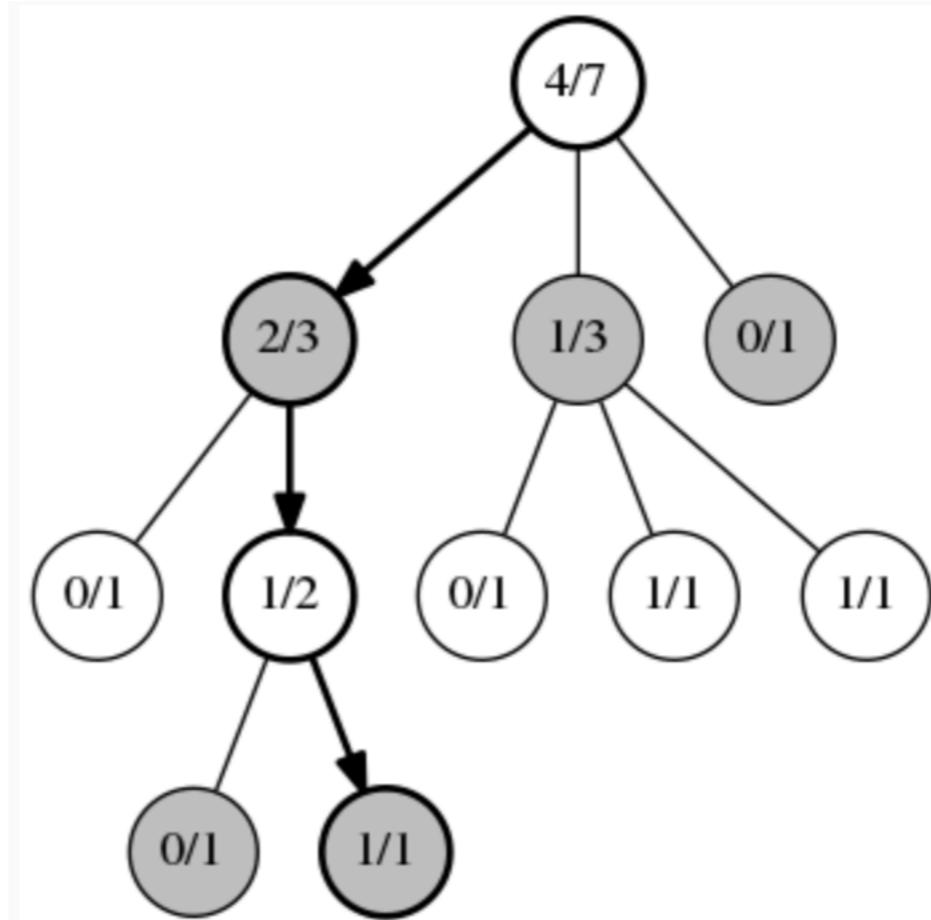
- In a standard Monte Carlo process, a large number of random simulations are run, in this case, from the board position that you want to find the best move for.
- The downside to this method, though, is that for any given turn in the simulation, there may be many possible moves, but only one or two that are good.
- Instead of doing many purely random simulations, UCB works by doing many multi-phase playouts where each position is itself viewed as a multi-armed bandit.

Monte Carlo Tree Search

$$\bar{x}_i \pm \sqrt{\frac{2 \ln n}{n_i}}$$

- \bar{x}_i : the mean payout for machine i
- n_i : the number of plays of machine i
- n : the total number of plays

UCB1



Selection

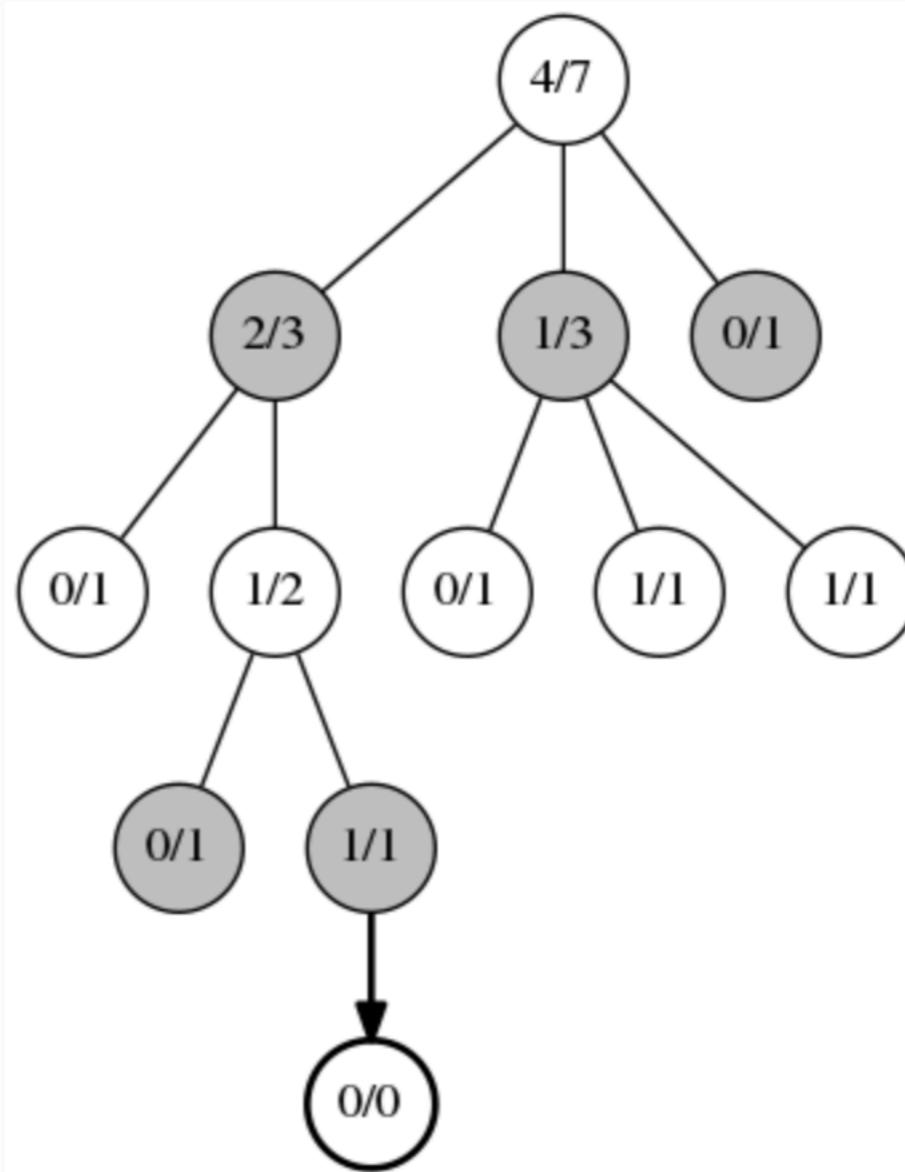
Here the positions and moves selected by the UCB1 algorithm at each step are marked in bold. Note that a number of playouts have already been run to accumulate the statistics shown. Each circle contains the number of wins / number of times played.

Monte Carlo Tree Search

$$\bar{x}_i \pm \sqrt{\frac{2 \ln n}{n_i}}$$

- \bar{x}_i : the mean payout for machine i
- n_i : the number of plays of machine i
- n : the total number of plays

UCB1



Expansion

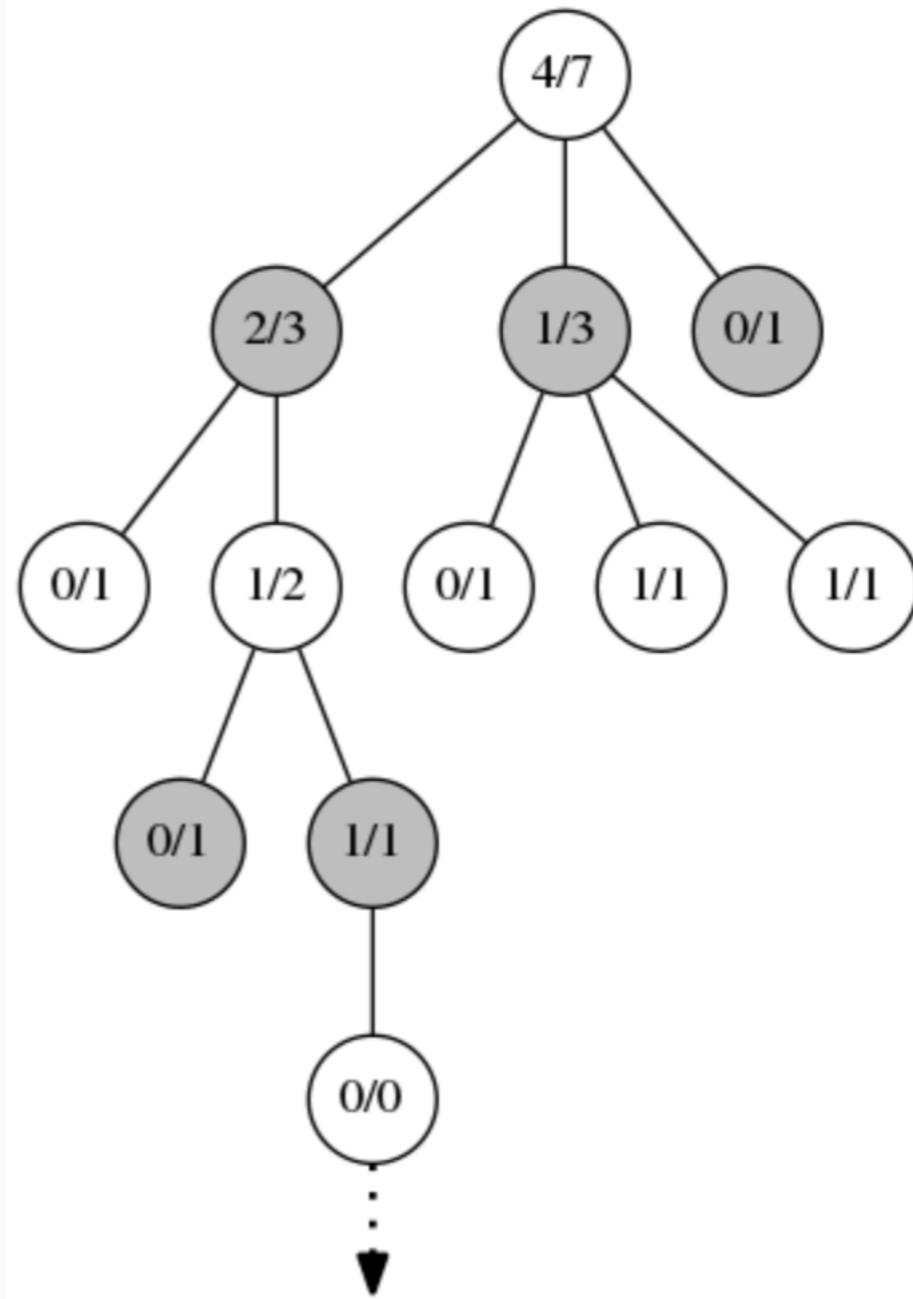
The position marked 1/1 at the bottom of the tree has no further statistics records under it, so we choose a random move and add a new record for it (bold), initialized to 0/0.

Monte Carlo Tree Search

$$\bar{x}_i \pm \sqrt{\frac{2 \ln n}{n_i}}$$

- \bar{x}_i : the mean payout for machine i
- n_i : the number of plays of machine i
- n : the total number of plays

UCB1



Simulation

Once the new record is added, the Monte Carlo simulation begins, here depicted with a dashed arrow. Moves in the simulation may be completely random, or may use calculations to weight the randomness in favor of moves that may be better.

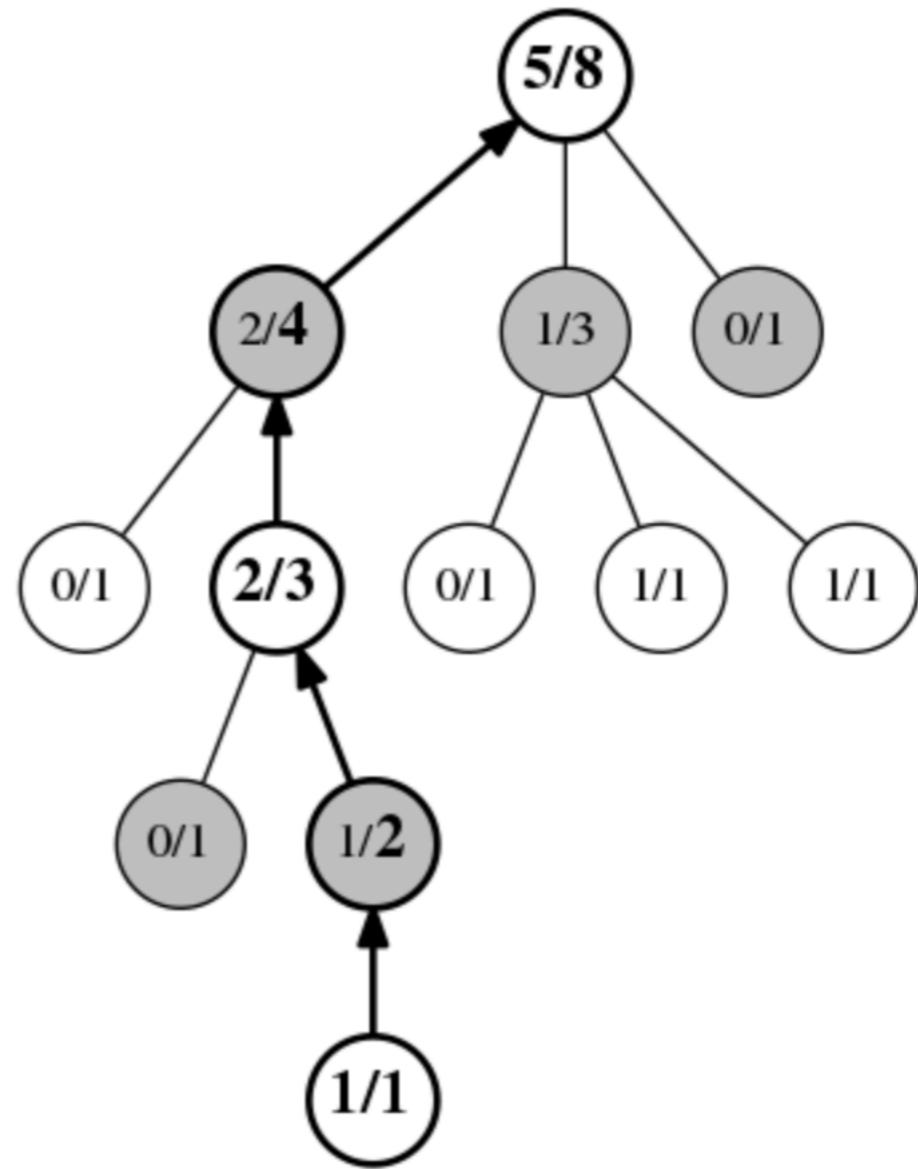
Monte Carlo Tree Search

- As more and more playouts are run, the tree of statistics grows in memory and the move that will finally be chosen will converge towards the actual optimal play, though that may take a very long time, depending on the game.

$$\bar{x}_i \pm \sqrt{\frac{2 \ln n}{n_i}}$$

- \bar{x}_i : the mean payout for machine i
 - n_i : the number of plays of machine i
 - n : the total number of plays

UCB1



Back-Propagation

After the simulation reaches an end, all of the records in the path taken are updated. Each has its play count incremented by one, and each that matches the winner has its win count incremented by one, here shown by the bolded numbers.

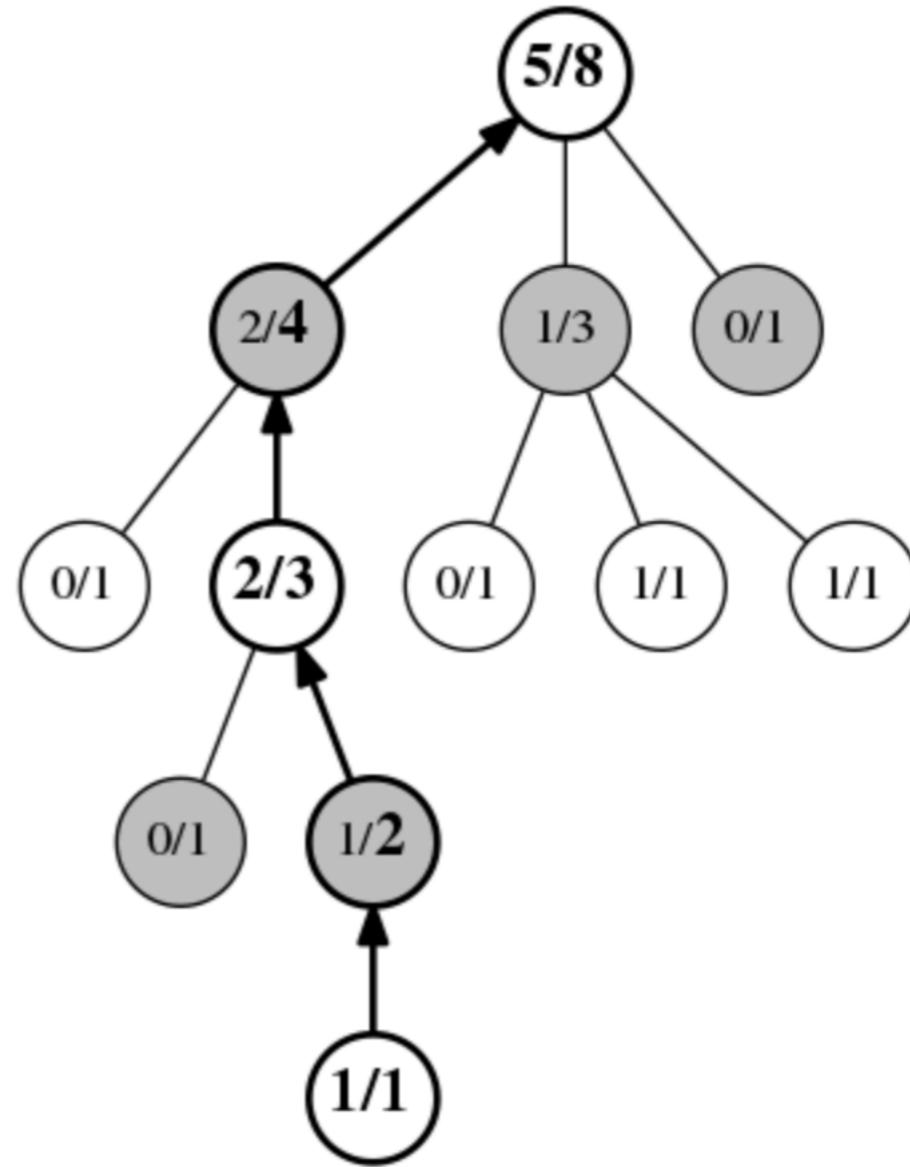
Monte Carlo Tree Search

- As more and more playouts are run, the tree of statistics grows in memory and the move that will finally be chosen will converge towards the actual optimal play, though that may take a very long time, depending on the game.

$$\bar{x}_i \pm \sqrt{\frac{2 \ln n}{n_i}}$$

- \bar{x}_i : the mean payout for machine i
- n_i : the number of plays of machine i
- n : the total number of plays

UCB1



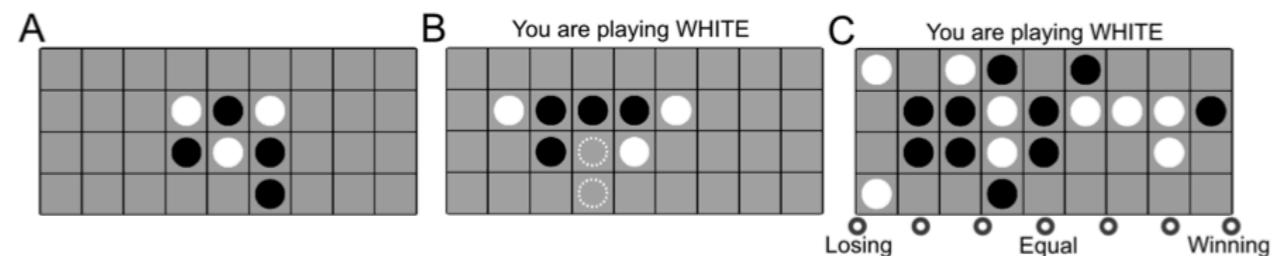
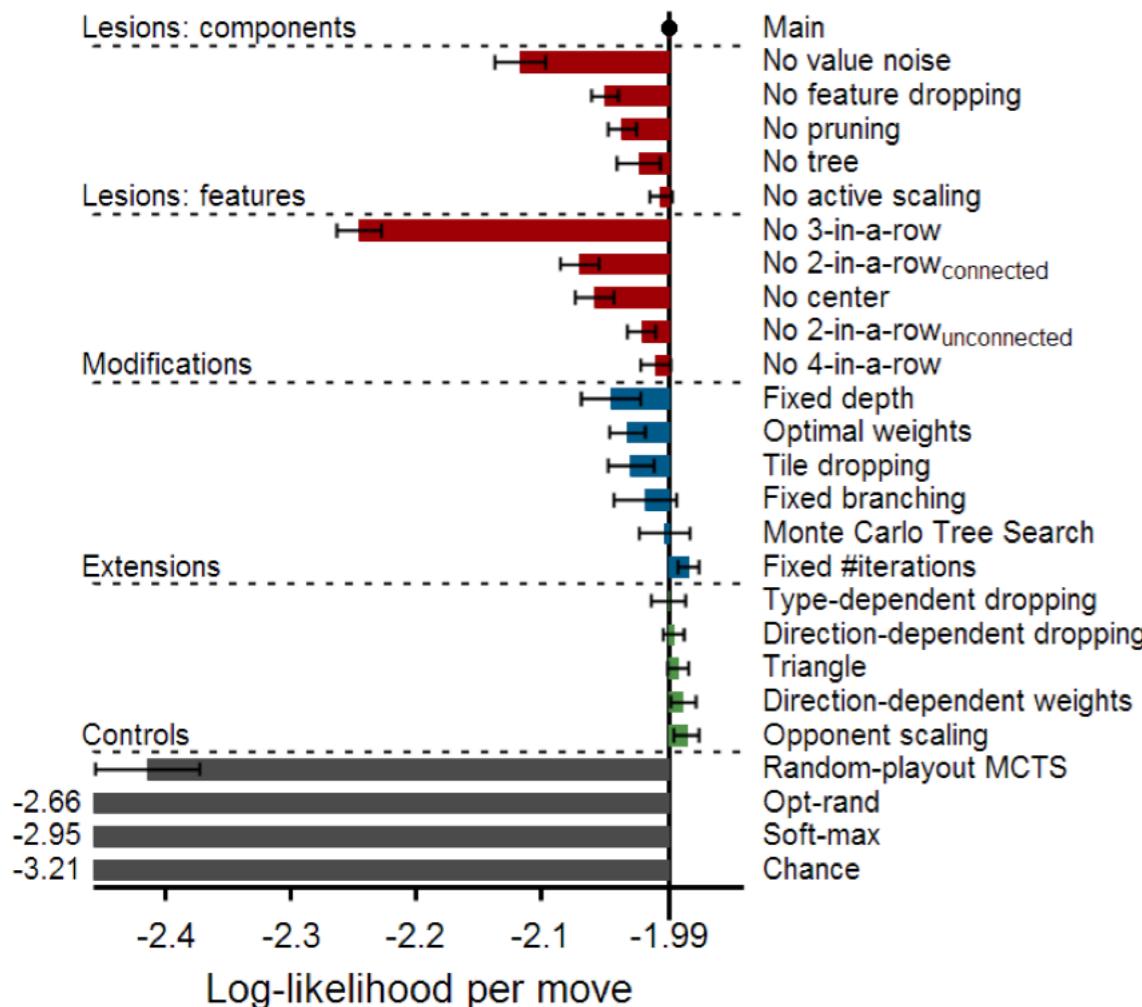
Back-Propagation

After the simulation reaches an end, all of the records in the path taken are updated. Each has its play count incremented by one, and each that matches the winner has its win count incremented by one, here shown by the bolded numbers.

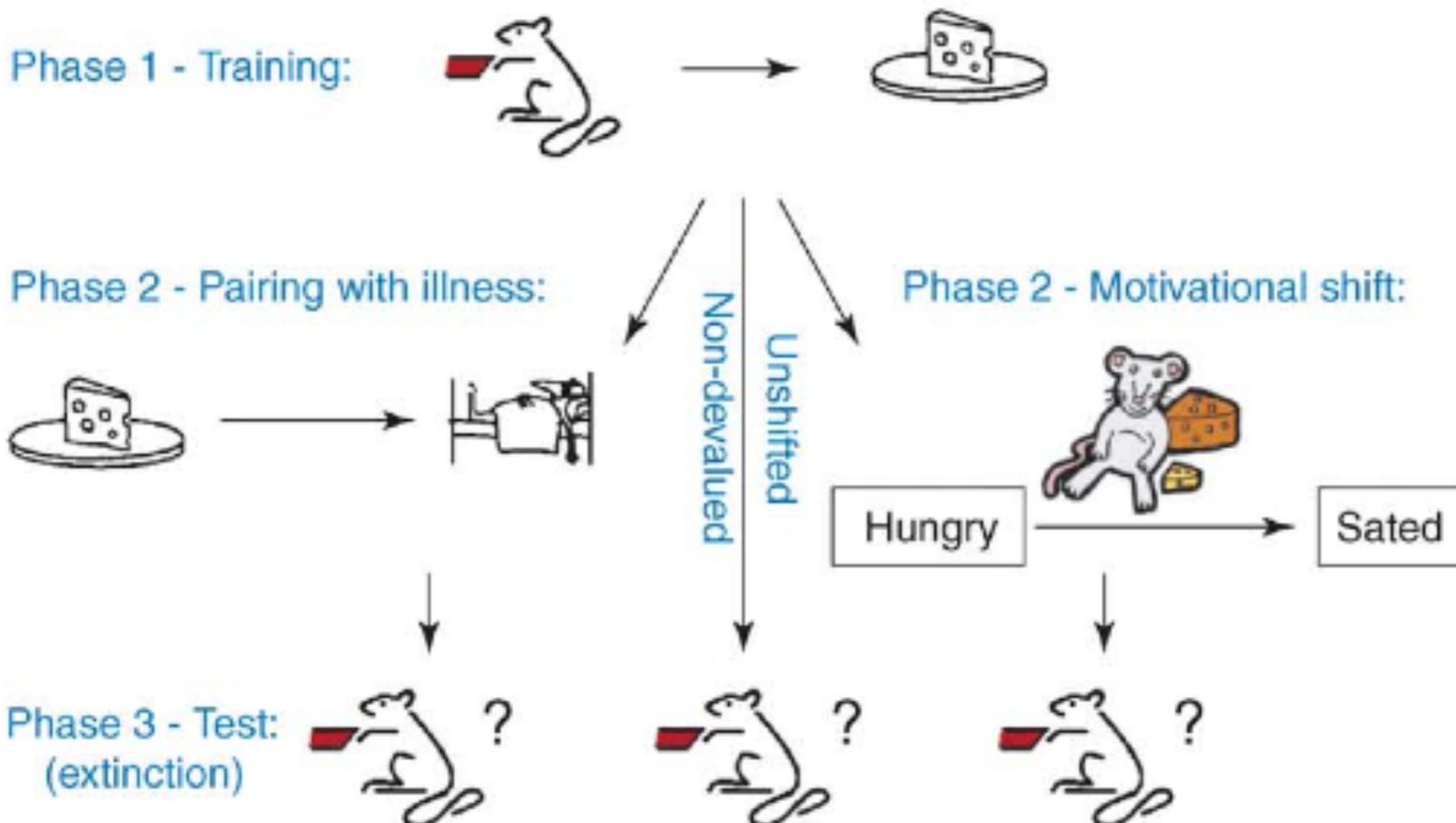
Human Tree Search?

Van Opheusden B, Galbiati G, Bnaya Z, Li Y, Ma WJ (2017),
A computational model for decision tree search.
Proceedings of the 39th Annual Meeting of the Cognitive Science Society, 1254-1259.

$$V(s) = c_{\text{self}} \sum_{i=0}^4 w_i f_i(s, \text{self}) - c_{\text{opp}} \sum_{i=0}^4 w_i f_i(s, \text{opponent})$$



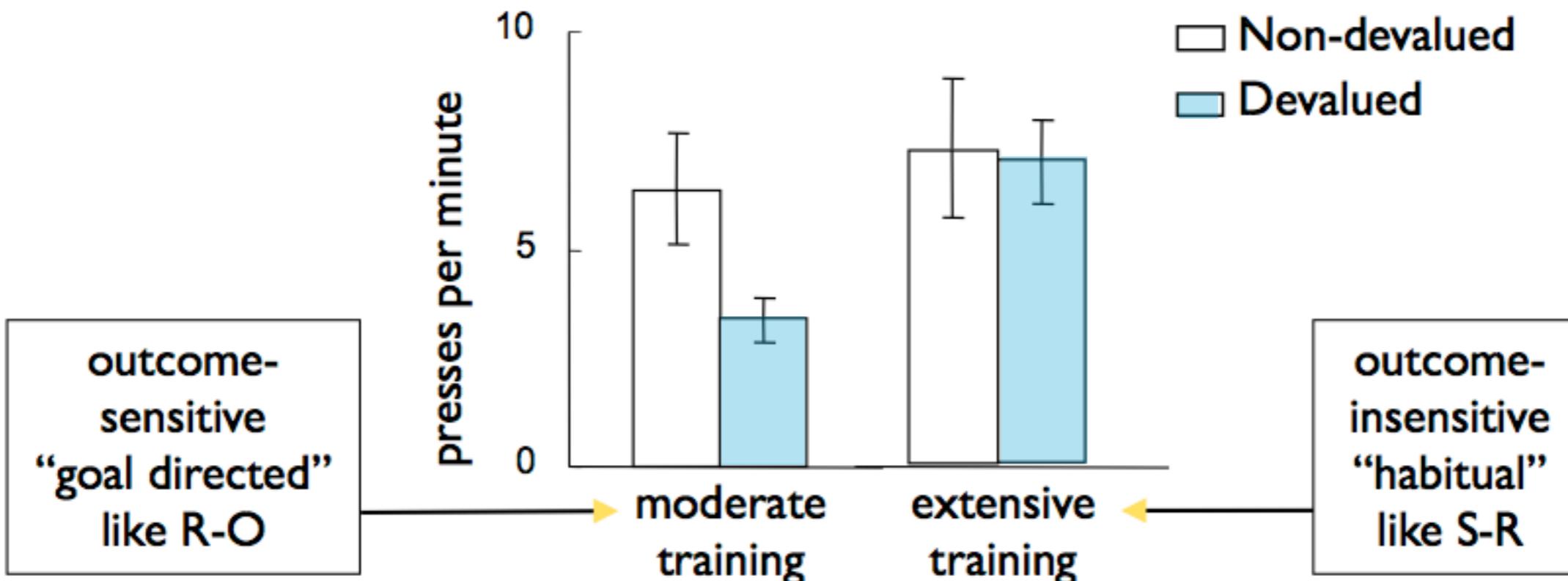
Do animals plan?



TRENDS in Cognitive Sciences

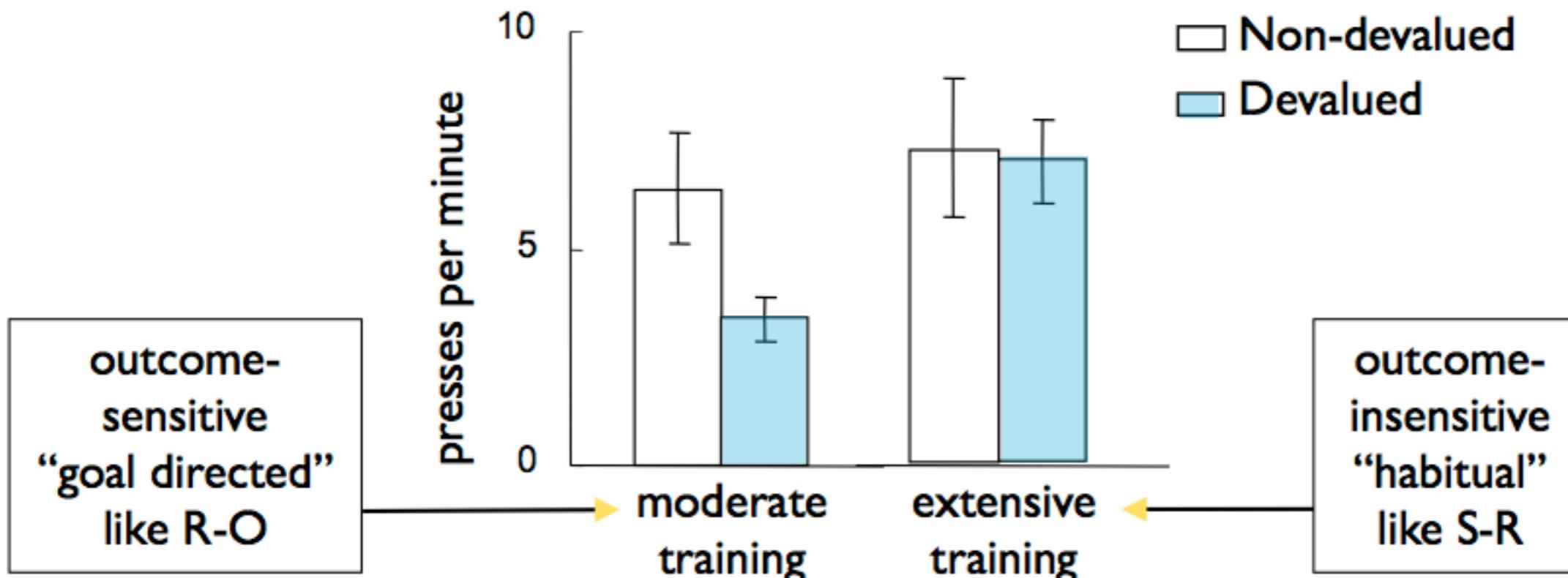
A good test: Outcome devaluation

Holland (2004) see also Dickinson article for examples



A good test: Outcome devaluation

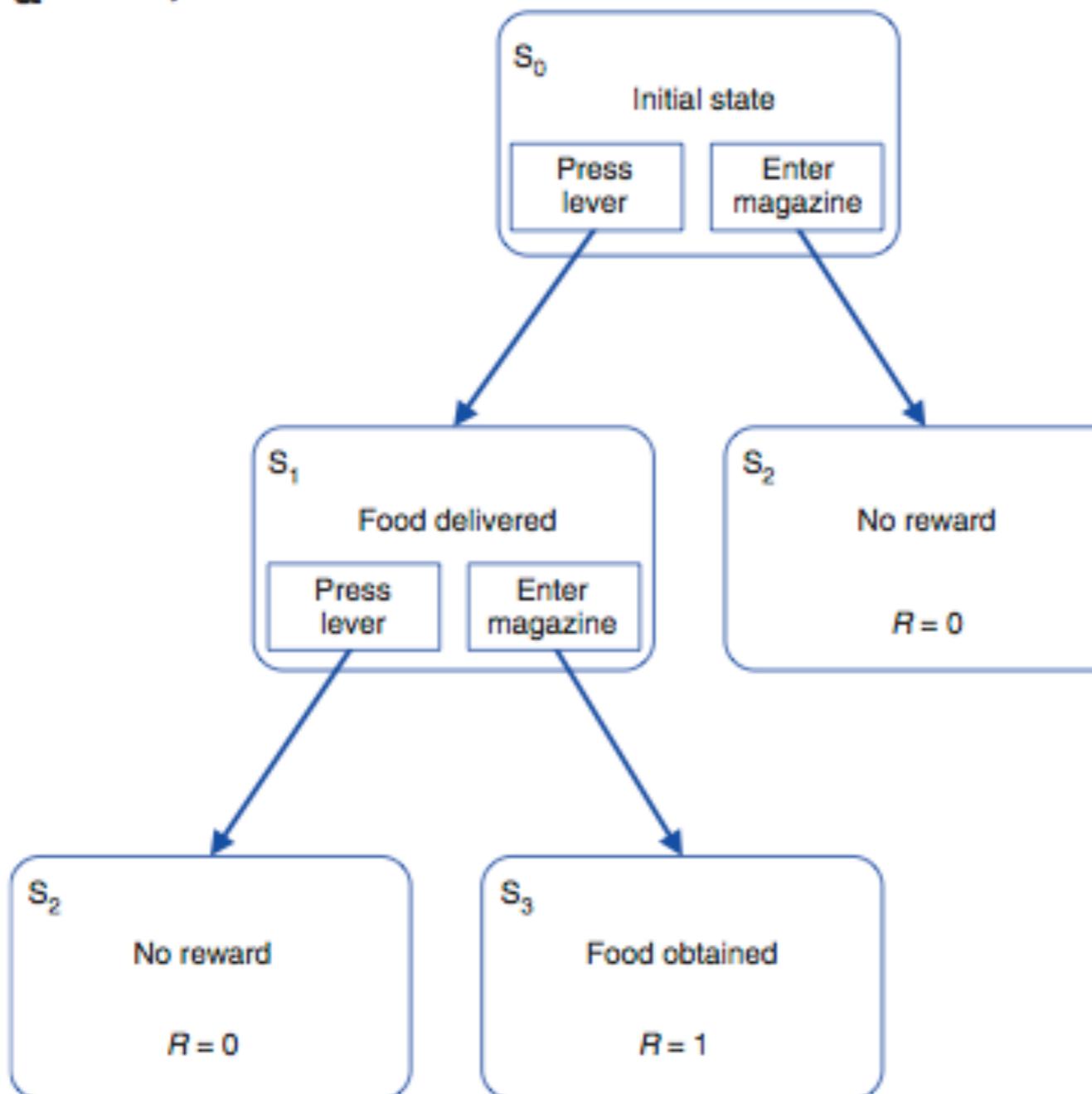
Holland (2004) see also Dickinson article for examples



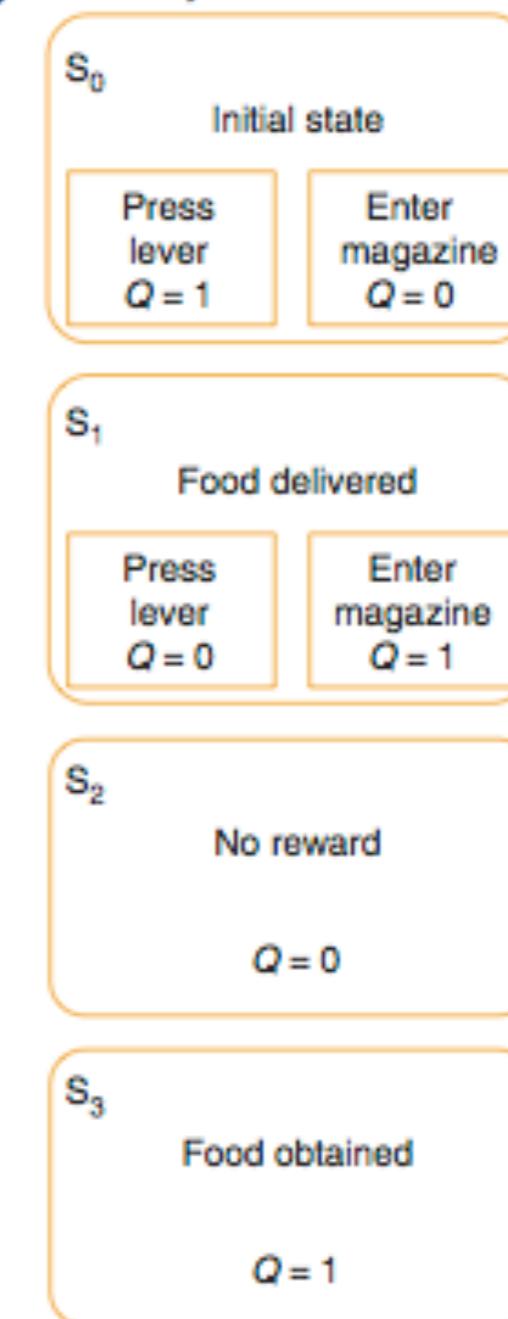
Suggests some measure of
reward or outcome **expectancy**
(shifts to habitual control with
extensive training)

Key idea

a Tree System



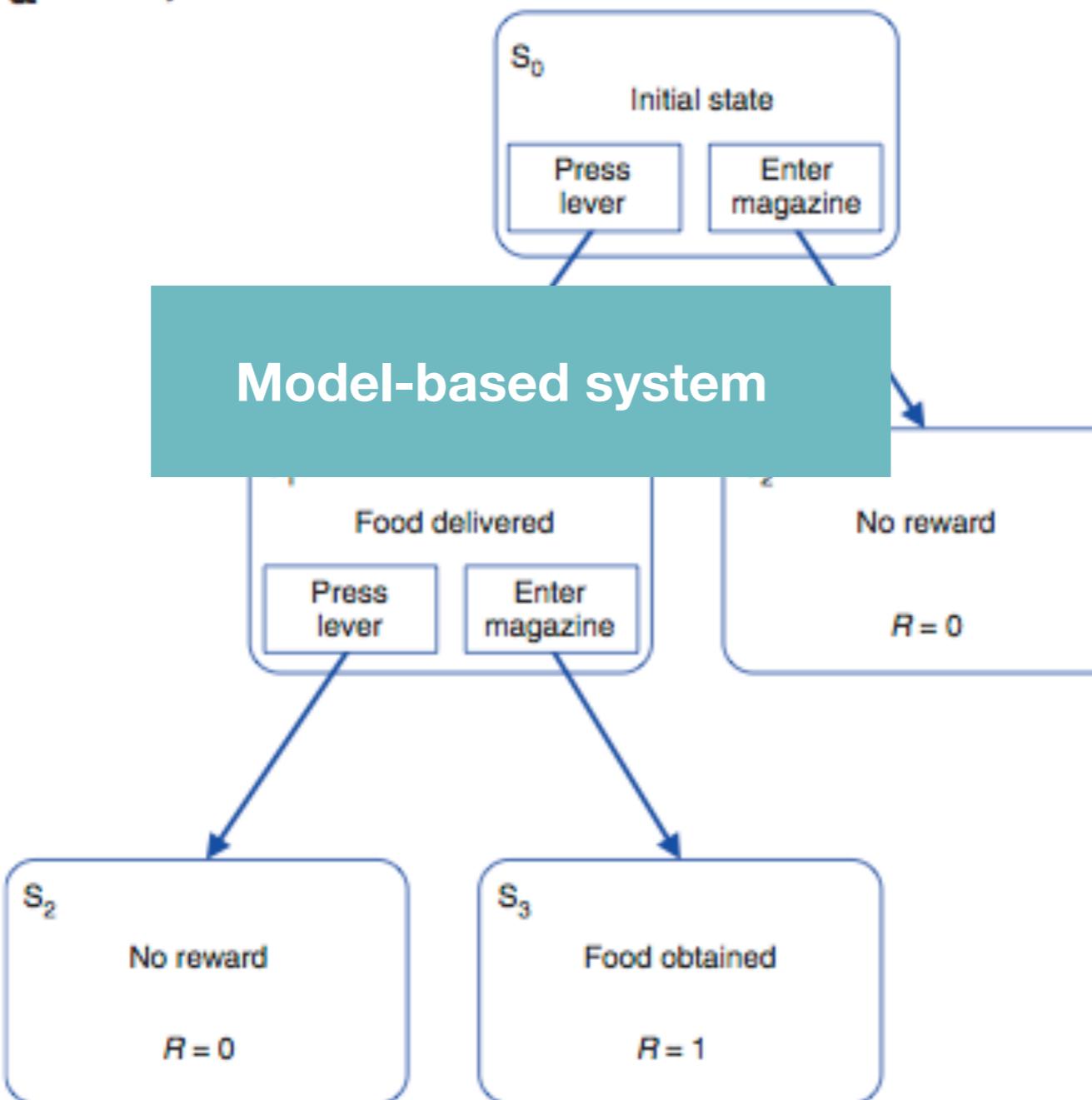
b Cache system



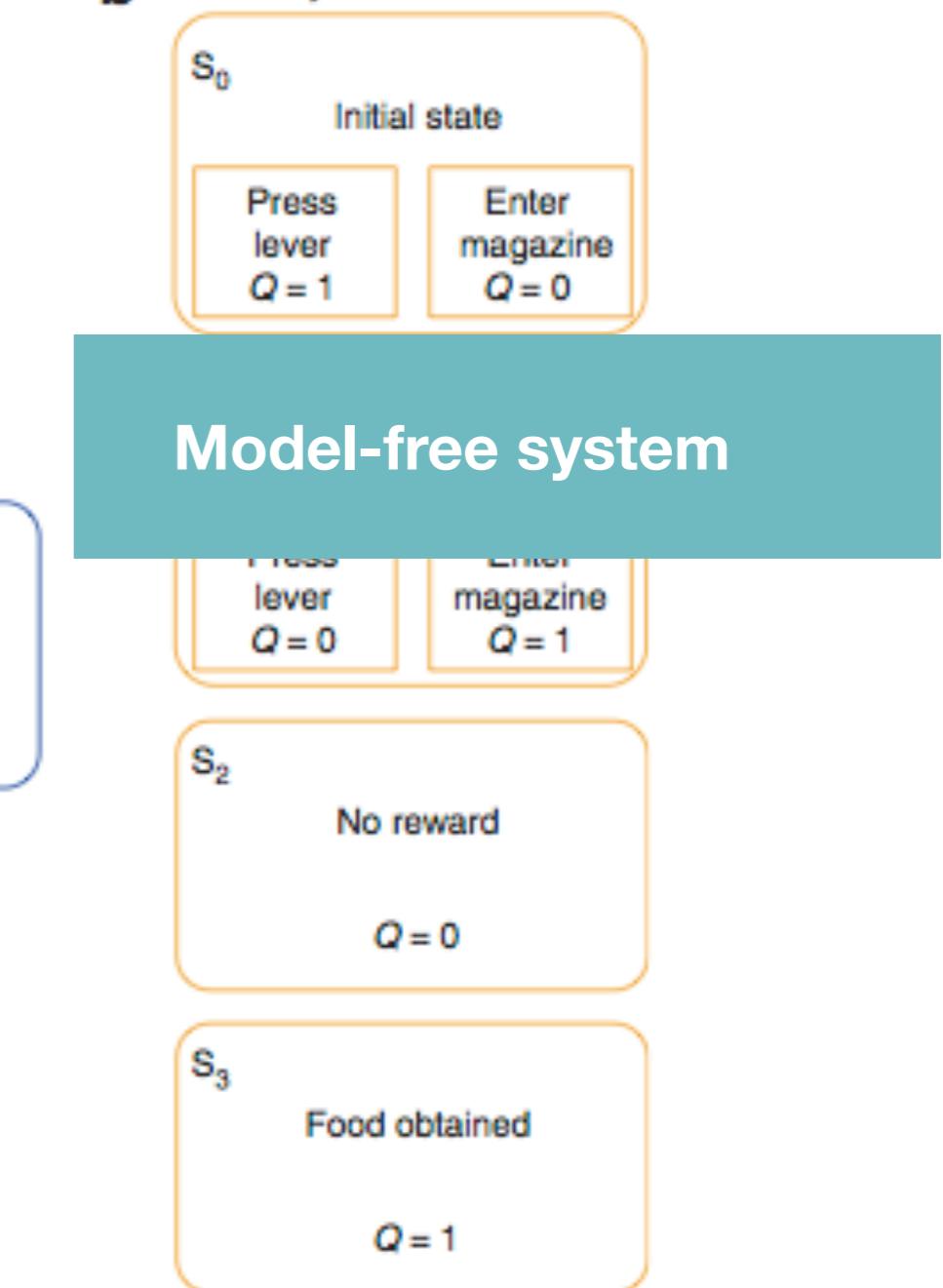
from Daw, Niv, Dayan, 2005

Key idea

a Tree System



b Cache system



from Daw, Niv, Dayan, 2005

two big questions

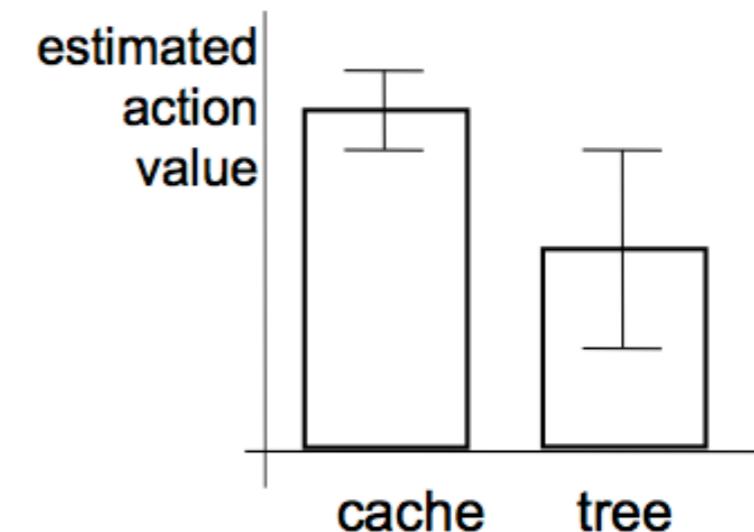
- Why should the brain use two different strategies/controllers in parallel?
- If it uses two: how can it arbitrate between the two when they disagree (*new decision making problem...*)



from Yael Niv

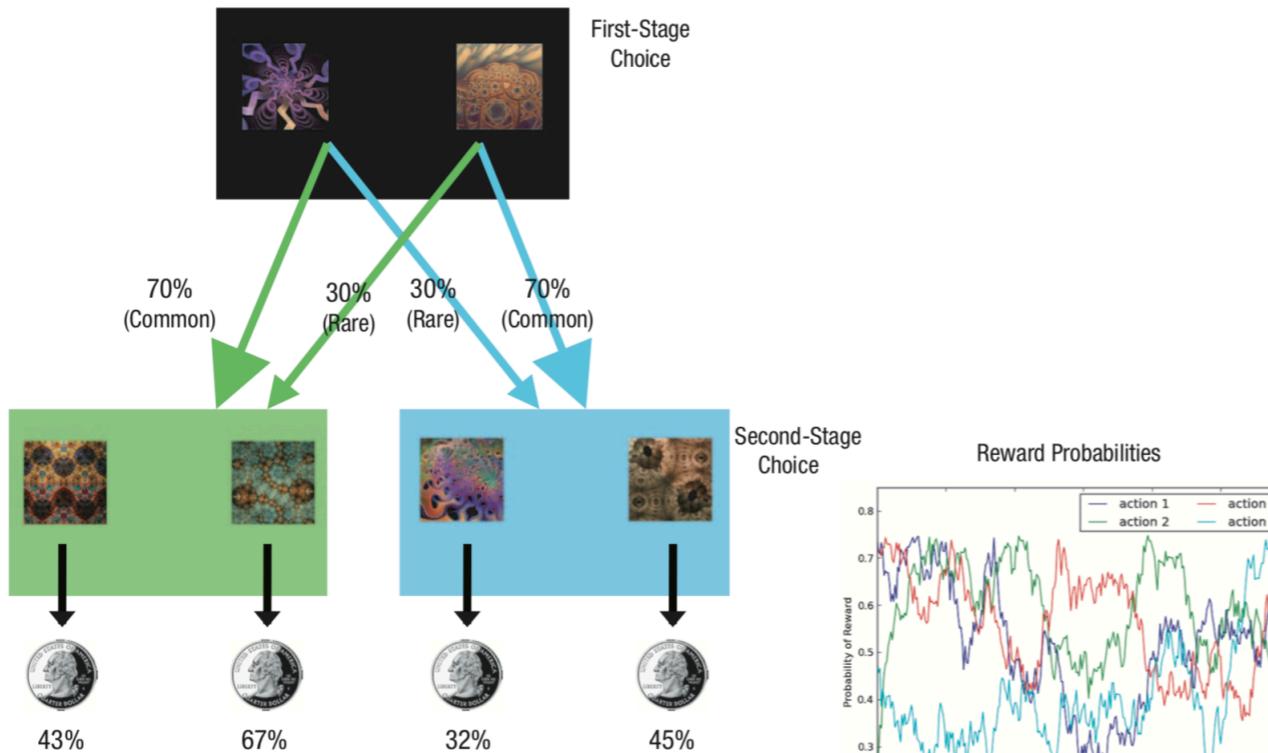
answers

- each system is best in different situations (use each one when it is most suitable/most accurate)
 - goal-directed (forward search) - good with limited training, close to the reward (don't have to search ahead too far)
 - habitual (cache) - good after much experience, distance from reward not so important
- arbitration: trust the system that is more confident in its recommendation
 - different sources of uncertainty in the two systems
 - don't always choose the highest value
 - uncertainty is different from risk

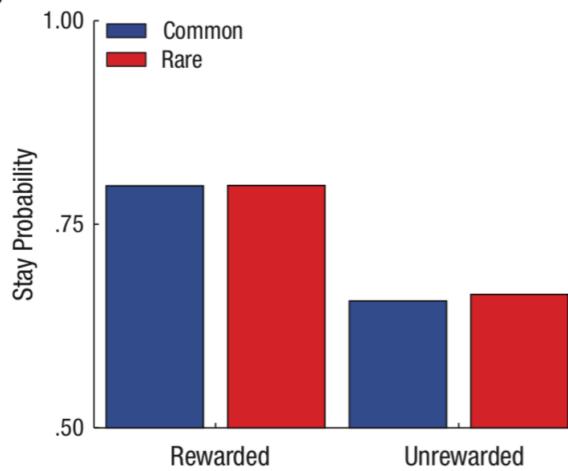


from Yael Niv

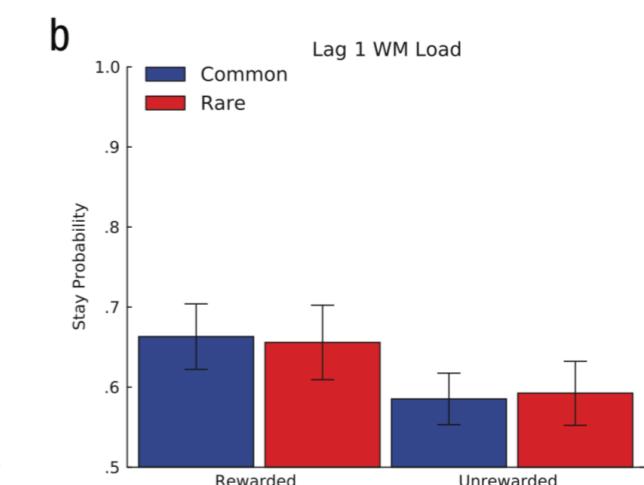
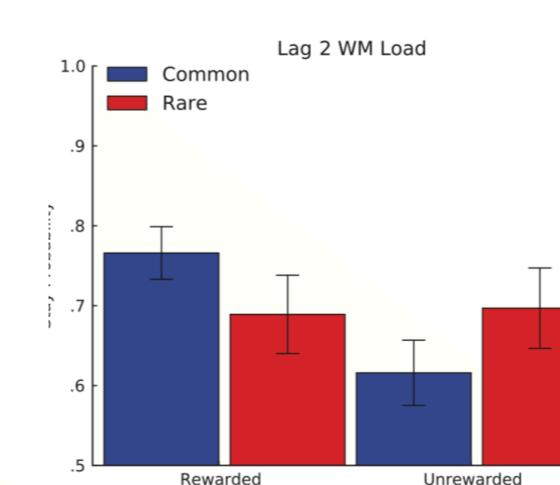
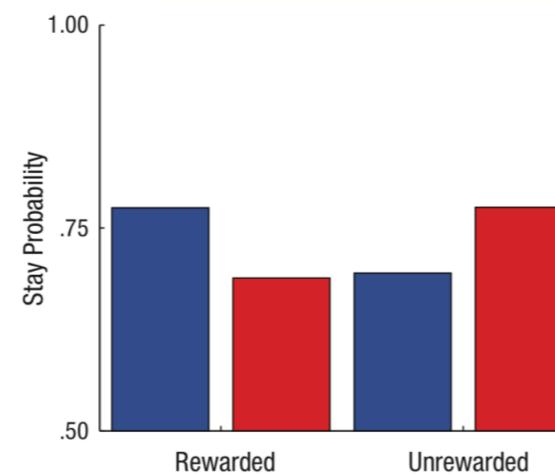
a



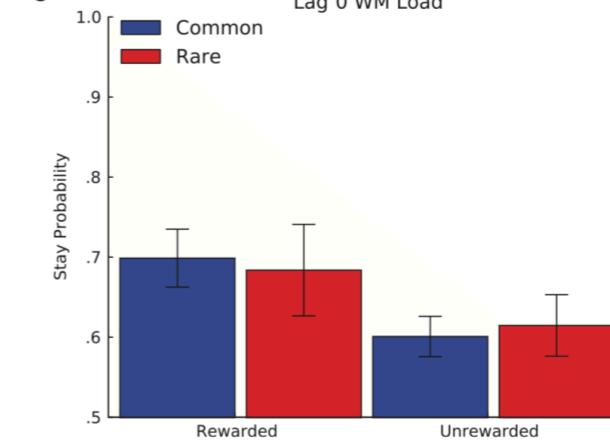
b



c



c



Working memory load interferes with online planning making human act more like model-free learners.

The Curse of Planning: Dissecting Multiple Reinforcement-Learning Systems by Taxing the Central Executive

A. Ross Otto¹, Samuel J. Gershman², Arthur B. Markman¹, and Nathaniel D. Daw³

¹Department of Psychology, University of Texas at Austin; ²Department of Psychology and Princeton Neuroscience Institute, Princeton University; and ³Department of Psychology and Center for Neural Science, New York University

Three levels of description (*David Marr, 1982*)

Computational

Why do things work the way they do?
What is the goal of the computation?
What are the unifying principles?



Algorithmic

What representations can implement such computations?
How does the choice of representations determine the algorithm?

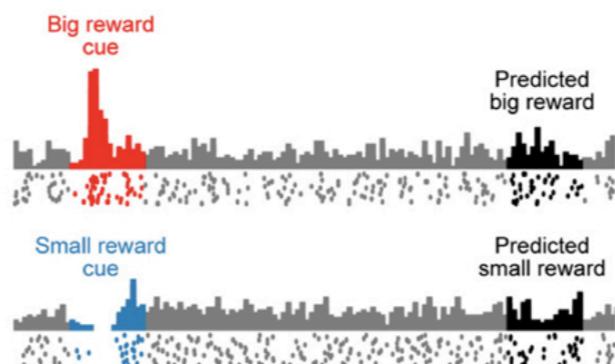
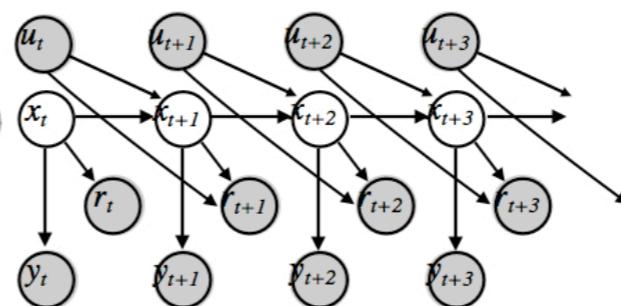
Implementational

How can such a system be built in hardware?
How can neurons carry out the computations?

maximize:

$$R_t = r_{t+1} + r_{t+2} + \dots + r_T$$

Bellman



Dynamic programming,
TD methods, Monte
Carlo

Neural firing patterns,
prediction errors,
system level
neuroscience

Playing around with RL...

Environments Documentation 

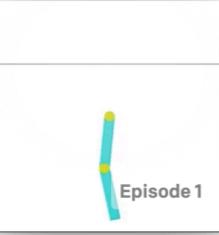
Algorithms

Atari

Box2D

Classic control

Control theory problems from the classic RL literature.

 Episode 1

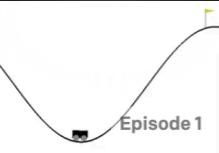
Acrobot-v1
Swing up a two-link robot.

 Episode 7

CartPole-v1
Balance a pole on a cart.

 Episode 1

MountainCar-v0
Drive up a big hill.

 Episode 1

MountainCarContinuous-v0
Drive up a big hill with continuous control.

 Episode 1

Pendulum-v0
Swing up a pendulum.

Third party environments 



Gym

Gym is a toolkit for developing and comparing reinforcement learning algorithms. It supports teaching agents everything from walking to playing games like Pong or Pinball.

[View documentation >](#)

[View on GitHub >](#)

Next time: Bayes



Slide Credits

Nathaniel Daw (exploration/gittins)

Alex Rich

Gillian Hayes (TD methods/explore)

Rich Sutton (general approach)

Andy Barto (general approach)