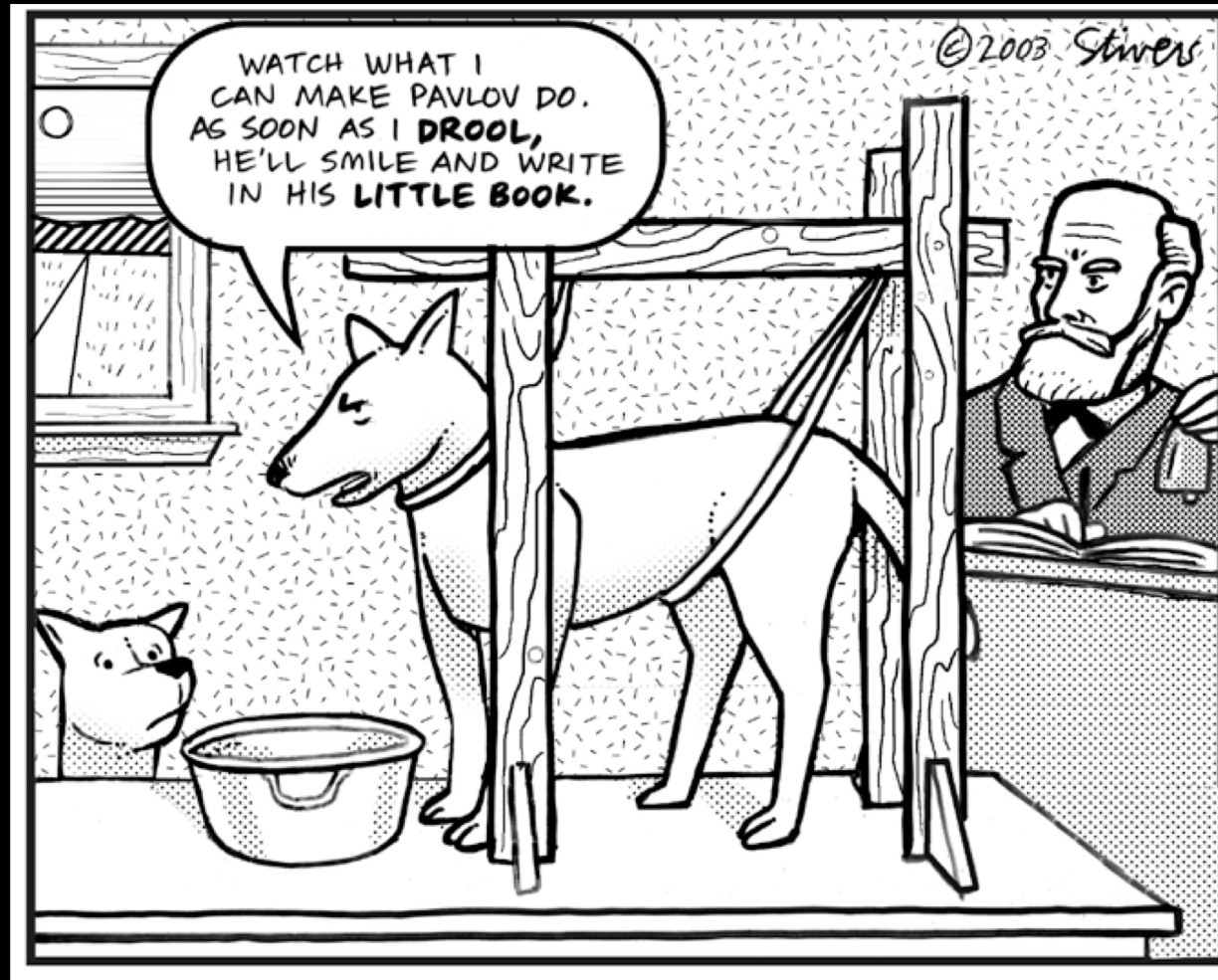


Temporal reinforcement learning



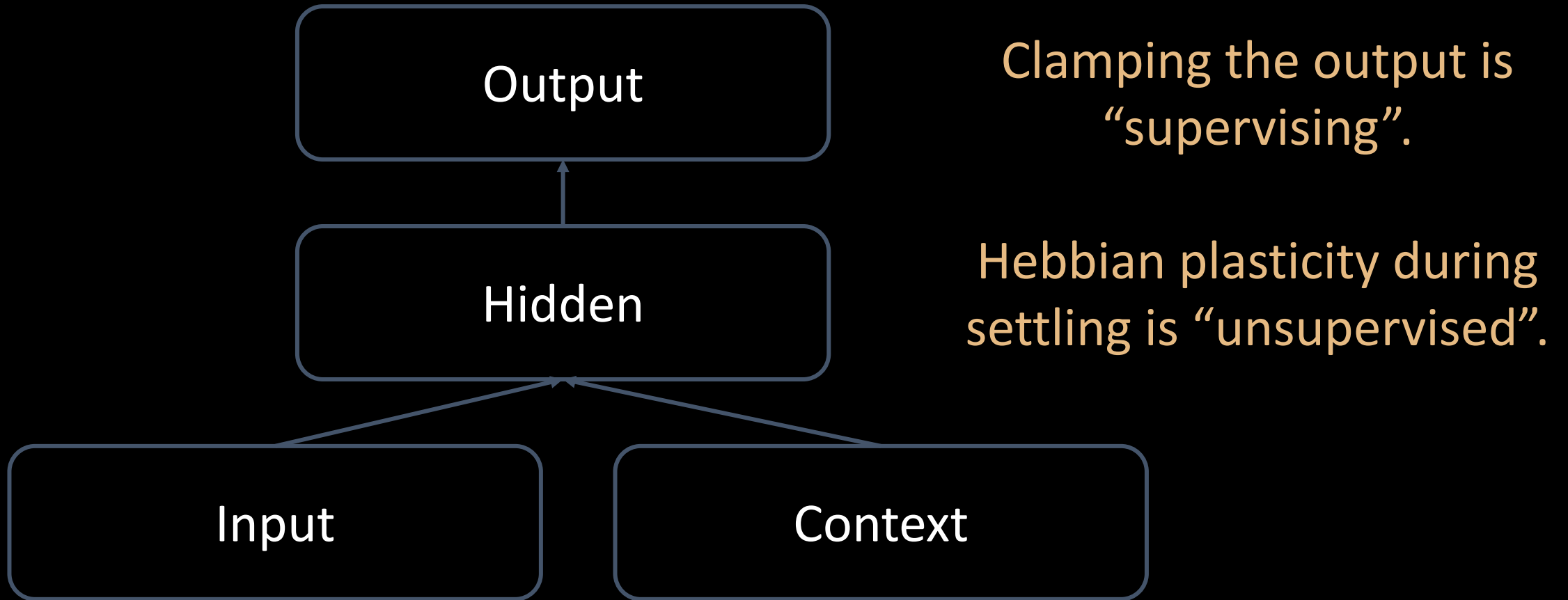
Slides adapted from:
Andra Geana
Jeff Cockburn
Michael Frank

Assume you're a brain

You need to eat things, you don't want to be eaten by other things, and you'd "like" to produce more brains.

- 1. What is learning?*
- 2. Why is learning important?*
- 3. What should you learn?*
- 4. When should you learn?*

Learning example: SRN



Reinforcement learning: In-between

Features:

- No 'teacher'
- No 'correct' answer given
- Only better/worse outcomes

What determines value?

Grounded in motivation

Some 'wants' are innate

Link events/actions to innate 'wants'



Supervised? Un?

Link spans time!

Reinforcement Learning

- Reinforcement (term from operant conditioning)
 - Something that makes it more likely that a certain response will re-occur

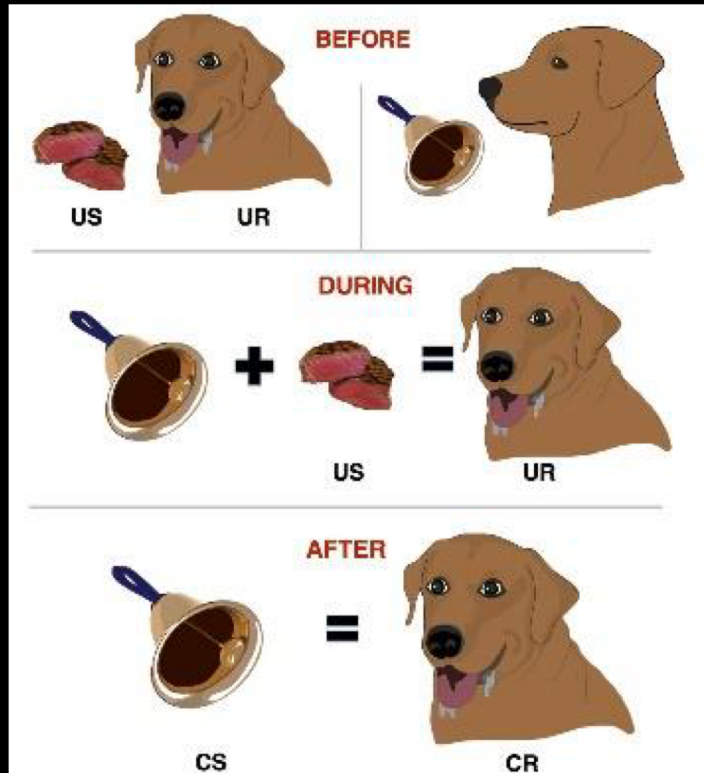


Reinforcement Learning

- Reinforcement (term from operant **conditioning**)
 - Something that makes it **more likely** that a certain response will re-occur



Conditioning: training an organism to respond in specific ways to certain stimuli (shaping behavior)



Classical (Pavlovian):
train involuntary responses (CR)



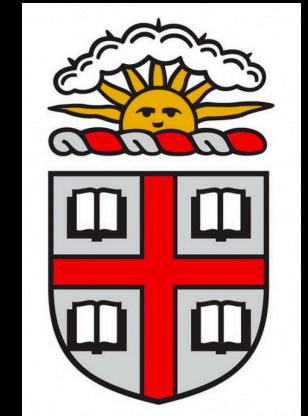
Instrumental (operant):
train voluntary responses (actions)

Reinforcement Learning

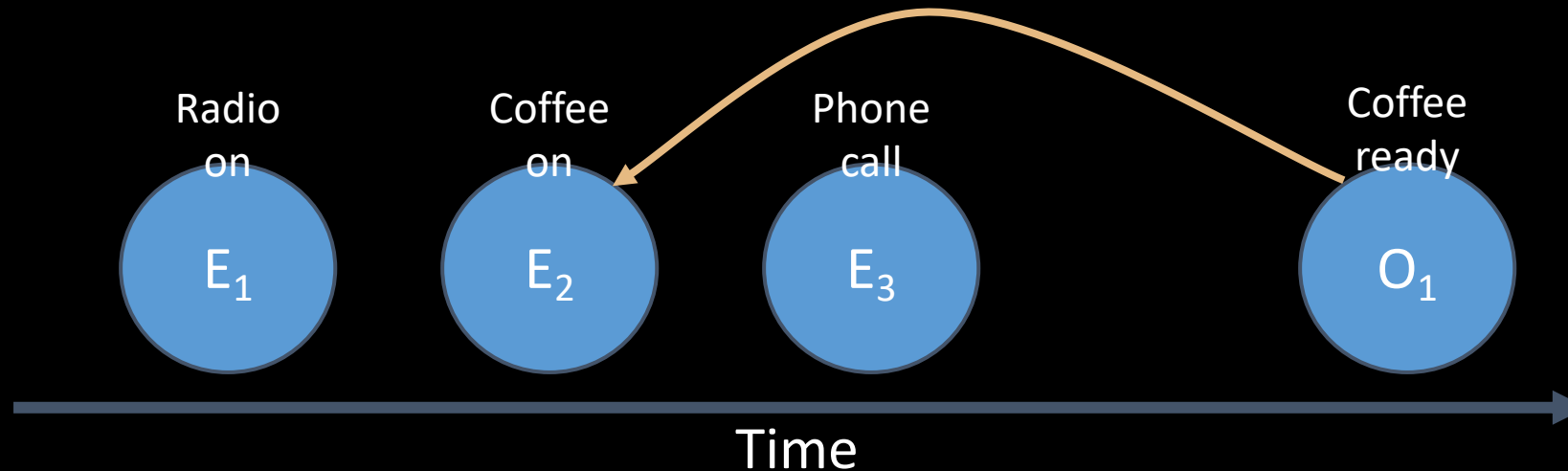
Using trial-and-error to learn how to map situations to actions so as to maximize a numerical reward signal

Learning about delayed outcomes

- The results of actions are often delayed
- Irrelevant information abounds
- Want “good” results

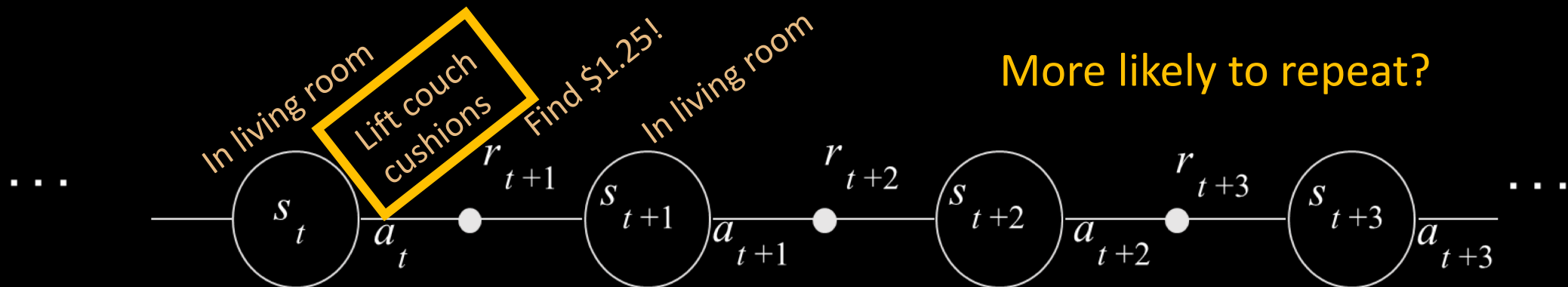
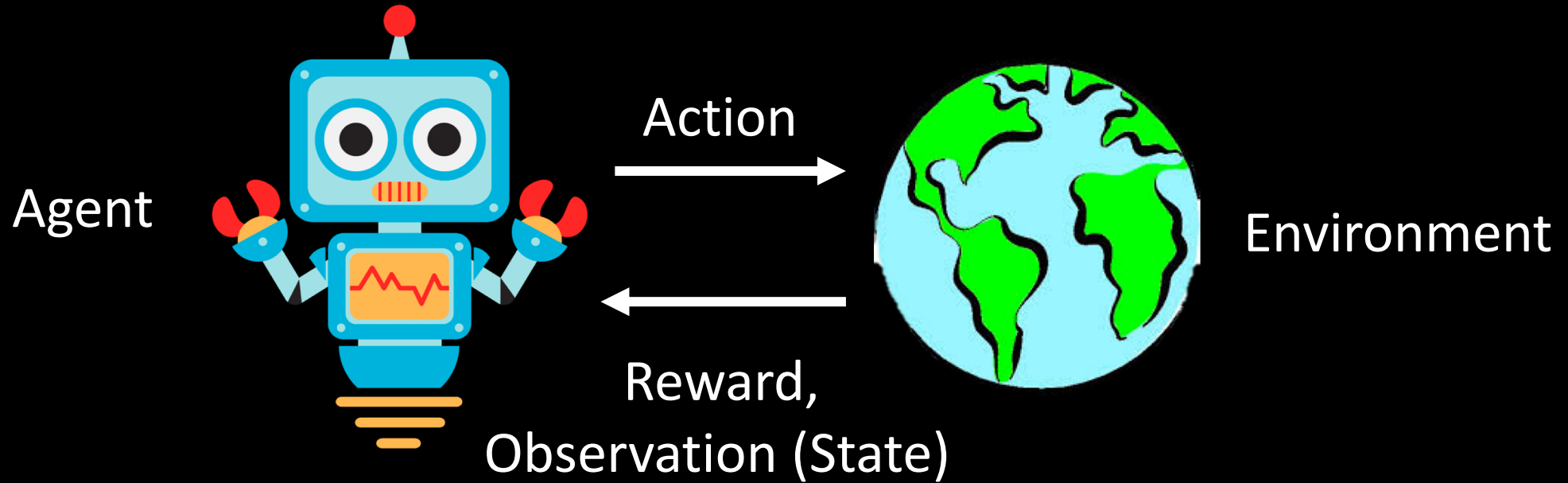


Goal: Learn to **predict** event outcomes



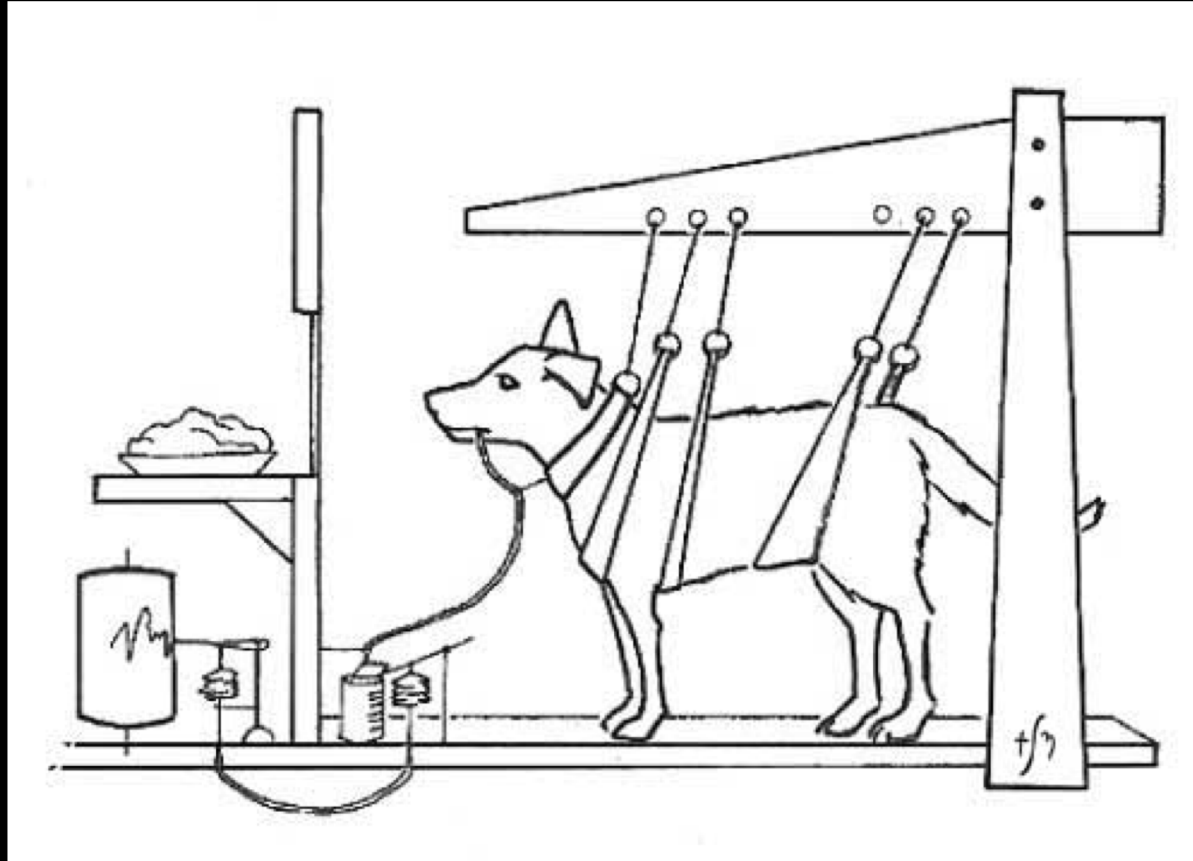
Referred to as: Temporal credit assignment

Temporal difference learning (informally)



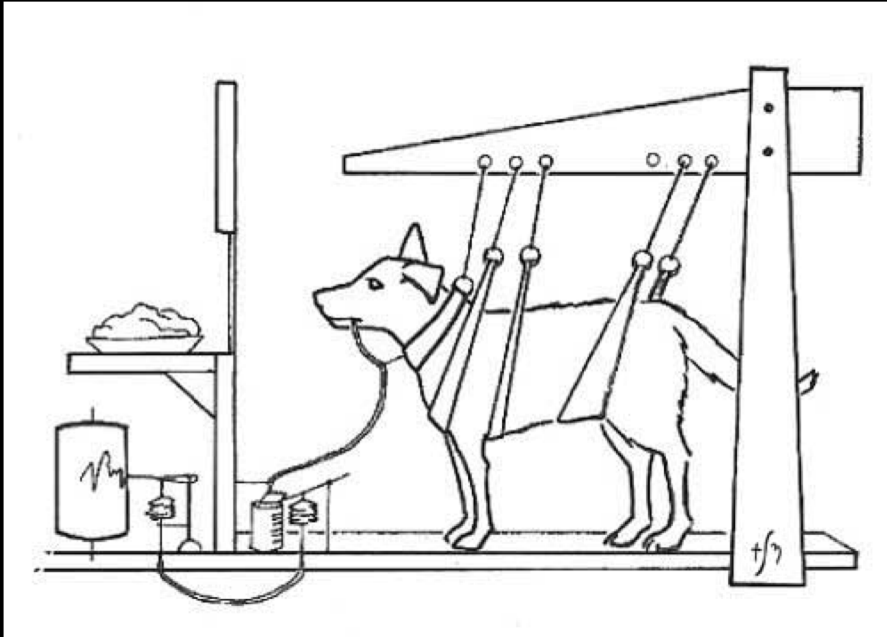
Classical Conditioning

What did Pavlov learn about what his dog learned?



Classical Conditioning

What did Pavlov learn about what his dog learned?



Day 1

A bell rings:

Was that: expected / unexpected ?

Is this: Better / Worse / Neutral to what you expected?

Is there anything to learn?

Food arrives:

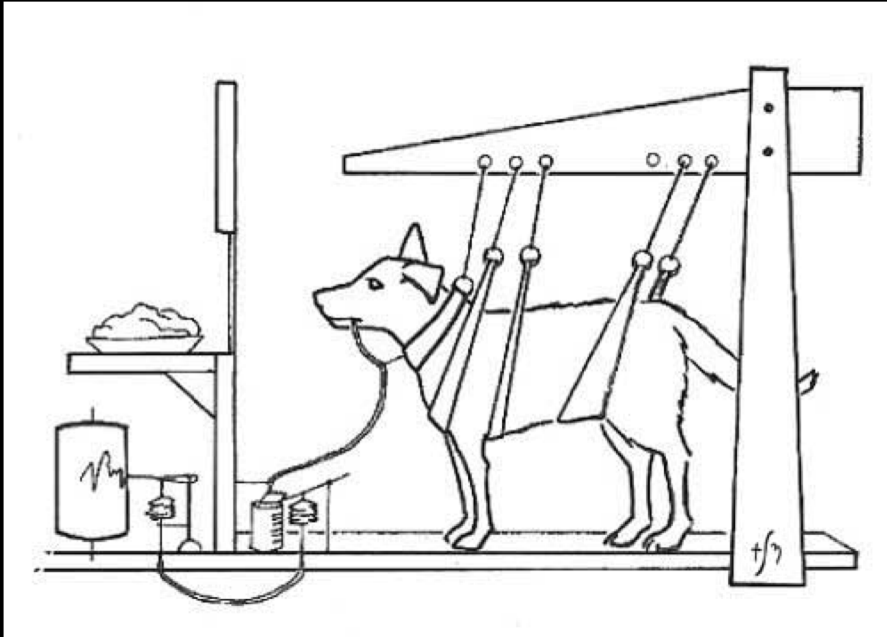
Was that: expected / unexpected ?

Is this: Better / Worse / Neutral to what you expected?

Is there anything to learn?

Classical Conditioning

What did Pavlov learn about what his dog learned?



Day 2

A bell rings:

Was that: expected / unexpected ?

Is this: Better / Worse / Neutral to what you expected?

Is there anything to learn?

Food arrives:

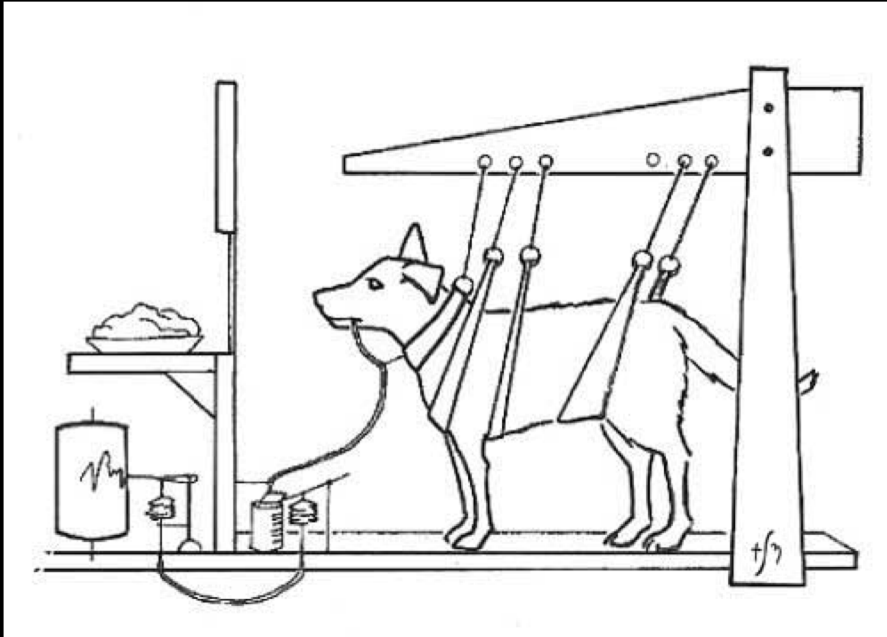
Was that: expected / unexpected ?

Is this: Better / Worse / Neutral to what you expected?

Is there anything to learn?

Classical Conditioning

What did Pavlov learn about what his dog learned?



Day 3

A bell rings:

Was that: expected / unexpected ?

Is this: Better / Worse / Neutral to what you expected?

Is there anything to learn?

Food arrives:

Was that: expected / unexpected ?

Is this: Better / Worse / Neutral to what you expected?

Is there anything to learn?

Rescorla & Wagner (1972)

- **Prediction error** or “surprise” driven learning rule

Prediction Error
difference b/n
observed and
predicted

Predicted Value
of state at time t

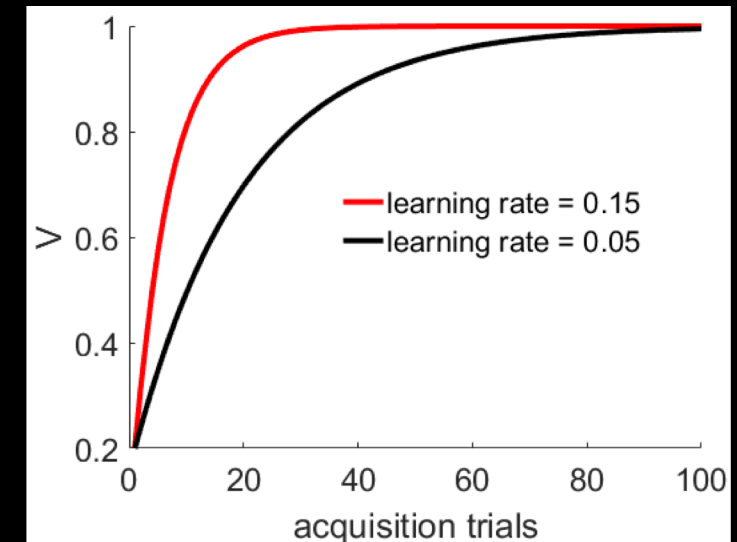
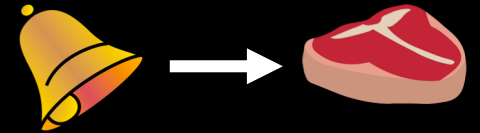
$$\delta_t = r_t - V_t(s)$$

$$V_{t+1}(s) = V_t(s) + \alpha \delta_t$$

Updated Value
at time t+1 of state

Learning rate: how much
do we care about each
new data point?

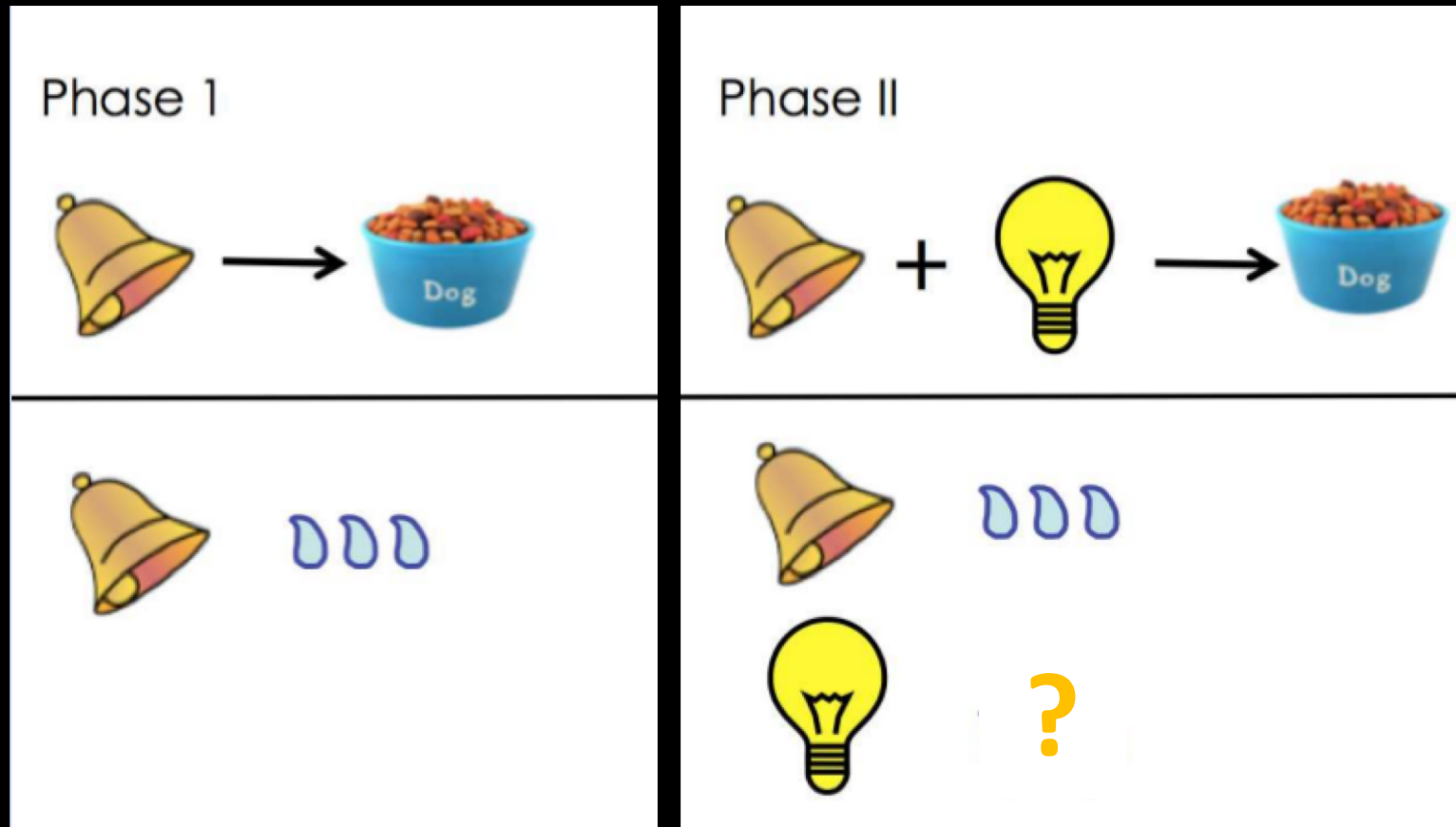
Why would you not always
want a learning rate of 1?



Higher learning rates
means we learn (*and*
forget) faster!

Blocking

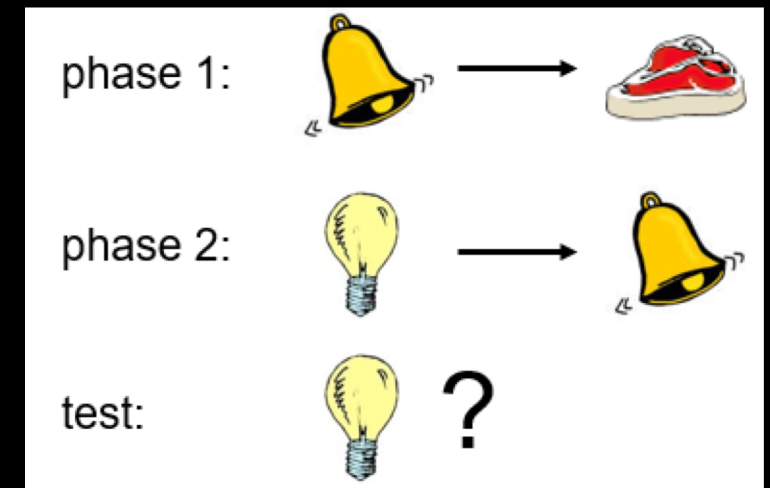
$$\delta_t = r_t - V_t(s)$$
$$V_{t+1}(s) = V_t(s) + \alpha \delta_t$$



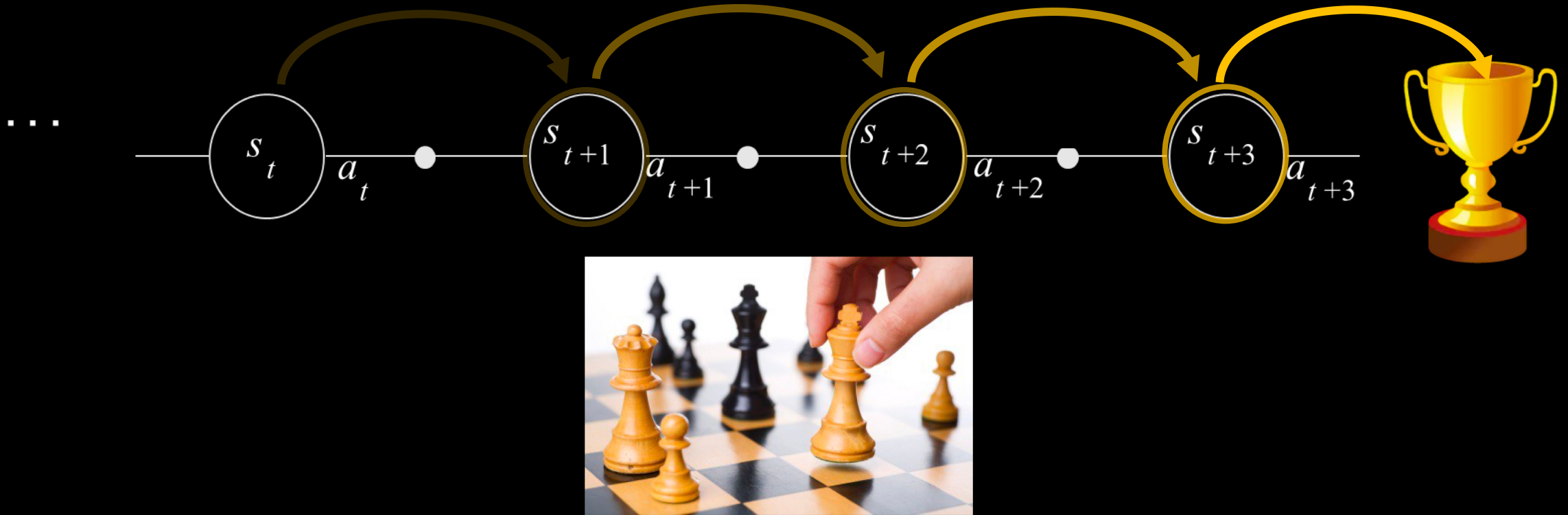
Problems with Rescorla Wagner

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

- Trial-level learning and prediction
- Scaling problem
- Can't explain effects like associative bias
 - Rats more likely to associate light and shock or flavor and poisoning
- Can't explain second order conditioning
 - *Why?*



How can RL handle sequences of states?



Rather than trial-level predicting based on immediate rewards,
learn to predict (discounted) future rewards!

Temporal difference learning (more formal)

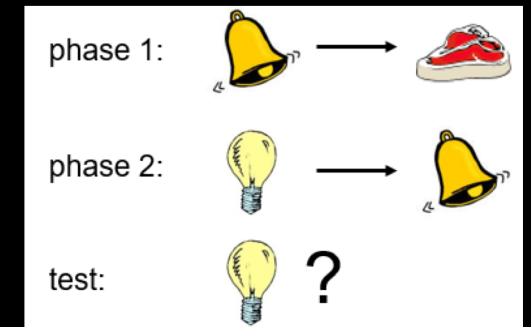
Prediction Error
difference b/n
outcome and
predicted

Outcome
Predicted value

Expectations
Predicted value

$$\delta_t = r_t + \underbrace{V(s_{t+1})}_{\text{Outcome}} - \underbrace{V(s_t)}_{\text{Expectations}}$$
$$V(s_t) = V(s_t) + \alpha \delta_t$$

$$V(s_{t+1}) = \text{sum}[R_{\text{future}}]$$

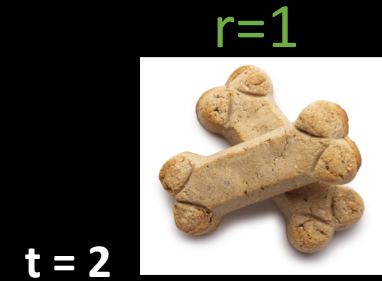
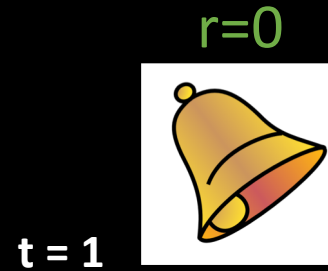
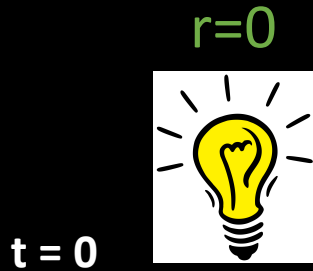
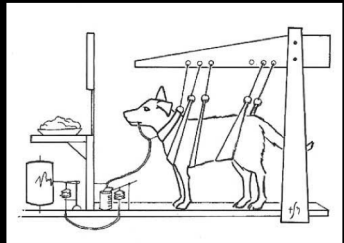


Pavlov's dog

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

$$V(s_t) \leftarrow V(s_t) + \alpha \delta_t$$

$$\alpha = 1$$



$V(\text{lightbulb})$	$V(\text{bell})$	$V(\text{bone})$
0	0	1

$$V(s_t) = V(\text{lightbulb}) = 0$$

$$V(s_t) = V(\text{bell}) = 0$$

$$V(s_t) = V(\text{bone}) = 0$$

$$V(s_{t+1}) = V(\text{bell}) = 0$$

$$V(s_{t+1}) = V(\text{bone}) = 0$$

$$V(s_{t+1}) = \text{None}$$

$$\delta = 0 + 0 - 0 = 0$$

$$\delta = 0 + 0 - 0 = 0$$

$$\delta = 1 - 0 = 1$$

$$V(\text{lightbulb}) = 0$$

$$V(\text{bell}) = 0$$

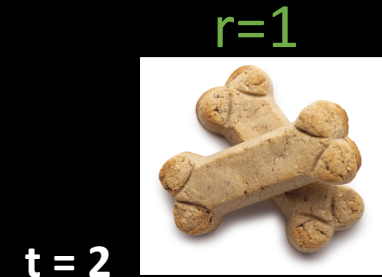
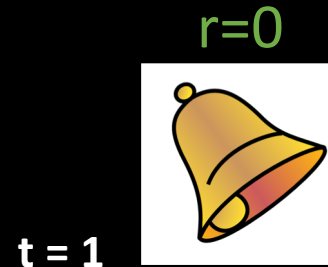
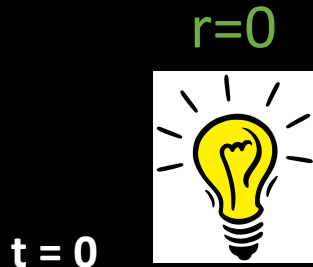
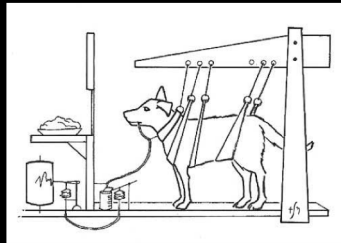
$$V(\text{bone}) = 1$$

Pavlov's dog

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

$$V(s_t) \leftarrow V(s_t) + \alpha \delta_t$$

$$\alpha = 1$$



$V(\text{lightbulb})$	$V(\text{bell})$	$V(\text{bone})$
0	0	1
0	1	1

$$V(S_t) = V(\text{lightbulb}) = 0$$

$$V(S_t) = V(\text{bell}) = 0$$

$$V(S_t) = V(\text{bone}) = 1$$

$$V(S_{t+1}) = V(\text{bell}) = 0$$

$$V(S_{t+1}) = V(\text{bone}) = 1$$

$$V(S_{t+1}) = \text{None}$$

$$\delta = 0 + 0 - 0 = 0$$

$$\delta = 0 + 1 - 0 = 1$$

$$\delta = 1 - 1 = 0$$

$$V(\text{lightbulb}) = 0$$

$$V(\text{bell}) = 1$$

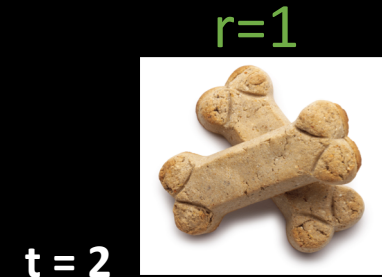
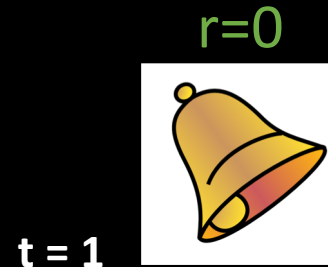
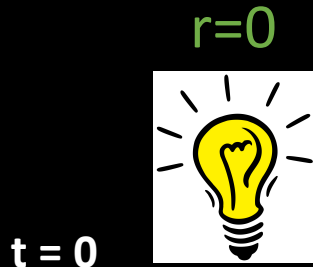
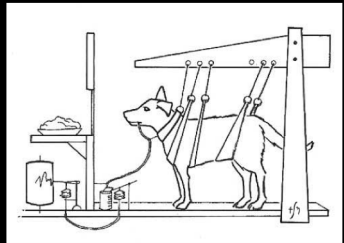
$$V(\text{bone}) = 1$$

Pavlov's dog

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

$$V(s_t) \leftarrow V(s_t) + \alpha \delta_t$$

$$\alpha = 1$$



$V(\text{lightbulb})$	$V(\text{bell})$	$V(\text{bone})$
0	0	1
0	1	1
1	1	1

$$V(S_t) = V(\text{lightbulb}) = 0$$

$$V(S_t) = V(\text{bell}) = 1$$

$$V(S_t) = V(\text{bone}) = 1$$

$$V(S_{t+1}) = V(\text{bell}) = 1$$

$$V(S_{t+1}) = V(\text{bone}) = 1$$

$$V(S_{t+1}) = \text{None}$$

$$\delta = 0 + 1 - 0 = 1$$

$$\delta = 0 + 1 - 1 = 0$$

$$\delta = 1 - 1 = 0$$

$$V(\text{lightbulb}) = 1$$

$$V(\text{bell}) = 1$$

$$V(\text{bone}) = 1$$

"bootstrapping"

Uses own estimation of future reward
to learn *even* in the absence of
immediate reward

Consider a long sequence



$$\delta_t = r_t + \mathbf{V(s_{t+1})} - V(s_t)$$
$$V(s_t) = V(s_t) + \alpha \delta_t$$

If $\mathbf{V(s_{t+1})} = \text{sum}[R_{future}]$

You are guaranteed all the kibbles...100 time points
from now

Consider a long sequence



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

Solution: **discount** the value of future reward!

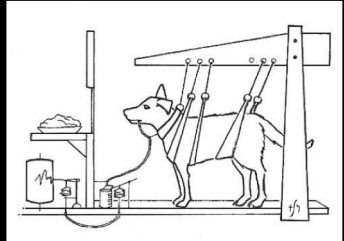
$V(s_t)$ = sum of discounted rewards

$$= r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \gamma^3 r_{t+4} + \dots$$

Pavlov's dog

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

$$V(s_t) \leftarrow V(s_t) + \alpha \delta_t$$



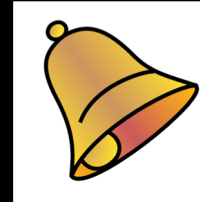
t = 0

r=0



t = 1

r=0



t = 2

r=1



$\alpha = 1$
 $\gamma = 0.5$

$V(\text{lightbulb})$	$V(\text{bell})$	$V(\text{bone})$
0	0	1

$$V(S_t) = V(\text{lightbulb}) = 0$$

$$V(S_t) = V(\text{bell}) = 0$$

$$V(S_t) = V(\text{bone}) = 0$$

$$V(S_{t+1}) = V(\text{bell}) = 0$$

$$V(S_{t+1}) = V(\text{bone}) = 0$$

$$V(S_{t+1}) = \text{None}$$

$$\delta = 0 + (0.5)0 - 0 = 0$$

$$\delta = 0 + (0.5)0 - 0 = 0$$

$$\delta = 1 - 0 = 1$$

$$V(\text{lightbulb}) = 0$$

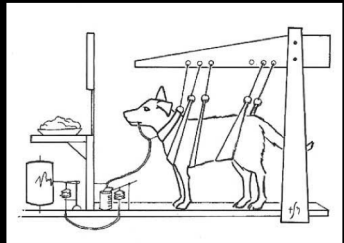
$$V(\text{bell}) = 0$$

$$V(\text{bone}) = 1$$

Pavlov's dog

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

$$V(s_t) \leftarrow V(s_t) + \alpha \delta_t$$



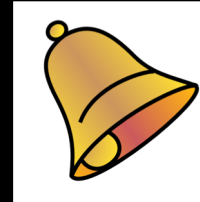
t = 0

r=0



t = 1

r=0



t = 2

r=1



$\alpha = 1$
 $\gamma = 0.5$

$V(\text{lightbulb})$	$V(\text{bell})$	$V(\text{bone})$
0	0	1
0	0.5	1

$$V(S_t) = V(\text{lightbulb}) = 0$$

$$V(S_t) = V(\text{bell}) = 0$$

$$V(S_t) = V(\text{bone}) = 1$$

$$V(S_{t+1}) = V(\text{bell}) = 0$$

$$V(S_{t+1}) = V(\text{bone}) = 1$$

$$V(S_{t+1}) = \text{None}$$

$$\delta = 0 + (.5)0 - 0 = 0$$

$$\delta = 0 + .5(1) - 0 = .5$$

$$\delta = 1 - 1 = 0$$

$$V(\text{lightbulb}) = 0$$

$$V(\text{bell}) = 0.5$$

$$V(\text{bone}) = 1$$

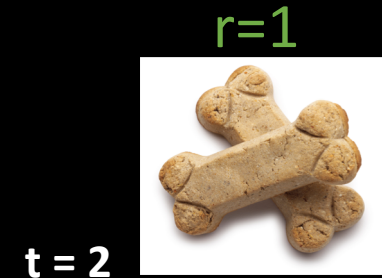
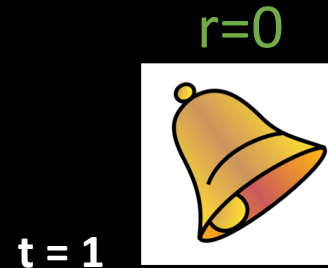
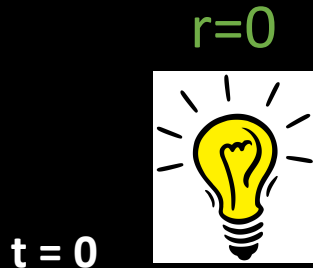
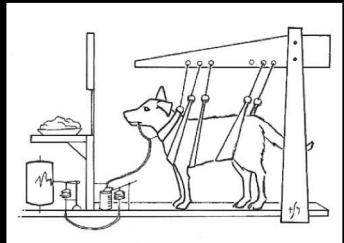
Pavlov's dog

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

$$V(s_t) \leftarrow V(s_t) + \alpha \delta_t$$

$$\alpha = 1$$

$$\gamma = 0.5$$



$V(\text{lightbulb})$	$V(\text{bell})$	$V(\text{bone})$
0	0	1
0	0.5	1
.25	.5	1

$$V(s_t) = V(\text{lightbulb}) = 0$$

$$V(s_t) = V(\text{bell}) = 0.5$$

$$V(s_t) = V(\text{bone}) = 1$$

$$V(s_{t+1}) = V(\text{bell}) = 0.5$$

$$V(s_{t+1}) = V(\text{bone}) = 1$$

$$V(s_{t+1}) = \text{None}$$

$$\delta = 0 + (.5).5 - 0 = .25$$

$$\delta = 0 + (.5)1 - .5 = 0$$

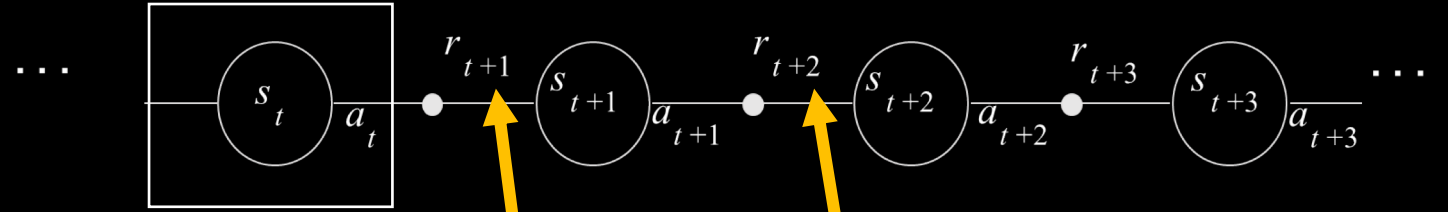
$$\delta = 1 - 1 = 0$$

$$V(\text{lightbulb}) = .25$$

$$V(\text{bell}) = 0.5$$

$$V(\text{bone}) = 1$$

TD-Learning



- Value of state expressed in future expected rewards

$$\begin{aligned} V_t &= E \left[\sum_{i=t}^T \gamma^{i-t} r_i \right] = E[\gamma^0 r_t + \gamma^1 r_{t+1} + \gamma^2 r_{t+2} \dots] \\ &= E[\gamma^0 r_t] + \gamma E[r_{t+1} + \gamma^1 r_{t+2} \dots] \\ &= E[r_t] + \gamma V_{t+1} \end{aligned}$$

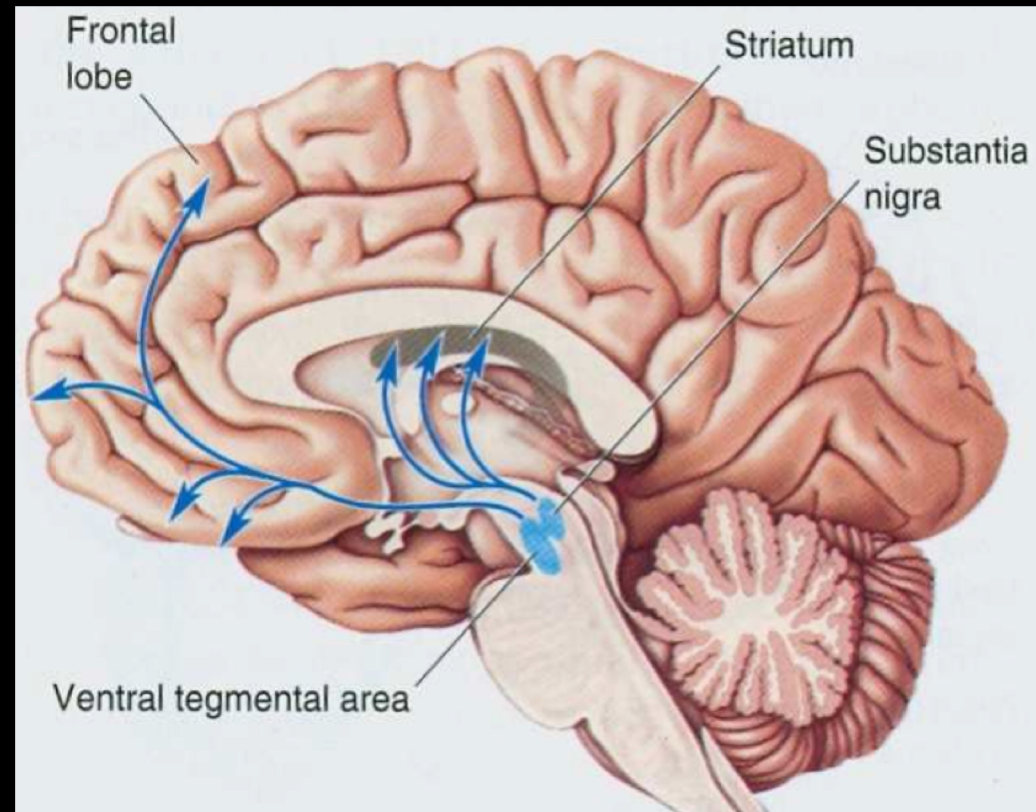
“Bellman Equation”

$$\begin{aligned} \delta_t &= r_t + \gamma \hat{V}_{t+1}(s) - \hat{V}_t(s) \\ \hat{V}_{t+1}(s) &= \hat{V}_t(s) + \alpha \delta_t \end{aligned}$$

“temporal difference”

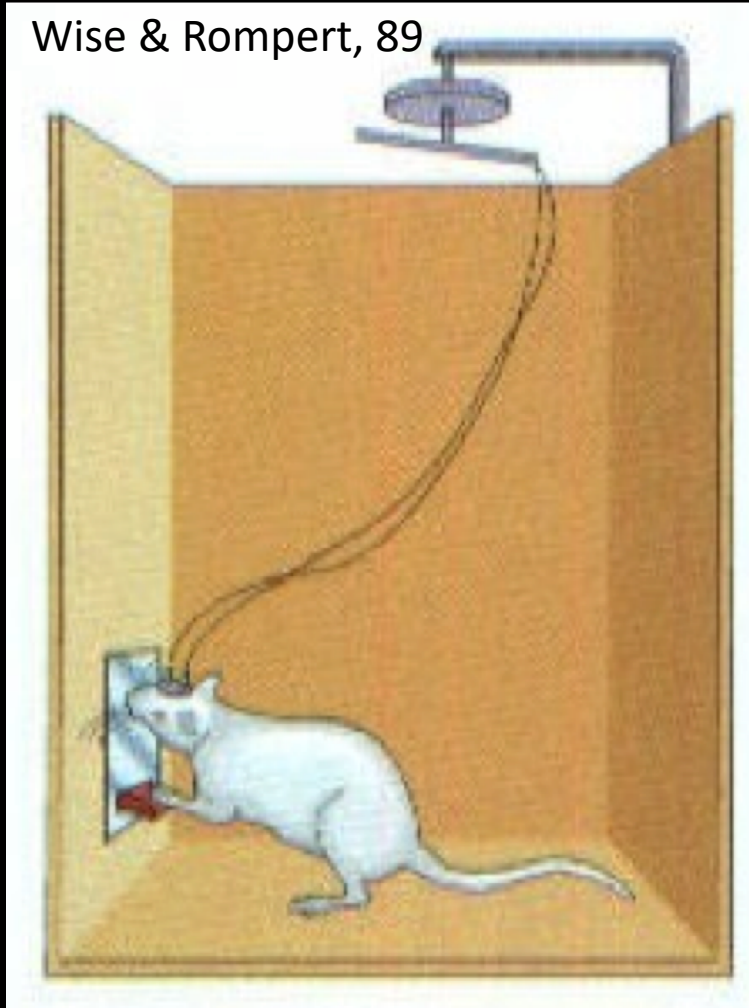
That's cool, Alana, but what about the brain?

RL in the brain: Dopamine system



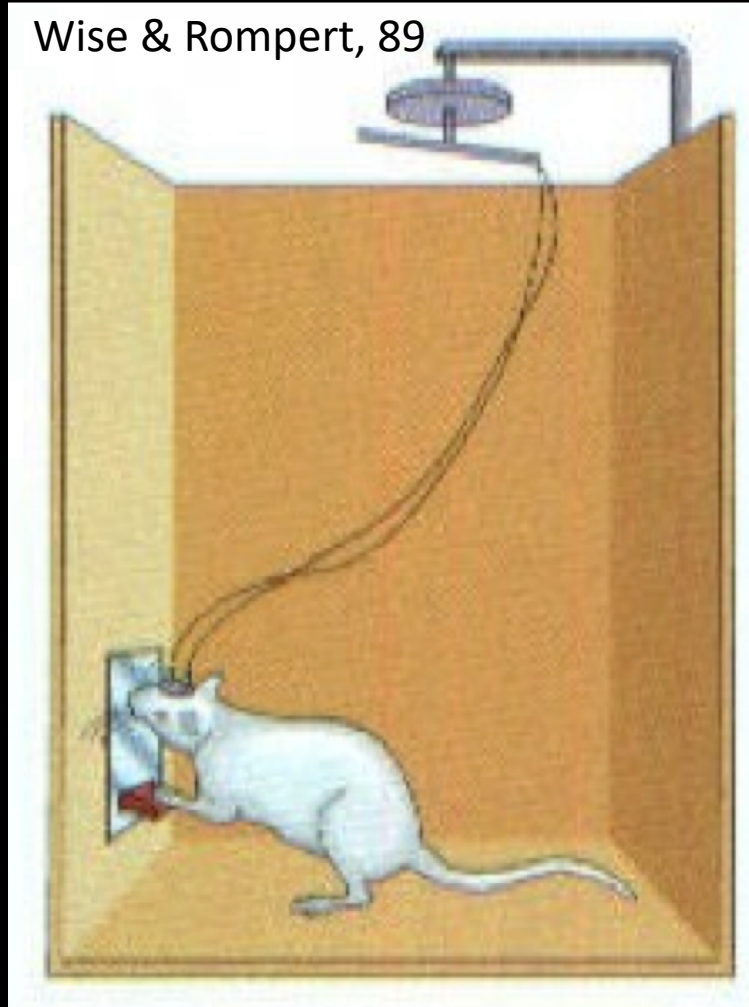
What is dopamine doing?

Dopamine carries the brain's reward signal



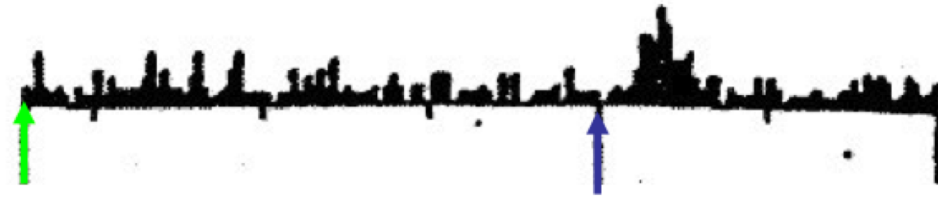
What is dopamine doing?

Dopamine carries the brain's reward signal



stimulus

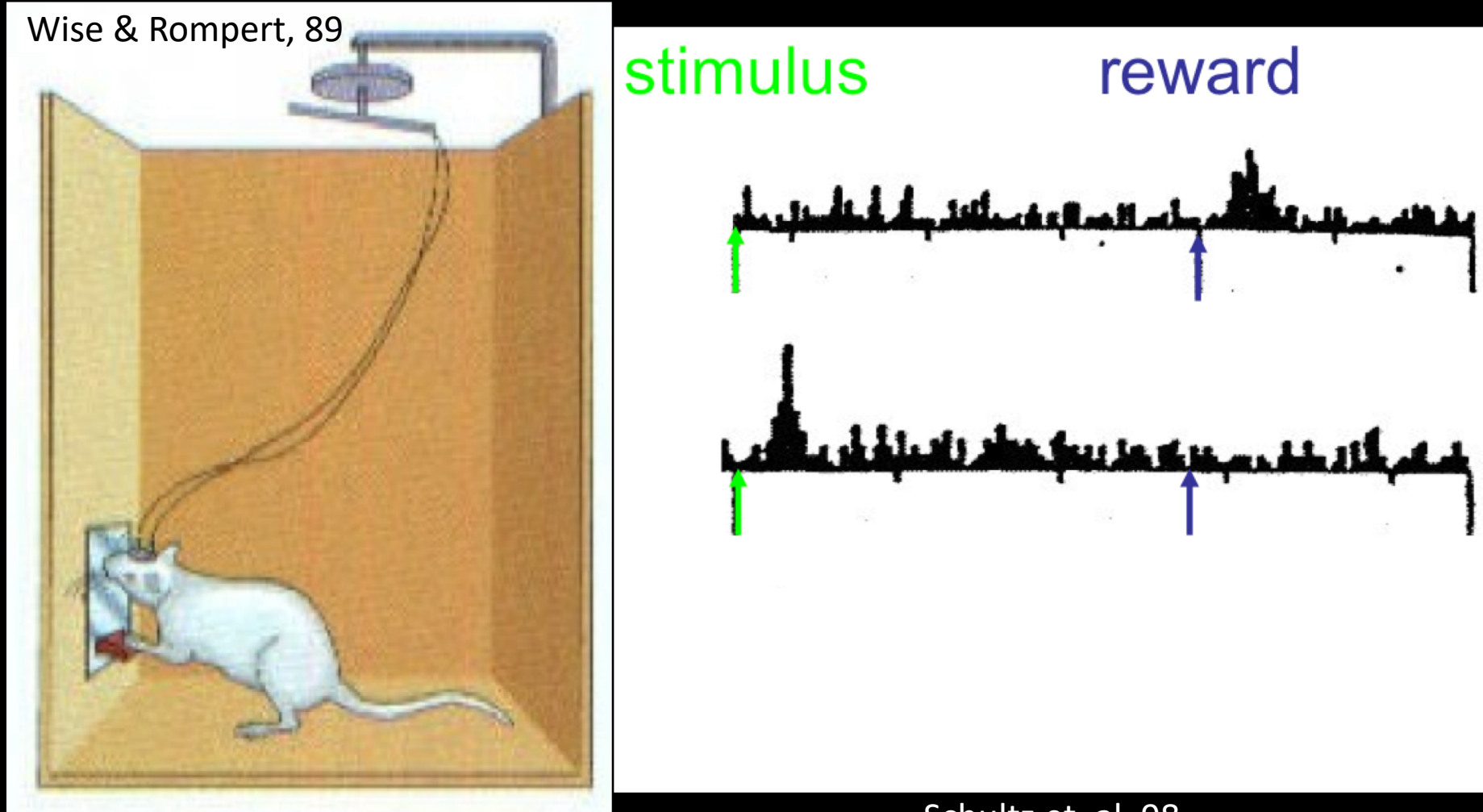
reward



Schultz et. al, 98

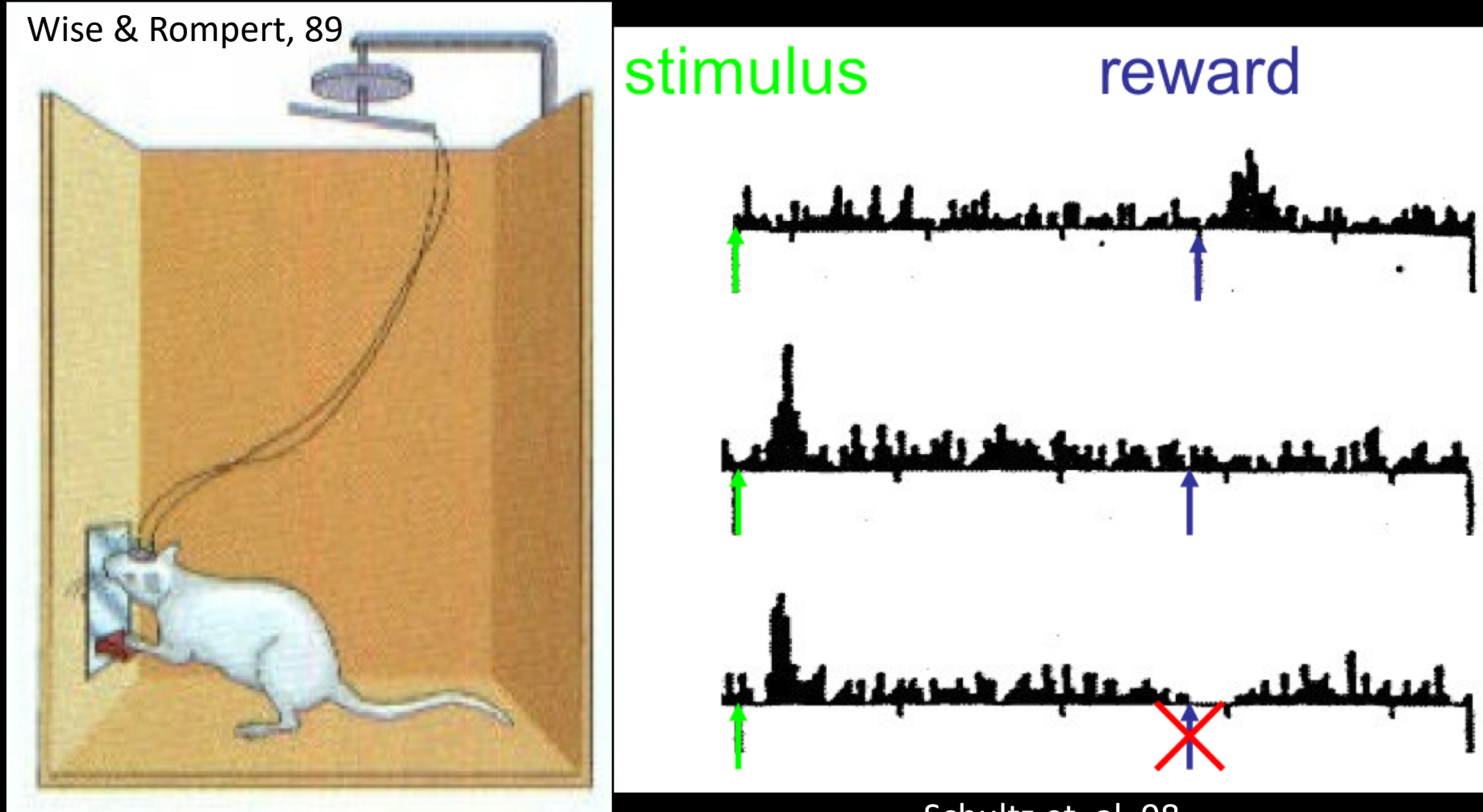
What is dopamine doing?

Dopamine carries the brain's reward signal



What is dopamine doing?

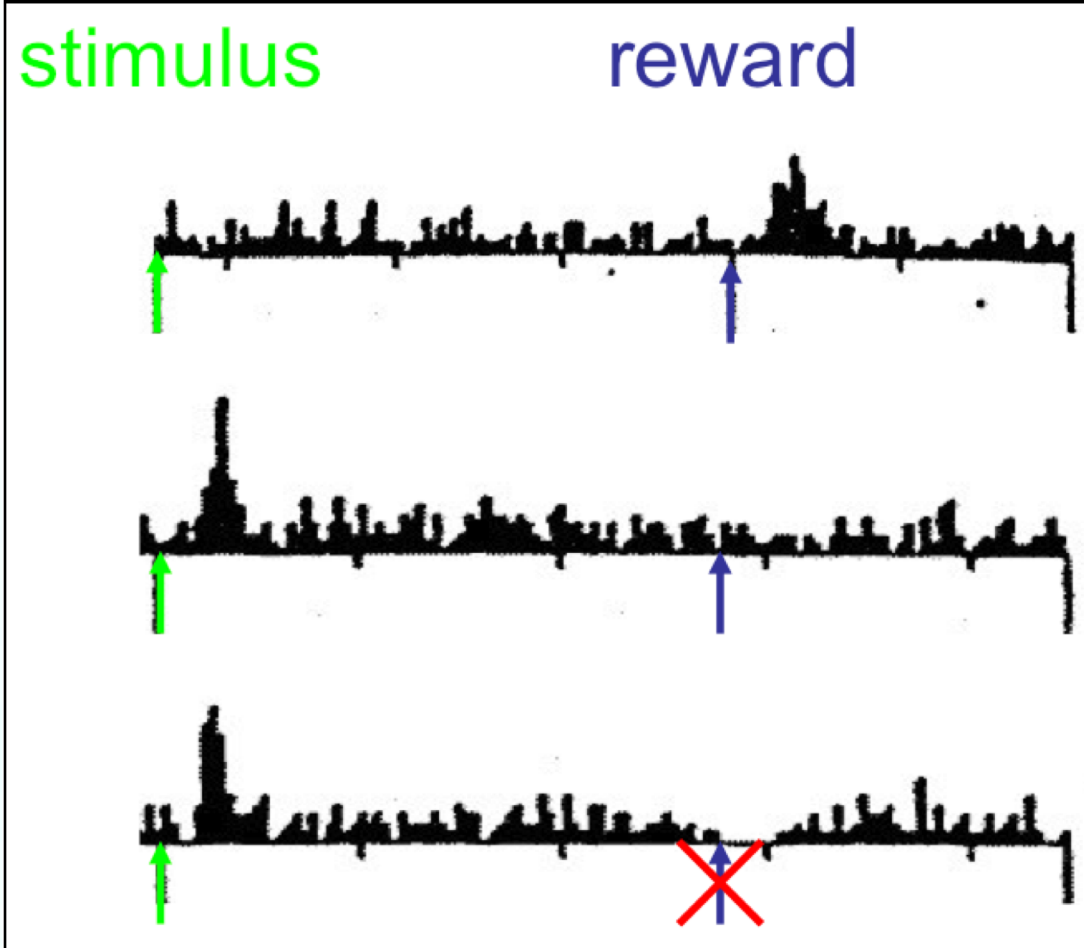
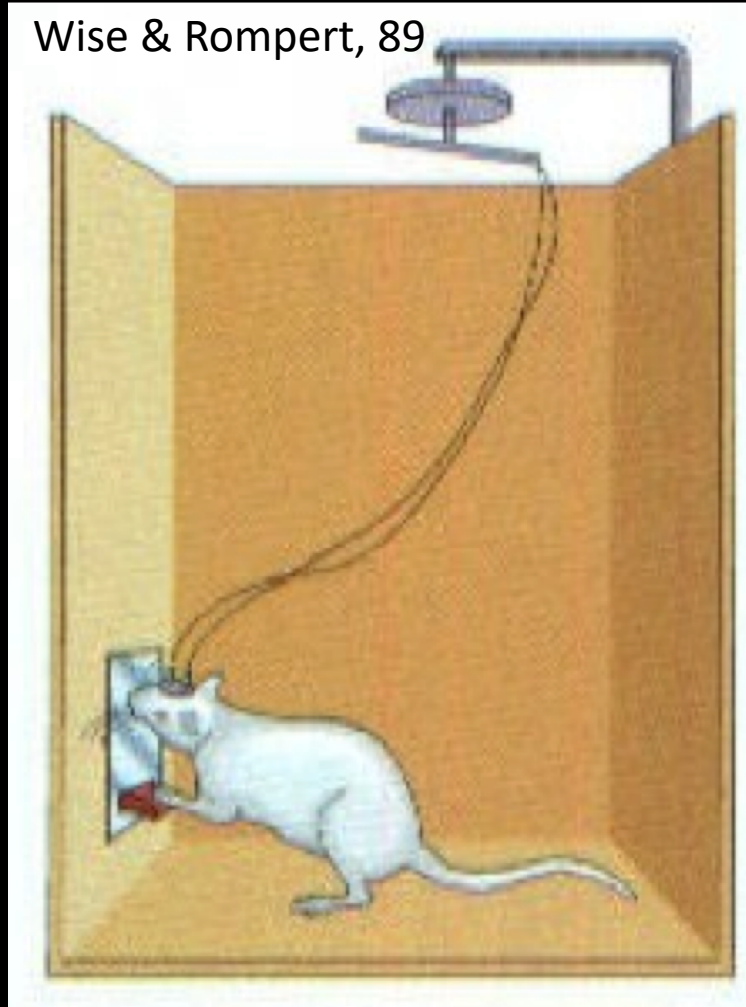
Dopamine carries the brain's reward signal



What is dopamine doing?

Dopamine carries the brain's ~~reward~~ signal

reward prediction error

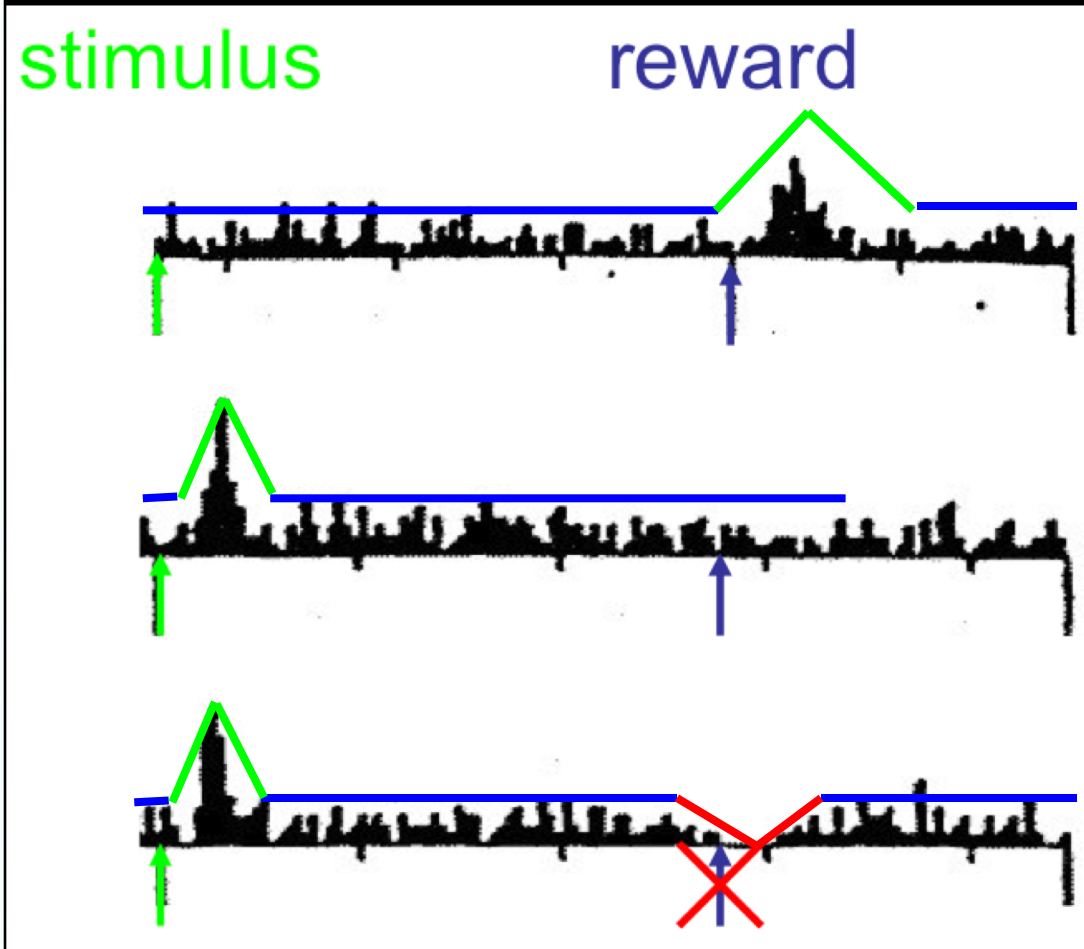
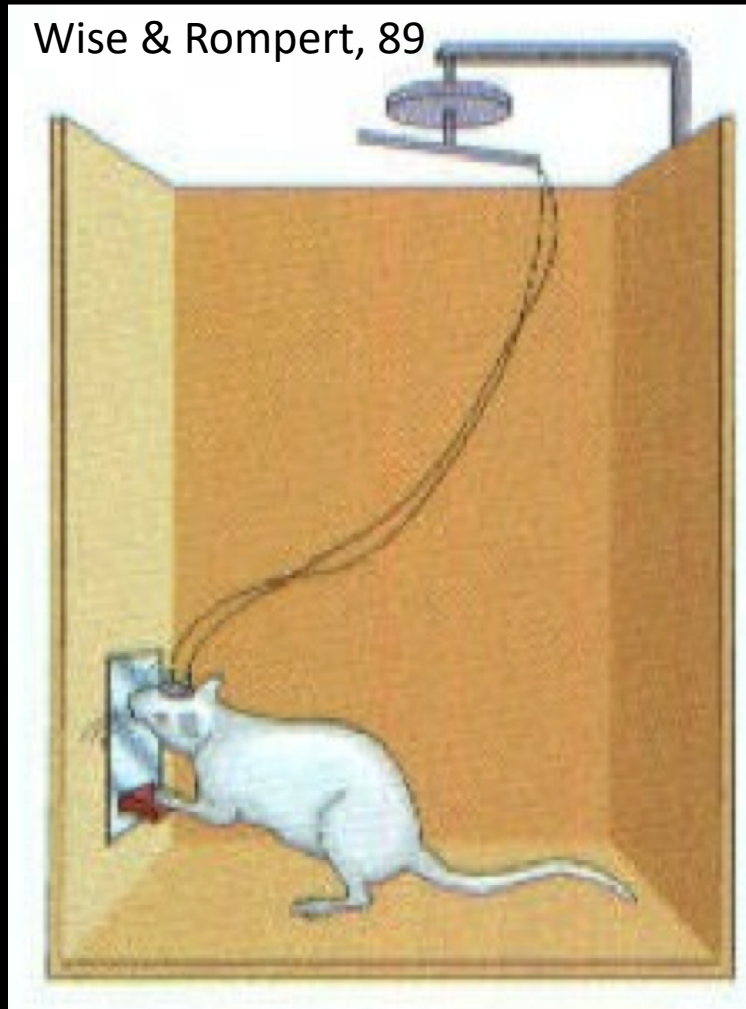


Schultz et. al, 98

What is dopamine doing?

Dopamine carries the brain's ~~reward~~ signal

reward prediction error



Schultz et. al, 98

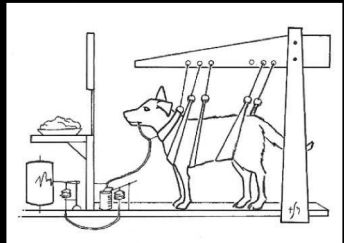
Pavlov's dog

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

$$V(s_t) \leftarrow V(s_t) + \alpha \delta_t$$

$$\alpha = 1$$

$$\gamma = 0.5$$



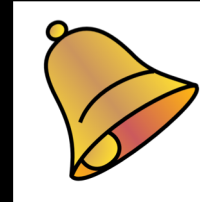
t = 0

r=0



t = 1

r=0



t = 2

r=1



$V(\text{lightbulb})$	$V(\text{bell})$	$V(\text{bone})$
0	0	1
0	0.5	1
.25	.5	1

$$V(S_t) = V(\text{lightbulb}) = 0$$

$$V(S_t) = V(\text{bell}) = 0.5$$

$$V(S_t) = V(\text{bone}) = 1$$

$$V(S_{t+1}) = V(\text{bell}) = 0.5$$

$$V(S_{t+1}) = V(\text{bone}) = 1$$

$$V(S_{t+1}) = \text{None}$$

$$\delta = 0 + (.5).5 - 0 = .25$$

$$\delta = 0 + (.5)1 - .5 = 0$$

$$\delta = 1 - 0 = 0$$

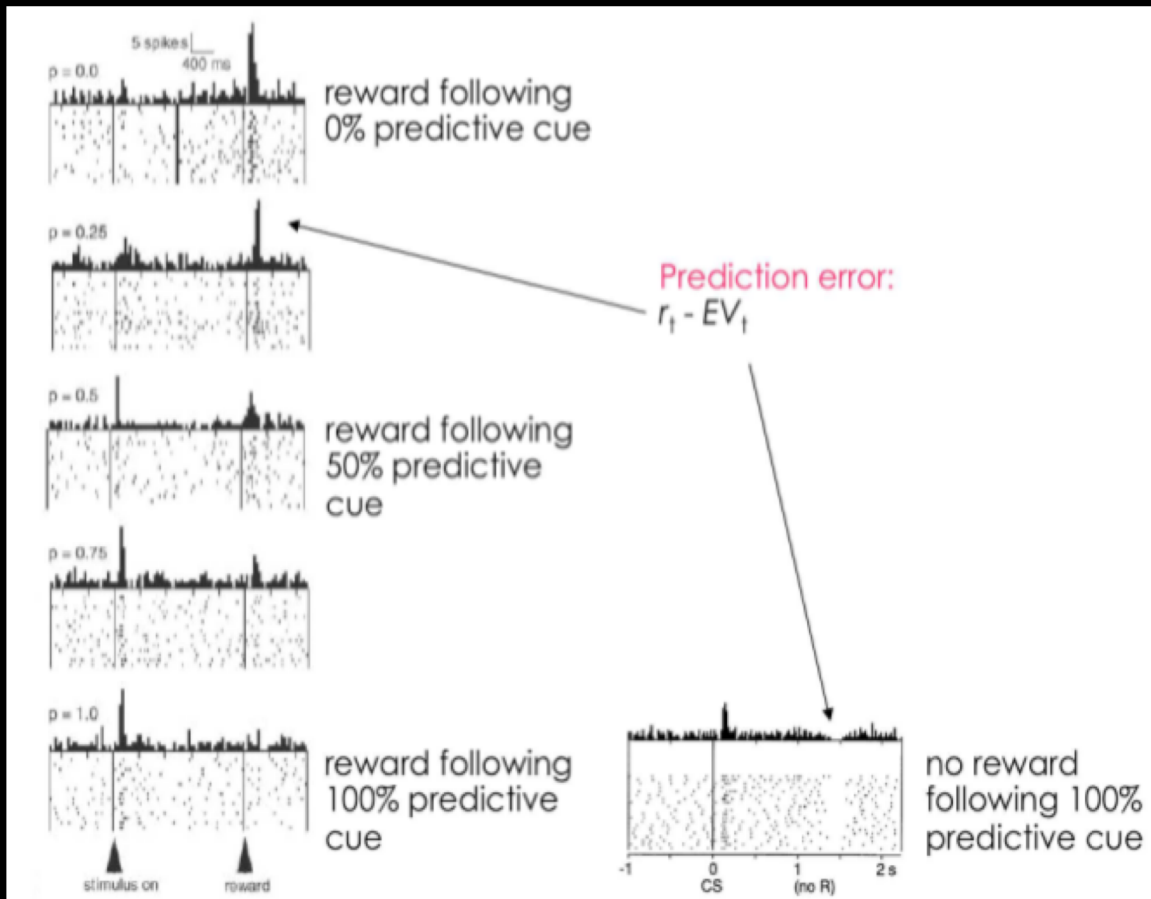
$$V(\text{lightbulb}) = .25$$

$$V(\text{bell}) = 0.5$$

$$V(\text{bone}) = 1$$

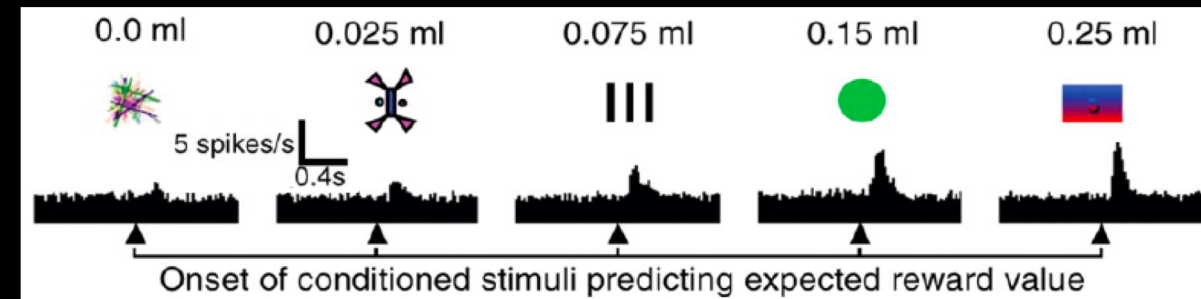
PE-hypothesis of DA

- Should be sensitive to **probability of reward** and **magnitude of reward**



Fiorillo et al, 2003

Tobler et al, 2005

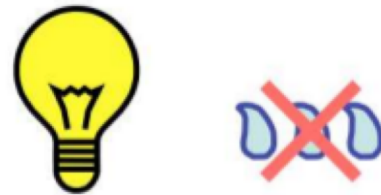
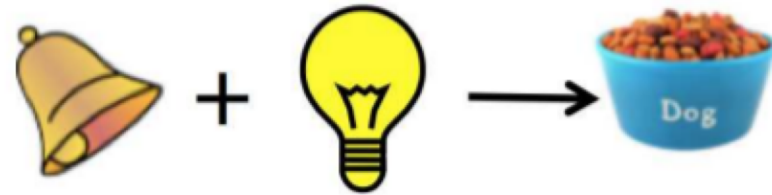


Blocking returns

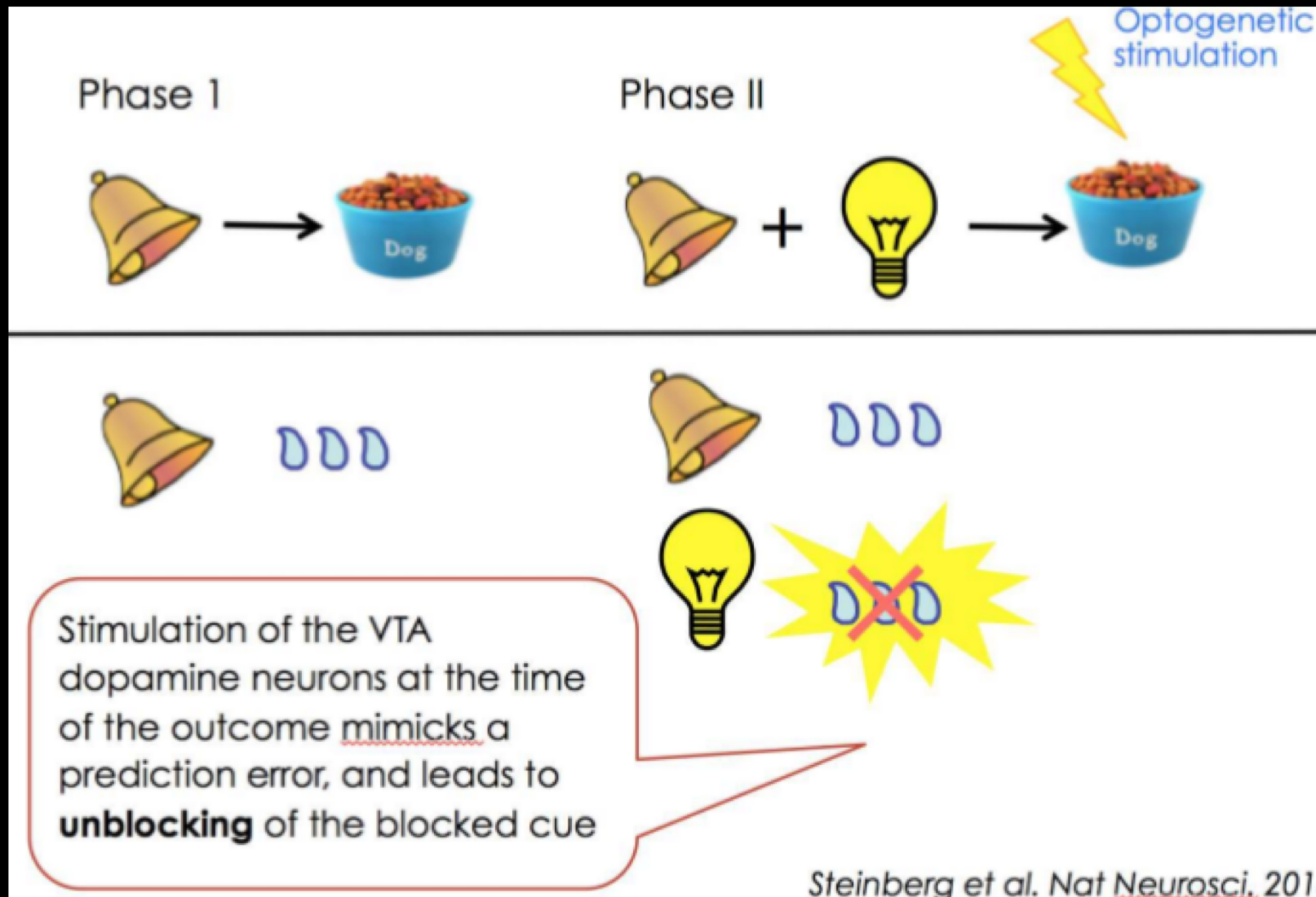
Phase 1



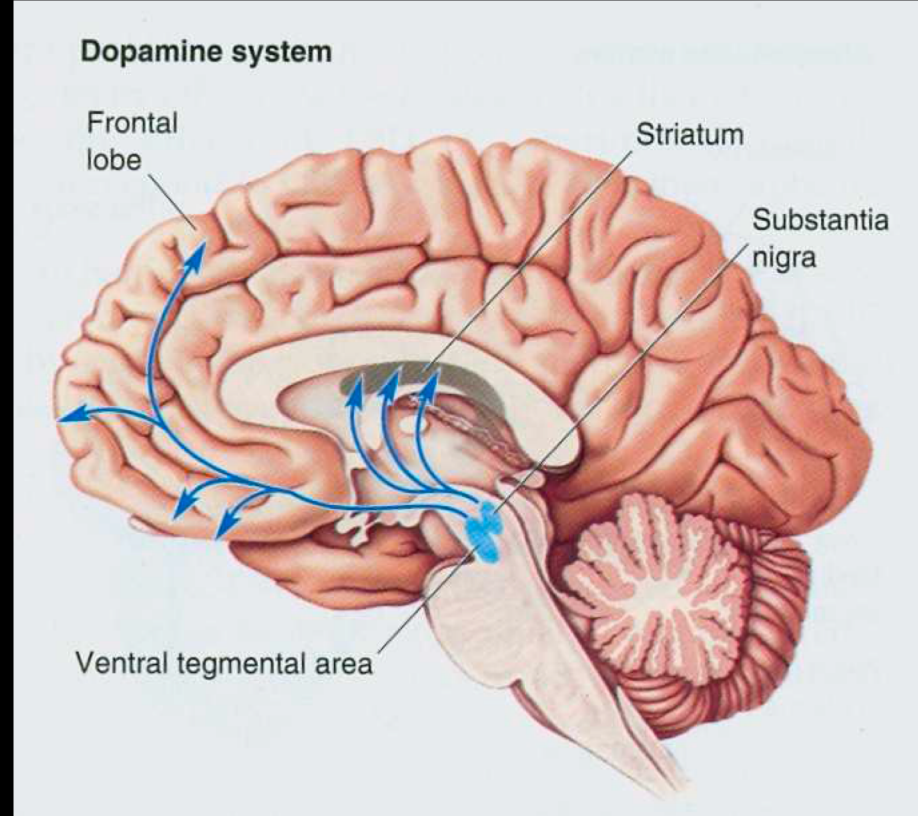
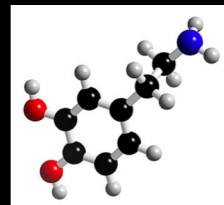
Phase II



Blocking: DA induces learning

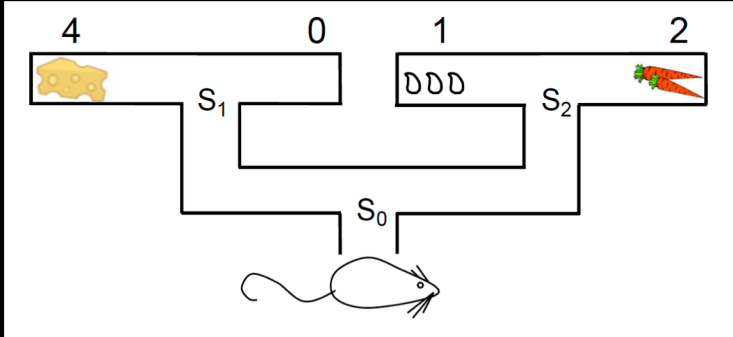


How are dopamine-based RPE signals used to select actions?



Will consider biological implementation in basal ganglia later

Q-Learning



The learning rate, i.e. that extent to which new information overrides old information. This is a number between 0 and 1.

The Q function we are updating, based on state s and action a at time t .

The reward earned when transitioning from time t to the next next turn, time $t+1$.

The value of the action that is estimated to return the largest (i.e. maximum) total future reward, based on all possible actions that can be made in the next state.

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha [r_{t+1} + \lambda \max_a Q(s_{t+1}, a) - Q(s_t, a_t)]$$

The arrow operator means update the Q function to the left. This is saying, add the stuff to the right (i.e. the difference between the old and the new estimated future reward) to the existing Q value. This is equivalent in programming to $A = A+B$.

The discount rate. Determines how much future rewards are worth, compared to the value of immediate rewards. This is a number between 0 and 1.

The existing estimate of the Q function, (a.k.a. current the action-value).

Alternative - SARSA: takes into account actual choice on next time step, “on policy”

Learning which actions to take

Critic

$V(s_t)$

expected value of
being in state s_t

Actor

$Q(s_t, a_t)$

preference (weight) for
taking action a_t in state s_t

Analogy: Player/coach

Learning which actions to take

Critic

$V(s_t)$

expected value of
being in state s_t

Actor

$Q(s_t, a_t)$

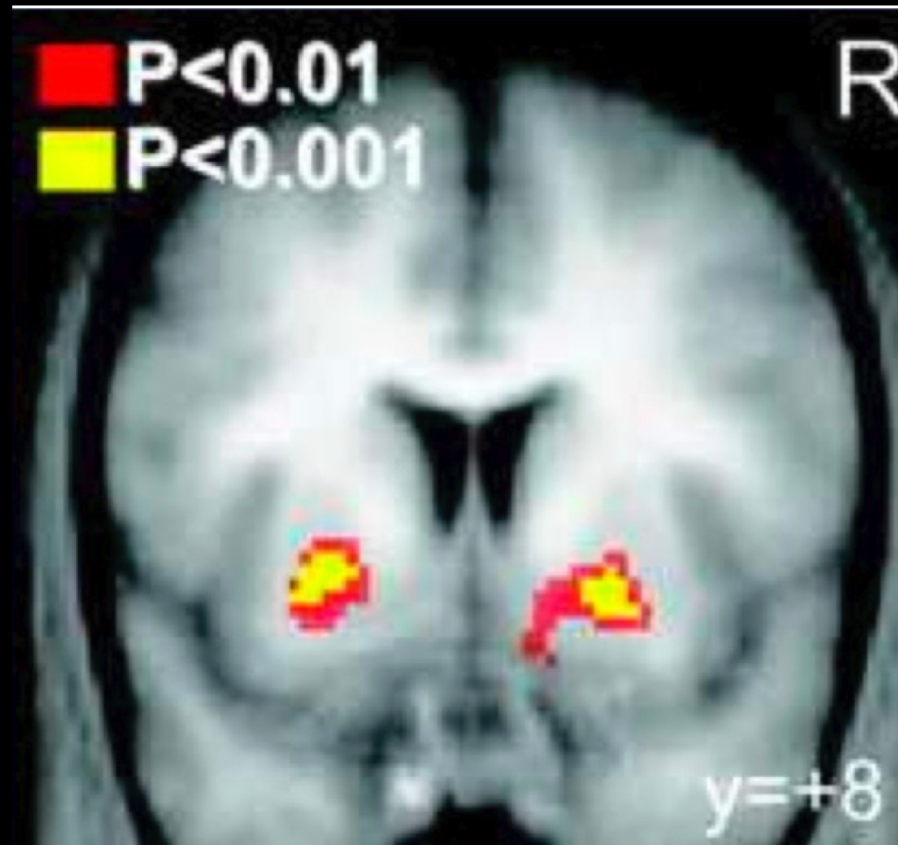
preference (weight) for
taking action a_t in state s_t

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

$$2) V(s_t) \leftarrow V(s_t) + \delta_t$$

$$3) Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \delta_t$$

fMRI & reward prediction errors



Ventral striatum correlates with reward prediction error: Critic!
Dorsal striatum correlates when actor involved

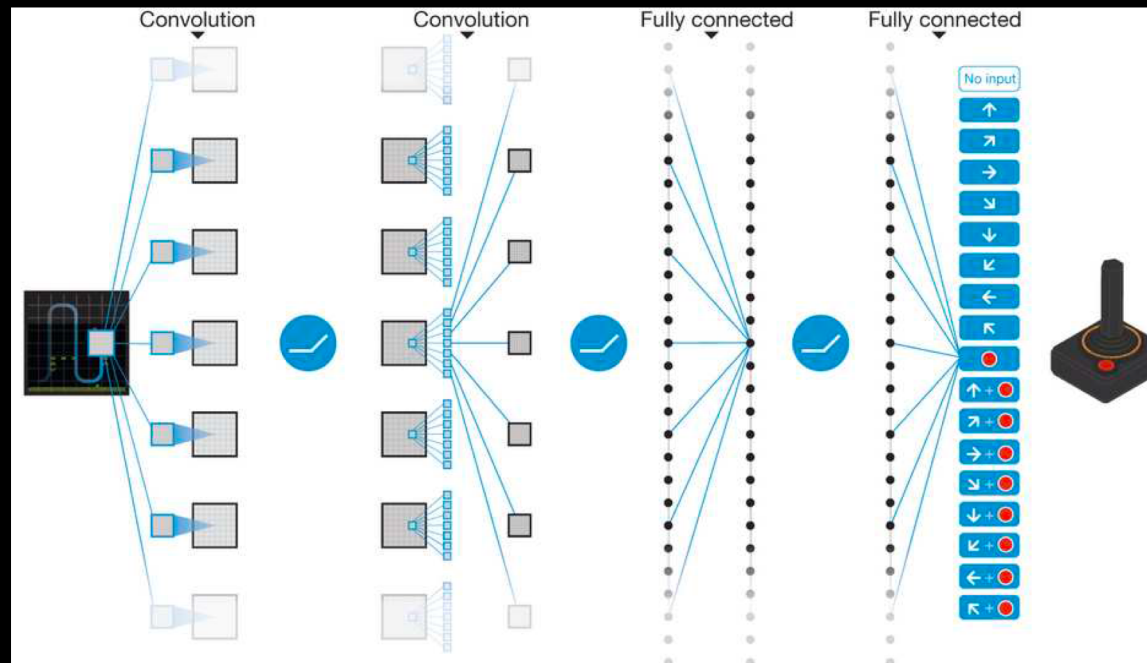
Google Deep Mind RL Network (“DQN”) Plays Atari

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. Fiedjeland¹, Georg Ostrovski¹, Stig Petersen¹, Charles Beattie¹, Amir Sadik¹, Ioannis Antonoglou¹, Helen King¹, Dharshan Kumaran¹, Daan Wierstra¹, Shane Legg¹ & Demis Hassabis¹



Assume you're a brain - revisited

You still need to eat things, you still don't want to be eaten by other things, and you'd still "like" to produce more brains.

- 1. What is learning?*
- 2. Why is learning important?*
- 3. What should you learn?*
- 4. When should you learn?*