Genela Morris
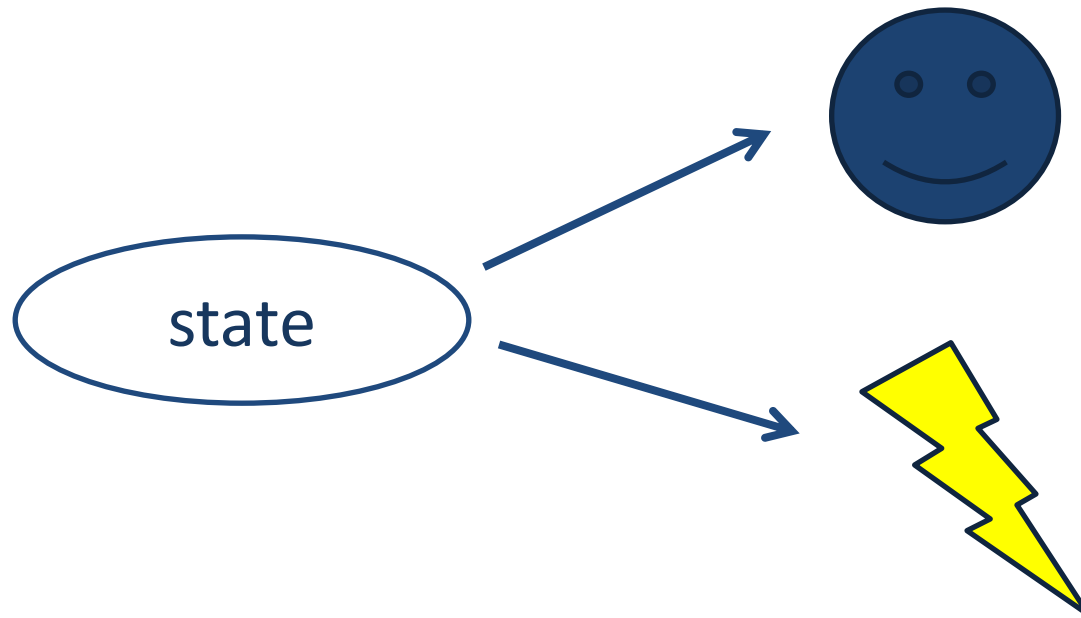Dept. of Neurobiology
Haifa University
gmorris@sci.haifa.ac.il

the role of dopamine in planning and action

# ON NEURAL CORRELATES OF REINFORCEMENT  LEARNING

# Reinforcement learning: finding correct action by trial and error

# Reinforcement learning
# the basics

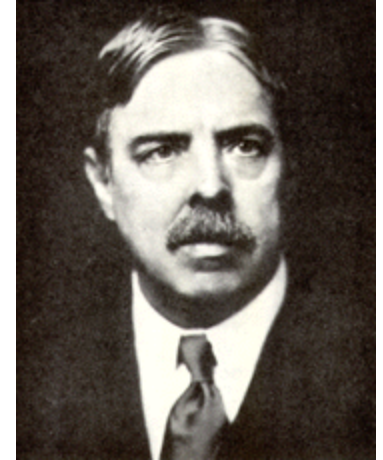Supervised learning –

all knowing teacher, detailed feedback

Reinforcement learning –

scalar (correct/incorrect) feedback

Unsupervised learning –

self organization
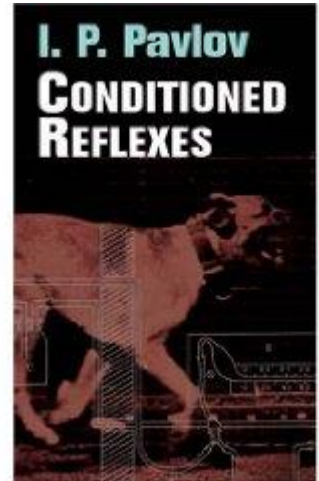
# Reinforcement learning: The law of effect

*"The Law of Effect is that: 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"*

Edward Lee Thorndike (1911)

# Early attempts at modeling

- By associative rules
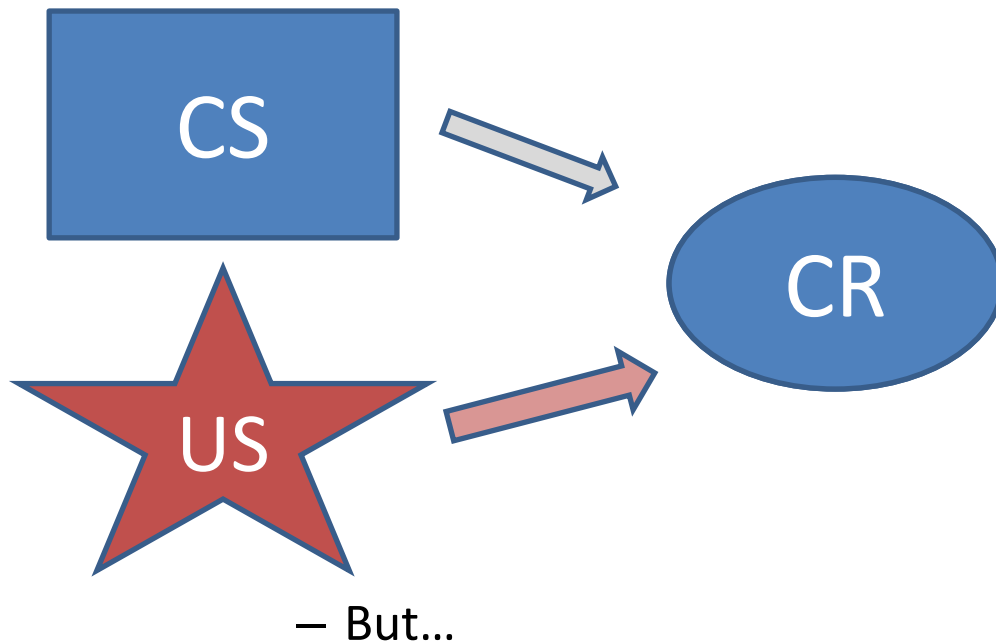- Classical conditioning

# Properties of classical conditioning

*(Pavlov 1927)*

- **Acquisition**.

- **Partial Reinforcement** (probabilistic).

- **Generalization**.

- **Interstimulus Interval (ISI) effects**.

- **Intertrial Interval (ITI) effects**.

# So far…

- A simple association (coincidence, Hebbian) model can explain the phenomenon.



— But…

- **Acquisition**.
- **Partial Reinforcemen**(probabilistic).
- **Generalization**.
- **Interstimulus Interval (ISI) effects**.
- **Intertrial Interval (ITI) effects**.

# Classical conditioning

**The Elements:**
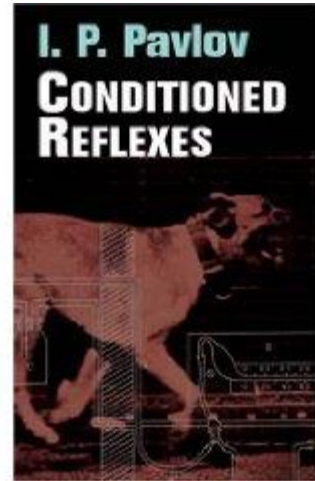**US**: Unconditioned stimulus
**UR**: Unconditioned response
**NS**: Neutral stimulus
**CS**: Conditioned stimulus
   **CS1**: Conditioned stimulus 1
    **CS2**: Conditioned stimulus 2
**CR**: Conditioned response

# Properties of classical conditioning

*(Cnt'd)*

- **Conditioned  Inhibition**

- **latent inhibition**

- **Relative validity** (Wagner 1968).

- **Blocking** (Kamin 1968)

- **…**

**CS must RELIABLY predict US**

# Which simple association can't explain

*Learning occurs not because two events co-occur, but because that co-occurrence is otherwise UNPREDICTED*

# Rescorla-Wagner rule (1972)

Learning to predict reward R given stimulus U=1

Goal: Form a prediction V of the reward of the form:

$V=\omega U$

And learn to change $\omega$ :

$\Delta \omega = \varepsilon(R-V)U$

After learning of consistent pairing: $\omega=R$

*Where:*
*U=CS availability (0,1);*
*V=reward prediction:*
*R=reward availability (0,1) :*
*$\omega$ = weight of the connection between U and V*
*$\varepsilon$ = learning rate*
*R-V = prediction error*

# Blocking with Rescorla Wagner

- Given U1, U2 and R, after U1 has been learnt:
- $\omega 1 = R$
- $V = \omega 1 U1 + \omega 2 U2$

  R          0

- Prediction error: $R - V = 0$

  And no learning occurs for $\omega 2$

# Critical problems, for control

1. Exploration/exploitation

# Solutions, for control

1. Variability in response policy

    1. Greedy ← → Random (gambling)

    2. Based on expected return

# Decision behaviour, theory and practice



C = choice

$C_{right}/(C_{right}+C_{left})$

1

0.5

0

**maximizing**

**probability-matching**

$$\frac{C_{right}}{C_{right}+C_{left}}(\infty) = \frac{R_{right}}{R_{right}+R_{left}\cdot\frac{\theta_{left}}{\theta_{right}}}$$

**monkeys?**

0.5

1

$R_{right}/(R_{right}+R_{left})$

R = reward

# Monkeys' decisions: probability matching

# … whether optimal or not

- Actions are related to their consequences

# Critical problems in reinforcement learning (and in Rescorla-Wagner)

2. Temporal credit assignment

# TD learning - solution for temporal credit assignment

1. Estimate value of current state $(V_t = r_t + \gamma' r_{t+1} + \cdots)$ : (discounted) sum of expected rewards

2. Measure 'truer' value of current state: reward at present state + estimated value of next state $(r_t + \gamma V_{t+1})$

3. TD error $\quad \delta_t = r_t + \gamma V_{t+1} - V_t$

4. Use TD error to improve 1 $(V_t^{k+1} = V_t^k + \eta\, \delta_t)$

*where:* $V_{t = value}$ *of the state reached at time t in iteration k*

$r_t$ = *reward given at time t;* $\eta$ = *learning rate,* $\delta$ = *prediction error*

# TD error: $\delta_t = r_t + \gamma V_{t+1} - V_t$



Before learning | After learning

stimulus (t)

reward (t)

value (t)

value (t+1)

TD error (t)

time

# TD error: $\delta_t = \gamma V_{t+1} - V_t + r_t$



Before learning | After learning

stimulus (t)

reward (t)

value (t)

value (t+1)

TD error (t)

time

# Basal ganglia - anatomy

# Intracranial self stimulation



ACTIVATES REWARD CIRCUITS

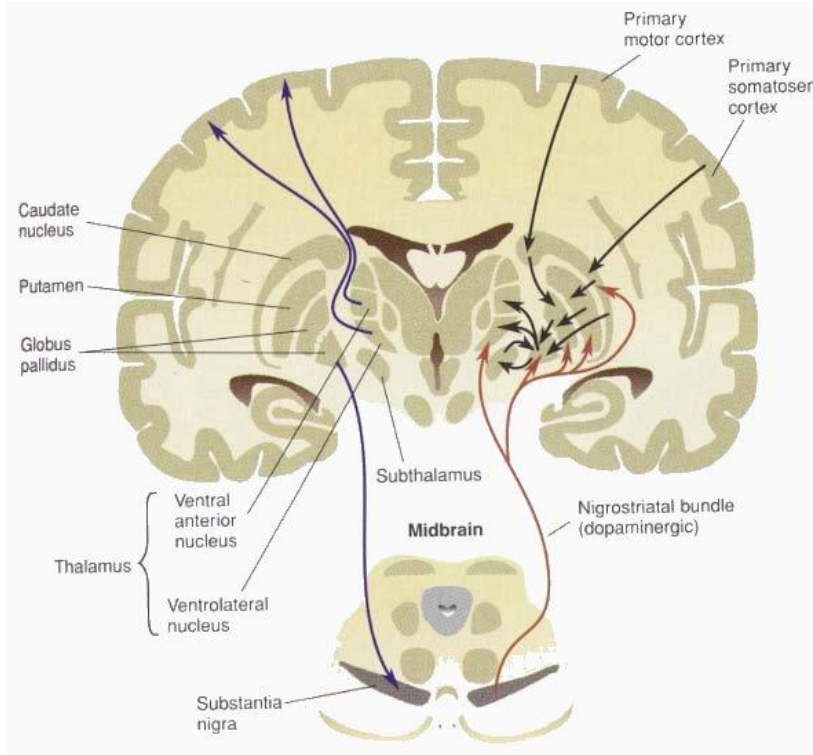# The midbrain dopamine system

# Dopamine and acetylcholine meet in the striatum



Monkey



Mouse

# Facts to remember (1)

- Basal ganglia receive cortical input
- Basal ganglia project to frontal cortex
- Dopamine and acetylcholine localization

# The midbrain dopamine system



Schultz et al,
JNS 13:
900-913, 1993

# Probabilistic instrumental conditioning task
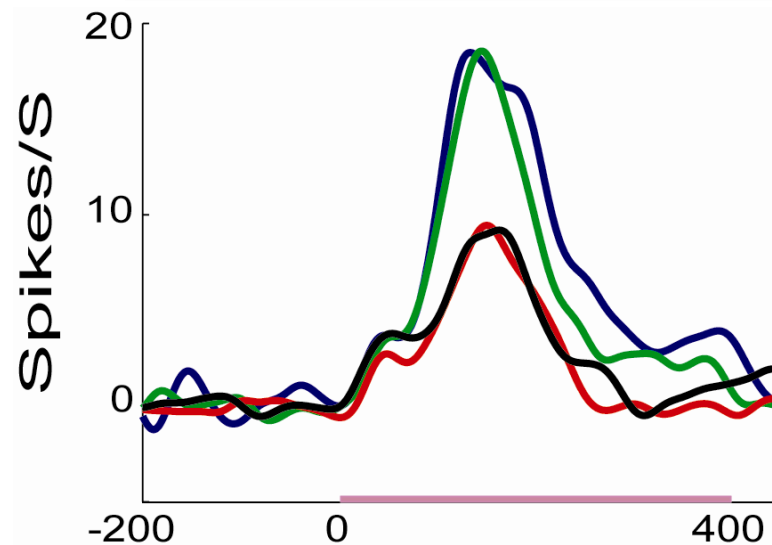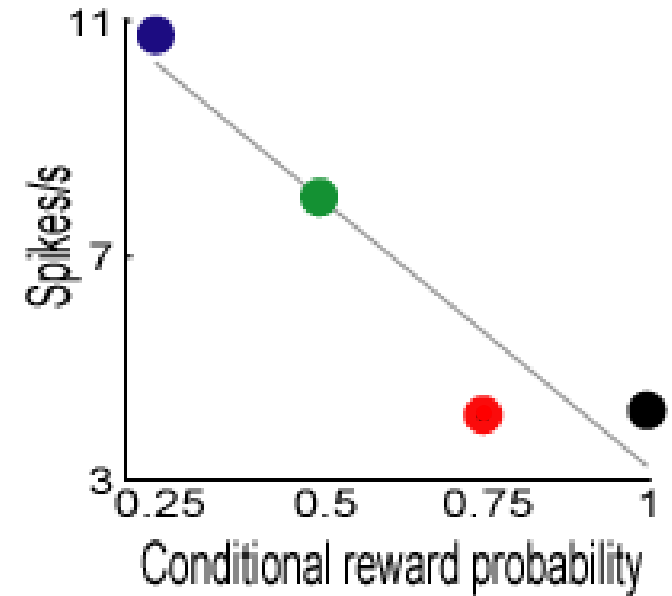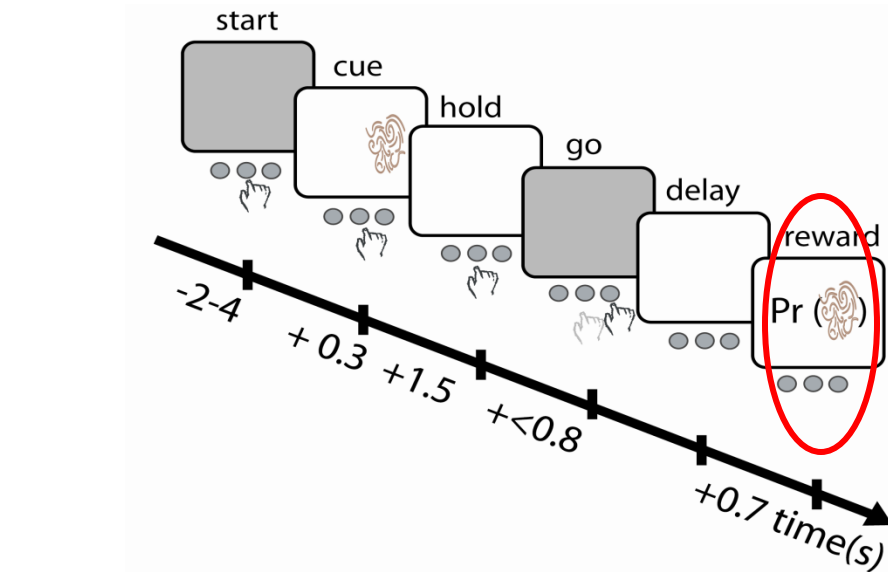


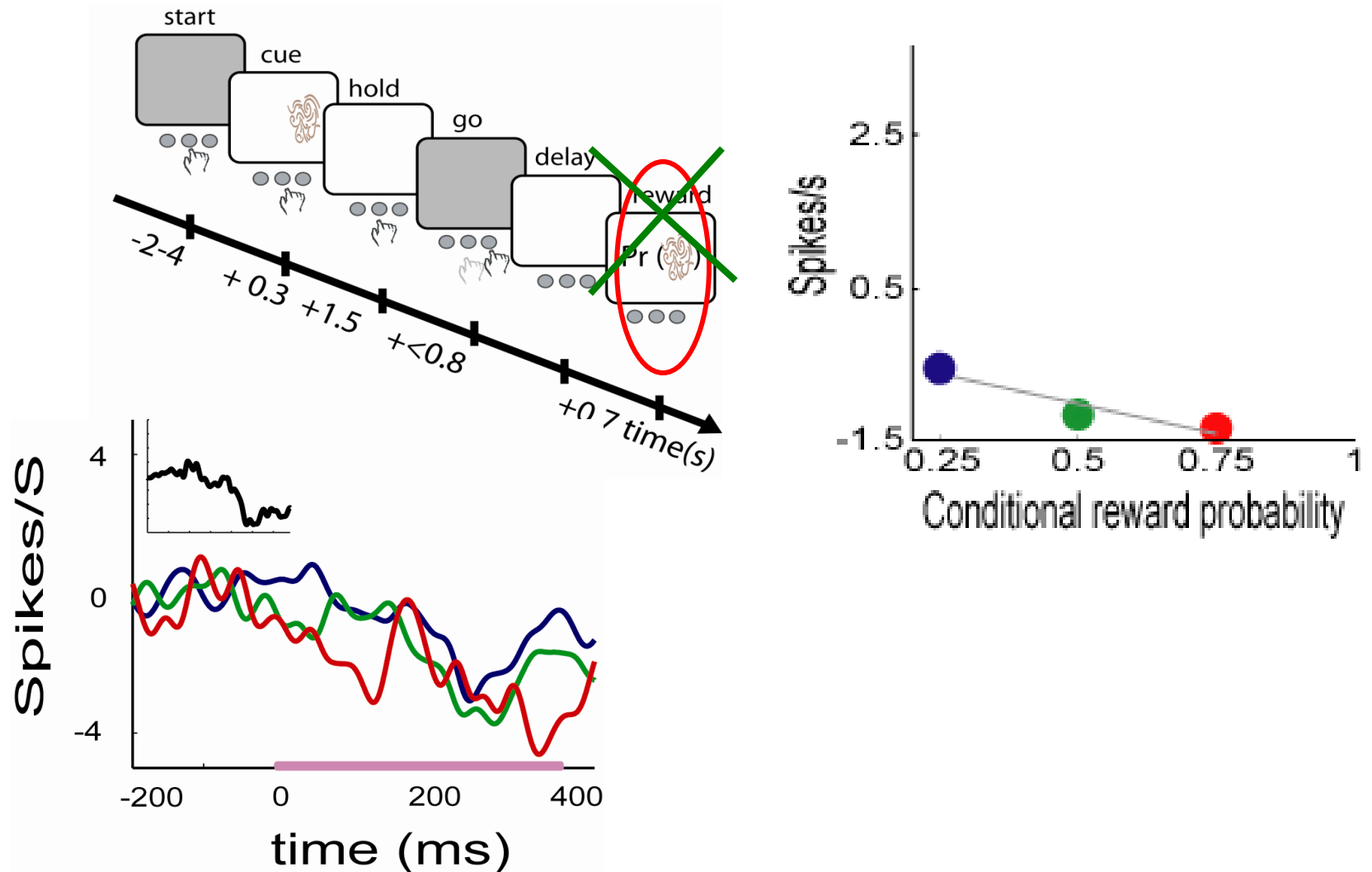$$\delta_t = \gamma V_{t+1} - V_t + r_t$$

# *DA response*

# Dopamine population response- cue

Weizmann systems

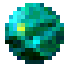# Dopamine population response- reward

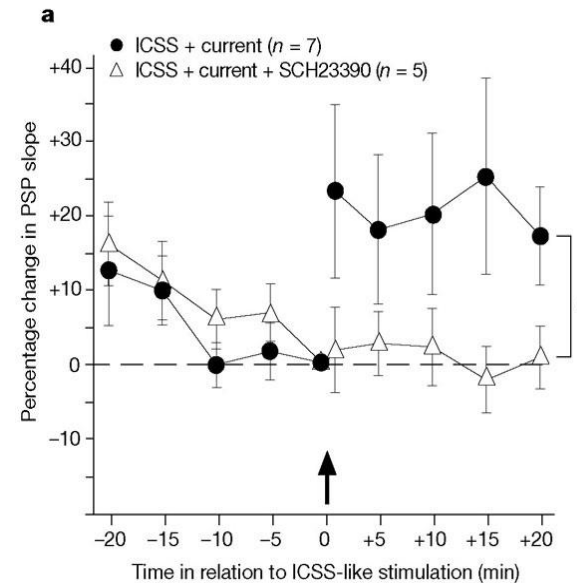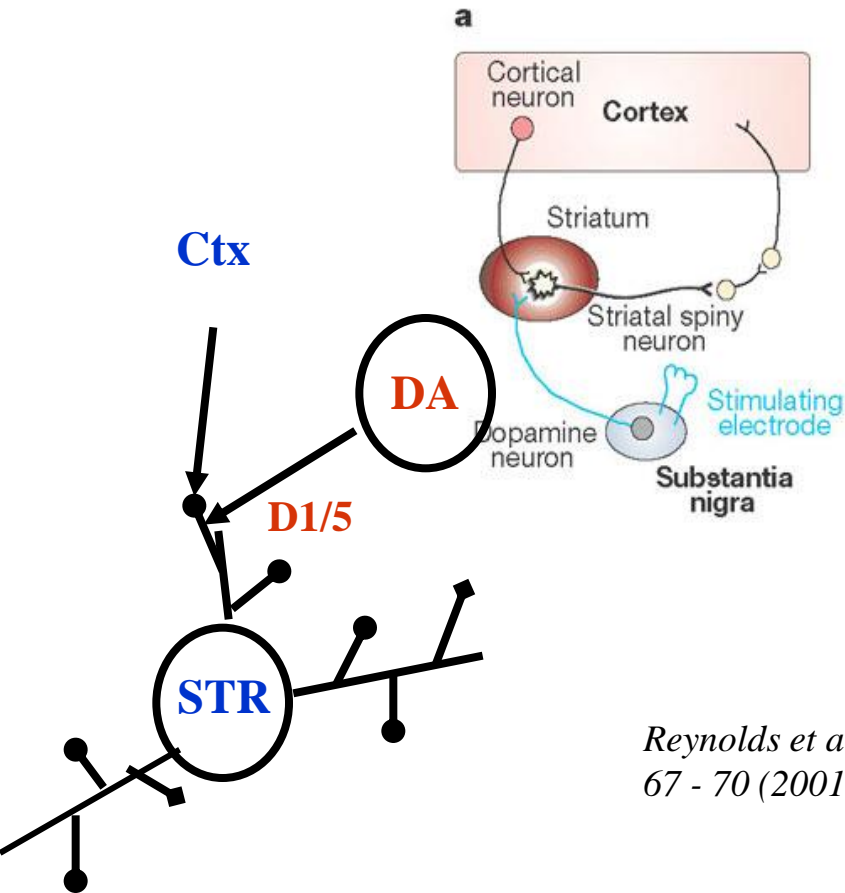# Dopamine population response – reward omission
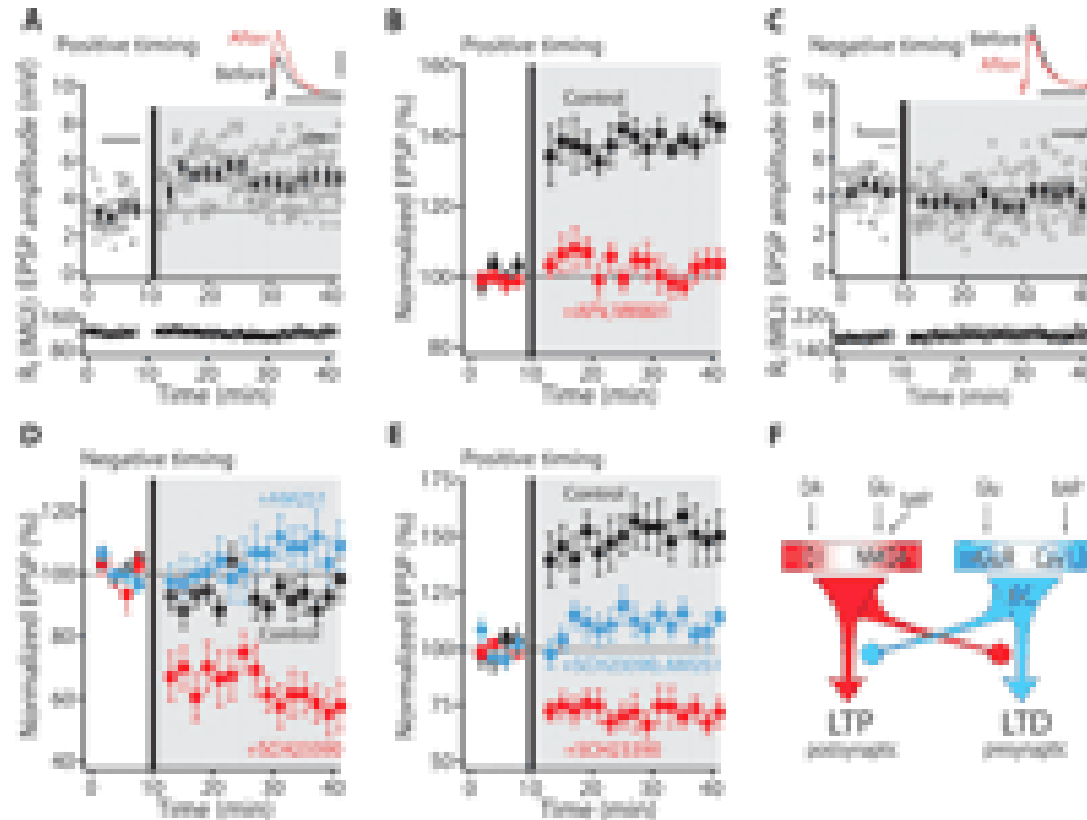
# Instrumental conditioning - results

- Responses to visual cue are correlated with future reward probability

- Responses to reward are inversely correlated with reward probability

- Responses to reward omission are indifferent to reward probability

- Dopamine neurons provide an accurate TD signal (but only in the positive domain)

# … and it can cause long term plasticity of cortico-striatal synapses



**Ctx**

**DA**

**D1/5**

**STR**

*Reynolds et al, A cellular mechanism of reward-related learning Nature 413, 67 - 70 (2001)*

# … and it can cause long term plasticity of cortico-striatal synapses

Weizmann systems *Shen et al., 2008*

# Facts to remember 2

- DA neurons provide a TD error signal
- To the cortico (state) striatal (action) synapses
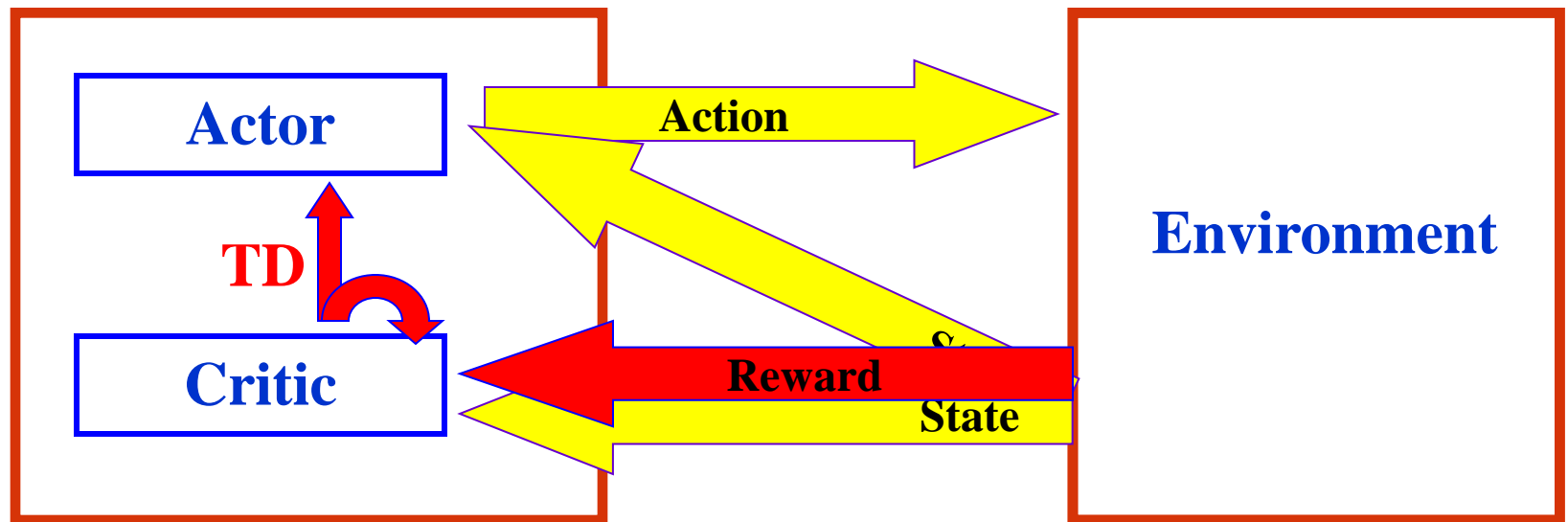- And DA modulates synaptic plasticity

# Control - Adding action



The agent has to:
  – Learn to predict reinforcement                    *state value*
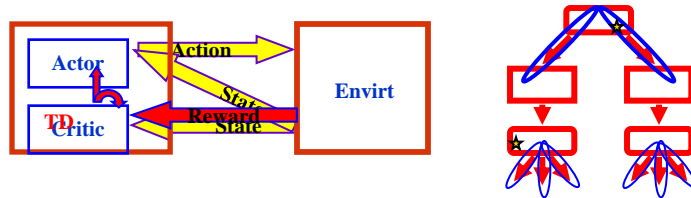  – Know the state-action-state transitions      *behavioural*
    *policy*

# Solution 1: actor/critic networks

# How can the dopamine signal contribute to decision behaviour?
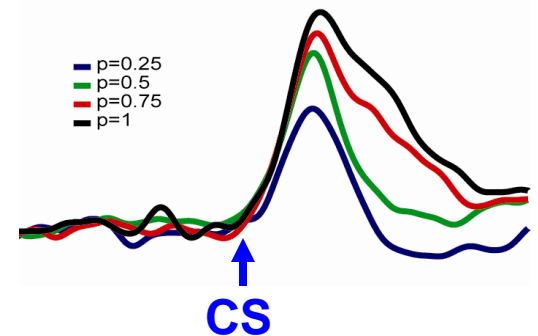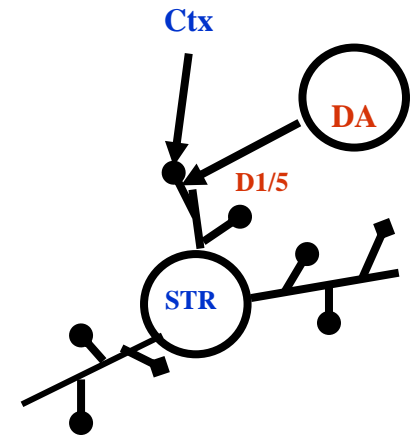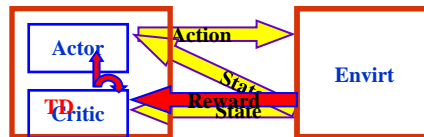
- Long term policy-shaping effect
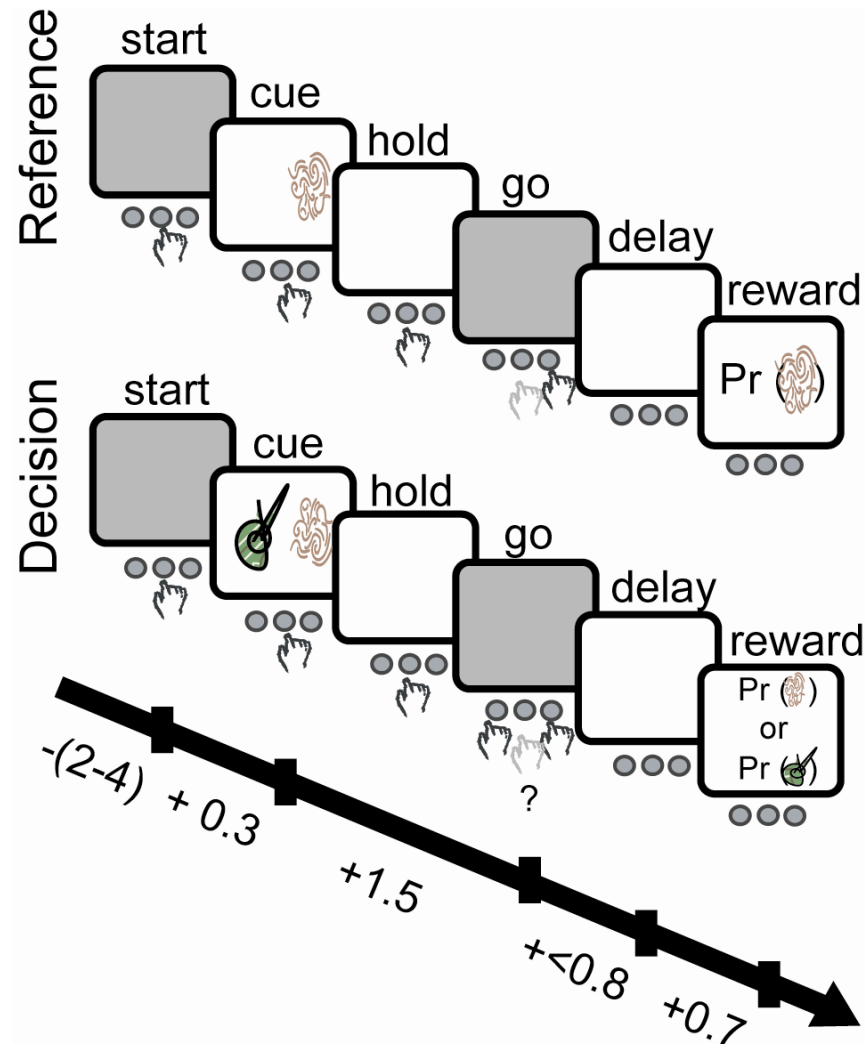
  through synaptic plasticity



- Immediate effect on action
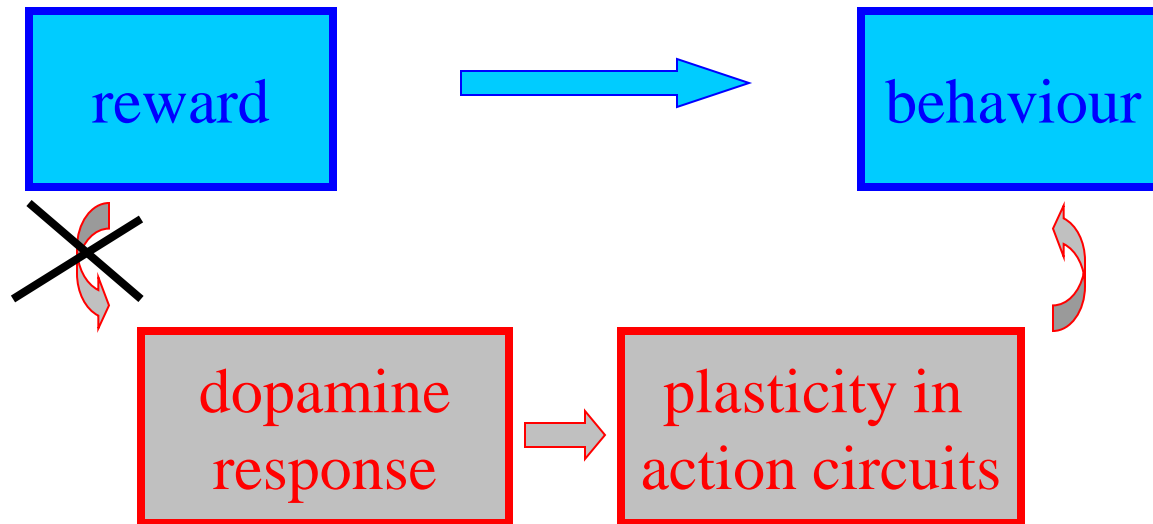
$$P_{action} = \frac{1}{1 + e^{-m\delta(t) + b}}$$
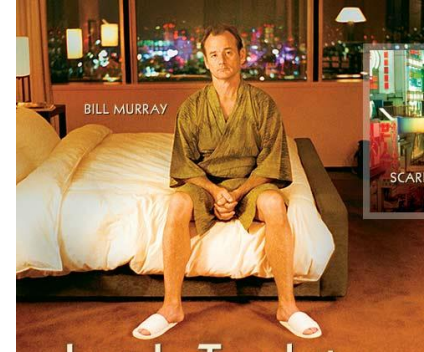
Weizmann systems

# The two armed bandit task

Weizmann systems

# Lost in translation?
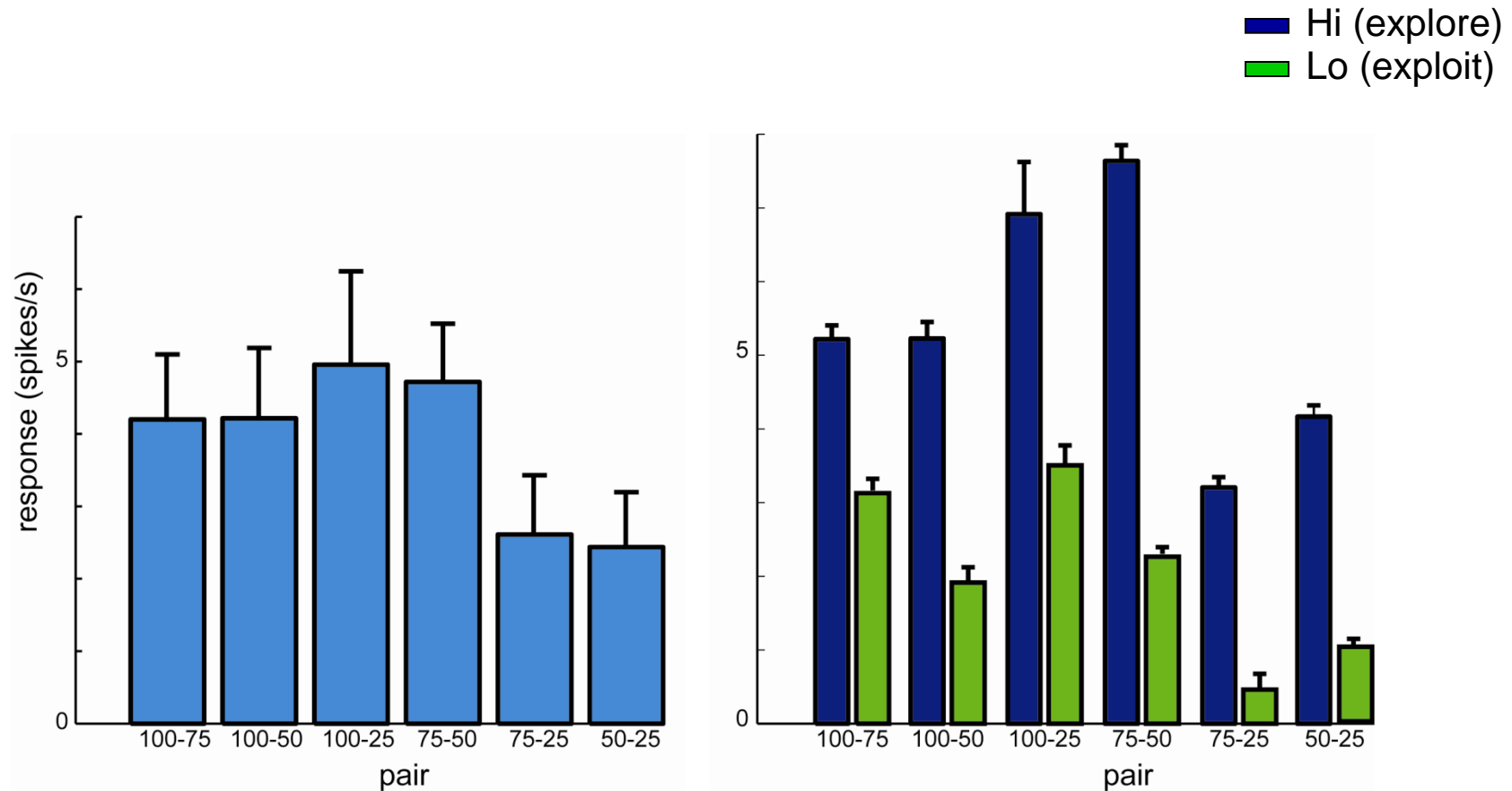


reward → behaviour

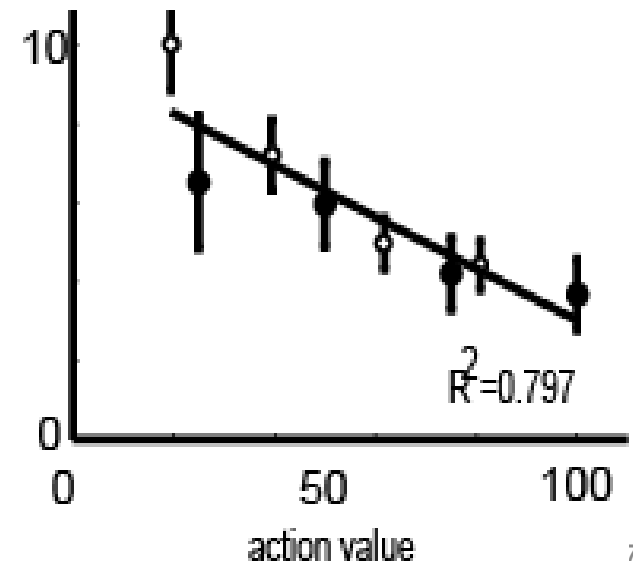dopamine response → plasticity in action circuits
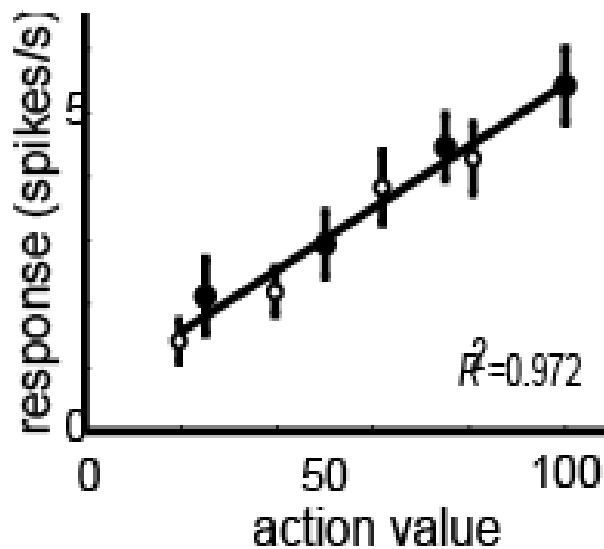
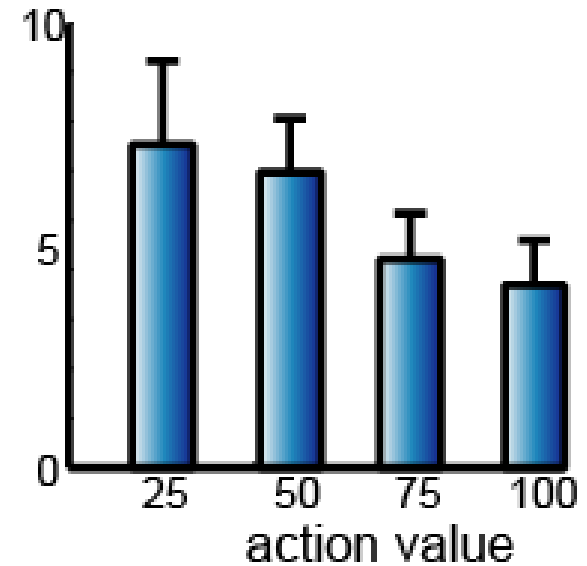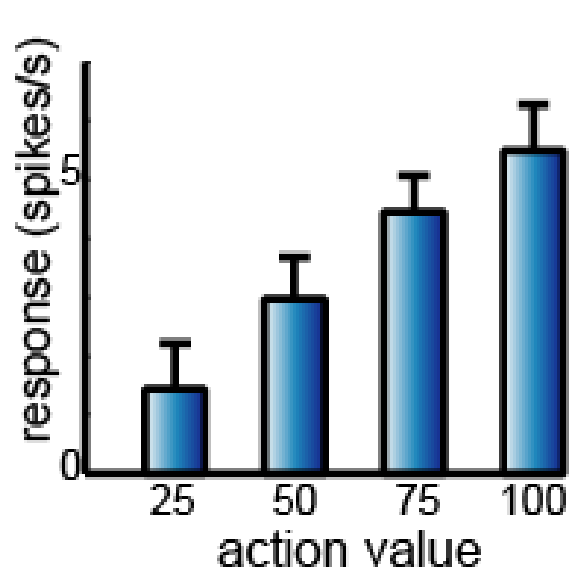# Monkeys' decisions: shaping by dopamine

# Dopamine neurons during decision

# Are DA neurons aware of future choice

# The learning is of state-action values



ystems

# Adding an internal model



Ctx

DA

D1/5

STR

p=0.25
p=0.5
p=0.75
p=1

CS
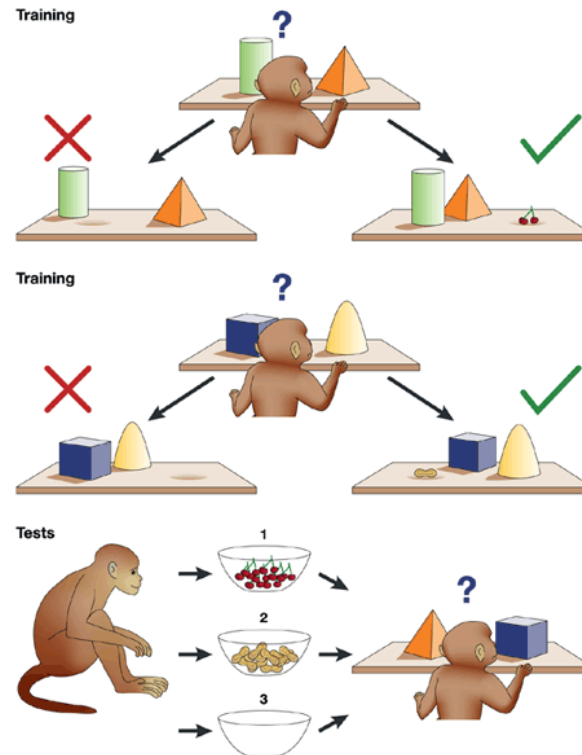
No explicit knowledge about the future

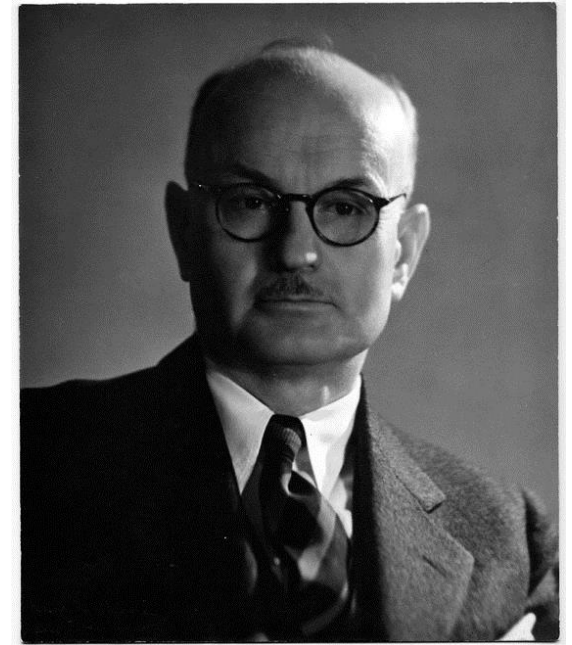# Reinforcement devaluation – evidence for model



Nature Reviews | Neuroscience

# Model based learning

We believe that in the course of learning something like a field map of the environment gets established in the rat's brain… Although we admit that the rat is bombarded by stimuli, we hold that his nervous system is surprisingly selective as to which of these stimuli it will let in at any given time… Rather, the incoming impulses are usually worked over and elaborated in the central control room into a tentative, cognitive-like map of the environment. And it is this tentative map, indicating routes and paths and environmental relationships, which finally determines what responses, if any, the animal will finally release.

*Cognitive maps of rats and men (1948)*
*The Psychological Review, 55(4),* 189-208

*Edward C. Tolman*

# Model based learning

We believe that in the course of learning something like a field map of the environment gets established in the rat's brain… **Although we admit that the rat is bombarded by stimuli, we hold that his nervous system is surprisingly selective as to which of these stimuli it will let in at any given time…** Rather, the incoming impulses are usually worked over and elaborated in the central control room into a tentative, cognitive-like map of the environment. And it is this tentative map, indicating routes and paths and environmental relationships, which finally determines what responses, if any, the animal will finally release.

*Cognitive maps of rats and men (1948)*
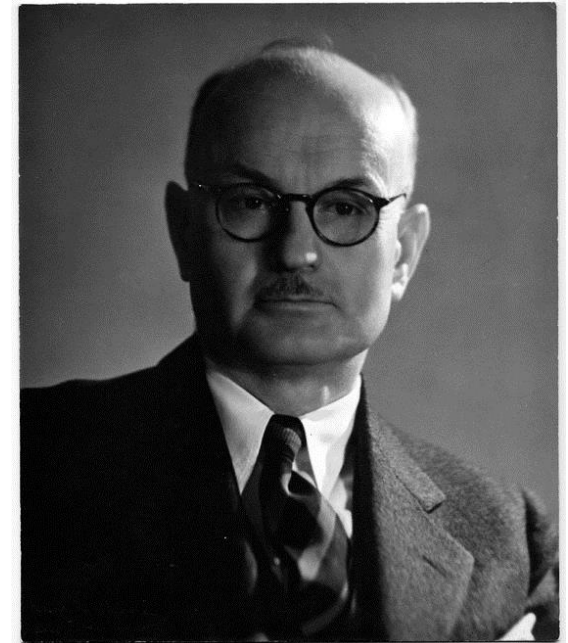*The Psychological Review, 55(4),* 189-208

*Edward C. Tolman*