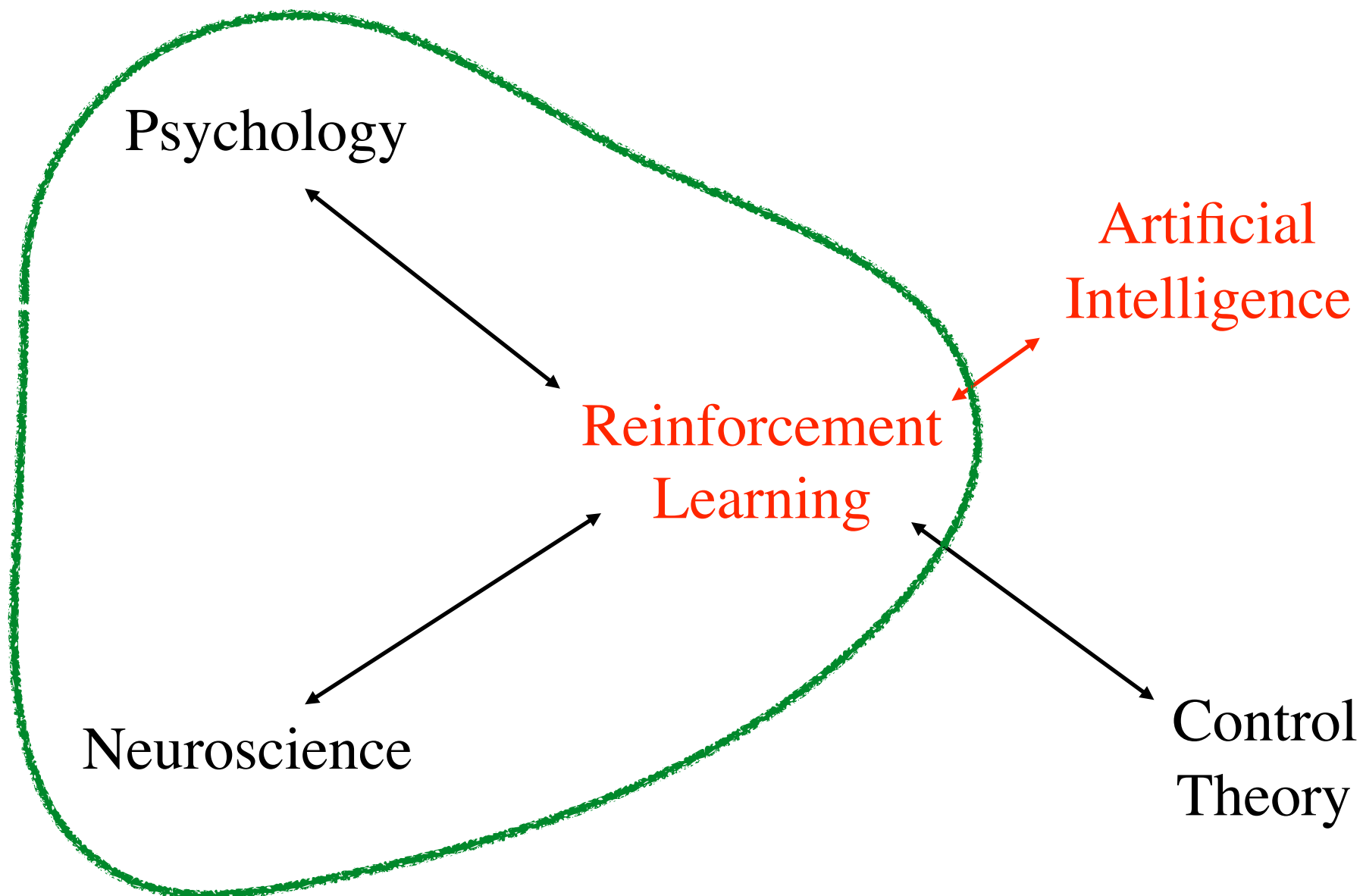


# Reinforcement Learning in Psychology and Neuroscience



with thanks to  
Elliot Ludvig  
University of Warwick



# Any information processing system can be understood at multiple “levels”

- The Computational Theory Level
  - *What* is being computed?
  - *Why* are these the right things to compute?
- Representation and Algorithm Level
  - *How* are these things computed?
- Implementation Level
  - How is this implemented physically?



# Goals for today's lecture

- To learn:
  - That psychology recognizes two fundamental learning processes, analogous to our prediction and control.
  - That all the ideas in this course are also important in completely different fields: psychology and neuroscience
  - That the details of the TD( $\lambda$ ) algorithm match key features of biological learning

# Psychology has identified two primitive kinds of learning

- *Classical* Conditioning
- *Operant* Conditioning (a.k.a. Instrumental learning)
- Computational theory:
  - ❖ *Classical* = Prediction
    - What is going to happen?
  - ❖ *Operant* = Control



# Classical Conditioning



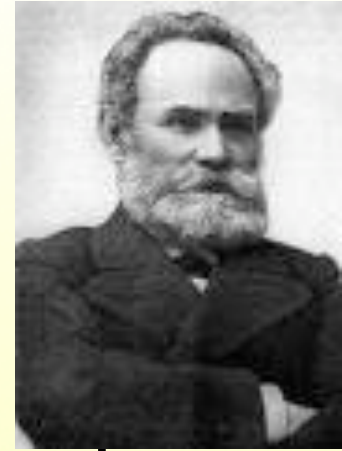


# Classical Conditioning as Prediction Learning

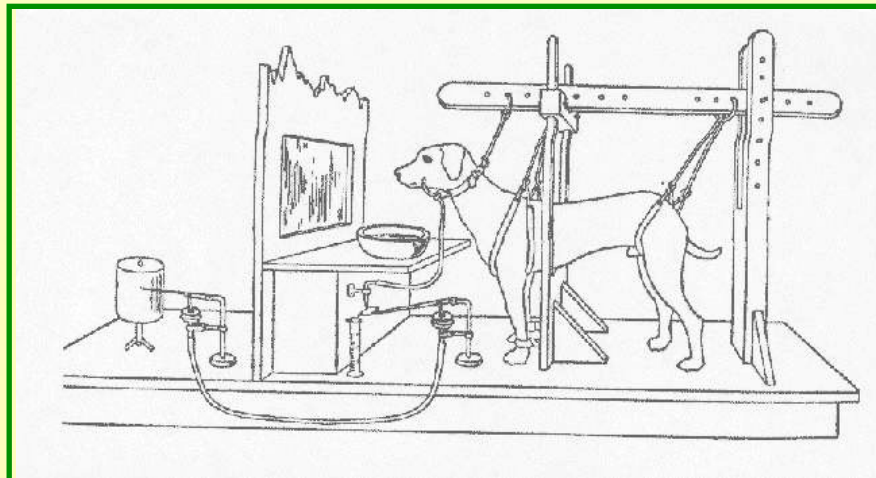
- Classical Conditioning is the process of learning to predict the world around you
  - ❖ Classical Conditioning concerns (typically) the subset of these predictions to which there is a hard-wired response



# Pavlov (1901)



- Russian physiologist
- Interested in how learning happened in the brain
- Condition<sup>al</sup> and Uncondition<sup>al</sup> Stimuli





Is it really  
predictions?



# Maybe Contiguity?

- Foundational principle of classical associationism (back to Aristotle)
  - ❖ Contiguity = Co-occurrence
  - ❖ Sufficient for association?

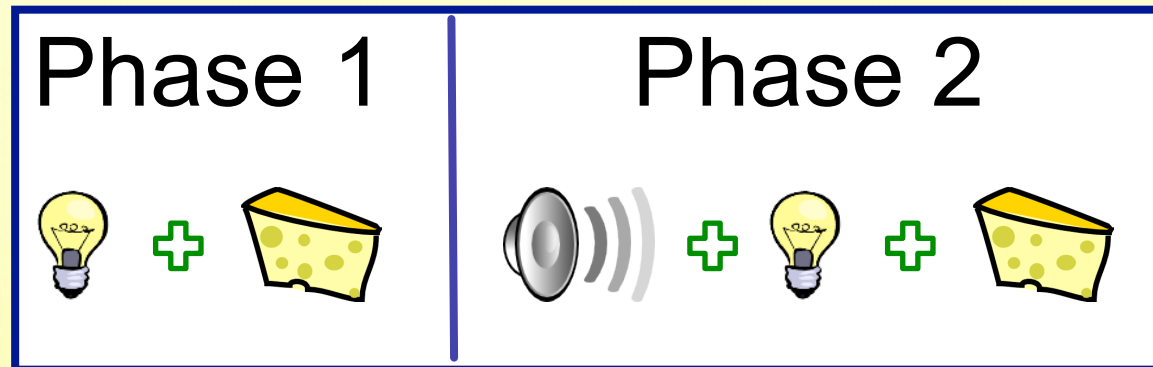


# Contiguity Problems

- Unnecessary:
  - ❖ Conditioned Taste Aversion
- Insufficient:
  - ❖ Blocking
  - ❖ Contingency Experiments



# Blocking





# Rescorla-Wagner Model (1972)



- Computational model of conditioning
  - ❖ Widely cited and used
- Learning as violation of expectations
  - ❖ As in linear supervised learning (p2)
  - ❖ TD learning is real-time extension



# Any information processing system can be understood at multiple “levels”

- The Computational Theory Level
  - *What* is being computed?
  - *Why* are these the right things to compute?
- Representation and Algorithm Level
  - *How* are these things computed?
- Implementation Level
  - How is this implemented physically?



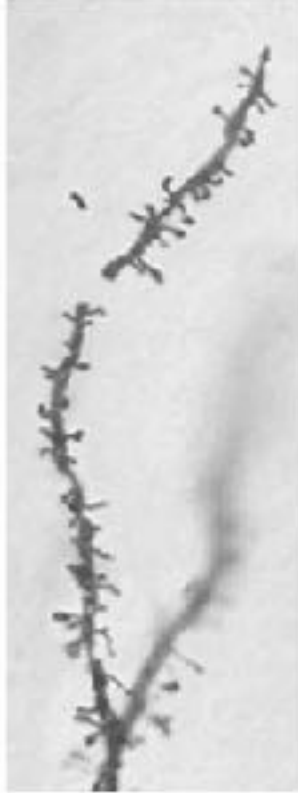


# The Basic TD Model

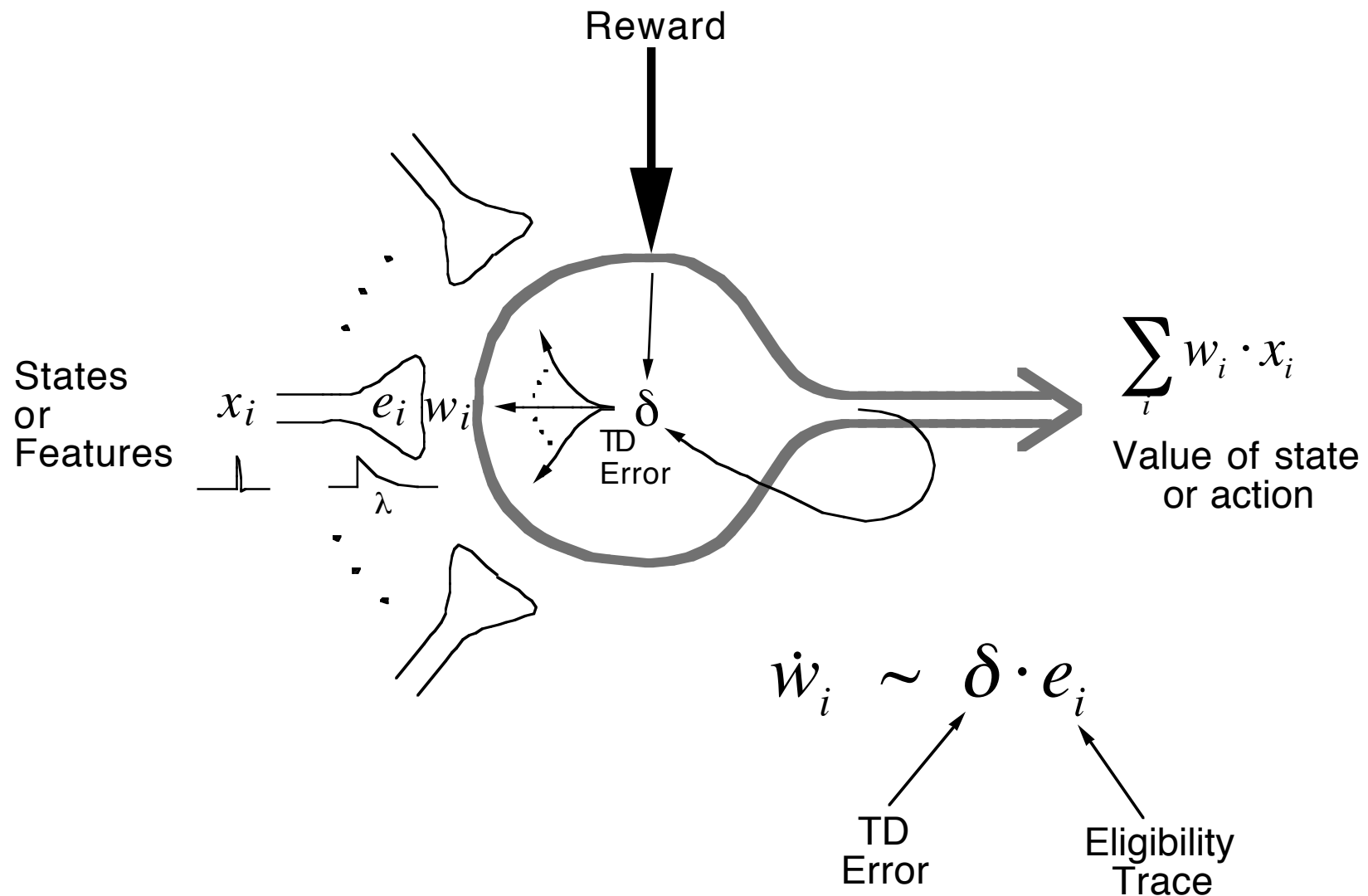
- Learn to predict discounted sum of upcoming reward through TD with linear function approximation
- The TD error is calculated as:

$$\delta_t = R_{t+1} + \gamma \hat{v}_{t+1} - \hat{v}_t$$



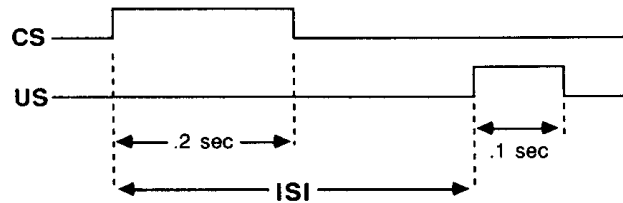


# TD( $\lambda$ ) algorithm/model/neuron

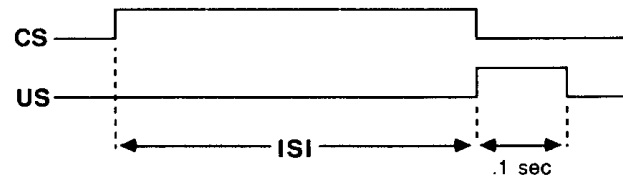


# Effect of inter-stimulus interval

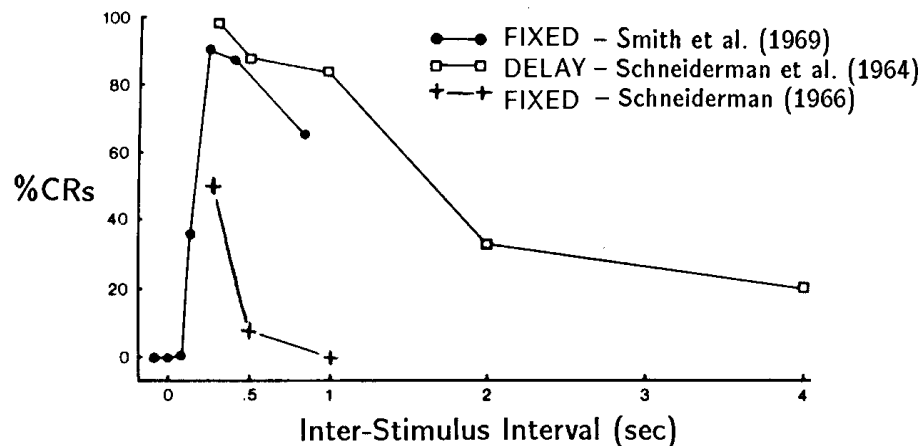
## FIXED-CS CONDITIONING



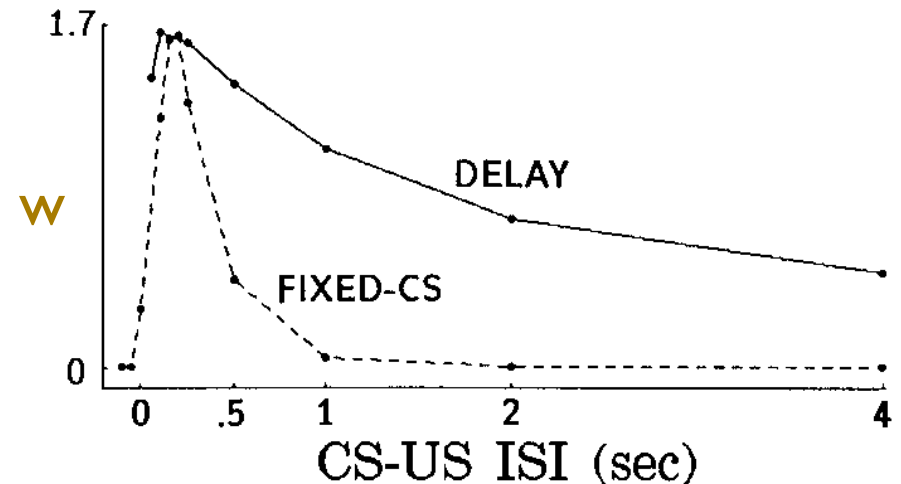
## DELAY CONDITIONING



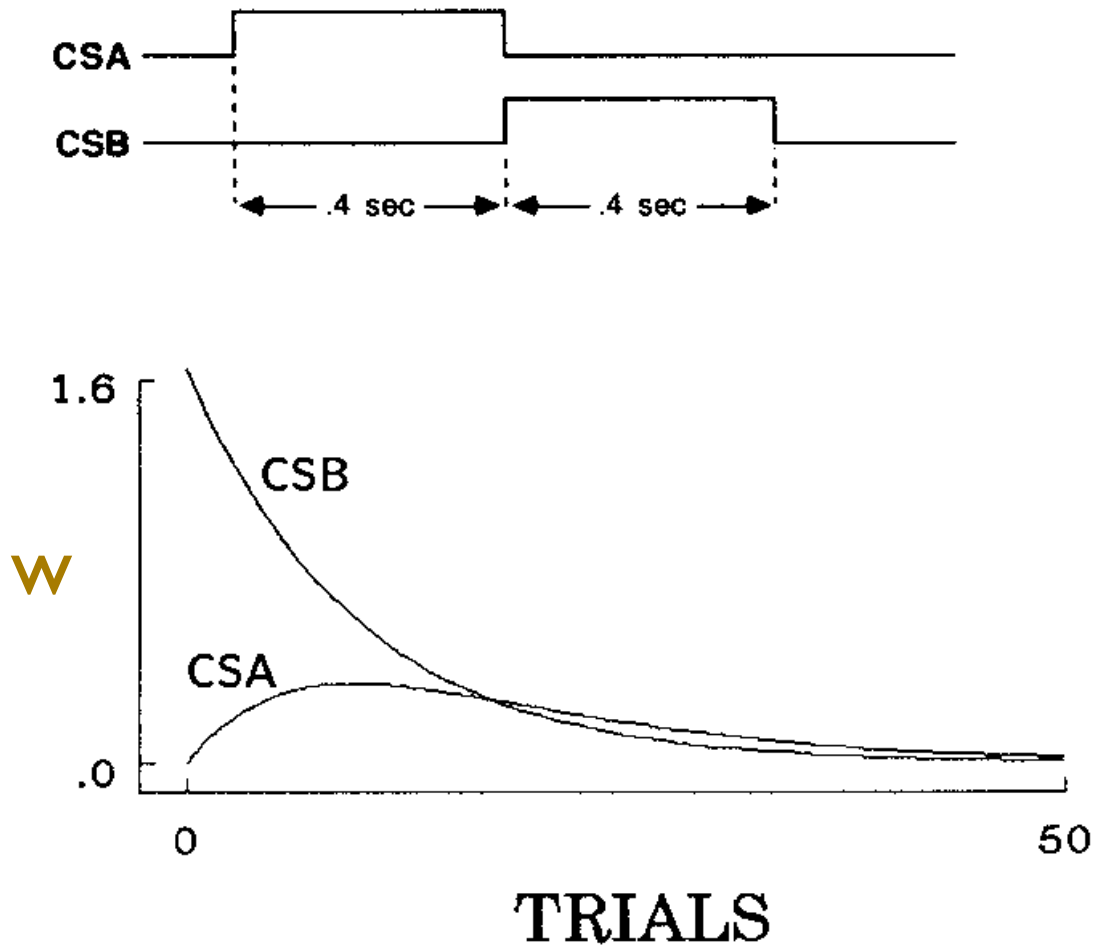
## Animal data (eyeblick)



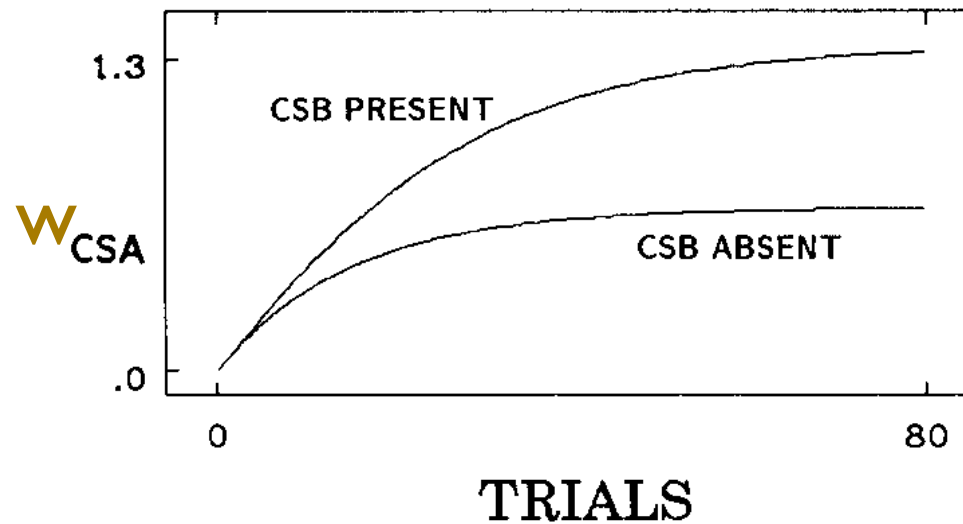
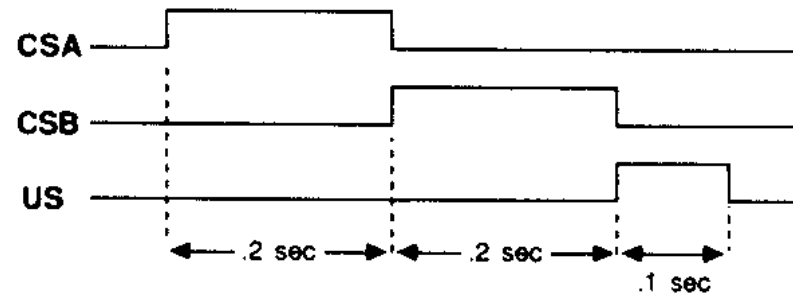
## TD model



# Second-order conditioning in the TD model



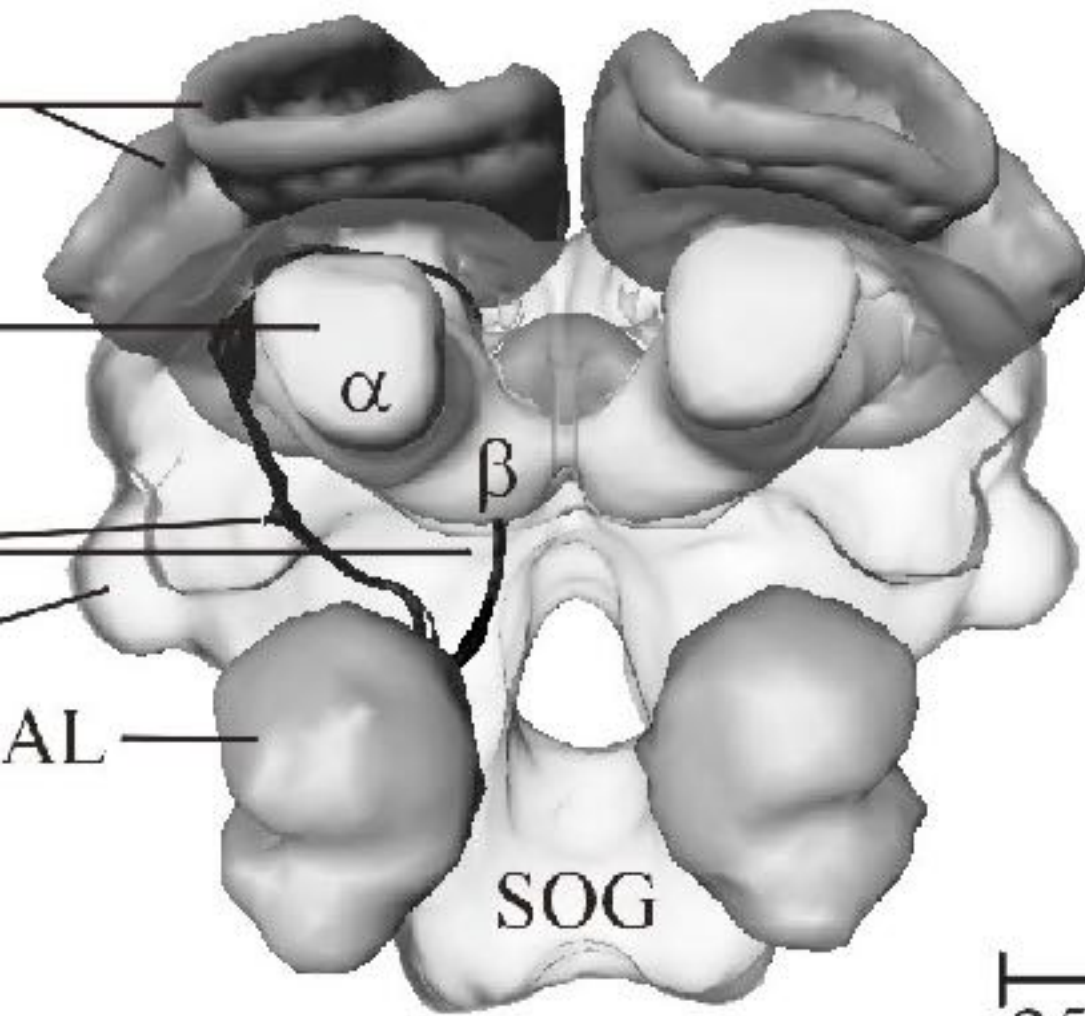
# Primacy effect in the TD model



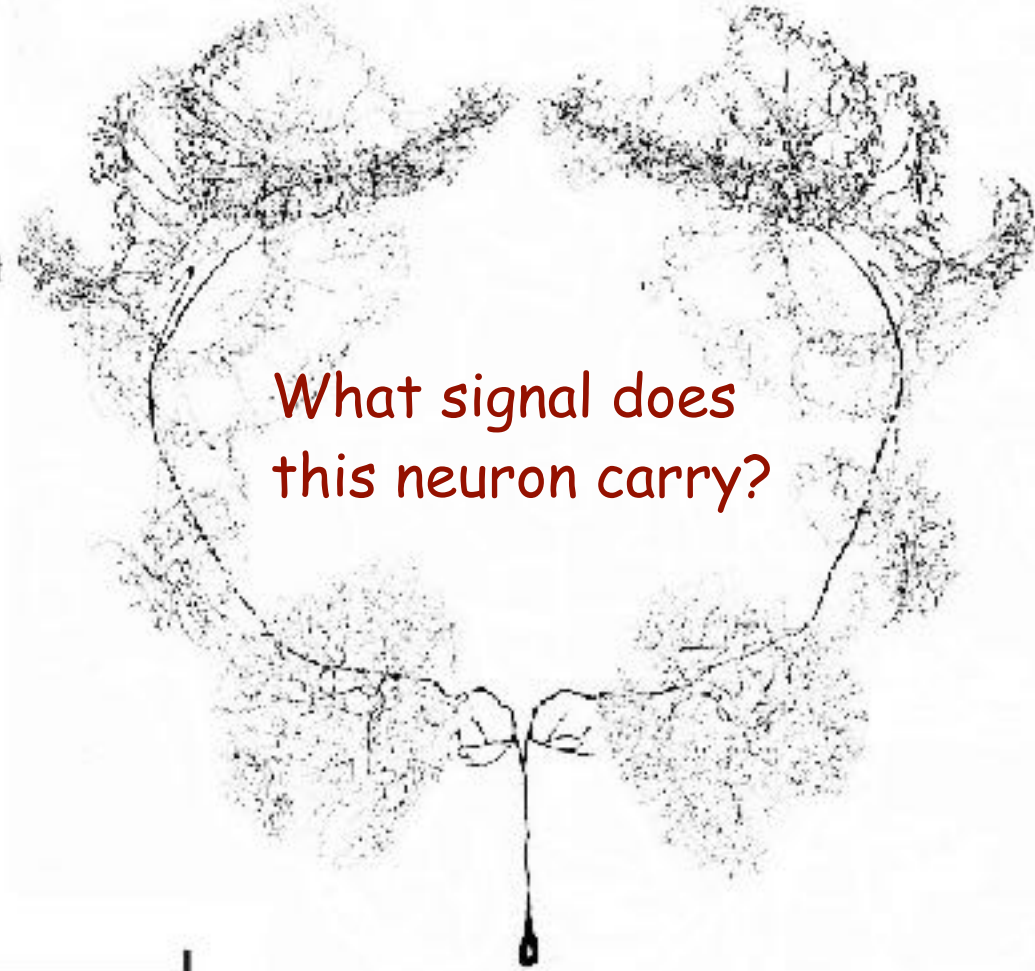
Facilitation of a remote association by an intervening stimulus



# Brain reward systems



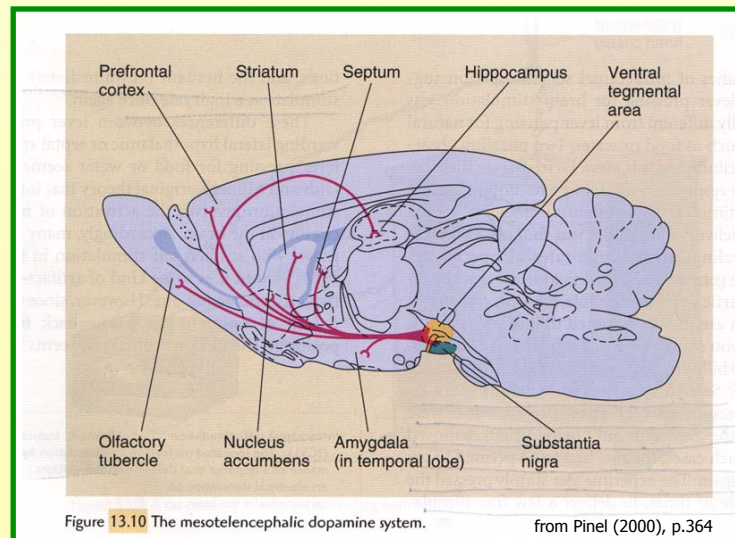
Honeybee Brain



VUM Neuron

# Dopamine

- Small-molecule Neurotransmitter
  - ❖ Diffuse projections from mid-brain throughout the brain



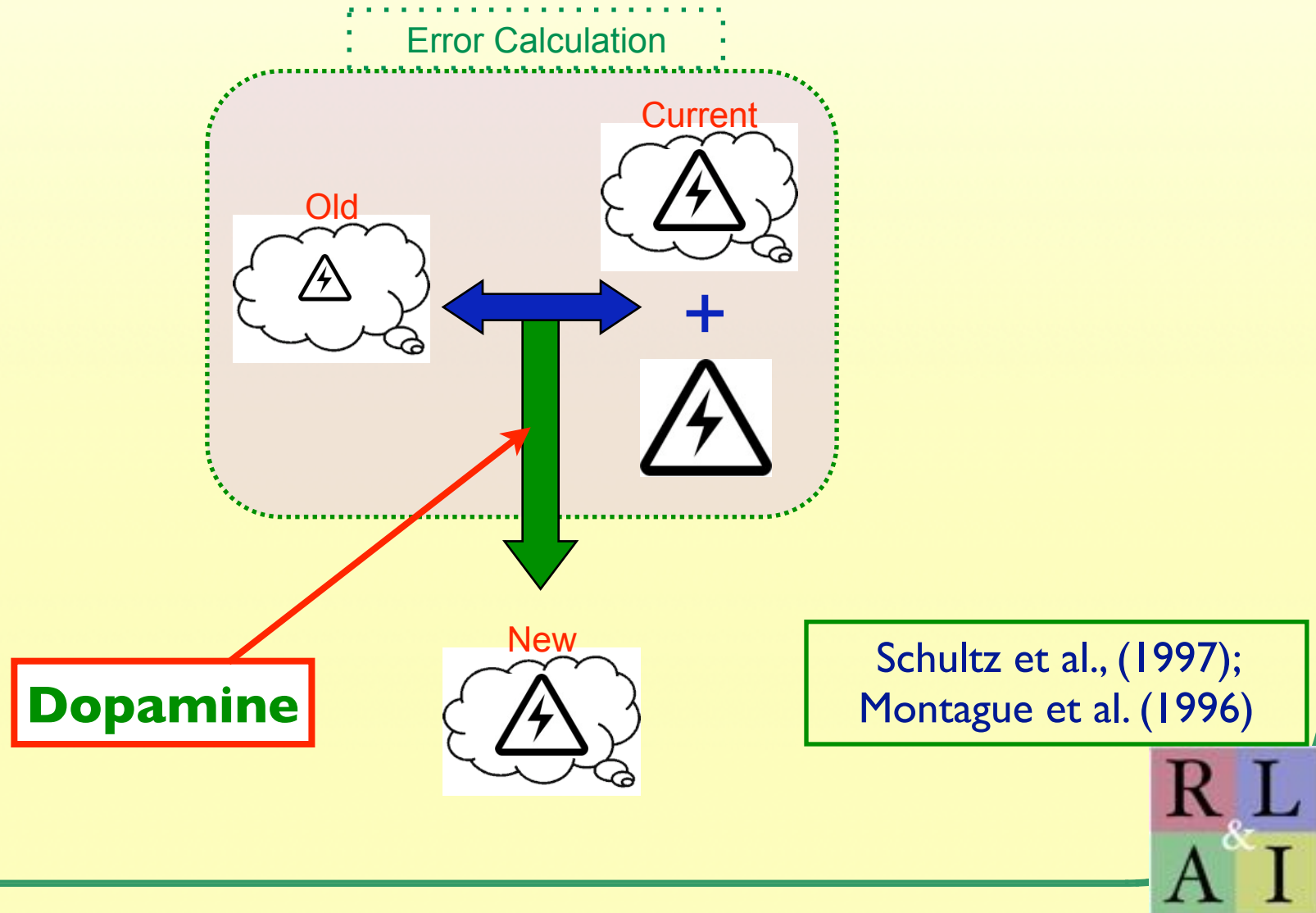
**Key Idea:** Phasic change in baseline dopamine responding = reward prediction error

# What does Dopamine Do?

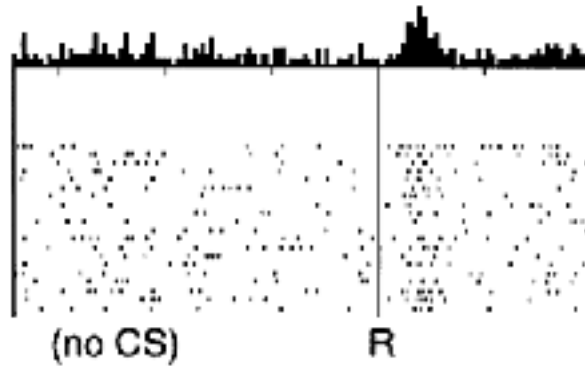
- Hedonic Impact
- Motivation
- Motor Activity
- Attention
- Novelty
- Learning



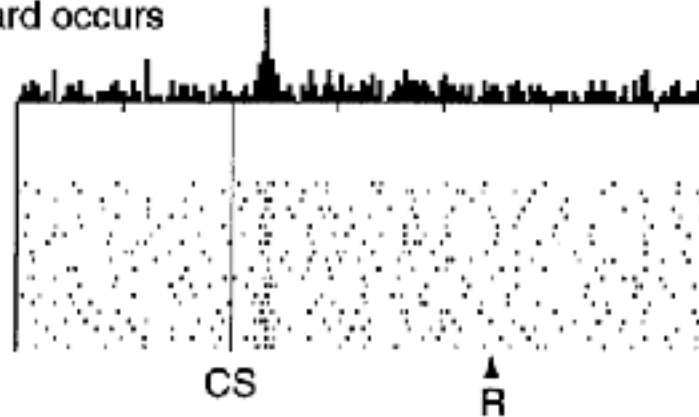
# TD Error = Dopamine



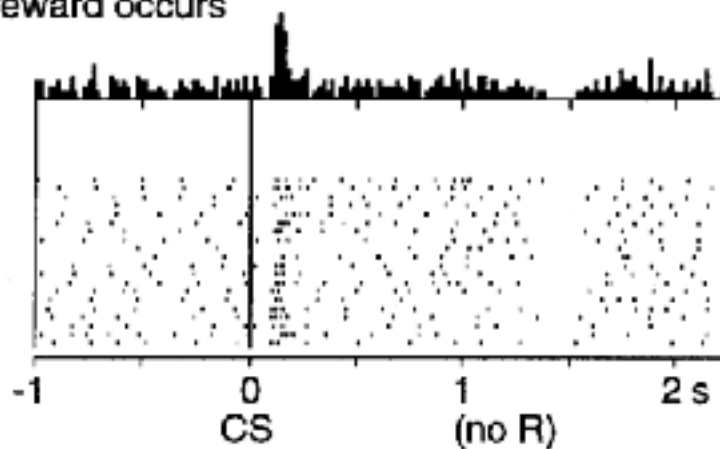
No prediction  
Reward occurs



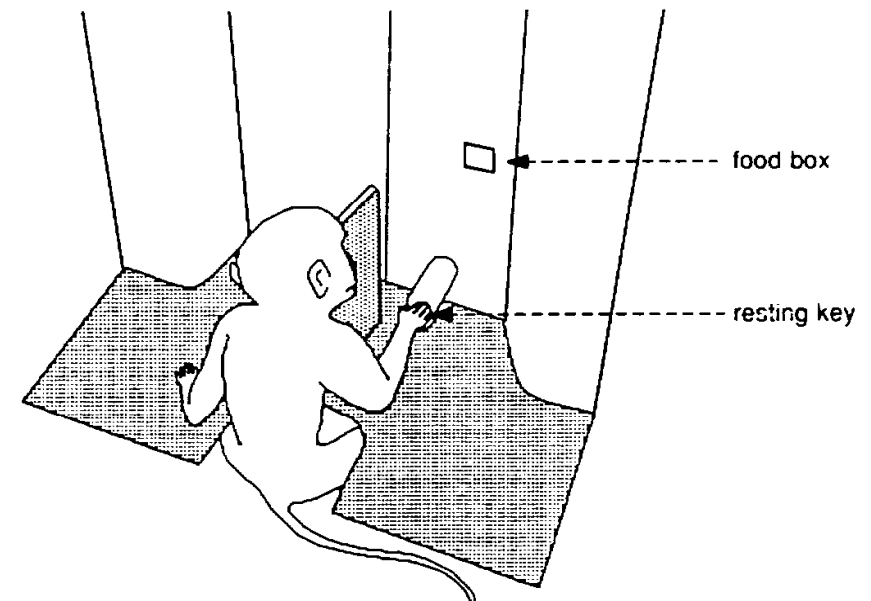
Reward predicted  
Reward occurs



Reward predicted  
No reward occurs



Dopamine neurons signal  
the error/change  
in prediction of reward

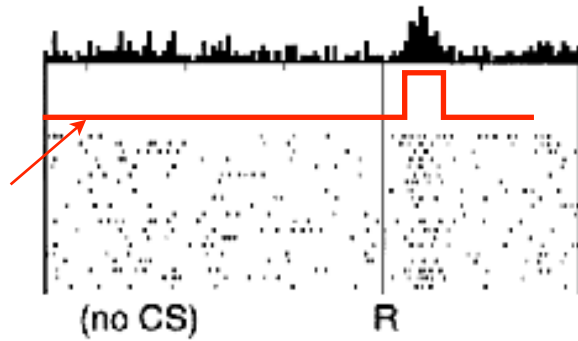


Wolfram Schultz, et al.



# Reward Unexpected

No prediction  
Reward occurs



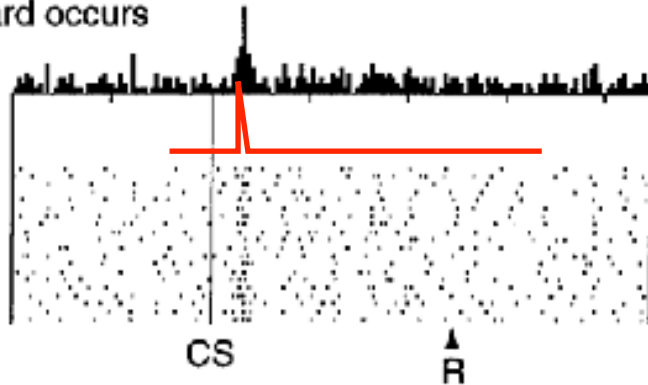
Reward

Value

TD error



Reward predicted  
Reward occurs

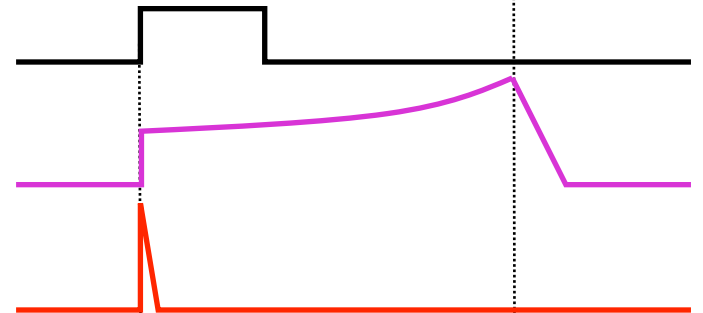


# Reward Expected

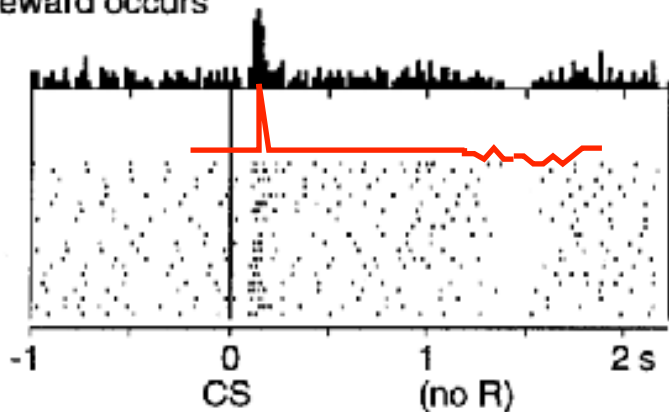
Cue

Value

TD error



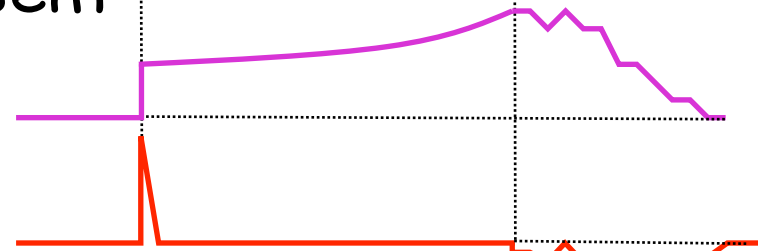
Reward predicted  
No reward occurs



# Reward Absent

Value

TD error



$$\delta_t = R_{t+1} + \gamma \hat{v}_{t+1} - \hat{v}_t$$



The theory that *Dopamine = TD error*  
is the *most important interaction ever*  
between AI and neuroscience

# Operant Learning

- The natural learning process directly analogous to reinforcement learning
- Control! What response to make when?

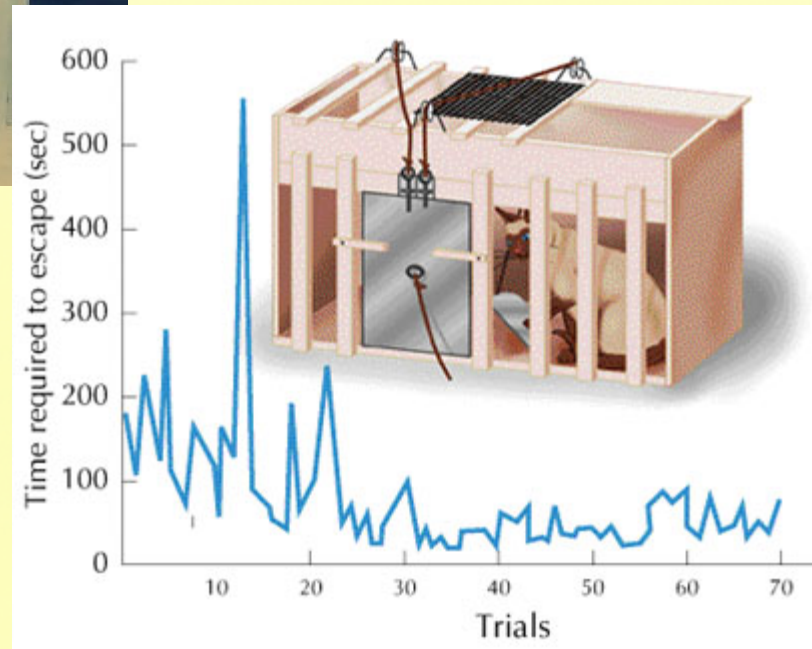
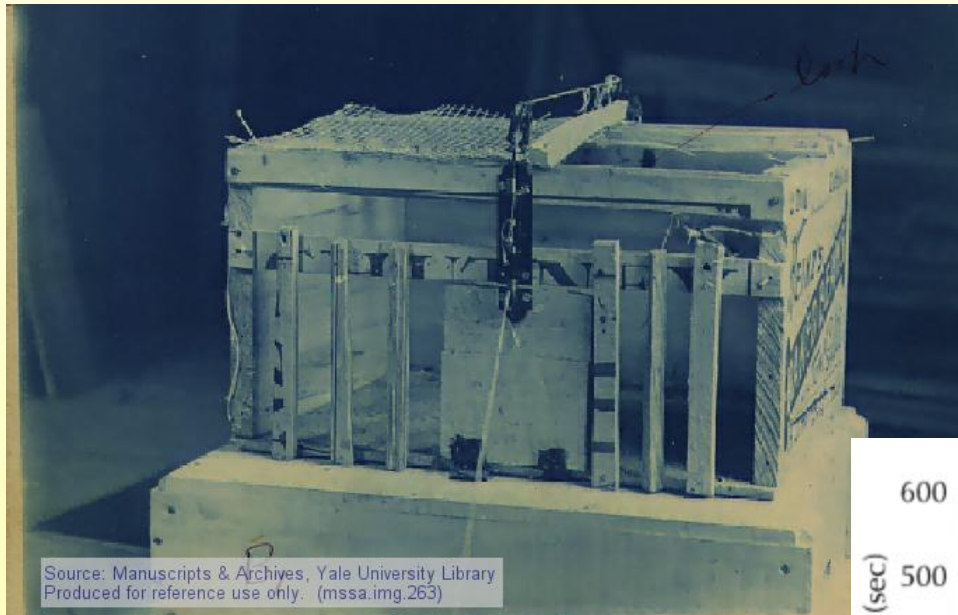


# Rat Basketball at Wofford College



QUINTESSENTIAL... A WOFFORD EDUCATION

# Thorndike's Puzzle Box (1910)

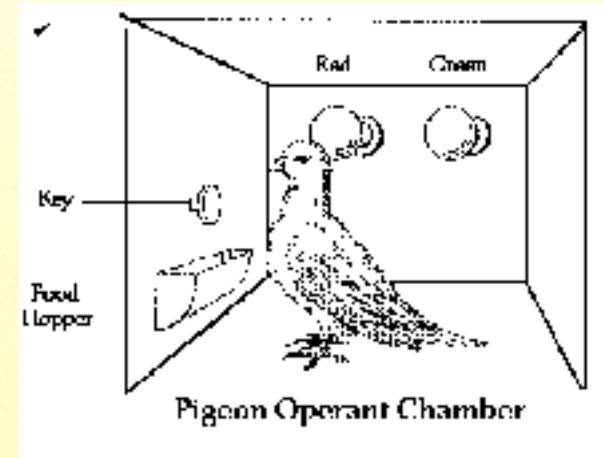


# Law of Effect



- “Of several **responses** made to the same situation, those which are accompanied by 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...” - Thorndike (1911), p. 244

# Operant Chambers





# Complex Cognition



# What have you learned about in this course (without buzzwords)?

- “Decision-making over time to achieve a long-term goal”
  - includes learning and planning
  - makes plain why value functions are so important
  - makes plain why so many fields care about these algorithms
    - AI
    - Control theory
    - Psychology and Neuroscience
    - Operations Research
    - Economics
  - all involve decision, goals, and time...
    - the essence of...

