

# Reinforcement Learning

# Reinforcement Learning

# Reinforcement Learning:

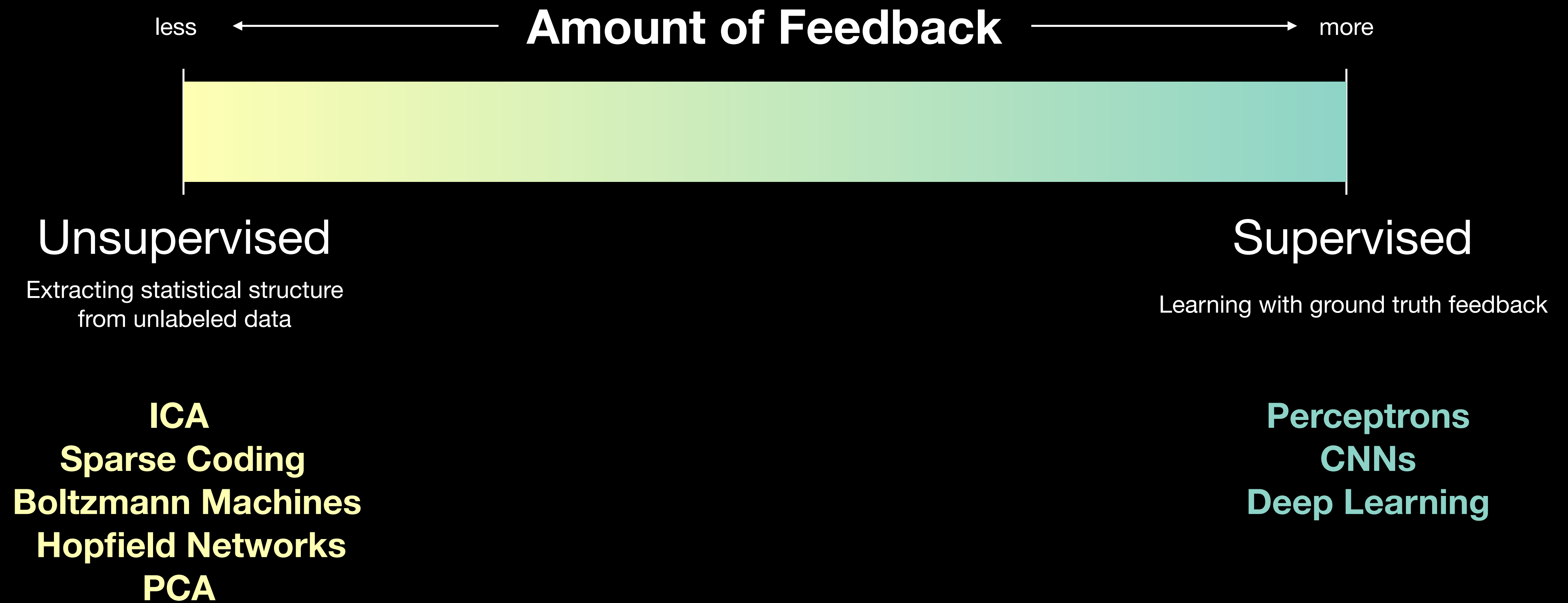
Learning behaviors from reward signals

# Types of Learning

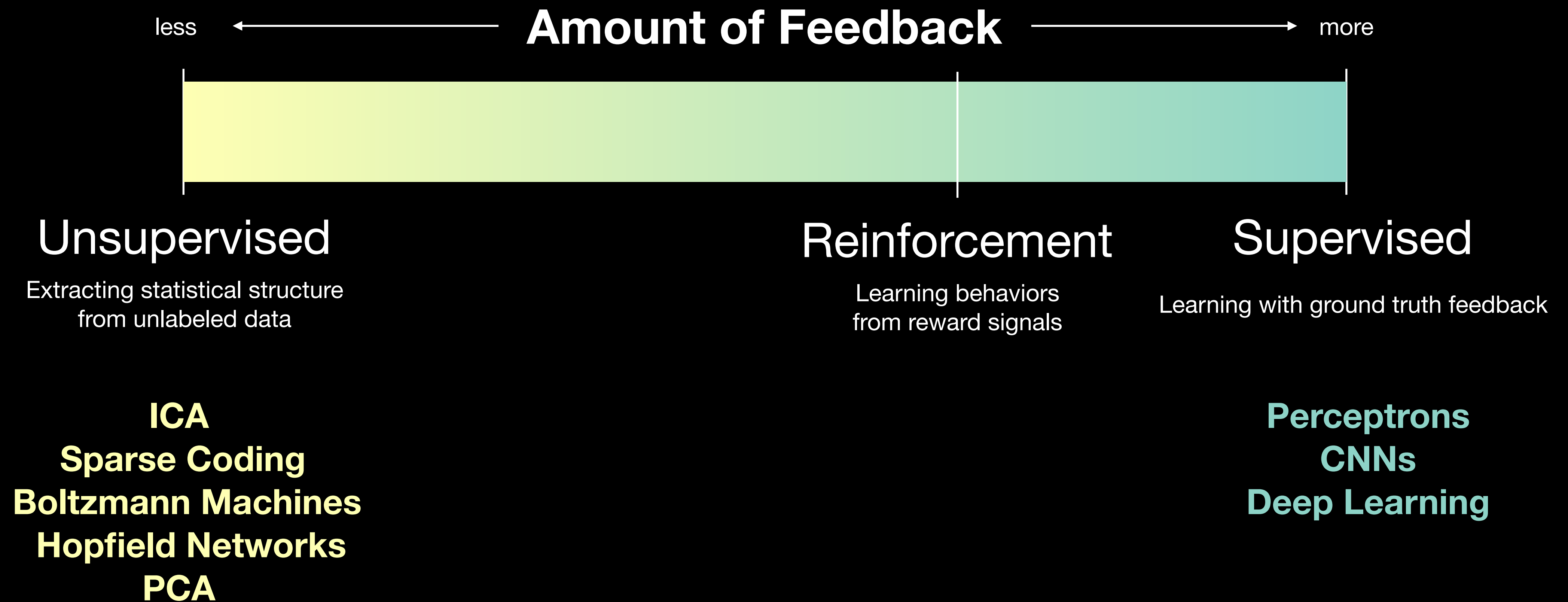




# Types of Learning

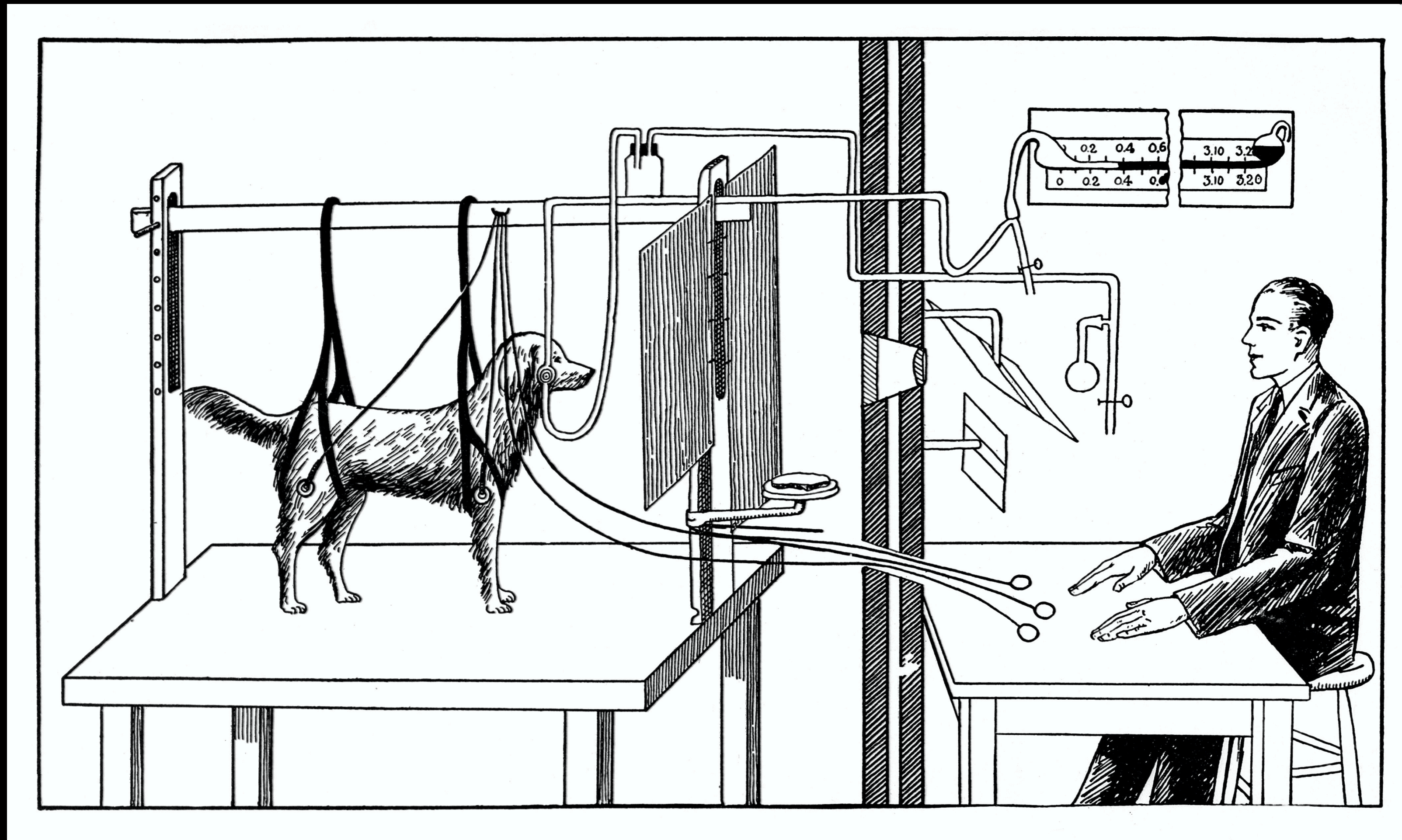


# Types of Learning



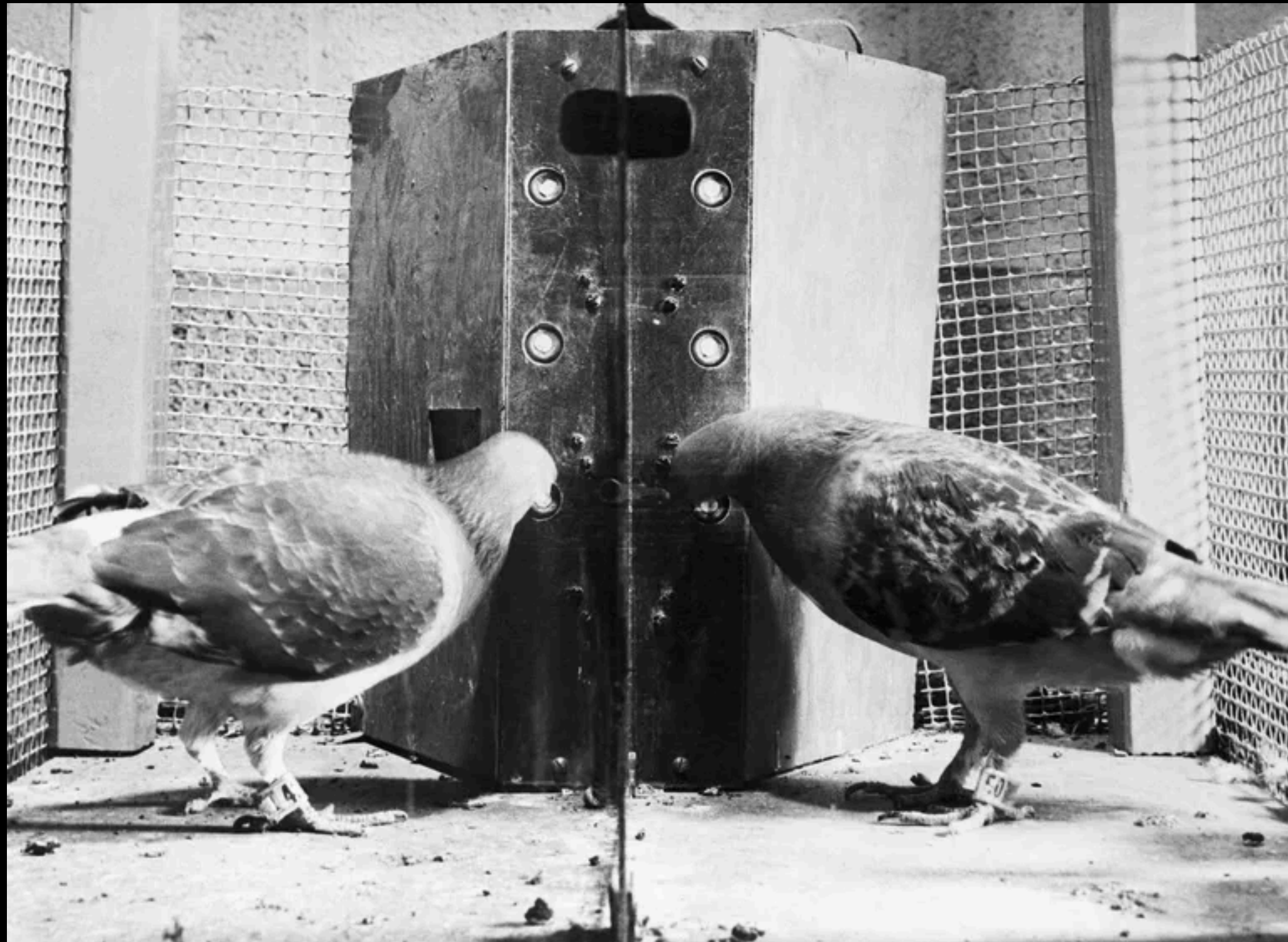


# Conditioning and Reinforcement





# Conditioning and Reinforcement





# Classical Conditioning Paradigms

Training	$A \rightarrow +$				
Test	$A \rightarrow +$				

Pavlovian

Extinction

Partial

Blocking

Overshadow

# Classical Conditioning Paradigms

Training	$A \rightarrow +$	$A \rightarrow +$ $A \rightarrow \cdot$			
Test	$A \rightarrow +$	$A \rightarrow \cdot$			

Pavlovian

Extinction

Partial

Blocking

Overshadow

# Classical Conditioning Paradigms

Training	$A \rightarrow +$	$A \rightarrow +$ $A \rightarrow \cdot$	$A \rightarrow +$ $AB \rightarrow +$		
Test	$A \rightarrow +$	$A \rightarrow \cdot$	$A \rightarrow \alpha +$		

Pavlovian

Extinction

Partial

Blocking

Overshadow

# Classical Conditioning Paradigms

Training	$A \rightarrow +$	$A \rightarrow +$ $A \rightarrow \cdot$	$A \rightarrow +$ $AB \rightarrow +$	$A \rightarrow +$ $AB \rightarrow +$
Test	$A \rightarrow +$	$A \rightarrow \cdot$	$A \rightarrow \alpha +$	$A \rightarrow +$ $B \rightarrow \cdot$

Pavlovian

Extinction

Partial

Blocking

Overshadow



# Classical Conditioning Paradigms

Training	$A \rightarrow +$	$A \rightarrow +$ $A \rightarrow \cdot$	$A \rightarrow +$ $AB \rightarrow +$	$A \rightarrow +$ $AB \rightarrow +$	$AB \rightarrow +$
Test	$A \rightarrow +$	$A \rightarrow \cdot$	$A \rightarrow \alpha +$	$A \rightarrow +$ $B \rightarrow \cdot$	$A \rightarrow \alpha +$ $B \rightarrow \beta +$
Pavlovian		Extinction		Partial	
Blocking		Overshadow			

# Rescorla-Wagner Model

$S$  : Stimulus Variable

$r$  : Actual Reward

$v$  : Expected Reward

$w$  : Weights

# Rescorla-Wagner Model

$S$  : Stimulus Variable

$r$  : Actual Reward

$v$  : Expected Reward

$w$  : Weights

Model

$$v = w^T S$$

# Rescorla-Wagner Model

$S$  : Stimulus Variable

$r$  : Actual Reward

$v$  : Expected Reward

$w$  : Weights

Model

$$v = w^T s$$

Objective

$$\min \sum \frac{1}{2} (r - v)^2$$

# Rescorla-Wagner Model

$S$  : Stimulus Variable

$r$  : Actual Reward

$v$  : Expected Reward

$w$  : Weights

Model

$$v = w^T s$$

Objective

$$\min \sum \frac{1}{2} (r - v)^2$$

Learning Rule

$$\begin{aligned} \delta &= r - v \\ w &\rightarrow w + \epsilon \delta s \end{aligned}$$

# Rescorla-Wagner Model

$S$  : Stimulus Variable

$r$  : Actual Reward

$v$  : Expected Reward

$w$  : Weights

Model

$$v = w^T s$$

Objective

$$\min \sum \frac{1}{2} (r - v)^2$$

Learning Rule

$$\delta = r - v$$
$$w \rightarrow w + \epsilon \delta s$$



Prediction Error

# Classical Conditioning Paradigms

Training	$A \rightarrow +$	$A \rightarrow +$ $A \rightarrow \cdot$	$A \rightarrow +$ $AB \rightarrow +$	$A \rightarrow +$ $AB \rightarrow +$	$AB \rightarrow +$
Test	$A \rightarrow +$	$A \rightarrow \cdot$	$A \rightarrow \alpha +$	$A \rightarrow +$ $B \rightarrow \cdot$	$A \rightarrow \alpha +$ $B \rightarrow \beta +$

Pavlovian

Extinction

Partial

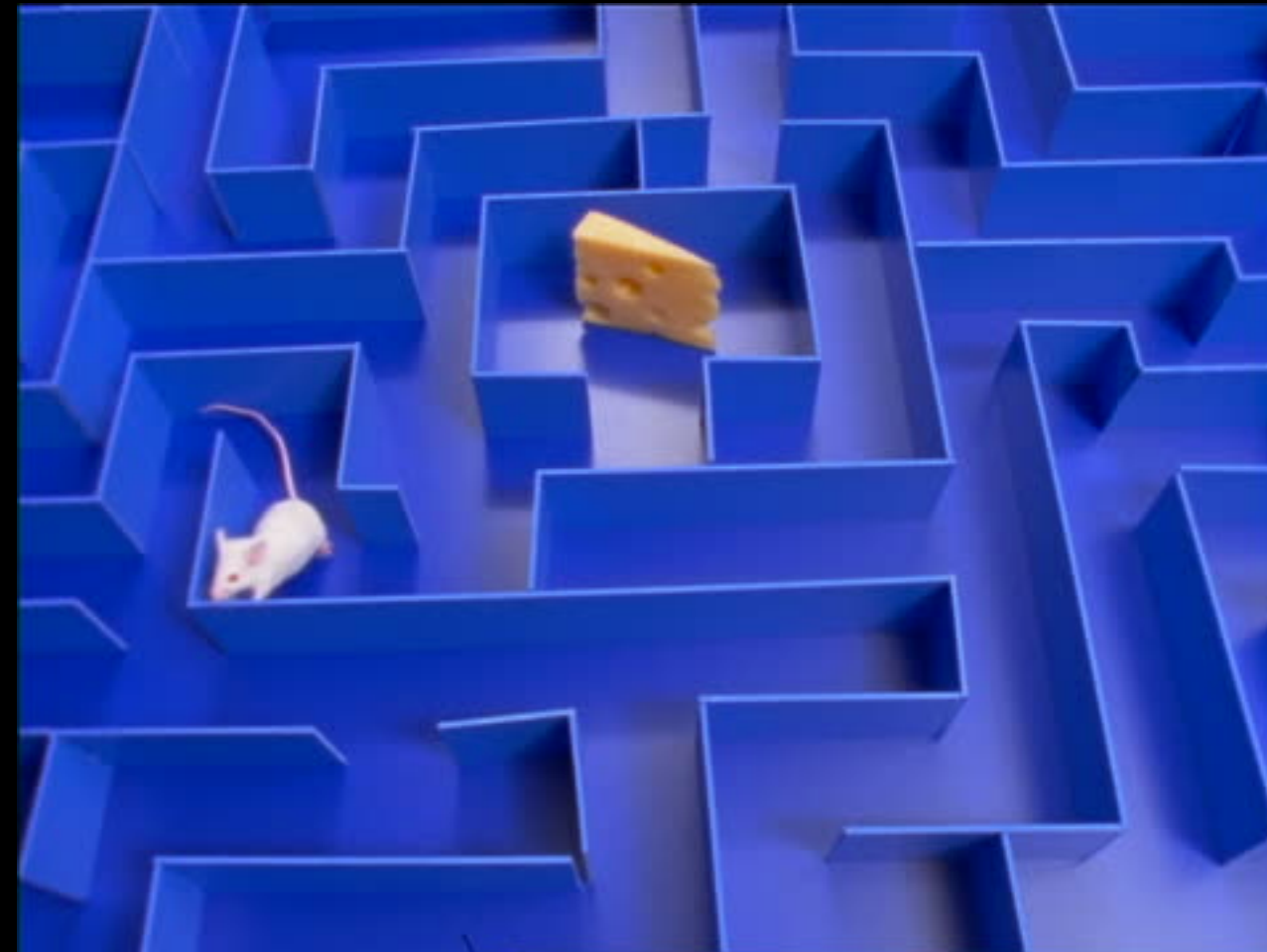
Blocking

Overshadow

Model	Objective	Learning Rule
$v = w^\top s$	$\min \sum \frac{1}{2}(r - v)^2$	$\delta = r - v$ $w \rightarrow w + \epsilon \delta s$

# Predicting Future Reward:

## Temporal Difference Learning



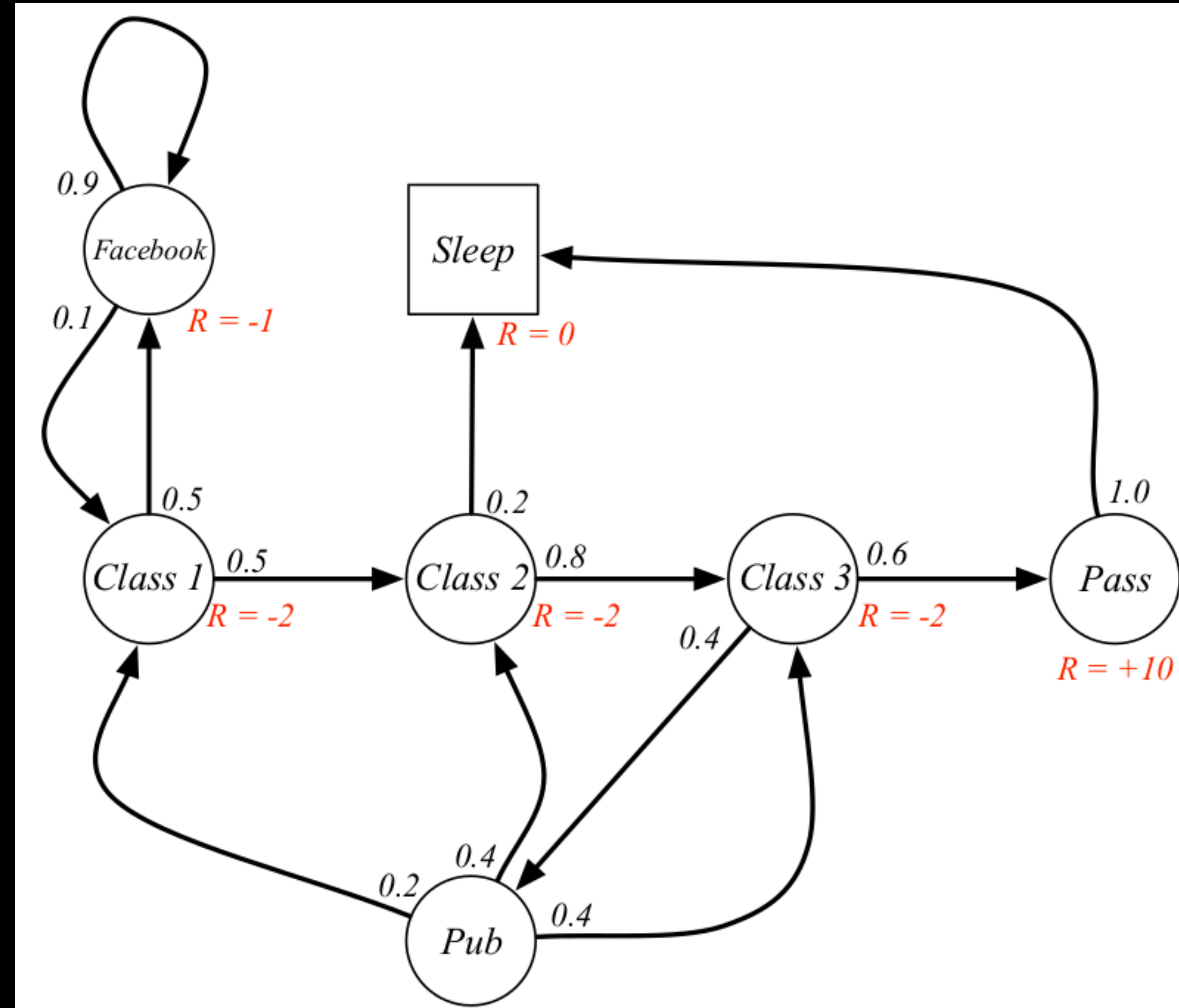


# Bringing an Agent in the Loop





# Bringing an Agent in the Loop



# Predicting Future Reward:

## Temporal Difference Learning

$\mathcal{S}$  : Stimulus Variable

$r$  : Actual Reward

$V$  : Expected Reward

$\mathcal{W}$  : Weights

# Predicting Future Reward:

## Temporal Difference Learning

$s_t$  : Stimulus Variable

$r_t$  : Actual Reward

$V_t$  : Expected Reward

$w_t$  : Weights

# Predicting Future Reward:

## Temporal Difference Learning

$s_t$  : Stimulus Variable

$r_t$  : Actual Reward

$V_t$  : Expected Reward

$w_t$  : Weights

Total Expected Reward

$$V_t = \mathbb{E} \left[ \sum_{i=0}^{\infty} \gamma^i r_{t+i} \right]$$

# Predicting Future Reward:

## Temporal Difference Learning

$s_t$  : Stimulus Variable

$r_t$  : Actual Reward

$V_t$  : Expected Reward

$w_t$  : Weights

Total Expected Reward

$$V_t = \mathbb{E} \left[ \sum_{i=0}^{\infty} \gamma^i r_{t+i} \right]$$

Bellman Equation

$$V_t = r_t + \gamma \mathbb{E}[V_{t+1}]$$

# Predicting Future Reward:

## Temporal Difference Learning

$s_t$  : Stimulus Variable

$r_t$  : Actual Reward

$V_t$  : Expected Reward

$w_t$  : Weights

Model

$$\hat{V}_t = w_t^\top s_t$$

Total Expected Reward

$$V_t = \mathbb{E} \left[ \sum_{i=0}^{\infty} \gamma^i r_{t+i} \right]$$

Bellman Equation

$$V_t = r_t + \gamma \mathbb{E}[V_{t+1}]$$

# Predicting Future Reward:

## Temporal Difference Learning

$s_t$  : Stimulus Variable

$r_t$  : Actual Reward

$V_t$  : Expected Reward

$w_t$  : Weights

Model

$$\hat{V}_t = w_t^\top s_t$$

Objective

$$\min \sum \frac{1}{2} (V_t - \hat{V}_t)^2$$

Total Expected Reward

$$V_t = \mathbb{E} \left[ \sum_{i=0}^{\infty} \gamma^i r_{t+i} \right]$$

Bellman Equation

$$V_t = r_t + \gamma \mathbb{E}[V_{t+1}]$$



# Predicting Future Reward:

## Temporal Difference Learning

$s_t$  : Stimulus Variable

$r_t$  : Actual Reward

$V_t$  : Expected Reward

$w_t$  : Weights

Model

$$\hat{V}_t = w_t^\top s_t$$

Objective

$$\min \sum \frac{1}{2} (V_t - \hat{V}_t)^2$$

Total Expected Reward

$$V_t = \mathbb{E} \left[ \sum_{i=0}^{\infty} \gamma^i r_{t+i} \right]$$

Bellman Equation

$$V_t = r_t + \gamma \mathbb{E}[V_{t+1}]$$

Prediction Error

$$\delta_t = r_t + \gamma \mathbb{E}[V_{t+1}] - \hat{V}_t$$

# Predicting Future Reward:

## Temporal Difference Learning

$s_t$  : Stimulus Variable

$r_t$  : Actual Reward

$V_t$  : Expected Reward

$w_t$  : Weights

Model

$$\hat{V}_t = w_t^\top s_t$$

Objective

$$\min \sum \frac{1}{2} (V_t - \hat{V}_t)^2$$

Update Rule

$$w_{t+1} = w_t + \epsilon s_t \delta_t$$

Total Expected Reward

$$V_t = \mathbb{E} \left[ \sum_{i=0}^{\infty} \gamma^i r_{t+i} \right]$$

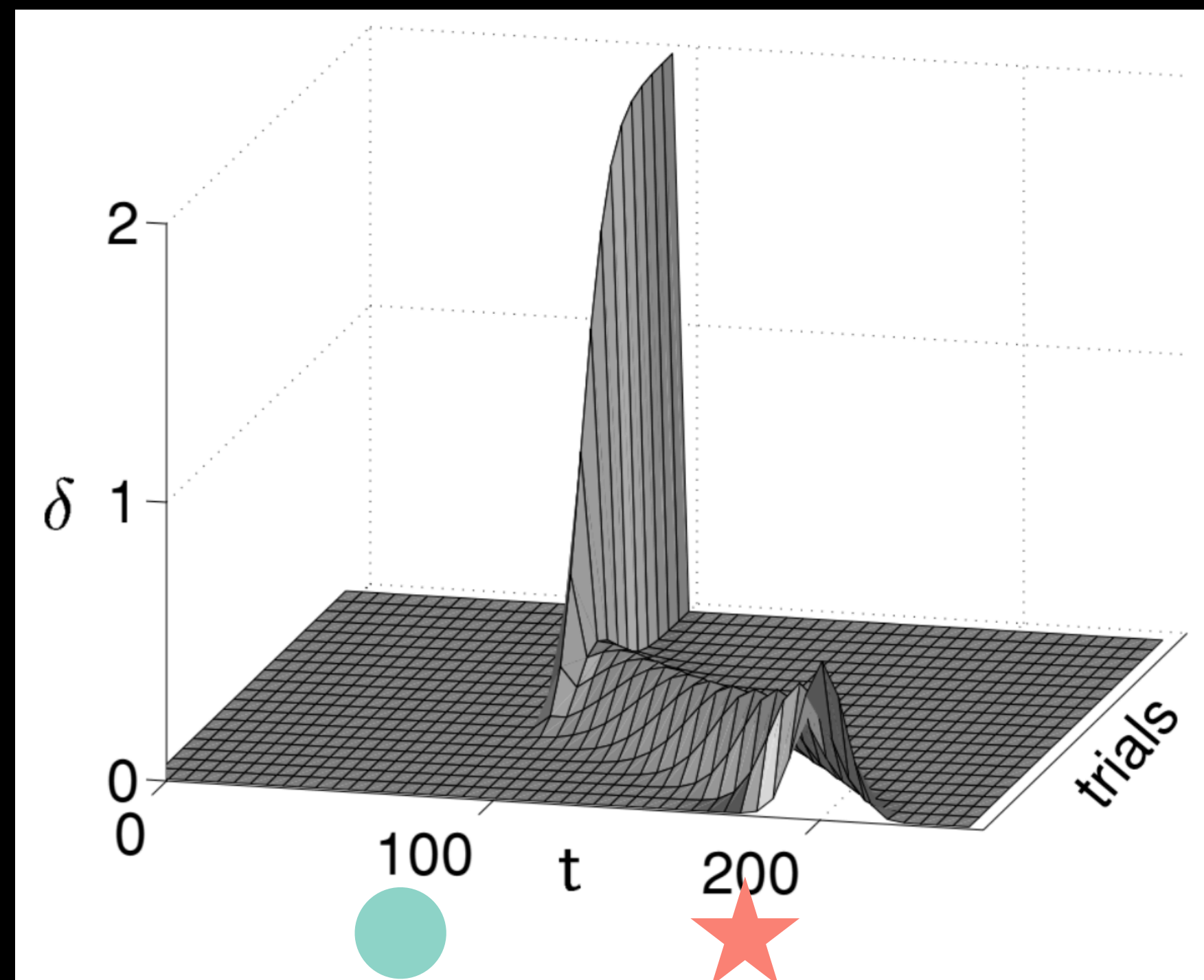
Bellman Equation

$$V_t = r_t + \gamma \mathbb{E}[V_{t+1}]$$

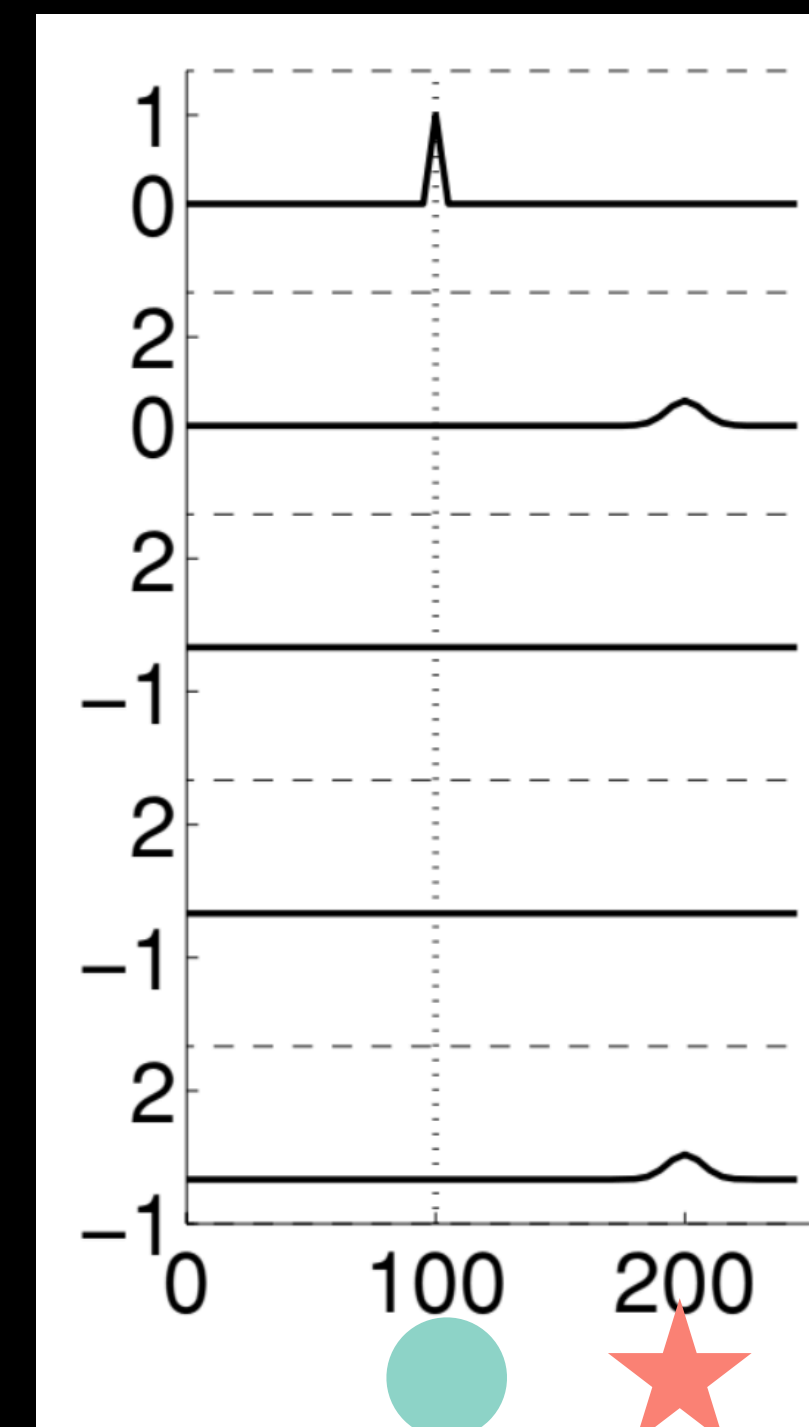
Prediction Error

$$\delta_t = r_t + \gamma \mathbb{E}[V_{t+1}] - \hat{V}_t$$

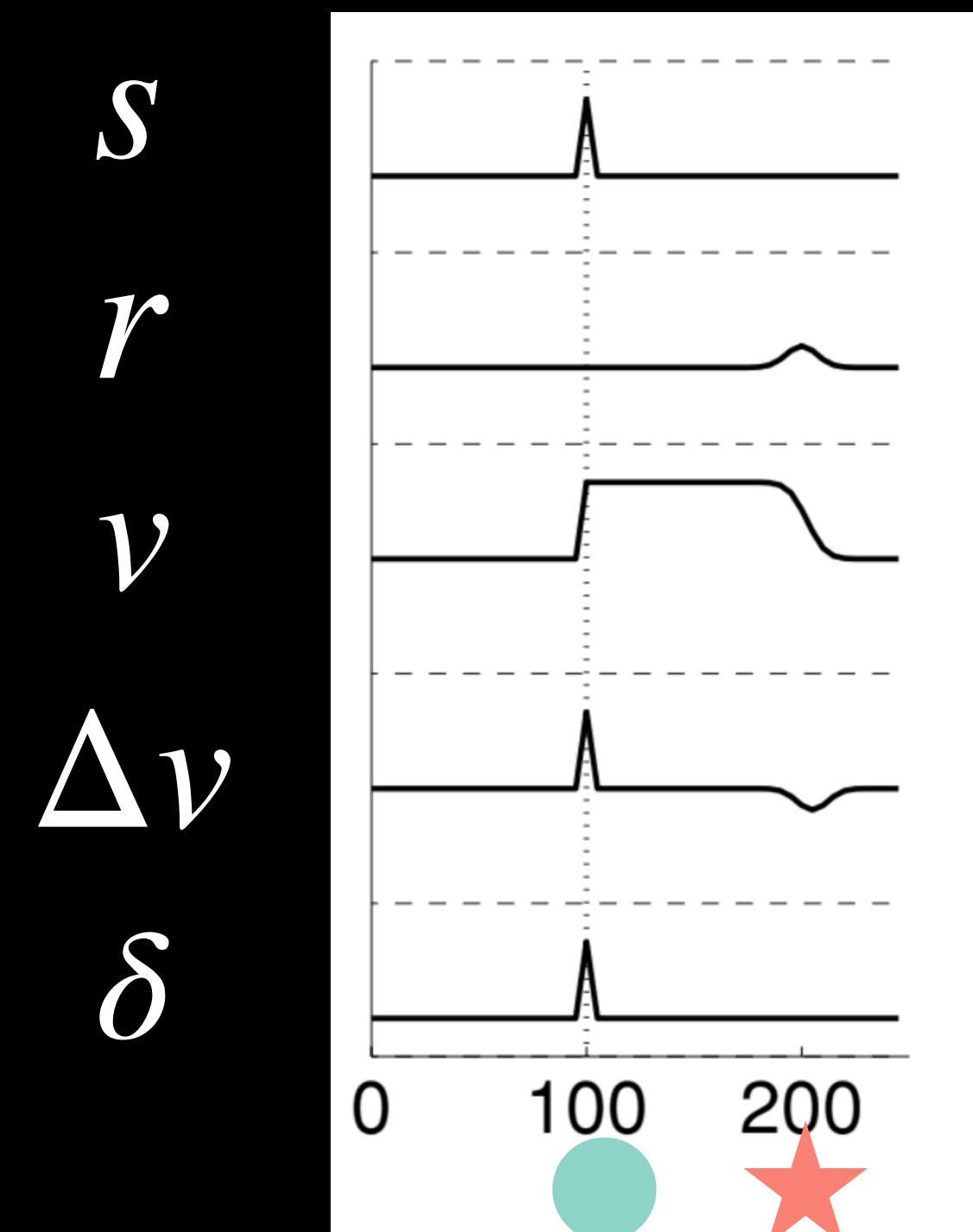
# Predicting Future Reward: Temporal Difference Learning



Prediction Error

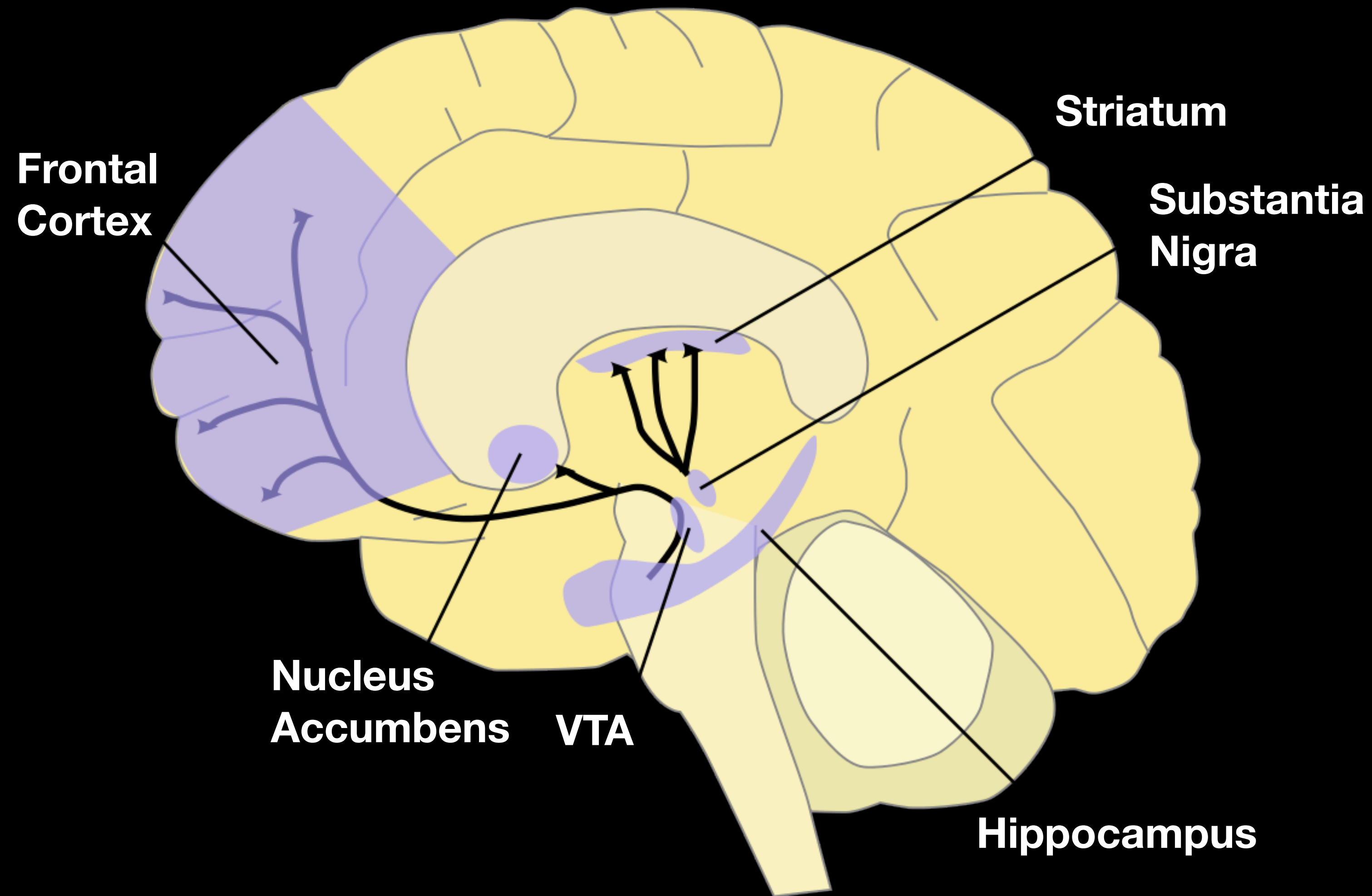


Before Training



After Training

# TD Learning in the Brain



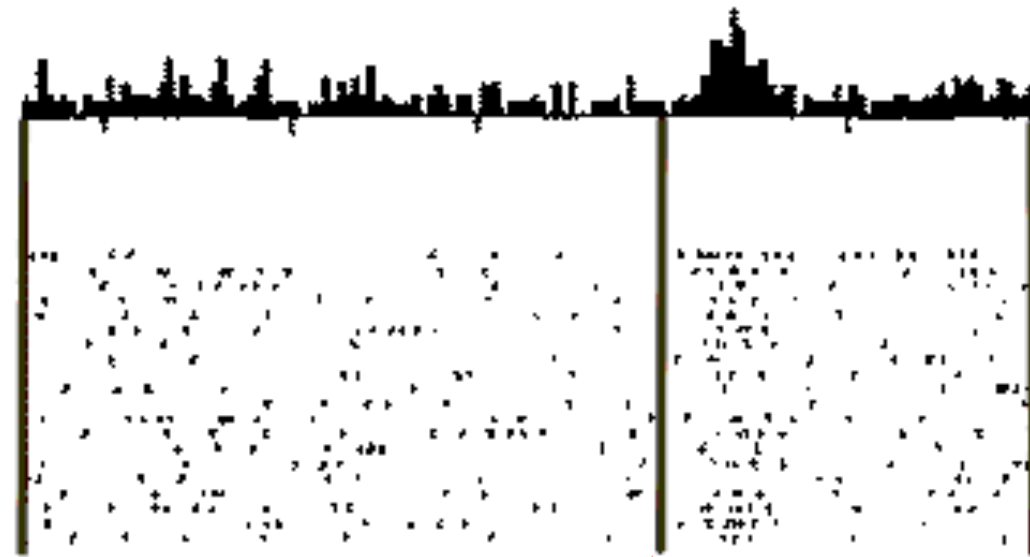


# TD Learning in the Brain

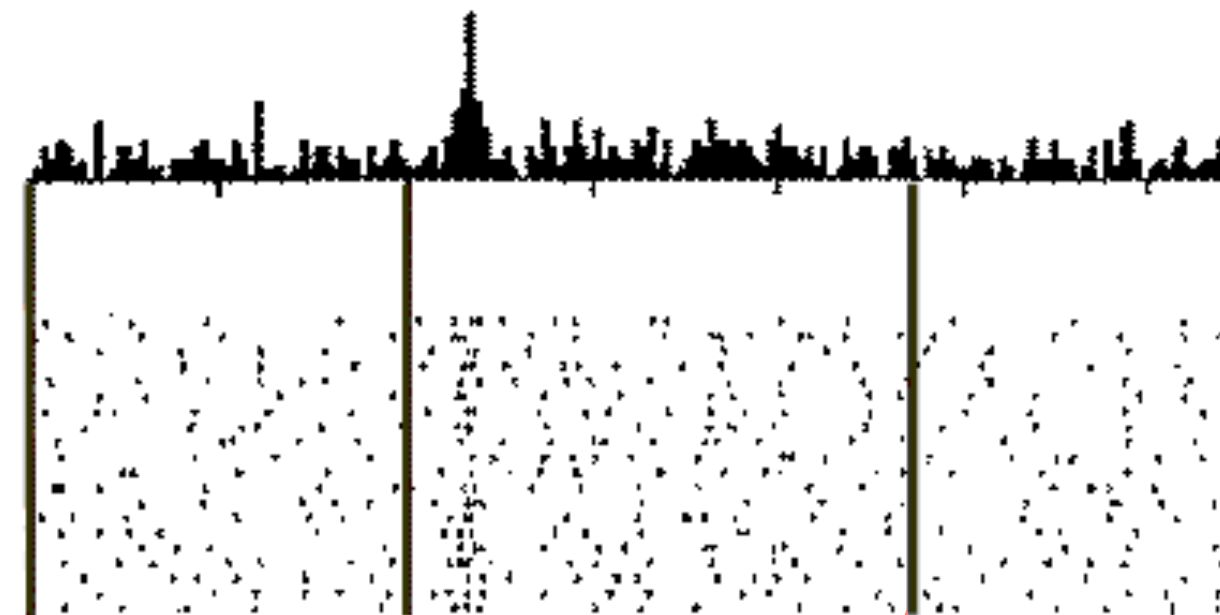


# TD Learning in the Brain

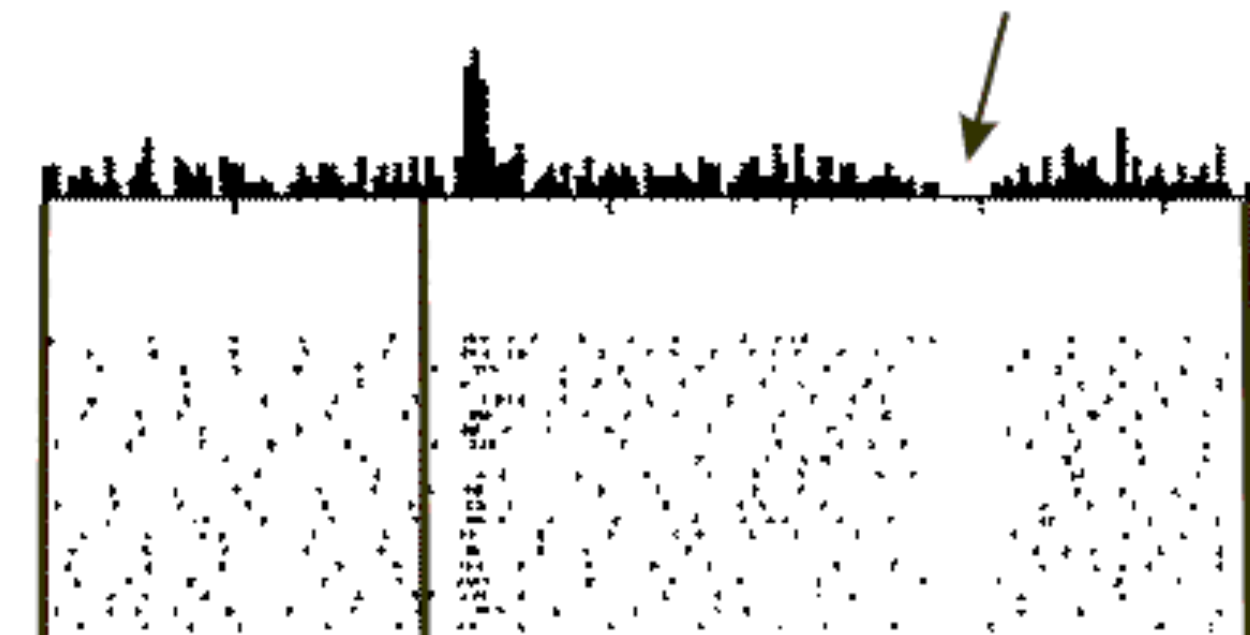
**Novelty Response**  
Reward, No Stimulus



**After Learning**  
Stimulus + Reward

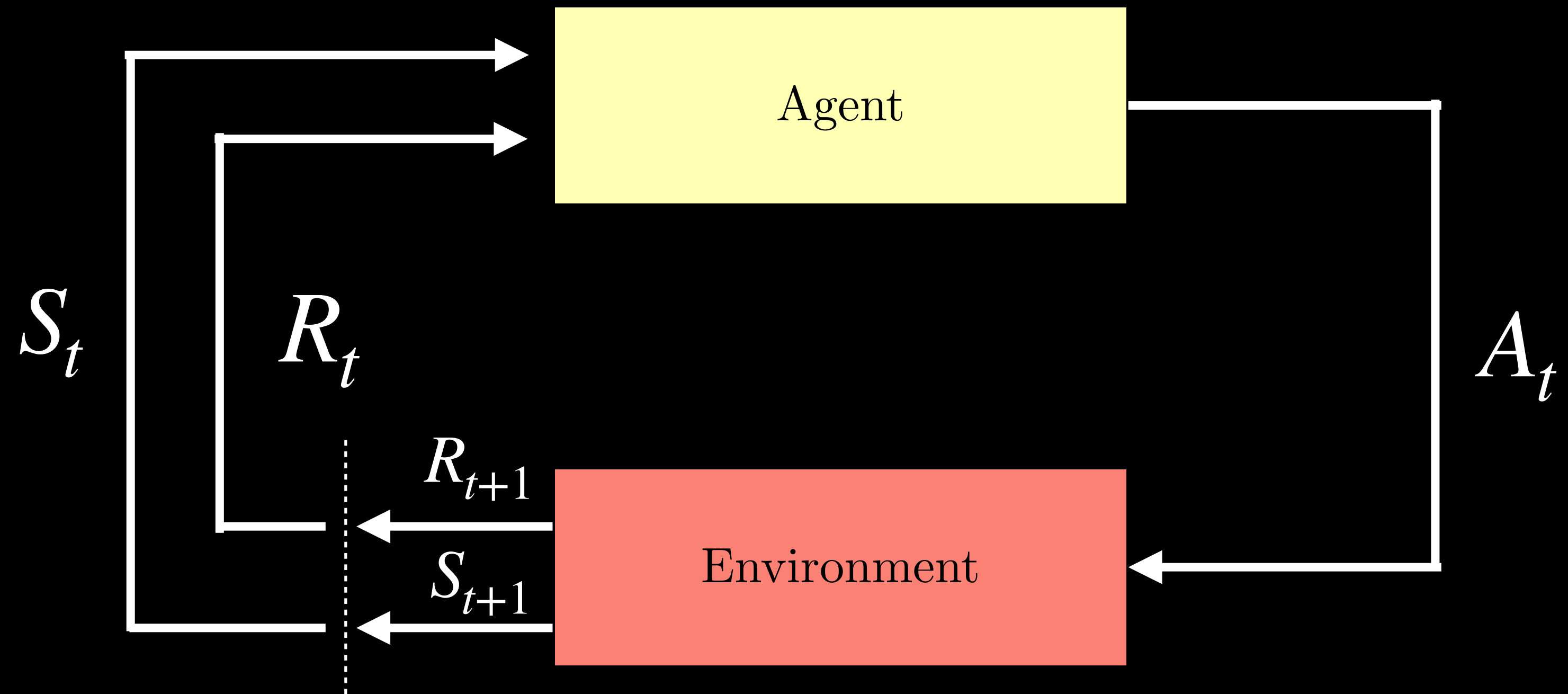


**After Learning**  
Stimulus, No Reward



# Bringing an Agent in the Loop

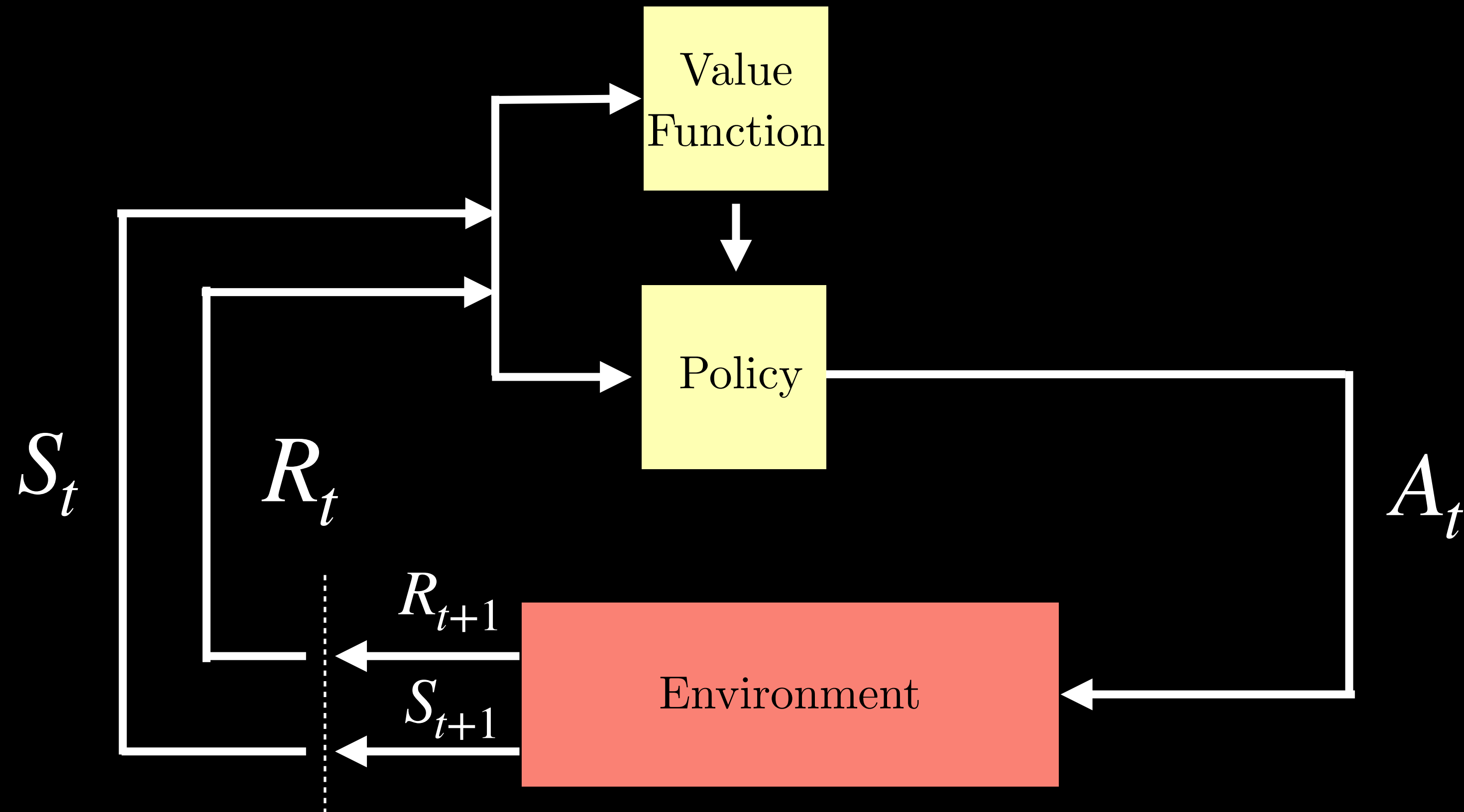
$S_t$  : State  
 $R_t$  : Reward  
 $A_t$  : Actions



# Bringing an Agent in the Loop

## The Actor-Critic Algorithm

$S_t$  : State  
 $R_t$  : Reward  
 $A_t$  : Actions

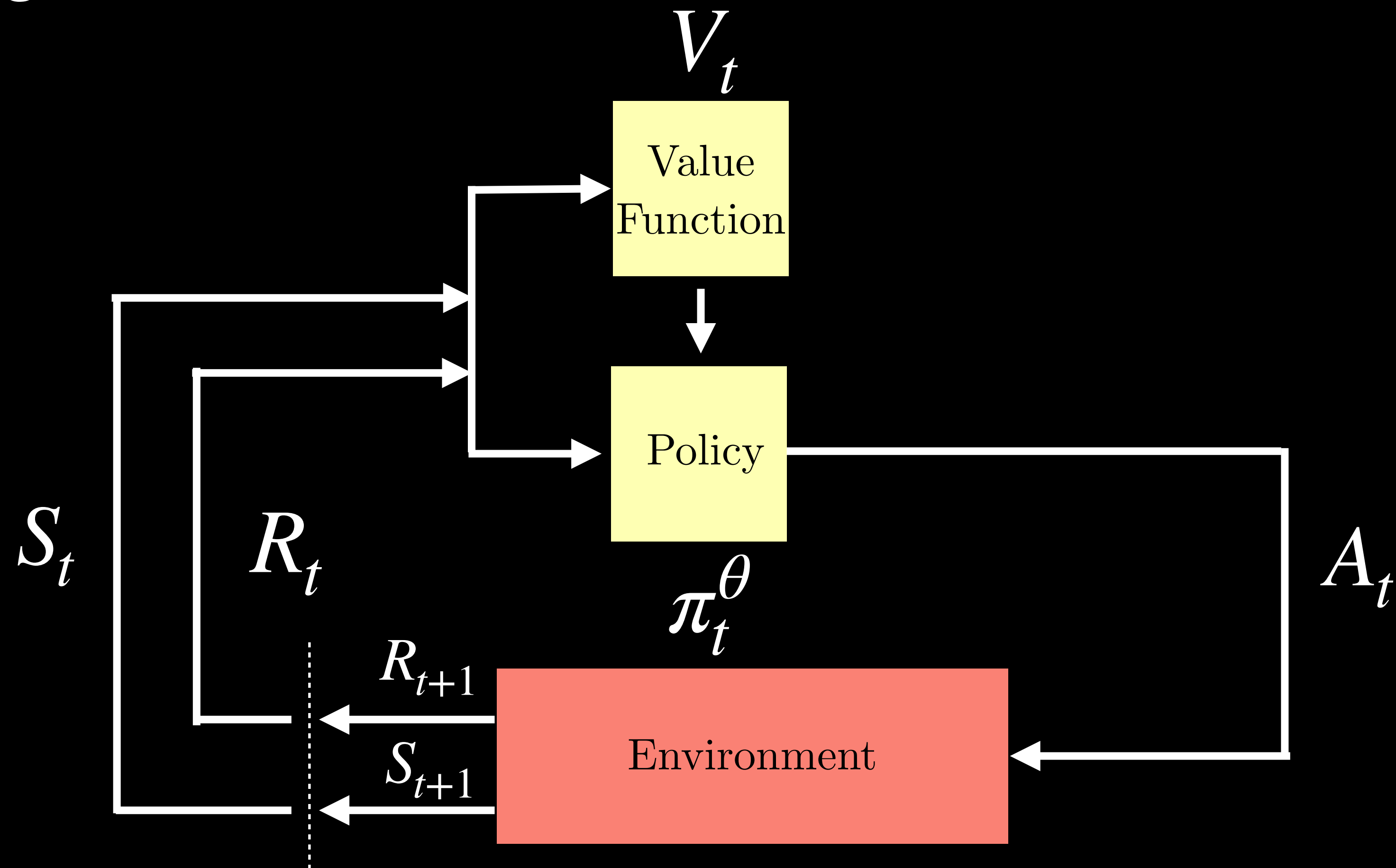




# Bringing an Agent in the Loop

## The Actor-Critic Algorithm

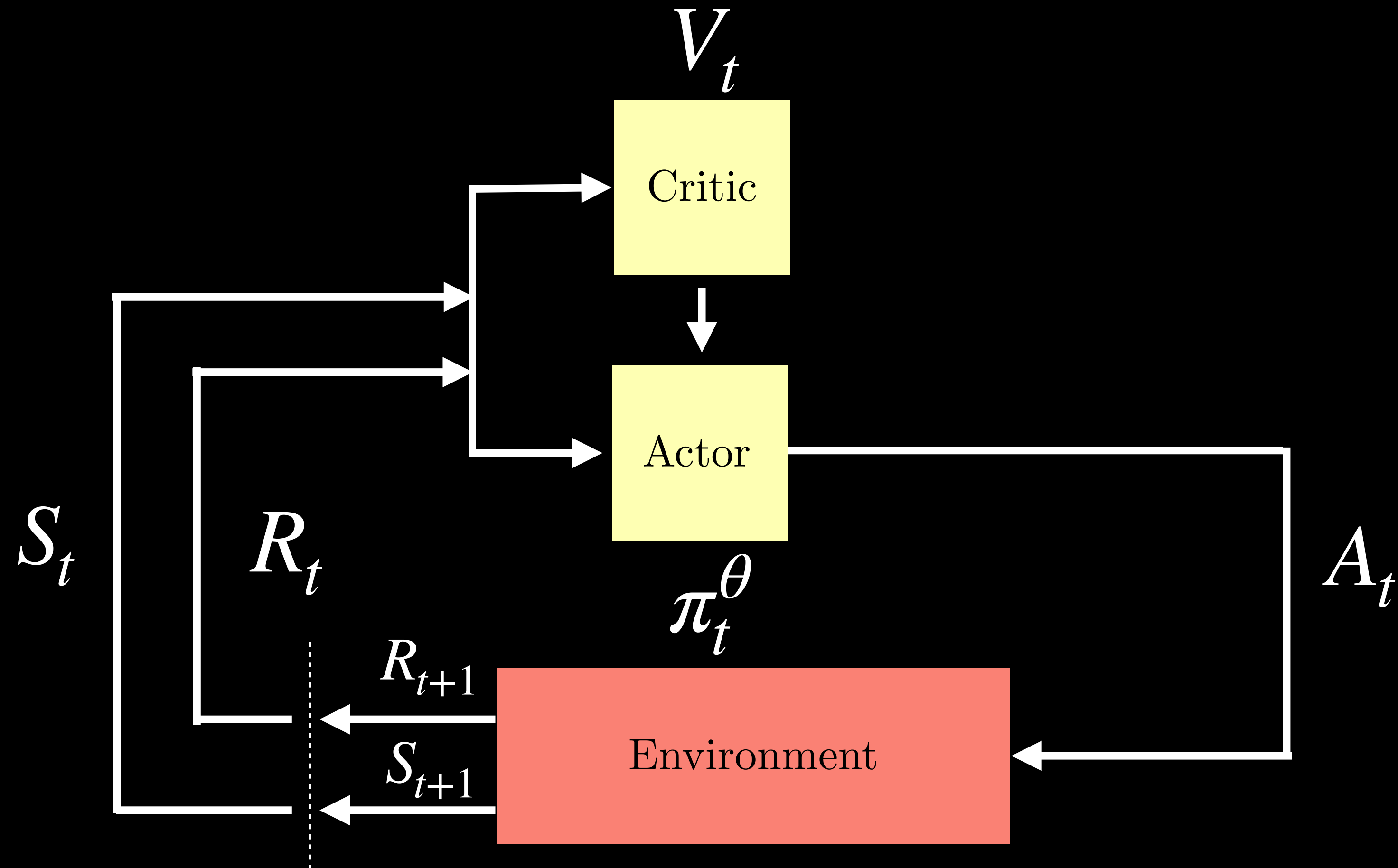
$S_t$  : State  
 $R_t$  : Reward  
 $A_t$  : Actions  
 $\pi_t^\theta$  : Policy  
 $V_t$  : Value Function



# Bringing an Agent in the Loop

## The Actor-Critic Algorithm

$S_t$  : State  
 $R_t$  : Reward  
 $A_t$  : Actions  
 $\pi_t^\theta$  : Policy  
 $V_t$  : Value Function



# Bringing an Agent in the Loop

## The Actor-Critic Algorithm

$S_t$  : State

$R_t$  : Reward

$A_t$  : Actions

$\pi_t^\theta$  : Policy

$V_t$  : Value Function

Policy

$$a_t \sim \pi_\theta$$

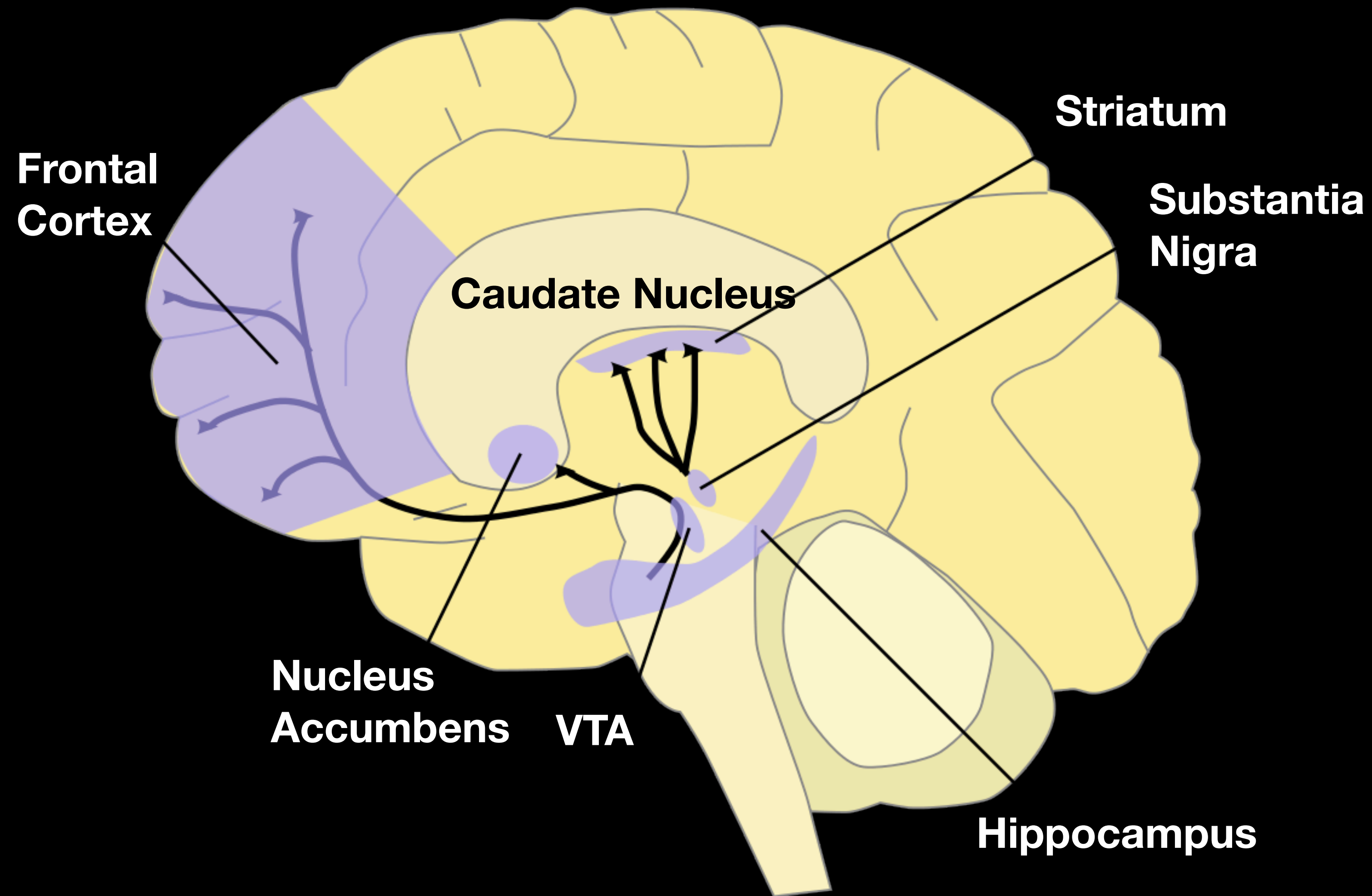
$$\pi_t^\theta = P(a_t = a \mid s_t = s)$$

$$P(a \mid s) = \frac{e^{\theta(s,a)}}{\sum_b e^{\theta(s,b)}}$$

Policy Update Rule

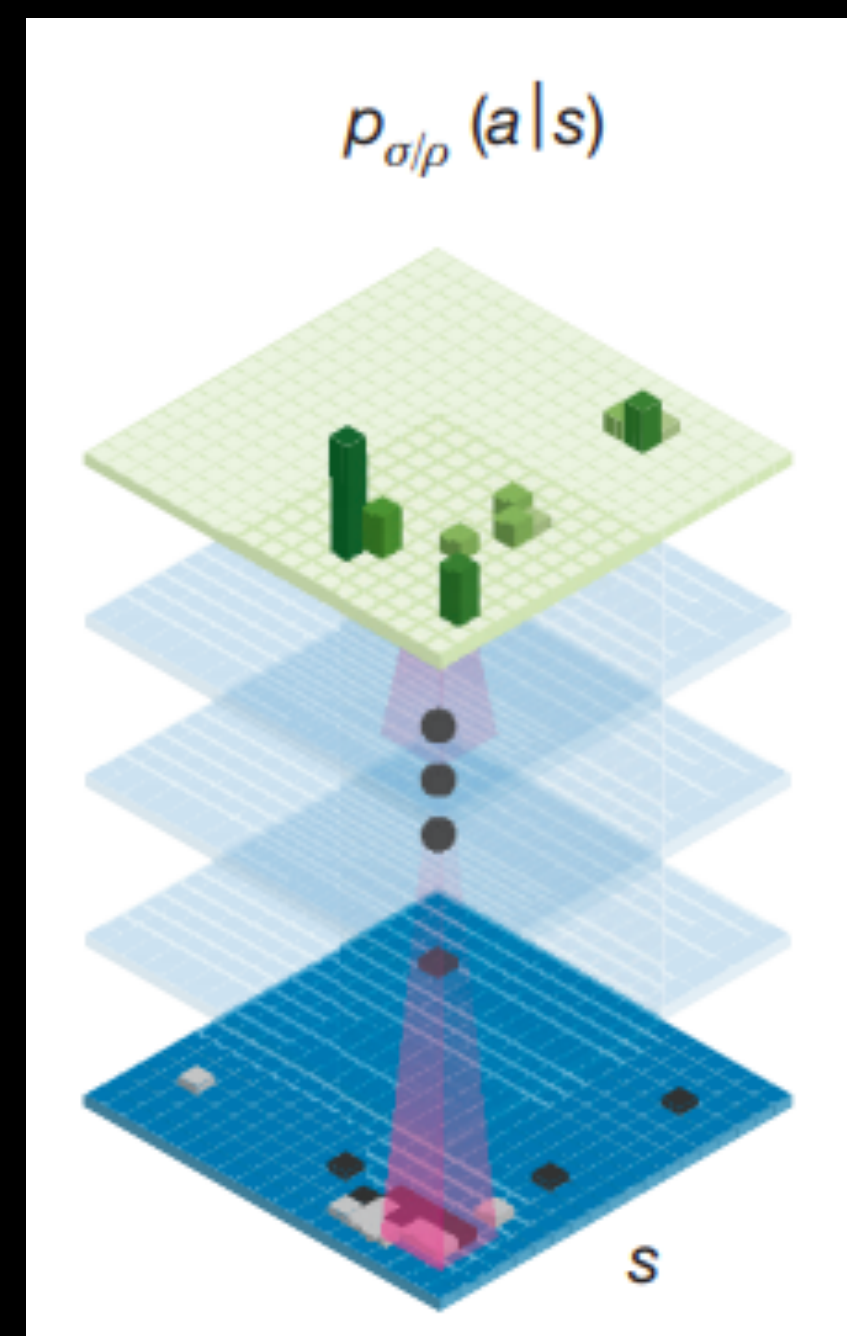
$$\theta(s_t, a_t) \leftarrow \theta(s_t, a_t) + \epsilon \delta_t$$

# Actor-Critic in the Brain

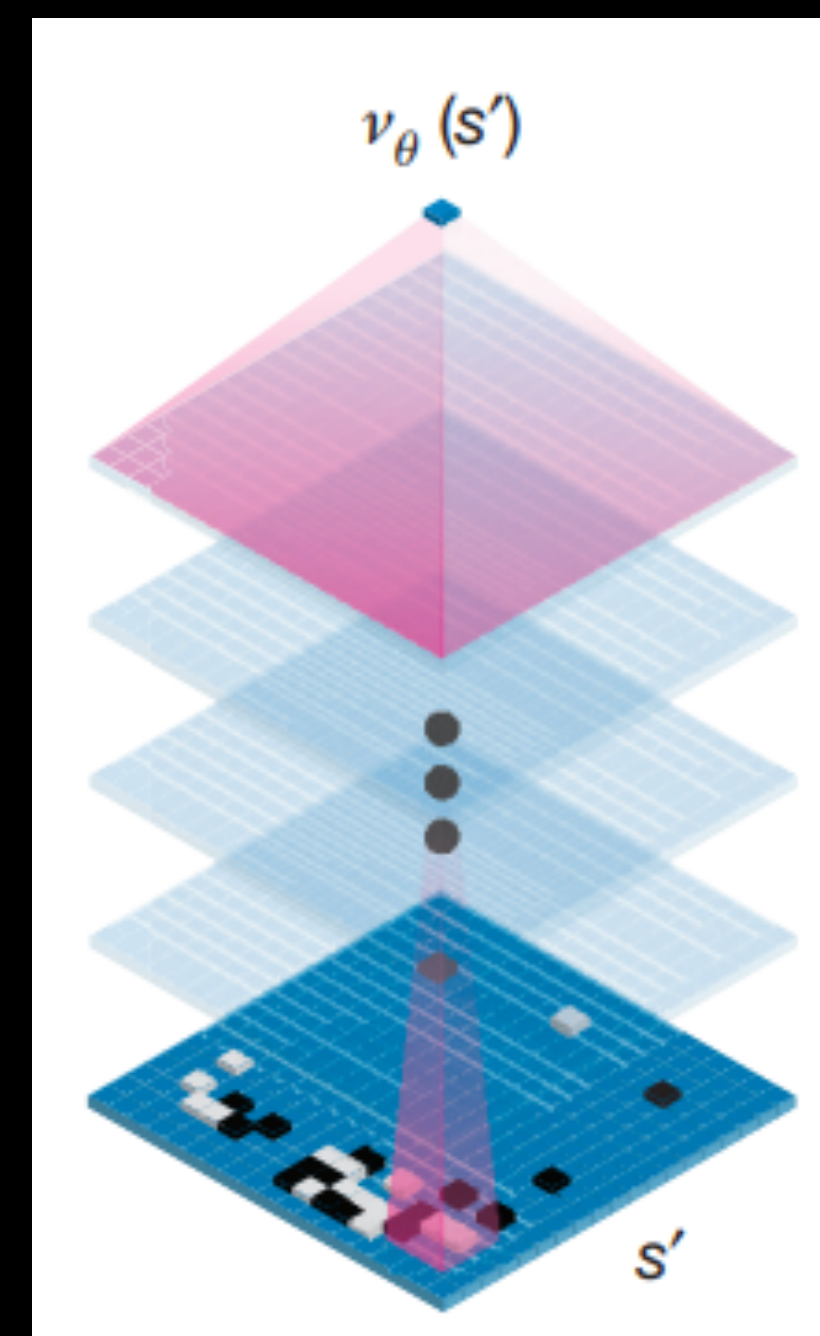


# Deep Reinforcement Learning

Policy Network

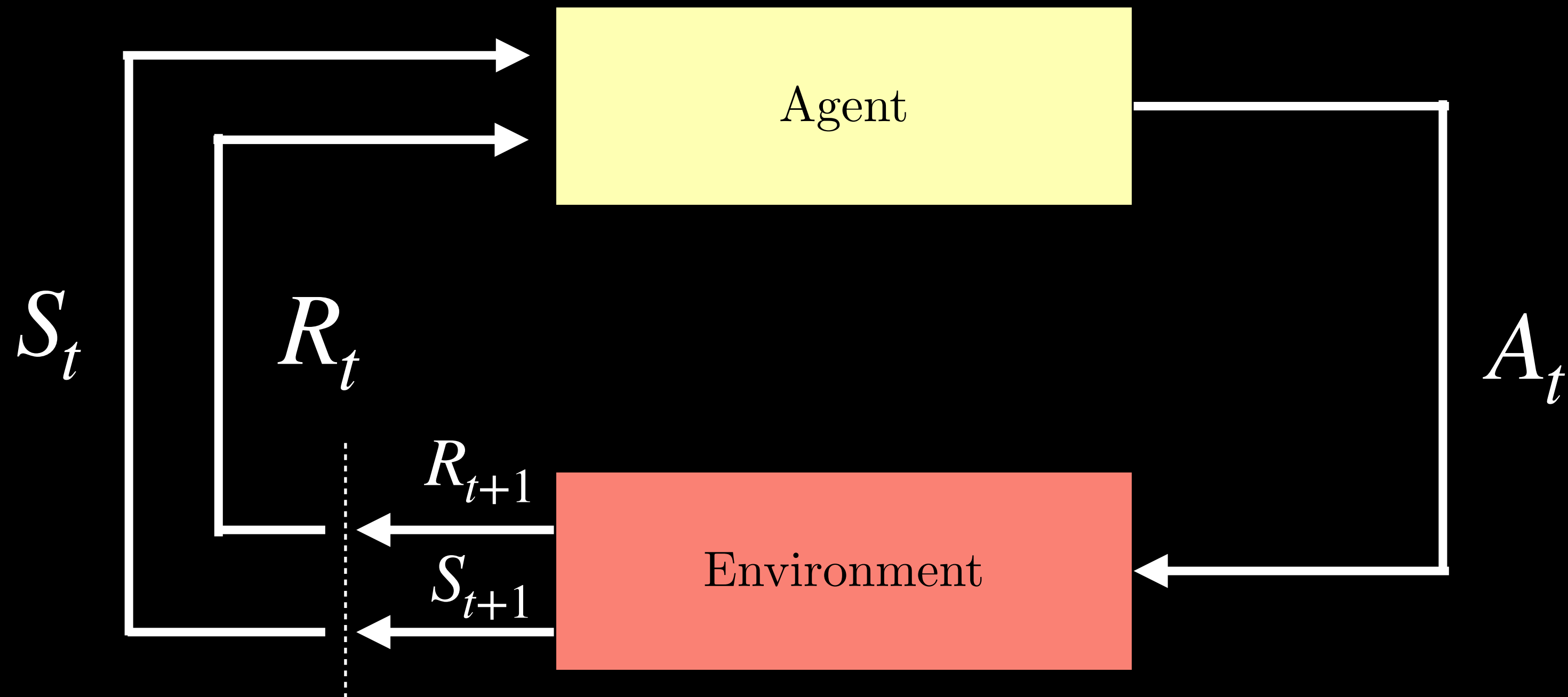


Value Network



# Model-Based Reinforcement Learning

$S_t$  : State  
 $R_t$  : Reward  
 $A_t$  : Actions



# Model-Based Reinforcement Learning

$S_t$  : State  
 $R_t$  : Reward  
 $A_t$  : Actions

