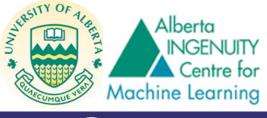
# RLAI

# Reinforcement Learning and Artificial Intelligence

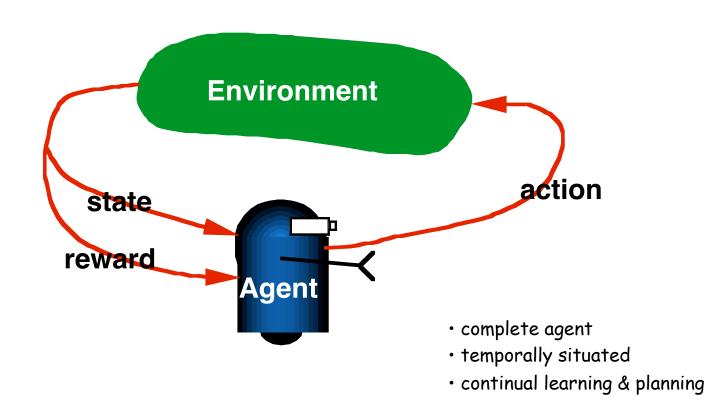


Pls:
Rich Sutton
Michael Bowling
Dale Schuurmans
Vadim Bulitko



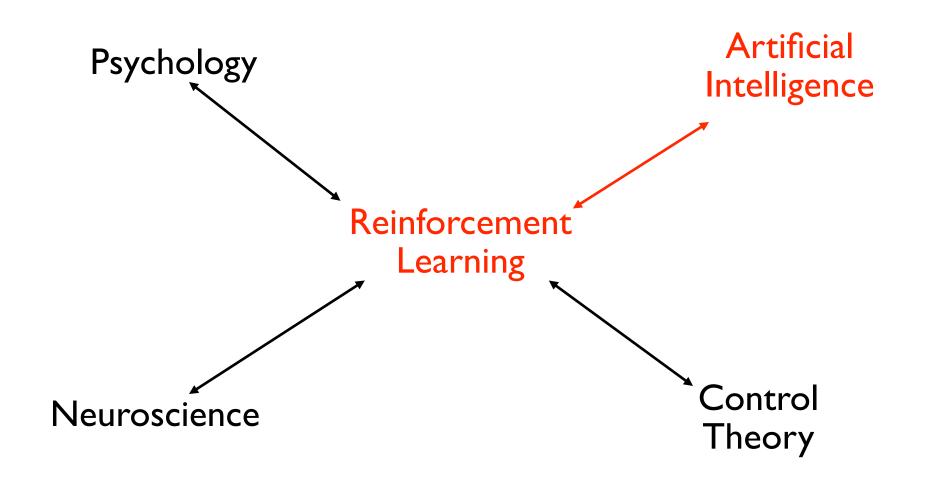


# Reinforcement learning is learning from interaction to achieve a goal



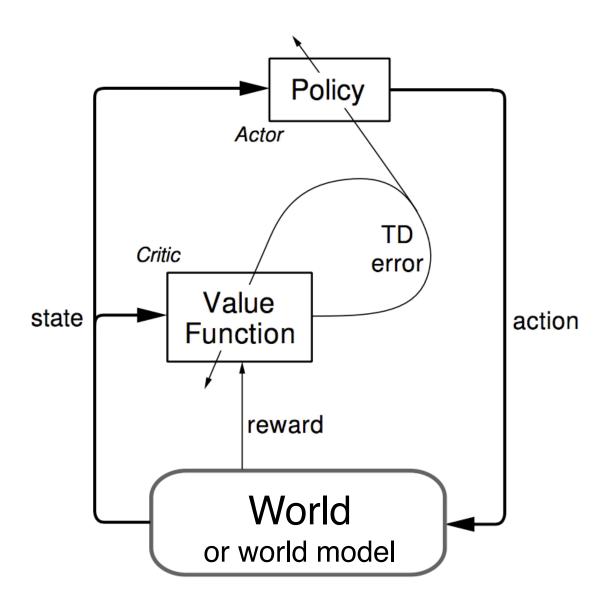
object is to affect environment

· environment stochastic & uncertain

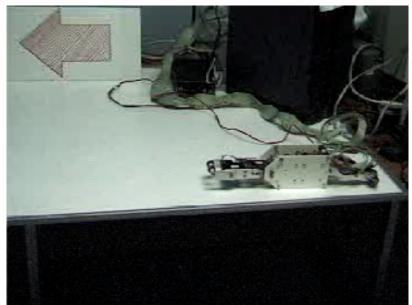


# Intro to RL with a robot bias

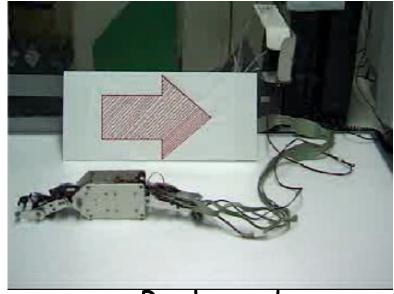
- policy
- reward
- value
- TD error
- world model
- generalized policy iteration



## Hajime Kimura's RL Robots



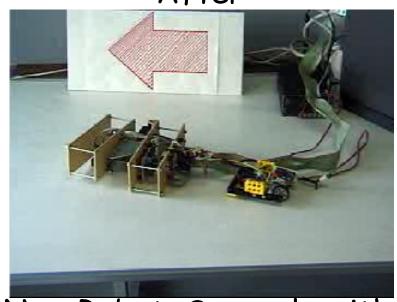
Before



Backward



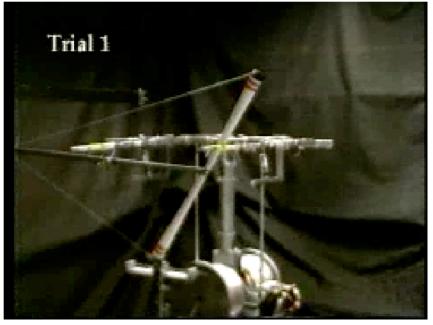
After



New Robot, Same algorithm

## Devilsticking

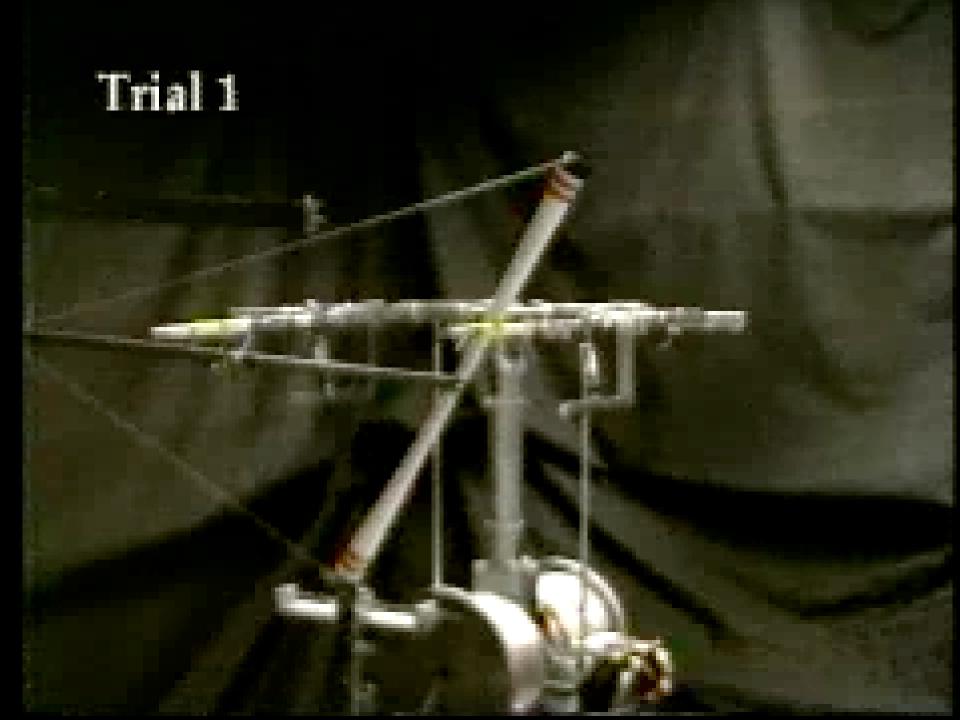




Finnegan Southey University of Alberta

Stefan Schaal & Chris Atkeson Univ. of Southern California "Model-based Reinforcement Learning of Devilsticking"





## The RoboCup Soccer Competition



### Autonomous Learning of Efficient Gait Kohl & Stone (UTexas) 2004









# **Policies**

- A policy maps each state to an action to take
  - Like a stimulus-response rule

We seek a policy that maximizes cumulative reward

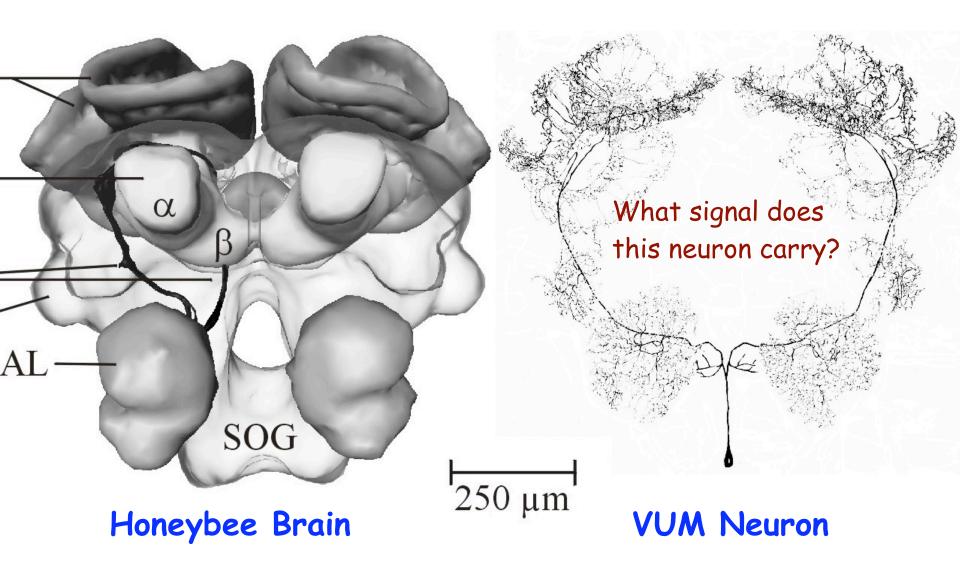
The policy is a subgoal to achieving reward

## The Reward Hypothesis

The goal of intelligence is to maximize the cumulative sum of a single received number: "reward" = pleasure - pain

Artificial Intelligence = reward maximization

## Brain reward systems



# Value

### Value systems are hedonism with foresight

We value situations according to how much reward we expect will follow them

All efficient methods for solving sequential decision problems determine (learn or compute) "value functions" as an intermediate step

Value systems are a means to reward, yet we care more about values than rewards

# Pleasure \neq good

"Even enjoying yourself you call evil whenever it leads to the loss of a pleasure greater than its own, or lays up pains that outweigh its pleasures. ... Isn't it the same when we turn back to pain? To suffer pain you call good when it either rids us of greater pains than its own or leads to pleasures that outweigh them."

-Plato, Protagoras

#### Reward

$$r: \mathsf{States} \to \mathsf{Pr}(\mathfrak{R})$$

e.g., 
$$r_t \in \mathfrak{R}$$

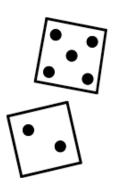
#### Policies

$$\pi$$
: States  $\rightarrow$  Pr(Actions)  $\pi^*$  is optimal

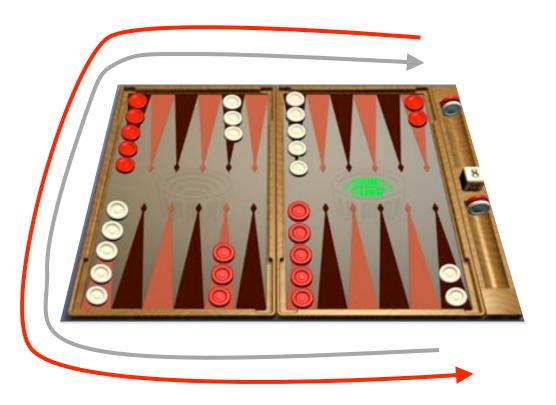
#### Value Functions

$$V^{\pi}: \mathsf{States} \to \Re$$

$$V^{\pi}(s) = E_{\pi} \left\{ \sum_{k=1}^{\infty} \gamma_{k}^{k-1} r_{t+k} \right\} \text{ given that } s_{t} = s$$
 discount factor  $\approx 1 \text{ but < 1}$ 



# Backgammon



STATES: configurations of the

playing board (≈10<sup>20</sup>)

**ACTIONS:** moves

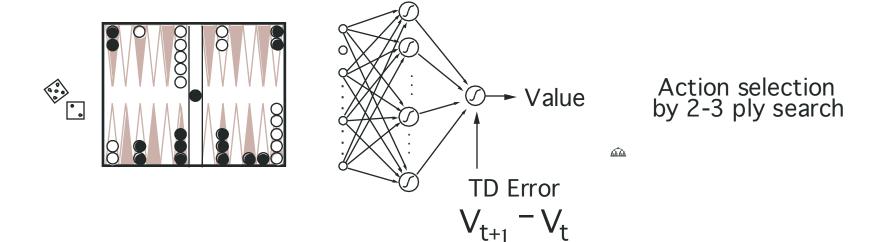
REWARDS: win: +1

lose: -I

else: 0

a "big" game

#### **TD-Gammon**



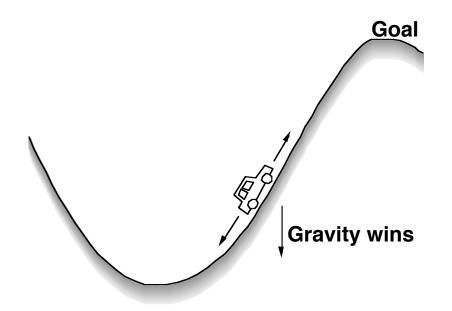
Start with a random Network

Play millions of games against itself

Learn a value function from this simulated experience

Six weeks later it's the best player of backgammon in the world

#### The Mountain Car Problem



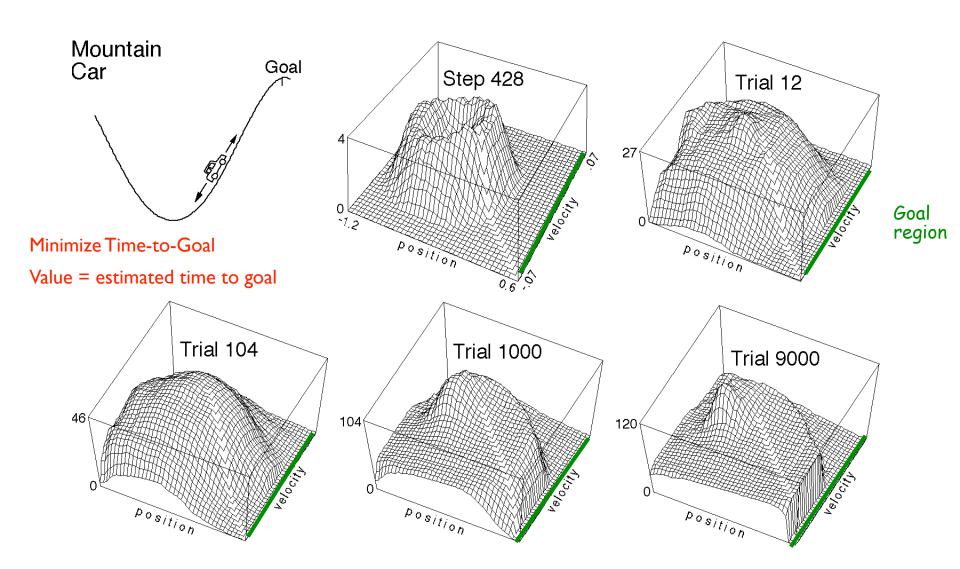
<u>SITUATIONS</u>: car's position and velocity

<u>ACTIONS</u>: three thrusts: forward, reverse, none

REWARDS: always -1 until car reaches the goal No Discounting

Minimum-Time-to-Goal Problem

# Value Functions Learned while solving the Mountain Car problem



# TD error

#### Reward

$$r: \mathsf{States} \to \mathsf{Pr}(\mathfrak{R})$$

e.g., 
$$r_t \in \mathfrak{R}$$

#### Policies

$$\pi$$
: States  $\rightarrow$  Pr(Actions)  $\pi^*$  is optimal

#### Value Functions

$$V^{\pi}: \mathsf{States} \to \Re$$

$$V^{\pi}(s) = E_{\pi} \left\{ \sum_{k=1}^{\infty} \gamma_{k}^{k-1} r_{t+k} \right\} \text{ given that } s_{t} = s$$
 discount factor  $\approx 1 \text{ but < 1}$ 

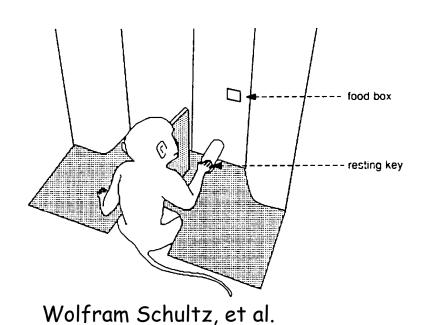
$$V_t = E\left\{\sum_{k=1}^{\infty} \gamma^{k-1} r_{t+k}\right\}$$

$$= E \left\{ r_{t+1} + \gamma \sum_{k=1}^{\infty} \gamma^{k-1} r_{t+1+k} \right\}$$

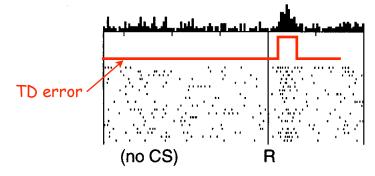
$$\approx r_{t+1} + \gamma V_{t+1}$$

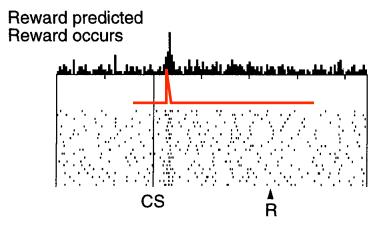
$$TD \ error_t = r_{t+1} + \gamma V_{t+1} - V_t$$

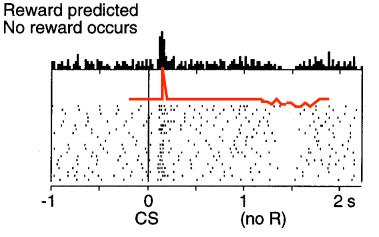
# Brain reward systems seem to signal TD error

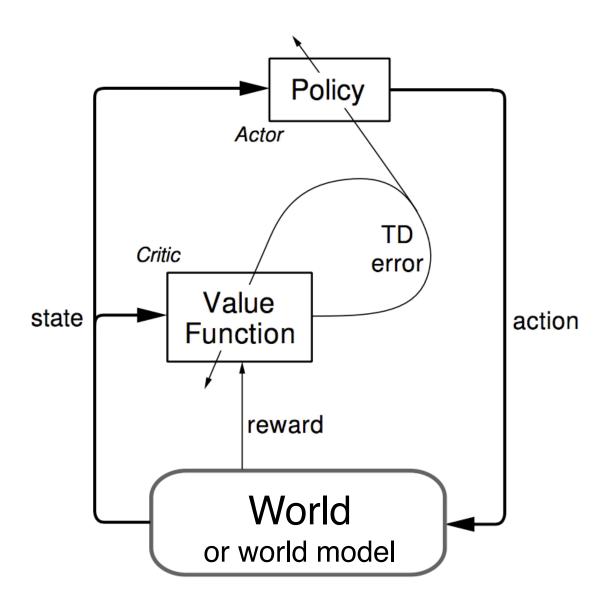


No prediction Reward occurs









# World models

# "Autonomous helicopter flight via Reinforcement Learning"

Ng (Stanford), Kim, Jordan, & Sastry (UC Berkeley) 2004



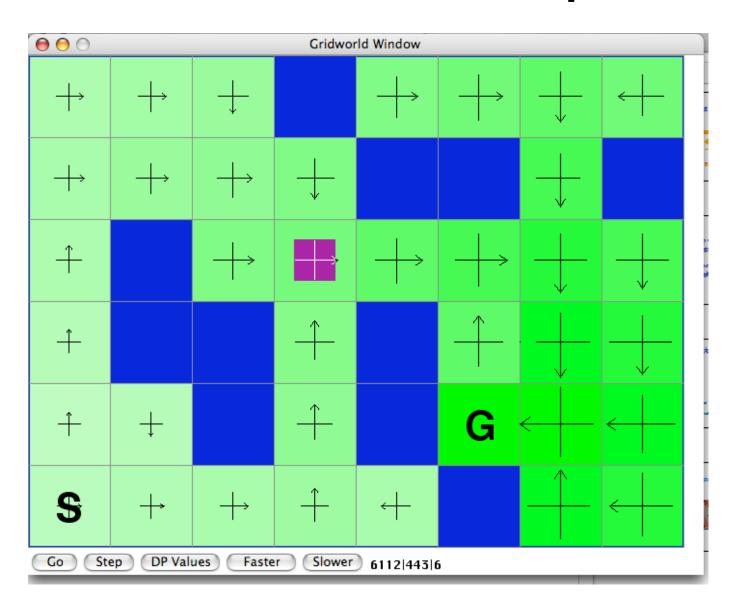




## Reason as RL over Imagined Experience

- I. Learn a predictive model of the world's dynamics transition probabilities, expected immediate rewards
- 2. Use model to generate imaginary experiences internal thought trials, mental simulation (Craik, 1943)
- 3. Apply RL as if experience had really happened vicarious trial and error (Tolman, 1932)

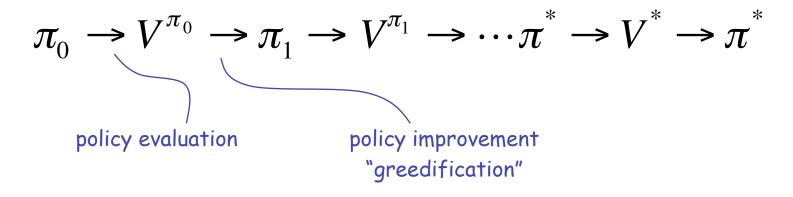
# GridWorld Example



# Intro to RL with a robot bias

- policy
- reward
- value
- TD error
- world model
- generalized policy iteration

### Policy Iteration



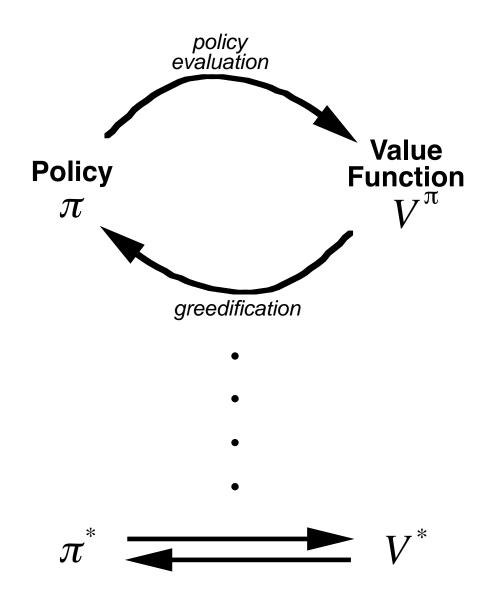
Improvement is monotonic

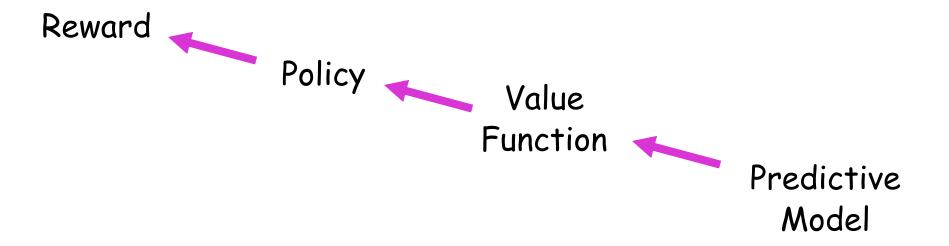
Converges is a finite number of steps

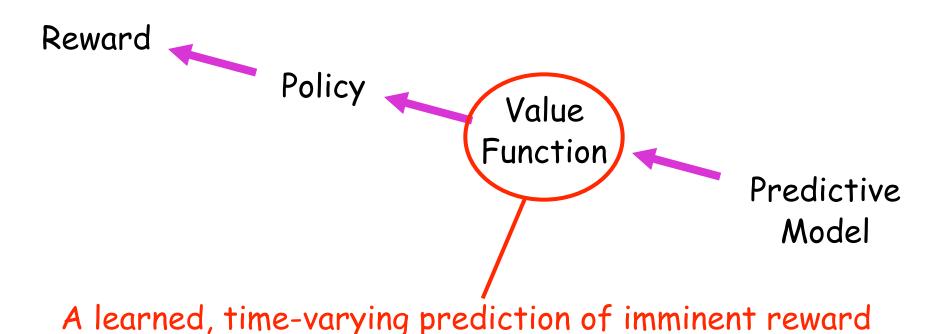
#### <u>Generalized</u> Policy Iteration:

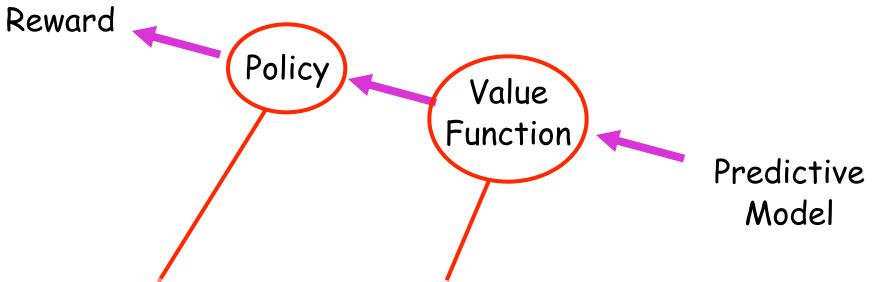
Intermix the two steps more finely, state by state, action by action, sample by sample

#### Generalized Policy Iteration

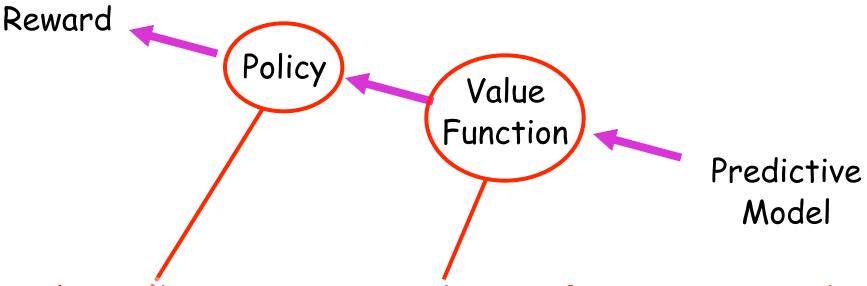






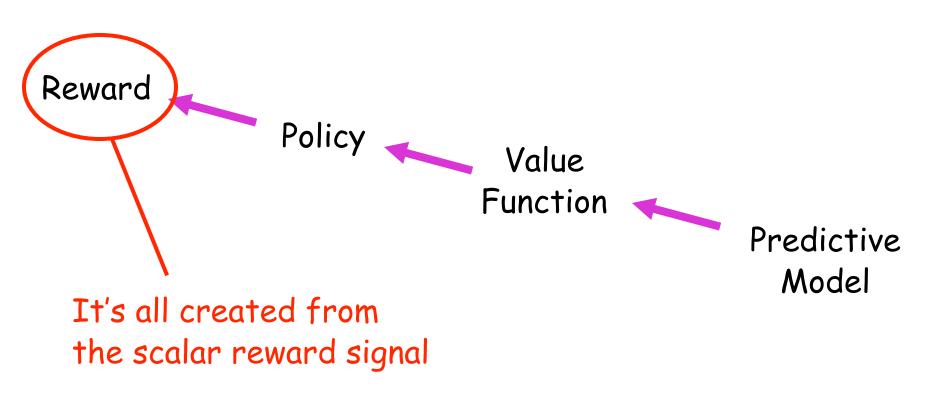


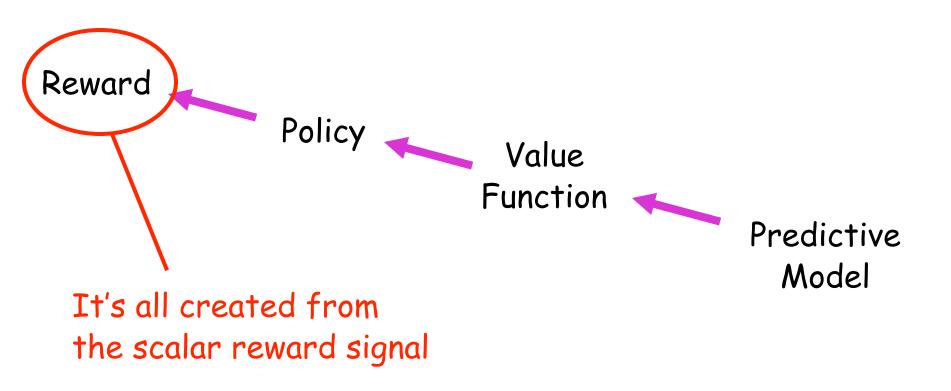
A learned, time-varying prediction of imminent reward Key to all efficient methods for finding optimal policies



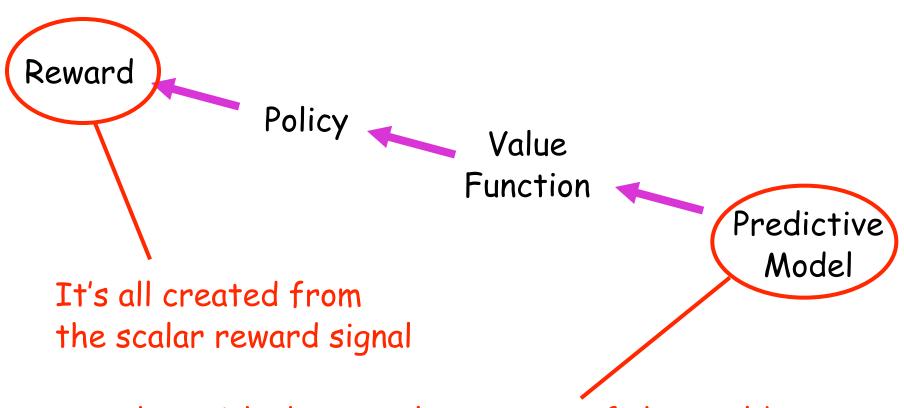
A learned, time-varying prediction of imminent reward Key to all efficient methods for finding optimal policies

This has nothing to do with either biology or computers





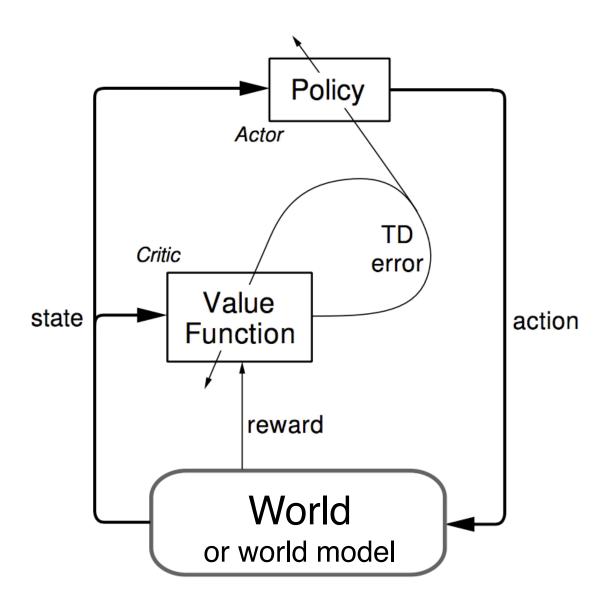
together with the causal structure of the world



together with the causal structure of the world

#### additions for next year

- sarsa
- how it might actually applied to something they might do



#### Current work

#### Knowledge is subjective

An individual's knowledge consists of predictions about its future (low-level) sensations and actions

Everything else (objects, houses, people, love) must be constructed from sensation and action

## Subjective/Predictive Knowledge

- "John is in the coffee room"
- "My car in is the South parking lot"
- What we know about geography, navigation
- What we know about how an object looks, rotates
- What we know about how objects can be used
- Recognition strategies for objects and letters
- "The portrait of Washington on the dollar in the wallet in my other pants in the laundry, has a mustache on it"

## RoboCup soccer keepaway

Stone, Sutton & Kuhlmann, 2005

