

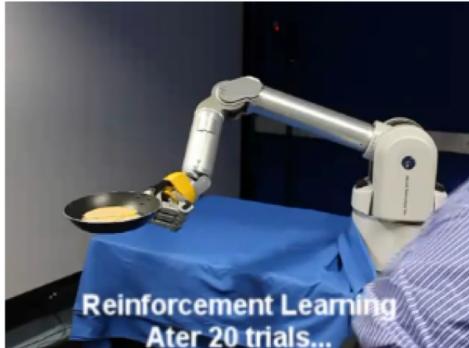
# CSC 665: Artificial Intelligence

## Reinforcement Learning II

Instructor: Pooyan Fazli  
San Francisco State University

# AI in the News

## ROBOT LEARNS TO FLIP PANCAKES



A ROBOT LEARNING to flip pancakes from Sylvain Calinon on Vimeo.

Flipping a pancake seems like one of those things you can do when you are just barely awake and still to get your morning caffeine.

Not so, if you are a robot. Then learning how to flip a pancake is quite a task and it can take 50 tries to get it right.

Two researchers at the Italian Institute of Technology—Petar Kormushev and Sylvain Calinon—taught a robot the

[https://www.youtube.com/watch?v=W\\_gxLKSsSIE](https://www.youtube.com/watch?v=W_gxLKSsSIE)

# Exploitation vs. Exploration

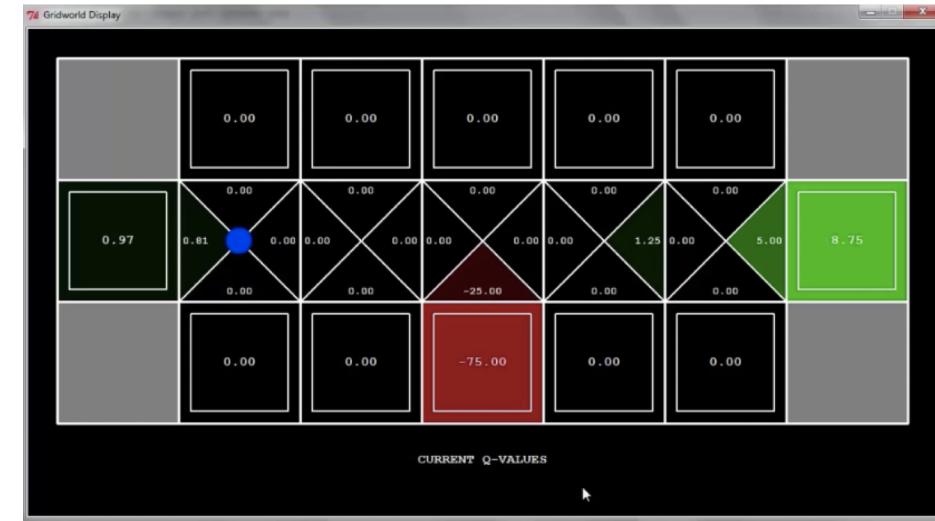
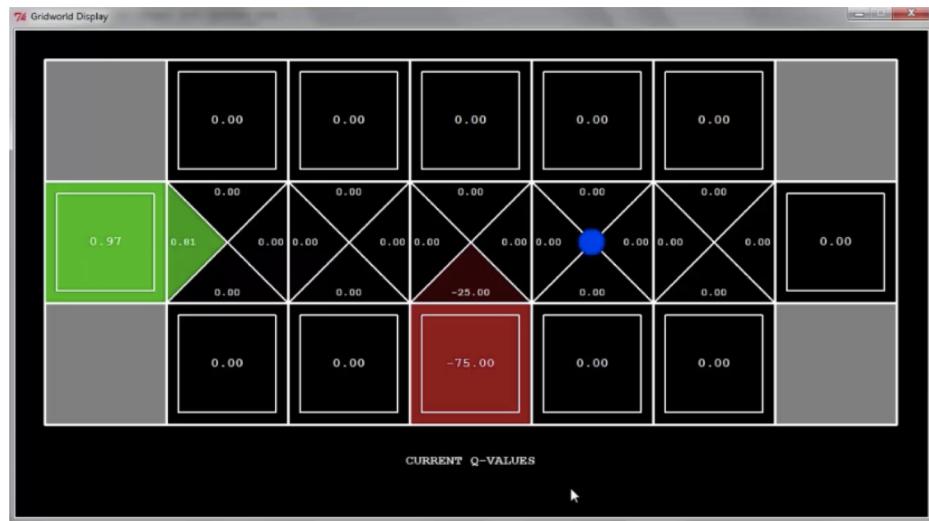
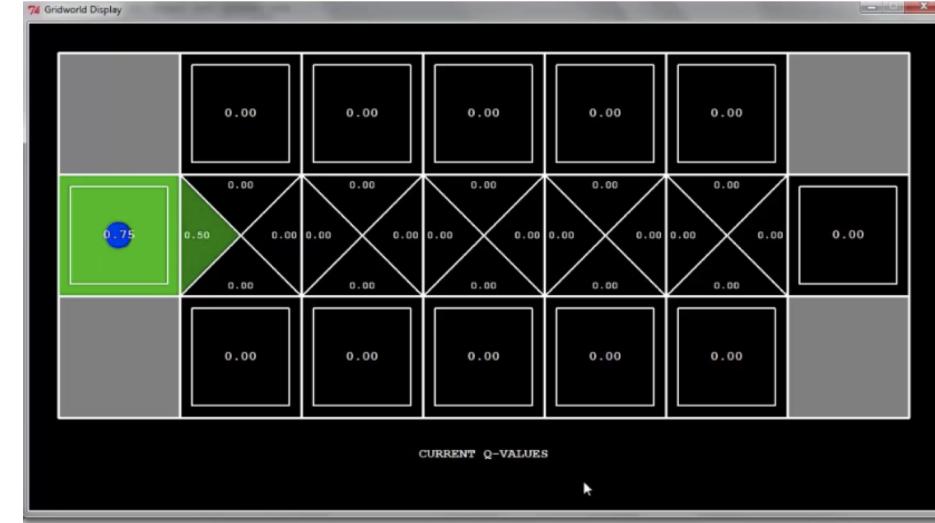
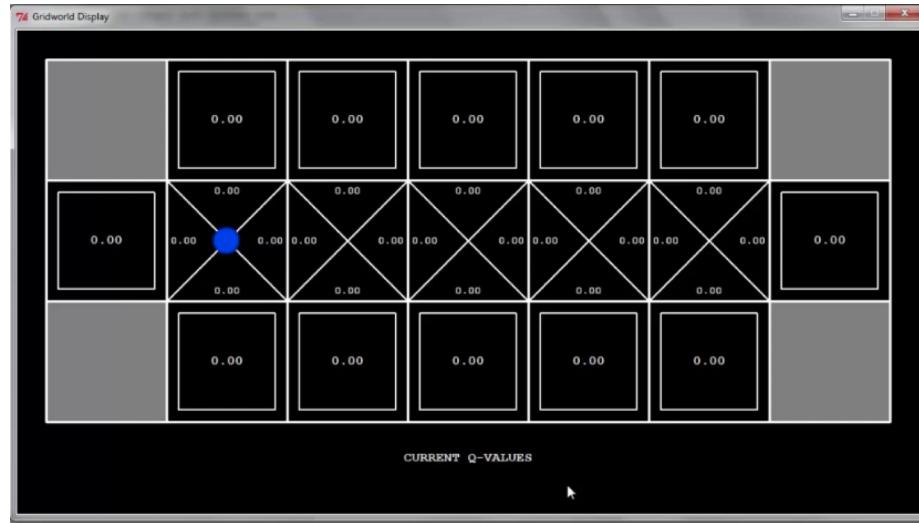
---

# Exploitation vs. Exploration

---

- **Exploitation:** Execute the current best (maybe optimal?) policy to get high payoff
- **Exploration:** Try new sequences of (possibly random) actions to improve the agent's knowledge of the environment even though current model doesn't believe they have high payoff

# Q-learning – Manual Exploration – Bridge Grid

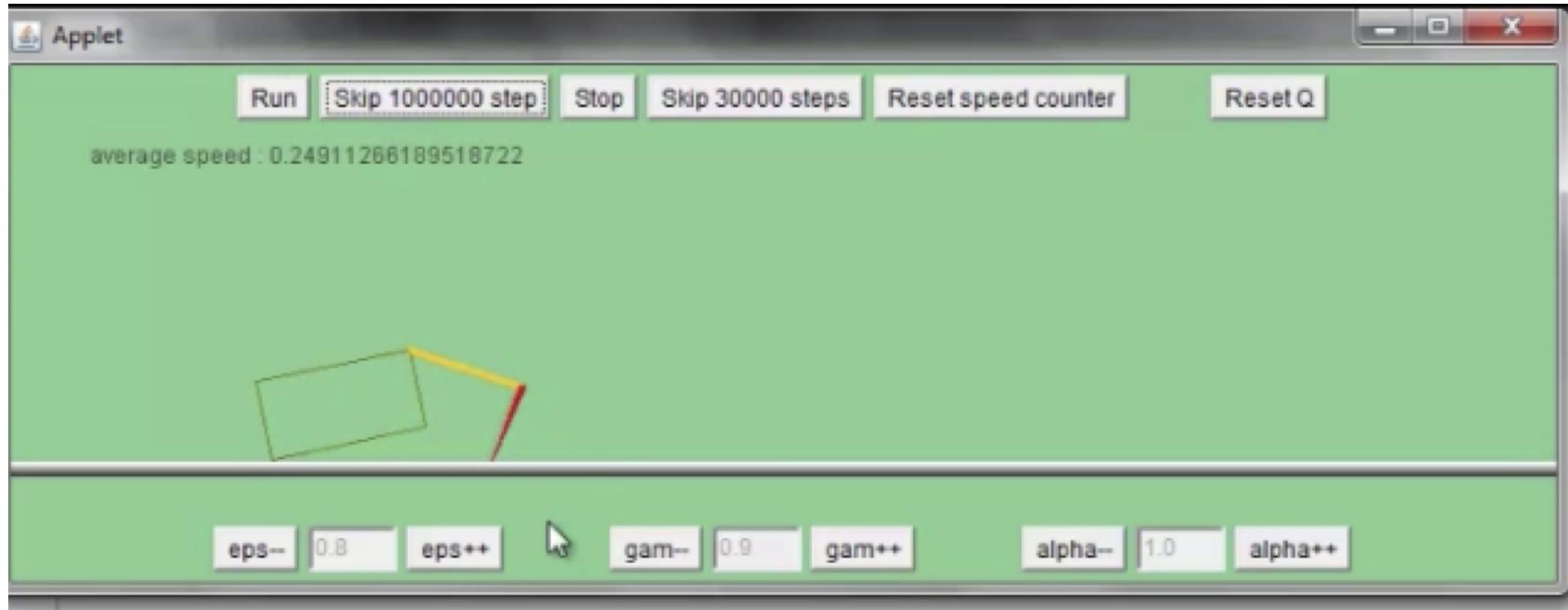


# How to Explore?

---

- Simplest: random actions ( $\epsilon$ -greedy)
  - Every time step, flip a coin
  - With (small) probability  $\epsilon$ , act **randomly**
  - With (large) probability  $1-\epsilon$ , act on current best policy
- Properties of  $\epsilon$ -greedy exploration
  - Every  $s,a$  pair is tried infinitely often
  - Does a lot of stupid things
    - Jumping off a cliff **lots of times** to make sure it hurts
  - Keeps doing stupid things for ever
    - Decay  $\epsilon$  towards 0

# Q-learning – Epsilon-Greedy – Crawler



# Exploration Functions

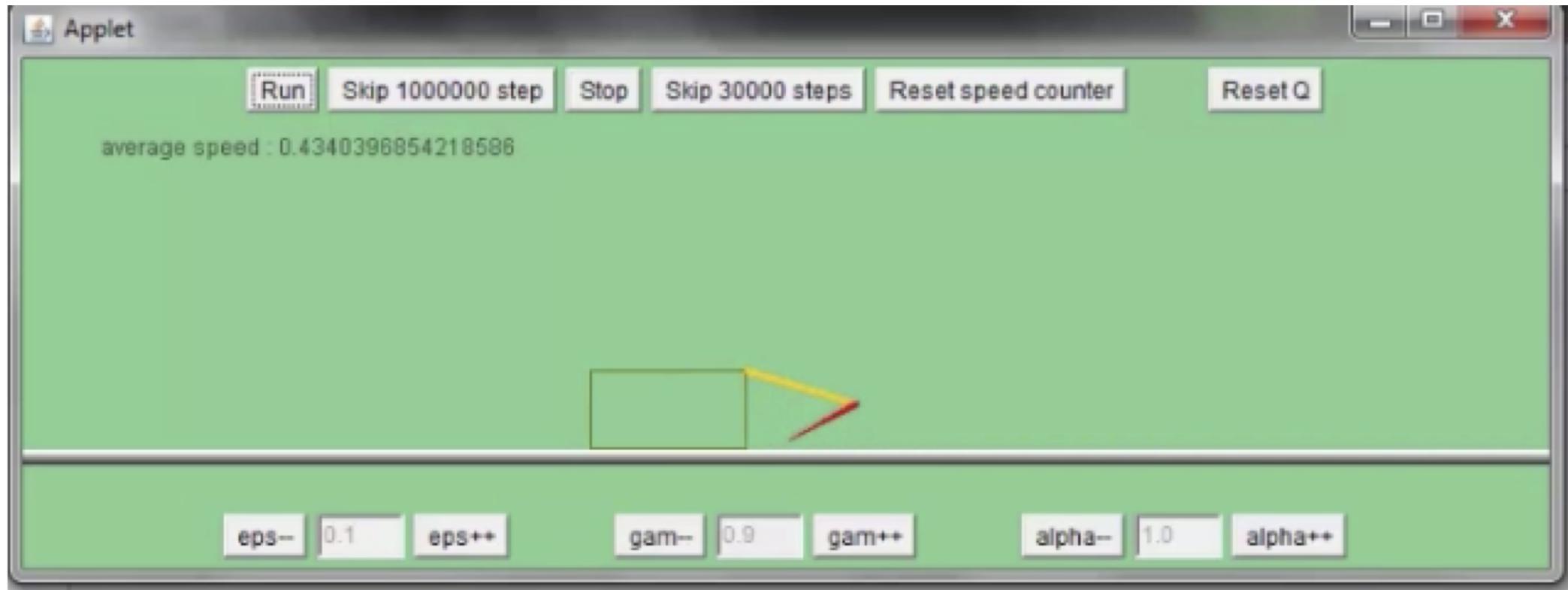
- What to explore?
  - Actions towards unexplored regions are encouraged (**much faster than  $\varepsilon$ -greedy!**)
- Exploration function
  - Takes a value estimate **u**, a visit count **n**, and a constant **k** returns:

$$f(u, n) = u + k/n$$

Regular Q-update:  $Q(s,a) \leftarrow (1-\alpha) Q(s,a) + \alpha [R(s,a,s') + \gamma \max_{a'} Q(s',a')]$

Modified Q-update:  $Q(s,a) \leftarrow (1-\alpha) Q(s,a) + \alpha [R(s,a,s') + \gamma \max_{a'} f(Q(s',a'), n(s',a'))]$

# Q-learning – Exploration Function – Crawler



# Approximate Q-Learning

---

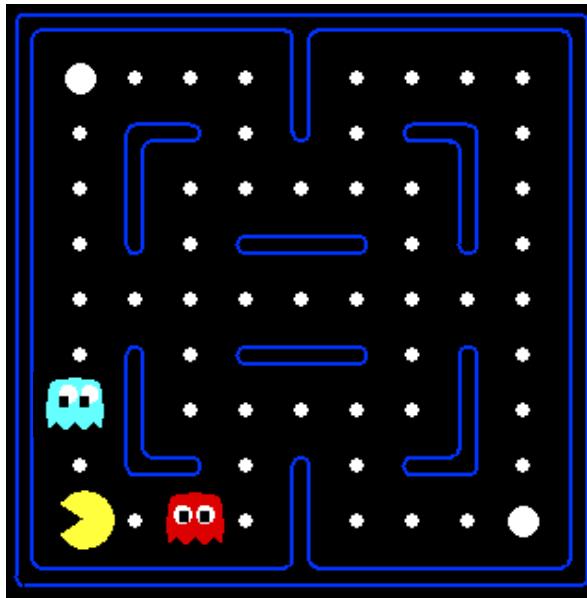
# Generalizing Across States

---

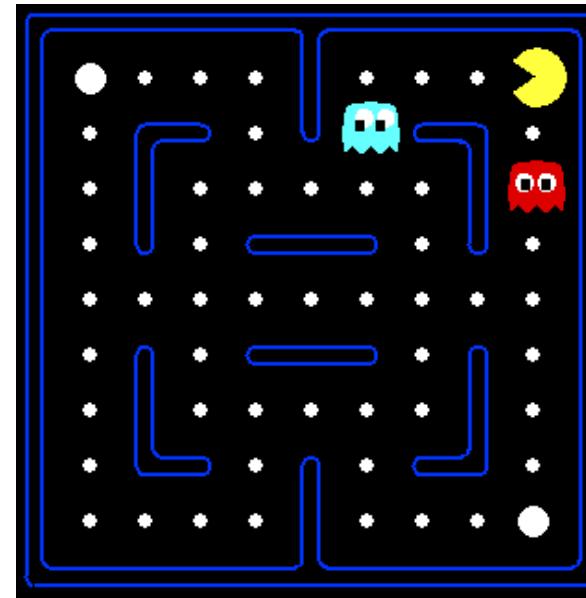
- Basic Q-Learning keeps a table of all Q-values
- In realistic situations, we cannot possibly learn about every single state!
  - Too many states to visit them all in training
  - Too many states to hold the Q-tables in memory
- Instead, we want to generalize:
  - Learn about some small number of states from experience
  - Generalize that experience to new, similar situations

# Example: Pacman

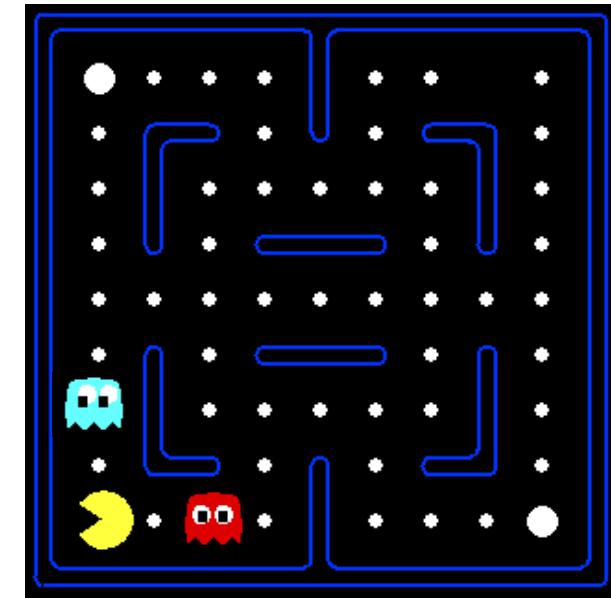
Let's say we discover through experience that this state is bad:



In naïve q-learning, we know nothing about this state:

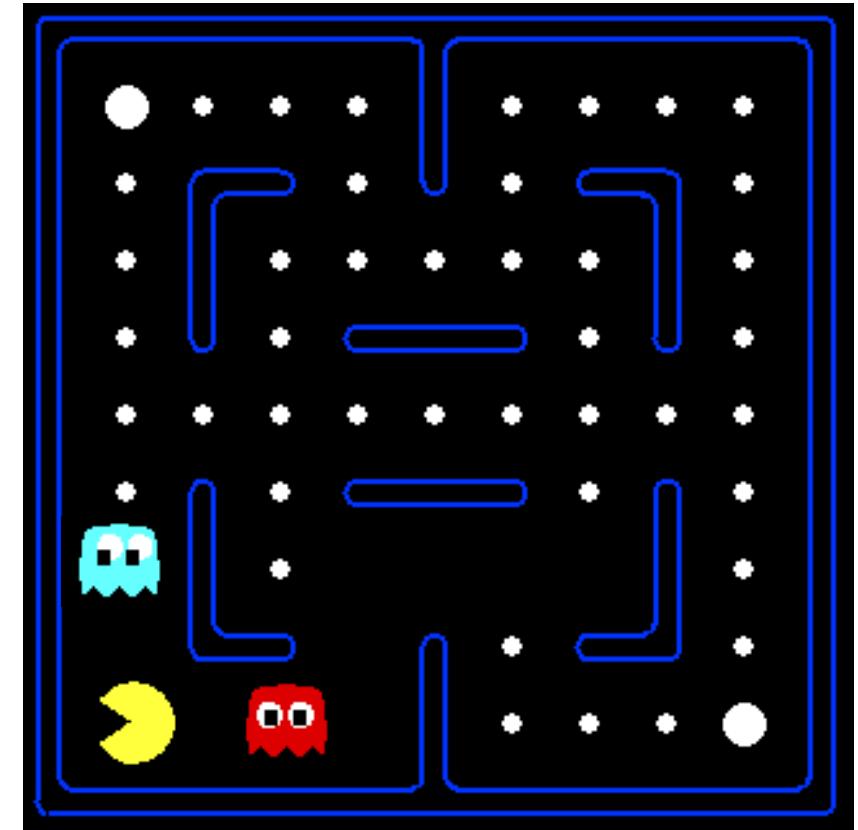


Or even this one!



# Feature-Based Representations

- **Solution:** describe a state using a vector of **features** (similar to the evaluation function in Assignment 2)
- **Features** are functions,  $f(s)$ , from states to real numbers that capture important properties of the state
  - Example features:
    - Distance to closest ghost
    - $1 / (\text{distance to closest ghost})$   $f_{GST}$
    - Number of ghosts
    - Distance to closest dot
    - $1 / (\text{dist to closest dot})$   $f_{DOT}$
    - Is Pacman in a tunnel? (0/1)
    - ..... etc.
- Can also describe a q-state  $(s, a)$  with features (e.g. action  $a$  in state  $s$  moves closer to food)



# Linear Value Functions

---

- Using a feature representation, we can write the Q-value (or the V-value) for any state using a few weights:

$$V(s) = w_1 f_1(s) + w_2 f_2(s) + \dots + w_n f_n(s)$$

What you used in Assignment 2

$$Q(s, a) = w_1 f_1(s, a) + w_2 f_2(s, a) + \dots + w_n f_n(s, a)$$

What you use in Approximate Q-learning

# Approximate Q-Learning

$$Q(s, a) = w_1 f_1(s, a) + w_2 f_2(s, a) + \dots + w_n f_n(s, a)$$

- Q-learning with linear functions:

transition =  $(s, a, r, s')$

difference =  $[r + \gamma \max_{a'} Q(s', a')] - Q(s, a)$

---

$Q(s, a) \leftarrow (1 - \alpha)Q(s, a) + (\alpha) \left[ r + \gamma \max_{a'} Q(s', a') \right]$

$Q(s, a) \leftarrow Q(s, a) + \alpha \text{ [difference]}$

---

What  
you saw  
before

Q-Learning

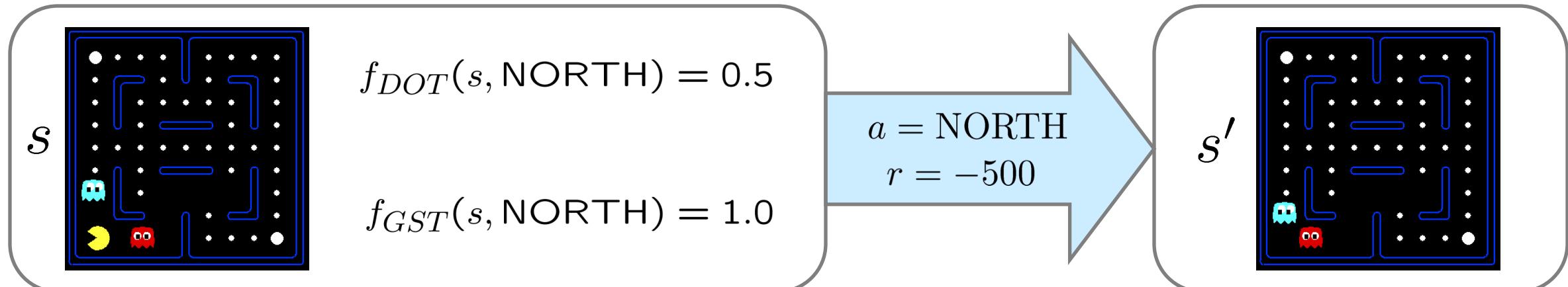
$w_i \leftarrow w_i + \alpha \text{ [difference]} f_i(s, a)$

Approximate Q-Learning

- Intuitive interpretation:
  - Adjust weights of features

# Example: Q-Pacman

$$Q(s, a) = 4.0f_{DOT}(s, a) - 1.0f_{GST}(s, a)$$



$$Q(s, \text{NORTH}) = +1$$

$$r + \gamma \max_{a'} Q(s', a') = -500 + 0$$

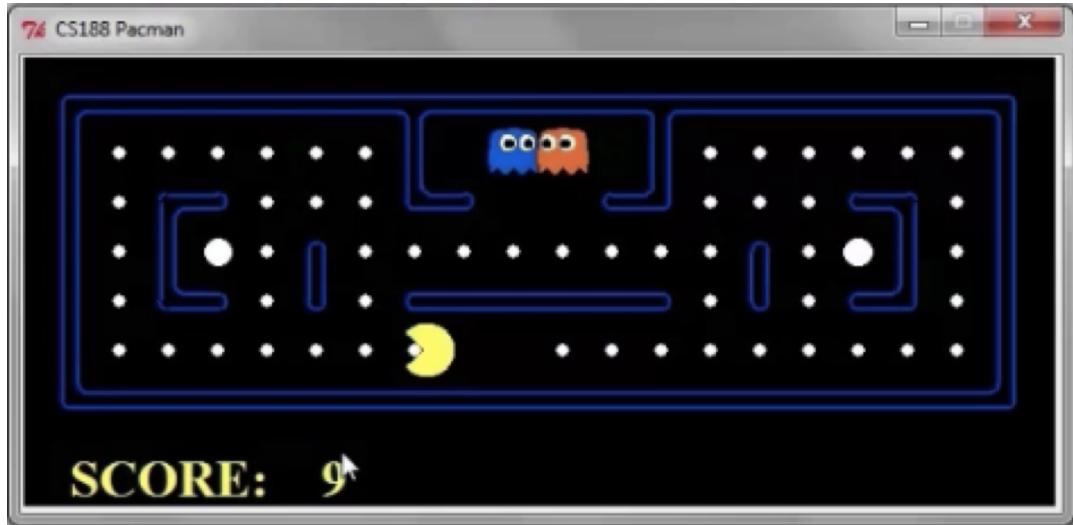
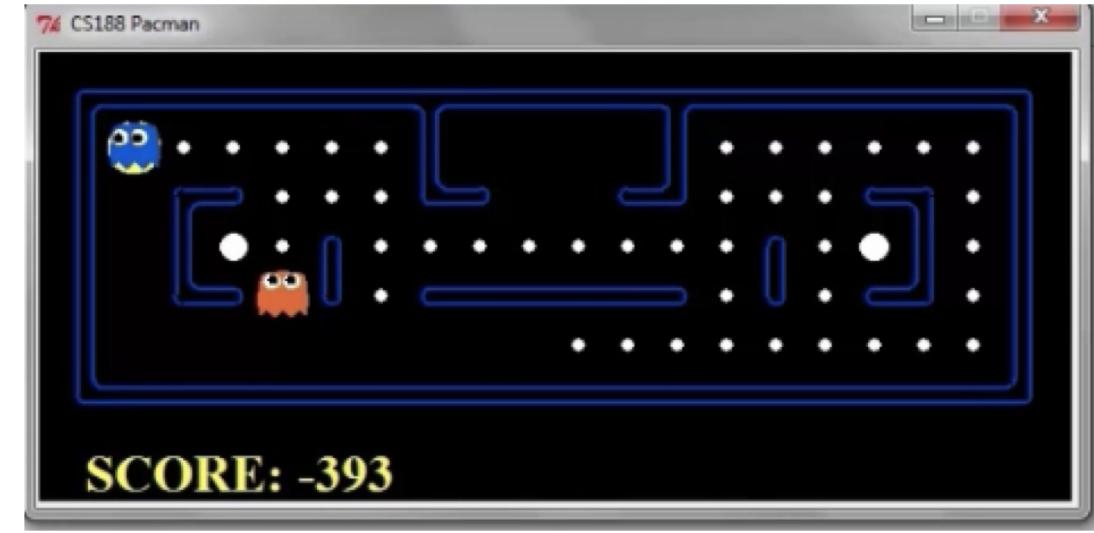
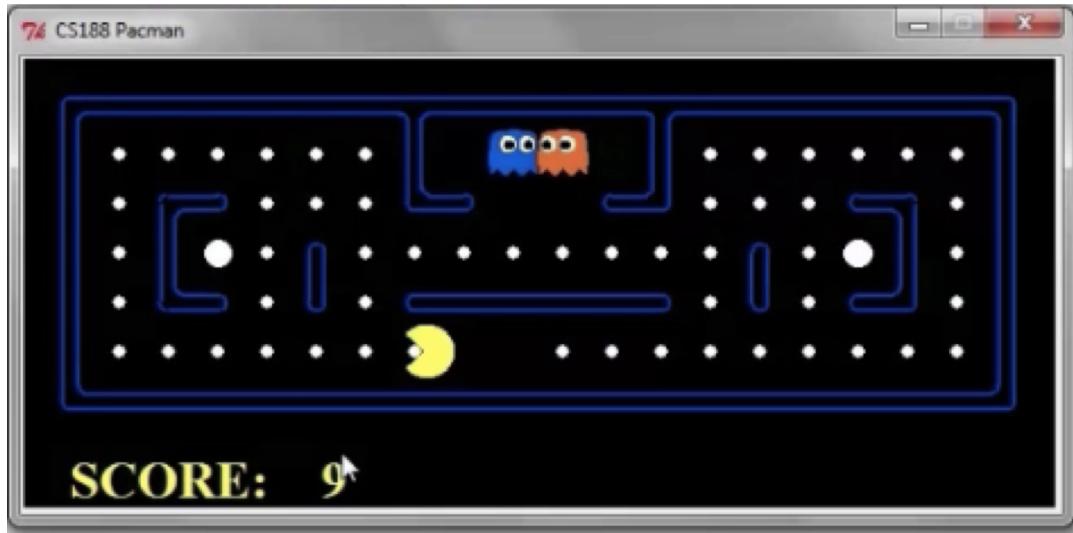
$$\text{difference} = -501$$

$$w_{DOT} \leftarrow 4.0 + \alpha [-501] 0.5$$
$$w_{GST} \leftarrow -1.0 + \alpha [-501] 1.0$$

$$Q(s, a) = 3.0f_{DOT}(s, a) - 3.0f_{GST}(s, a)$$

alpha = 0.004

# Approximate Q-Learning -- Pacman



# Conclusion

---

- We've seen how AI methods can solve problems in:
  - Search
  - Constraint Satisfaction Problems
  - Games
  - Markov Decision Problems
  - Reinforcement Learning
- Next up: Uncertainty and Learning!

# Reading

---

- Read Sections 21.4, 21.5 in the AIMA textbook