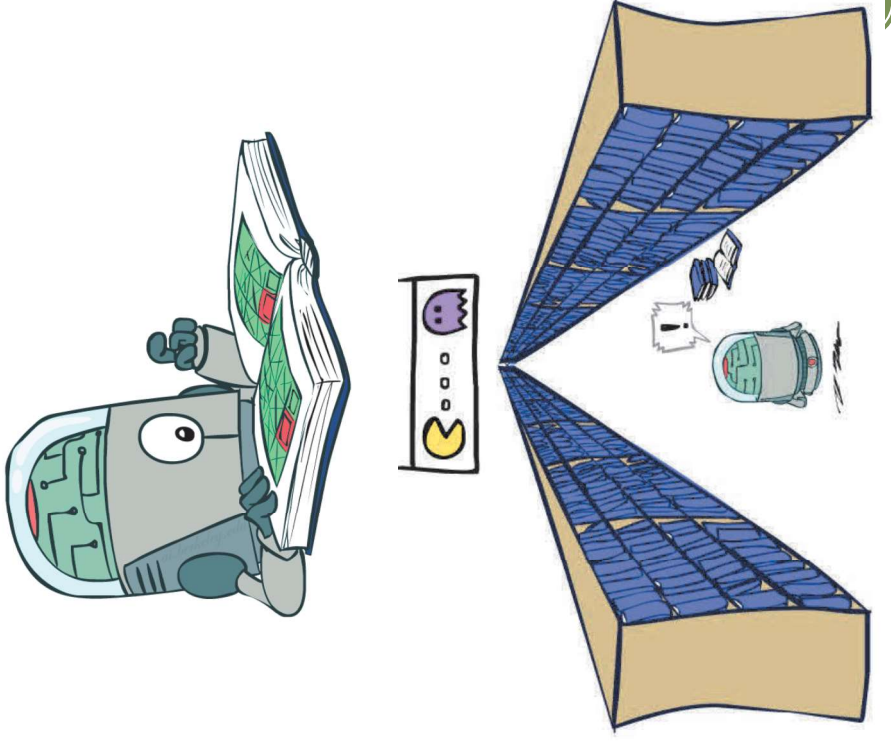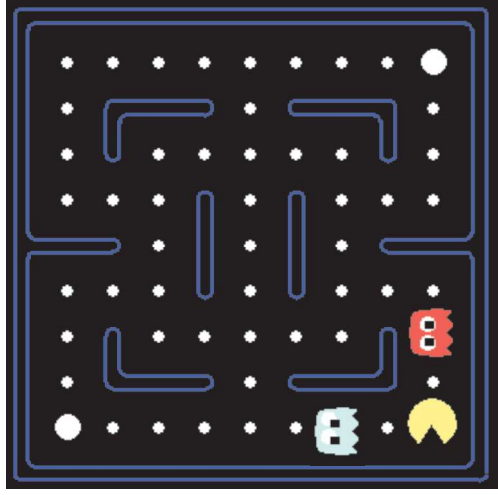# Generalizing Across States

- 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 training states from experience
  - Generalize that experience to new, similar situations
  - This is a fundamental idea in machine learning, and we'll see it over and over again
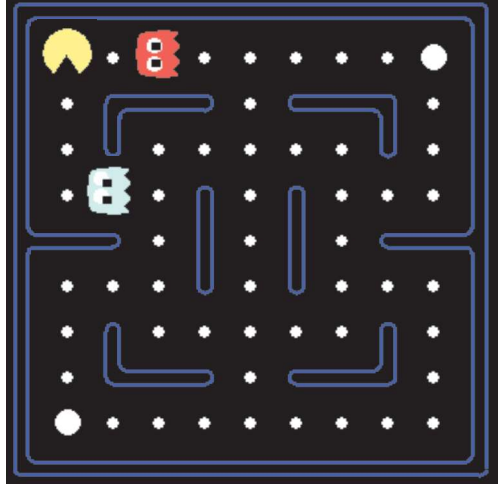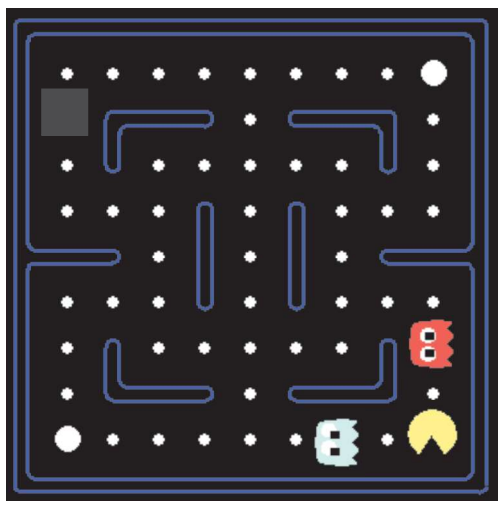
# Example: Pacman

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

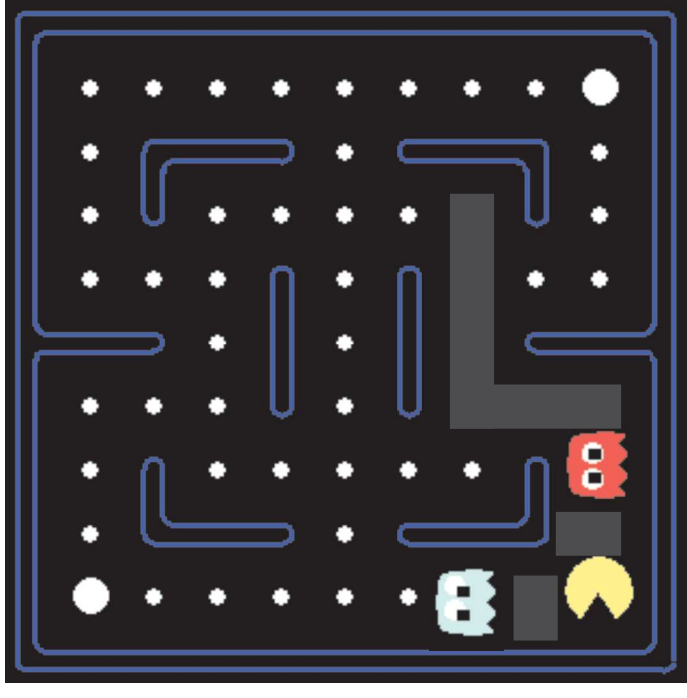In naïve q-learning, we know nothing about this state:

Or even this one!

# Feature–Based Representations

- Solution: describe a state using a vector of features (properties)
  - Features are functions from states to real numbers (often 0/1) that capture important properties of the state
  - Example features:
    - Distance to closest ghost
    - Distance to closest dot
    - Number of ghosts
    - $1 / (\text{dist to dot})^2$
    - Is Pacman in a tunnel? (0/1)
    - …… etc.
    - Is it the exact state on this slide?
  - Can also describe a q-state (s, a) with features (e.g. action moves closer to food)

# Linear Value Functions

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

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

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

- Advantage: our experience is summed up in a few powerful numbers

- Disadvantage: states may share features but actually be very different in value!

# Approximate Q-Learning

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

- Q-learning with linear Q-functions:

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

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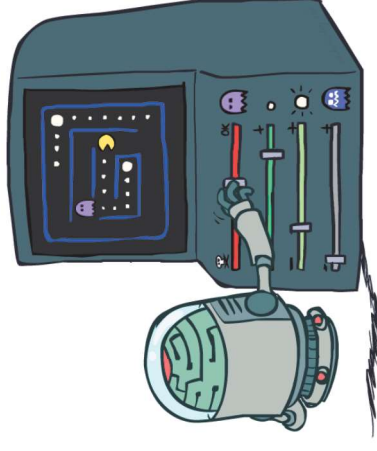
  $Q(s,a) \leftarrow Q(s,a) + \alpha \, [\text{difference}]$      Exact Q's

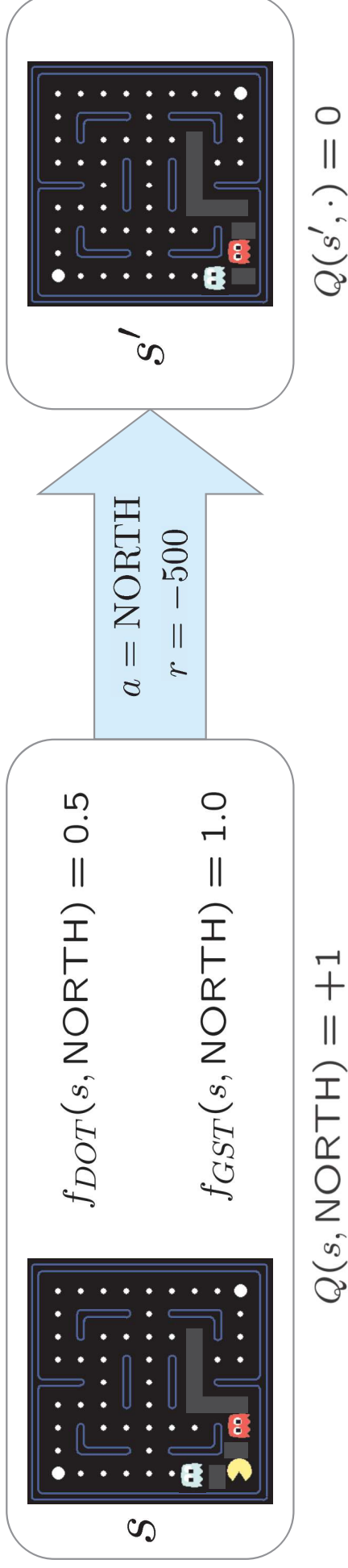  $w_i \leftarrow w_i + \alpha \, [\text{difference}] \, f_i(s,a)$      Approximate Q's

- Intuitive interpretation:

- Adjust weights of active features

- E.g., if something unexpectedly bad happens, blame the features that were on: disprefer all states with that state's features

- Formal justification: online least squares

# Example: Q-Pacman

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



$s$

$f_{DOT}(s, \text{NORTH}) = 0.5$

$f_{GST}(s, \text{NORTH}) = 1.0$

$a = \text{NORTH}$
$r = -500$



$s'$

$Q(s', \cdot) = 0$

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

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

difference $= -501$

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

$w_{GST} \leftarrow -1.0 + \alpha\,[-501]\,1.0$

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

# Q-learning with Linear Approximation

**Algorithm 4:** Q-learning with linear approximation.

1 Initialize q-value function $Q$ with random weights $w$: $Q(s,a;w) = \sum_m w_m f_m(s,a)$;

2 **for** $episode = 1 \to M$ **do**

3    Get initial state $s_0$;

4    **for** $t = 1 \to T$ **do**

5       With prob. $\epsilon$, select a random action $a_t$;

6       With prob. $1 - \epsilon$, select $a_t \in \text{argmax}_a\, Q(s_t, a; w)$;

7       Execute selected action $a_t$ and observe reward $r_t$ and next state $s_{t+1}$;

8       Set target $y_t = \begin{cases} r_t & \text{if episode terminates at step } t+1 \\ r_t + \gamma \max_{a'} Q(s_{t+1}, a'; w) & \text{otherwise} \end{cases}$;

9       Perform a gradient descent step to update $w$: $w_m \leftarrow w_m + \alpha\left[y_t - Q(s_t, a_t; w)\right] f_m(s,a)$;

Thanh H. Nguyen

18