

MDP and methods do solve it

With cute videos

This presentation, as well as some practical tasks,
mostly borrowed from SHAD course

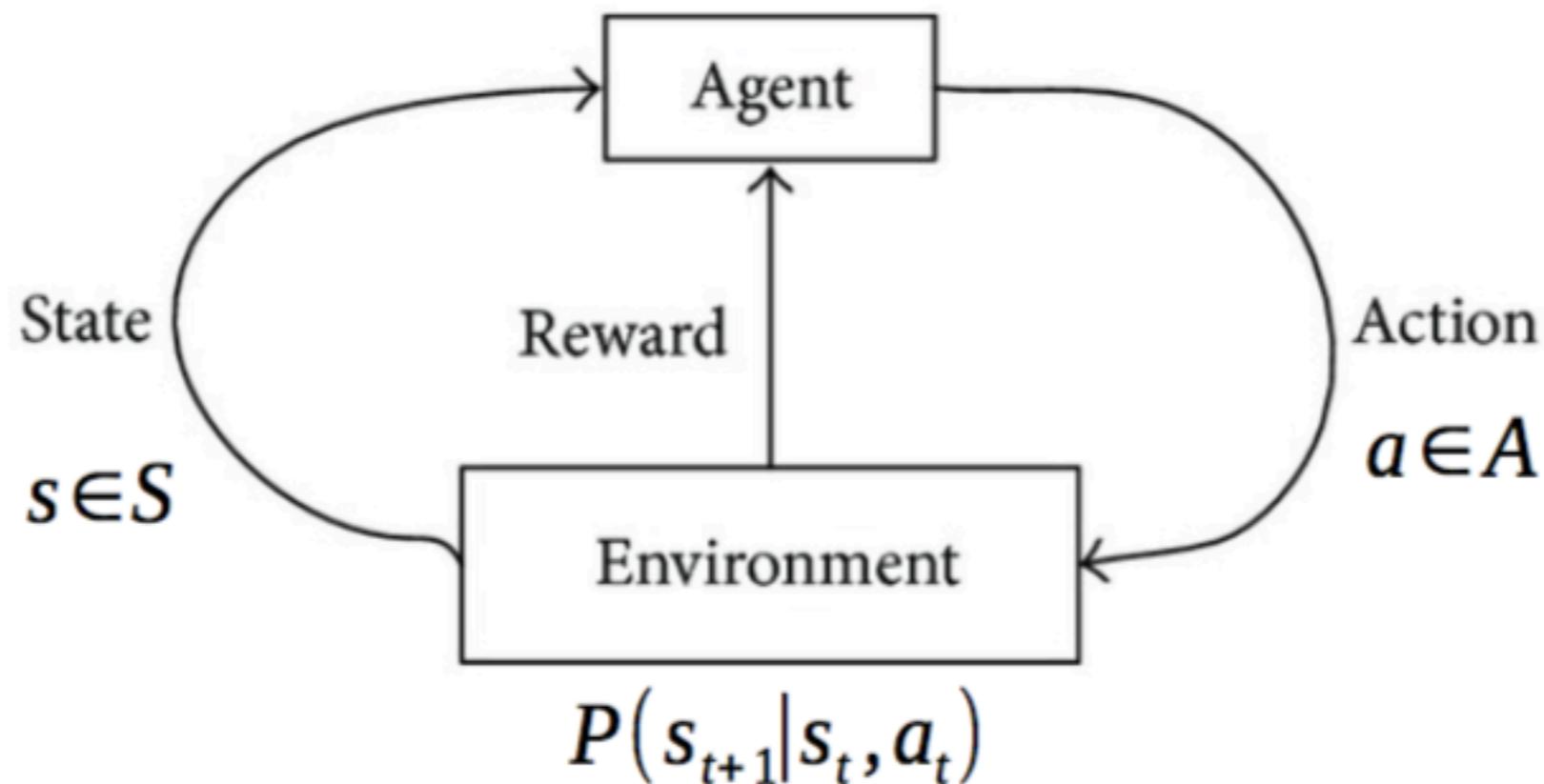
https://github.com/yandexdataschool/Practical_RL/

MDP

Markov Decision Process



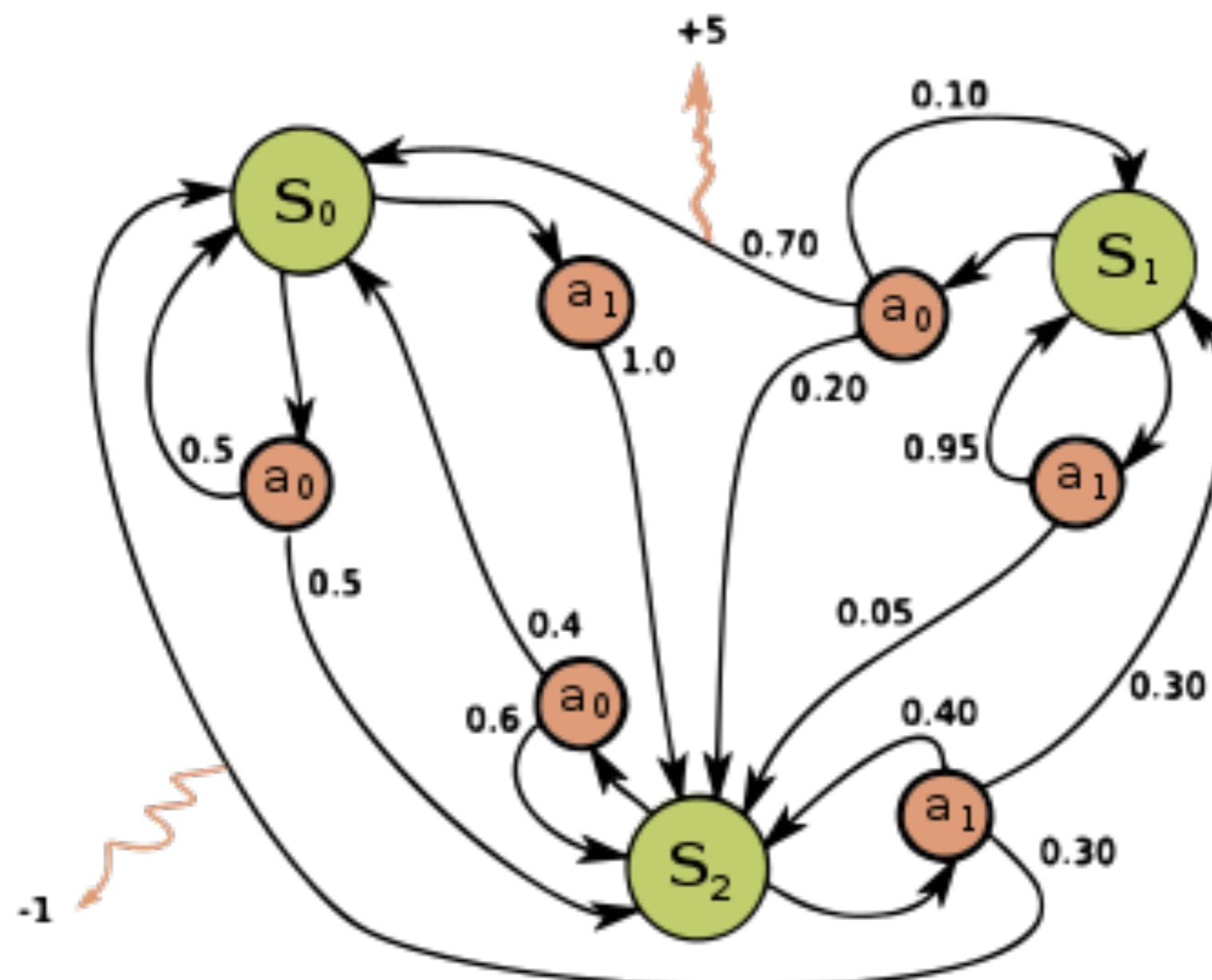
Markov Decision Process



Markov assumption

$$P(s_{t+1}|s_t, a_t, s_{t-1}, a_{t-1}) = P(s_{t+1}|s_t, a_t)$$

Markov Decision Process



$$P(s_{t+1}|s_t, a_t, s_{t-1}, a_{t-1}) = P(s_{t+1}|s_t, a_t)$$

Total reward



Total reward for session:

$$R = \sum_t r_t$$

Agent's policy:

$$\pi(a|s) = P(\text{take action } a \text{ in state } s)$$

Problem: find policy with highest reward:

$$\pi(a|s) : E_{\pi}[R] \rightarrow \max$$

Objective

The easy way:

$E_{\pi} R$ is an expected sum of rewards
that agent with policy π earns per session

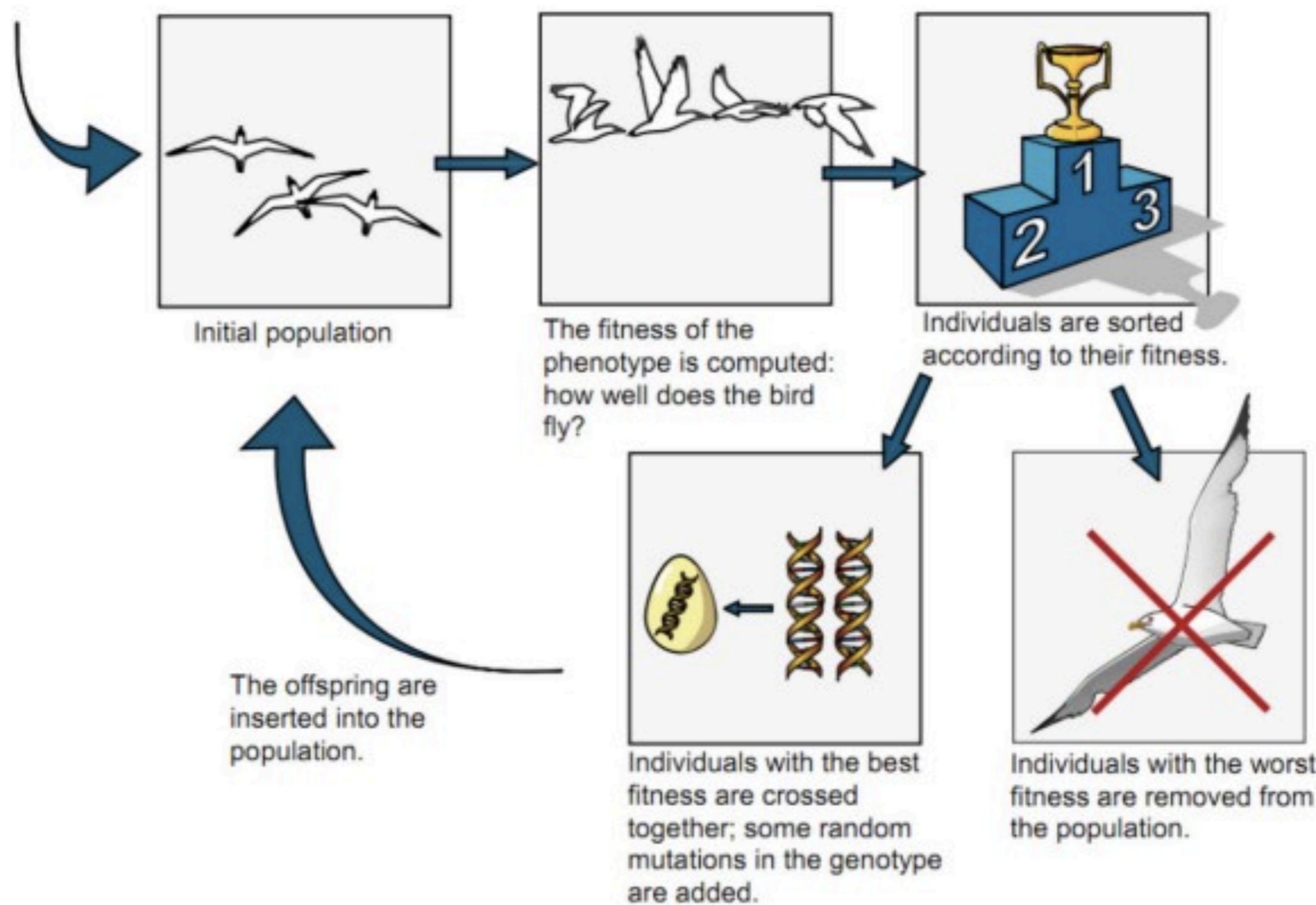
The hard way:

$$E \quad E \quad E \quad \dots \quad E \quad [r_0 + r_1 + r_2 + \dots + r_T]$$
$$s_0 \sim p(s_0), a_0 \sim \pi(a|s_0), s_1, r_0 \sim P(s', r | s, a) \quad s_T, r_T \sim P(s', r | s_{T-1}, a_{t-1})$$

How to solve?

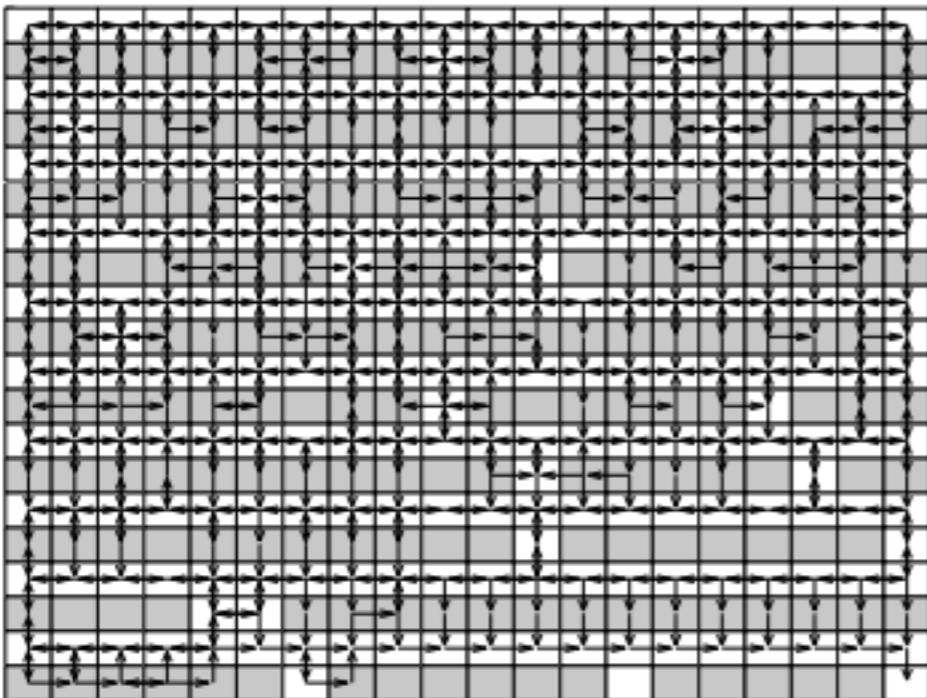
Ideas?

1. Genetic algorithms

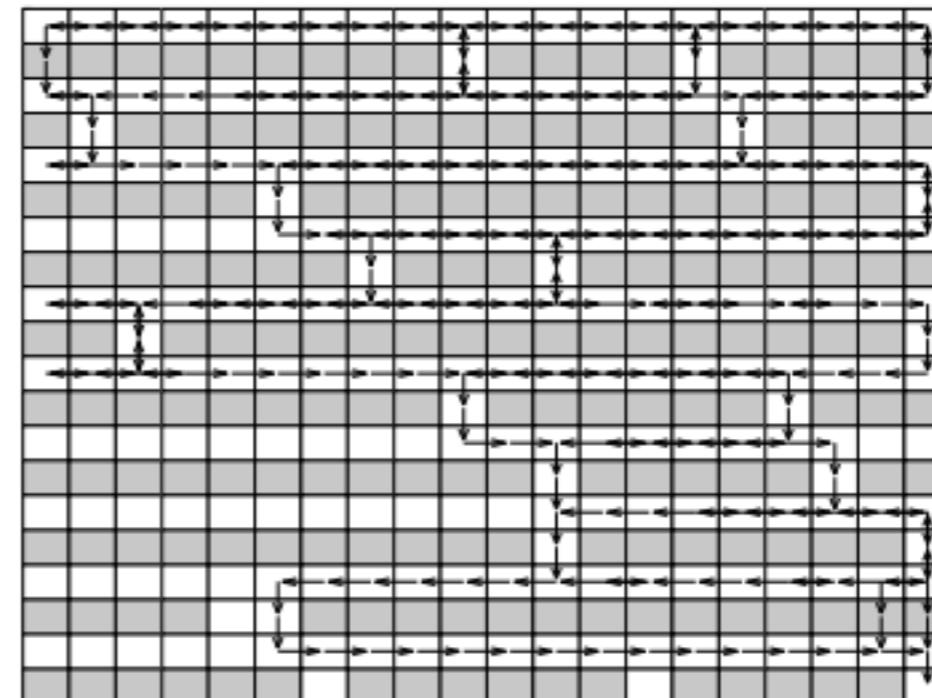


2. Cross-entropy method

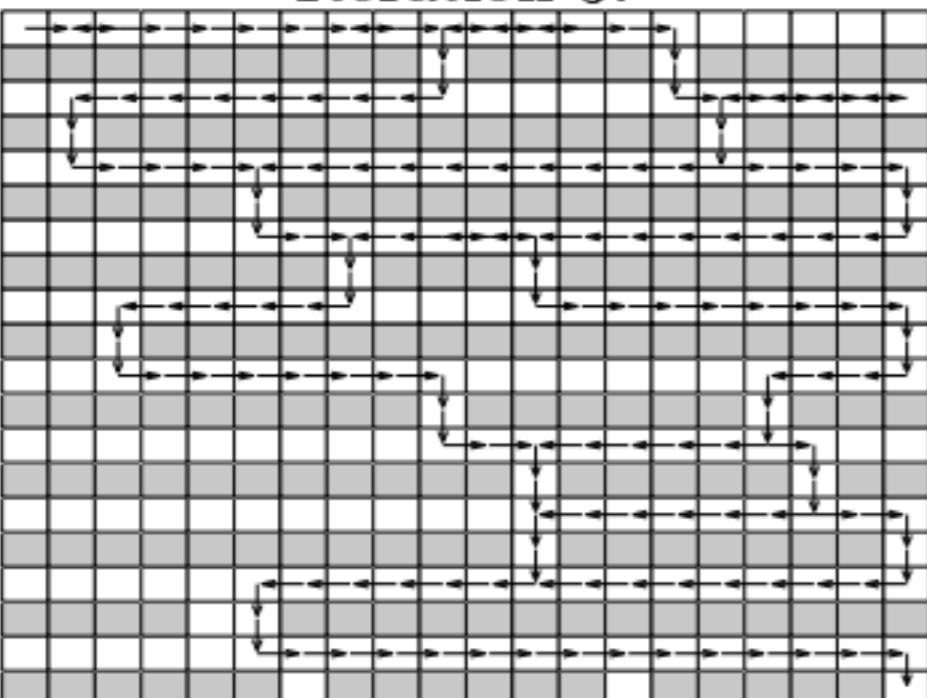
Iteration 1:



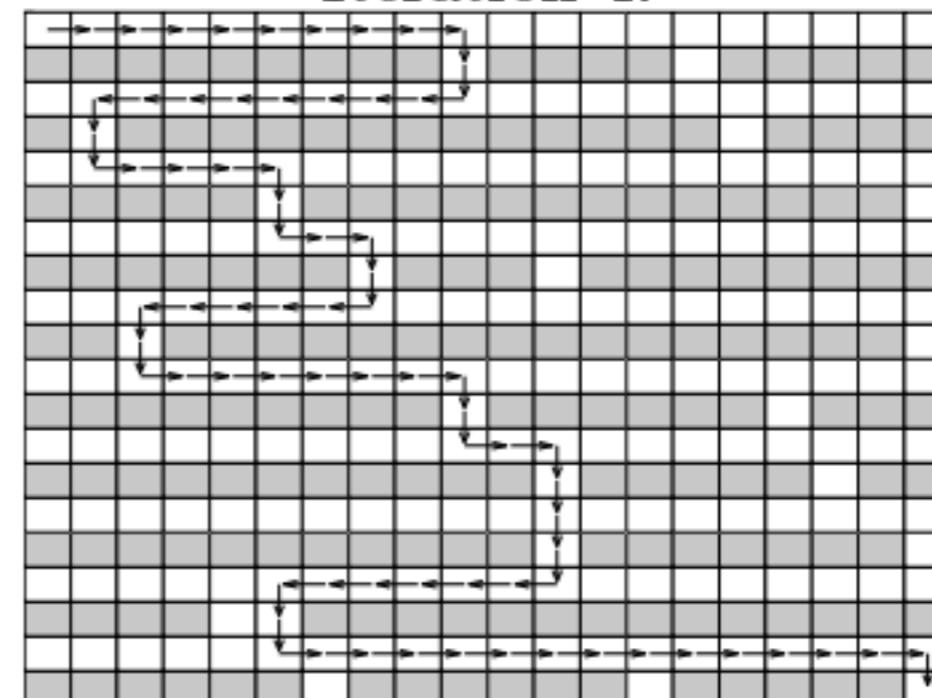
Iteration 2:



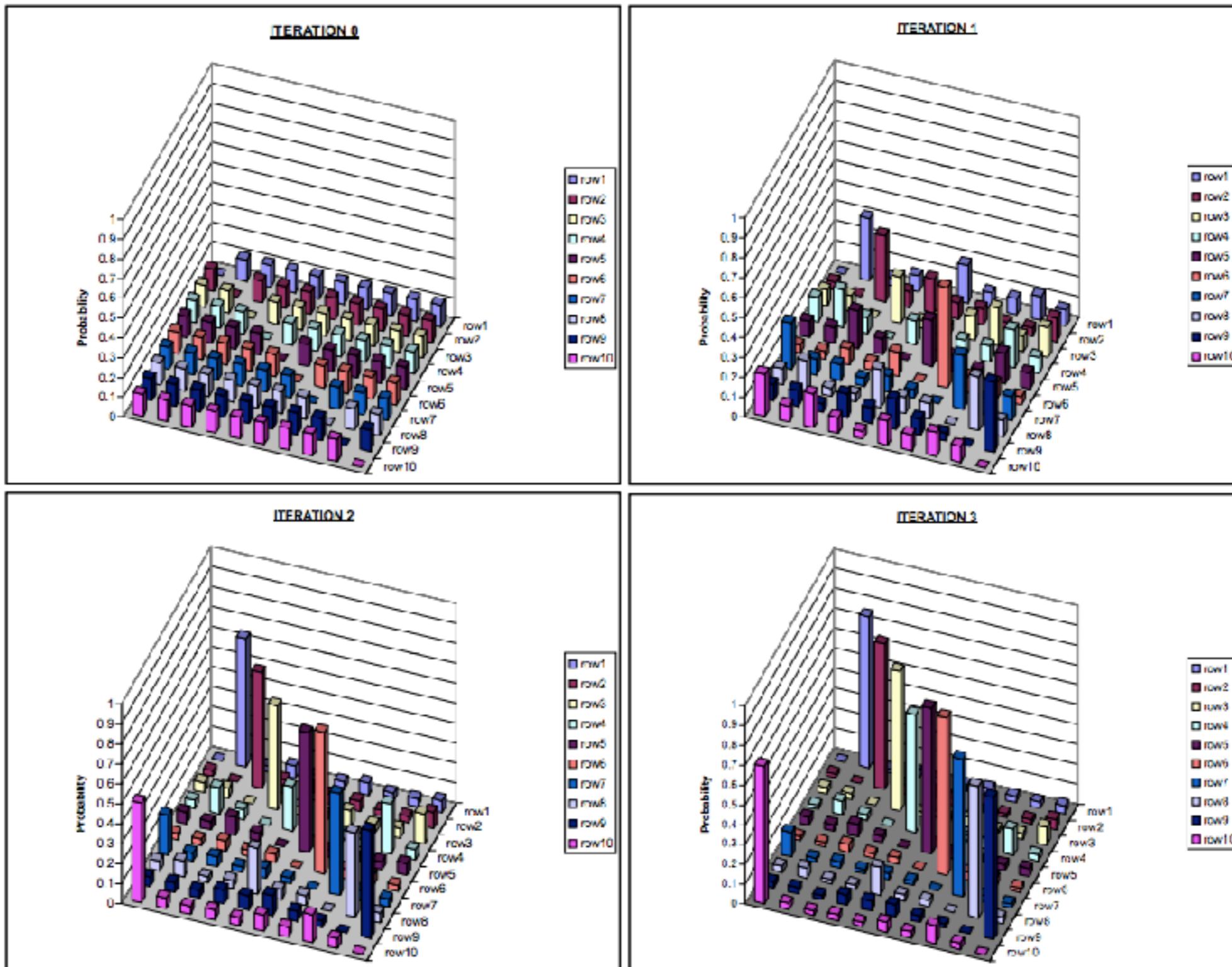
Iteration 3:



Iteration 4:



2. Cross-entropy method

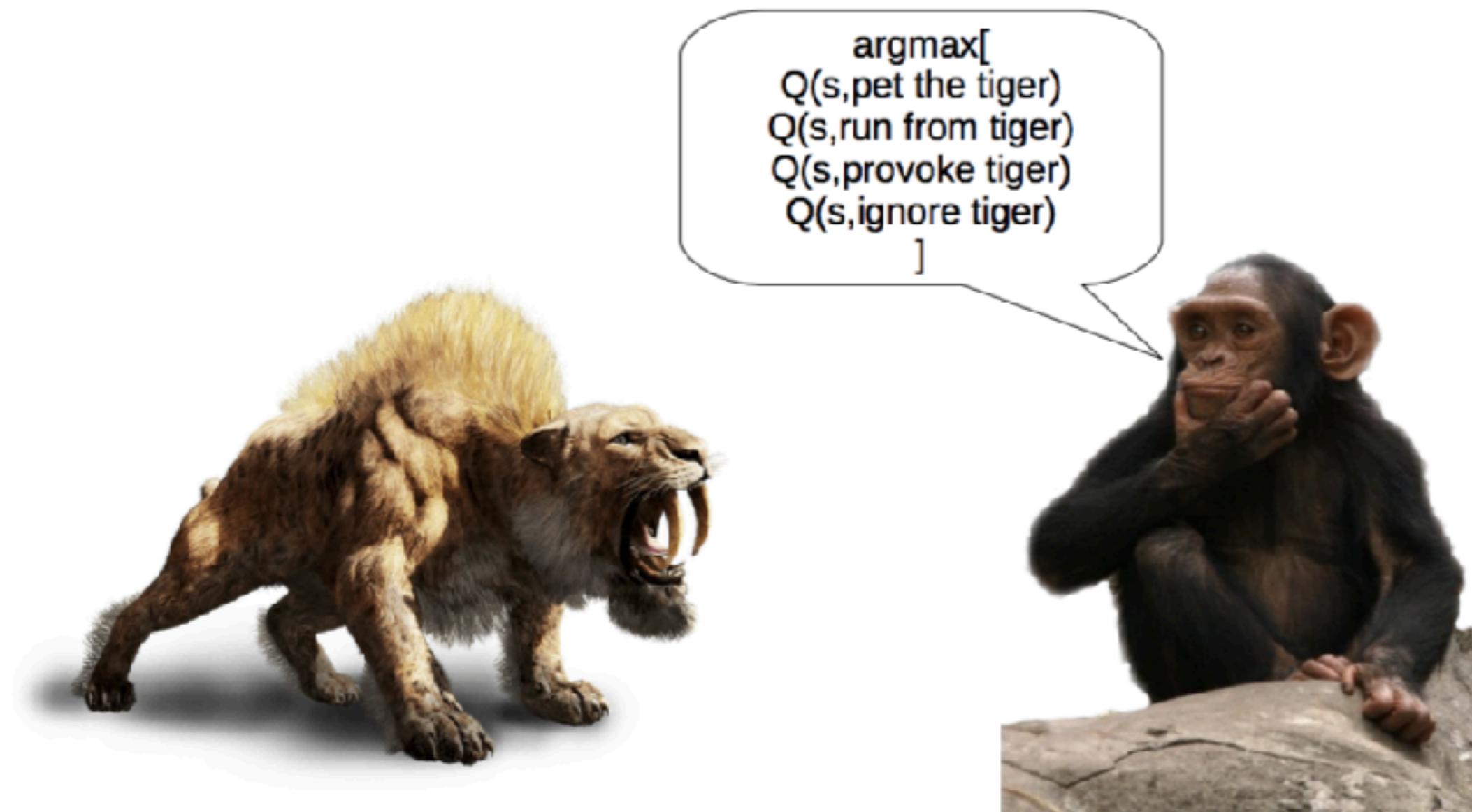


3. Q-learning



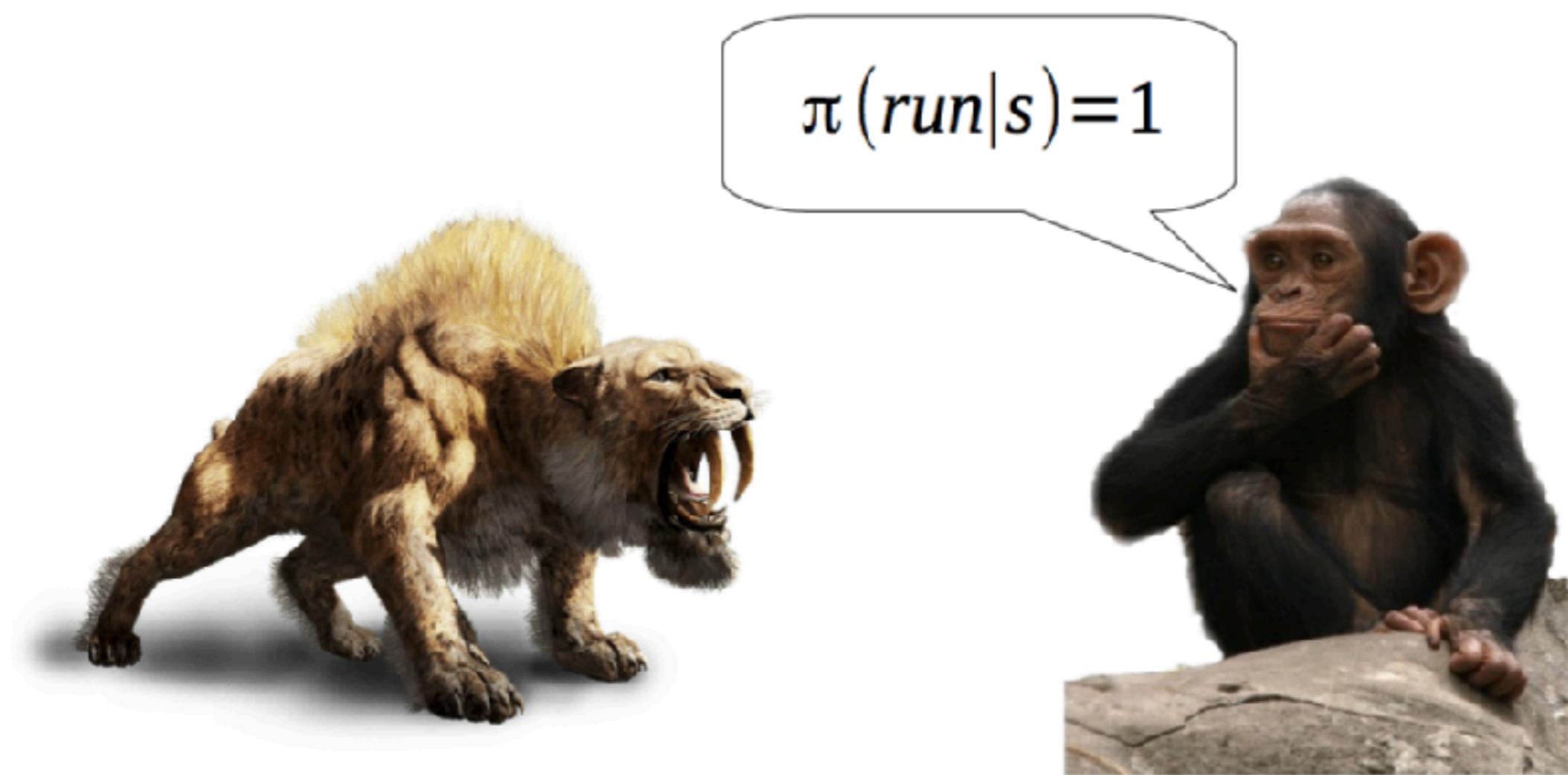
3. Policy-based

NOT how humans survived



3. Policy-based

how humans survived



Genetic algorithms

Genetic algorithms

Context: FrozenLake

- A grid world with a goal tile and ice holes

SFFF (S: starting point, safe)

FHFH (F: frozen surface, safe)

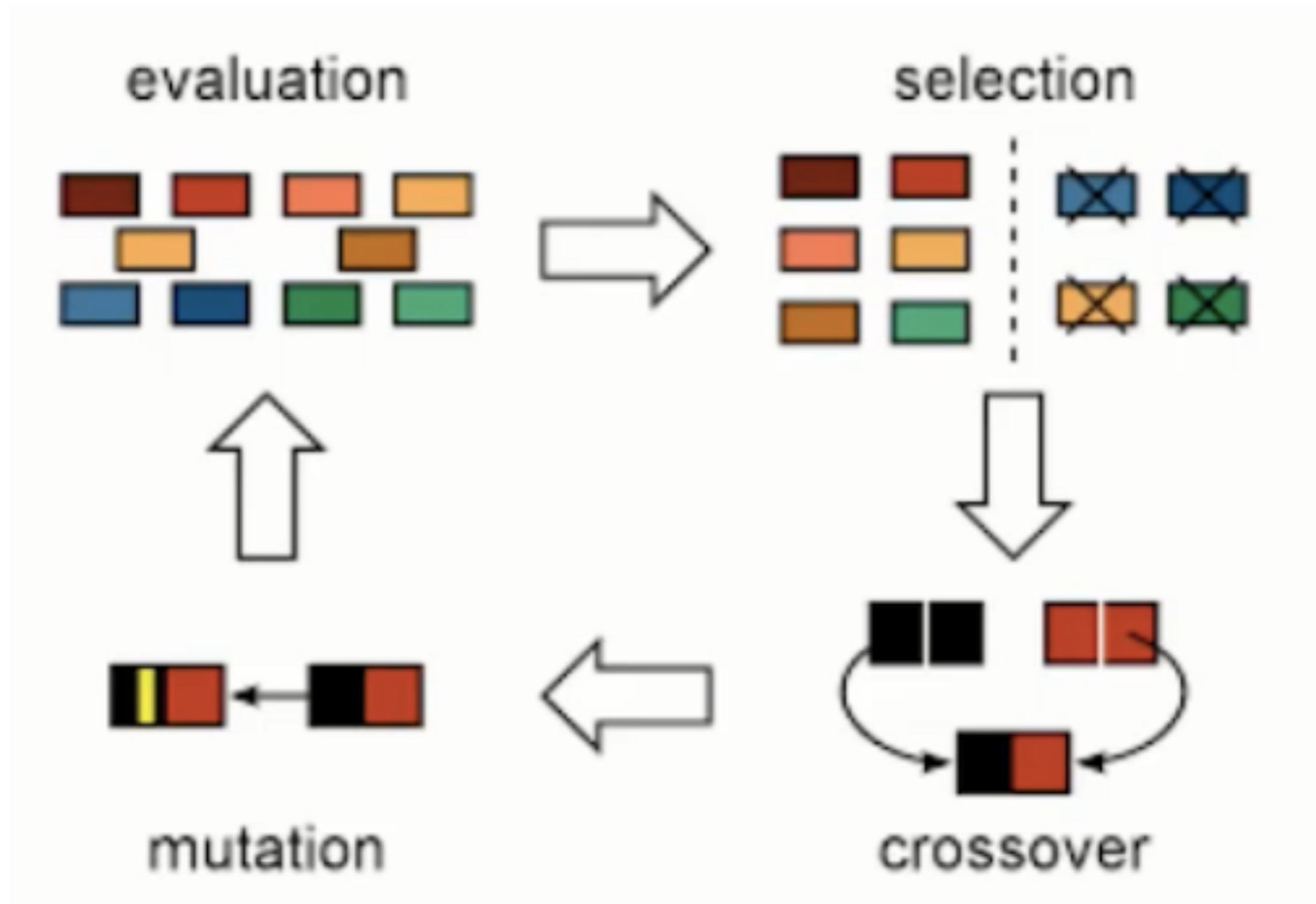
FFFH (H: hole, fall to your doom)

HFFG (G: goal, where the frisbee is located)

Quiz: what states, actions and rewards are used?

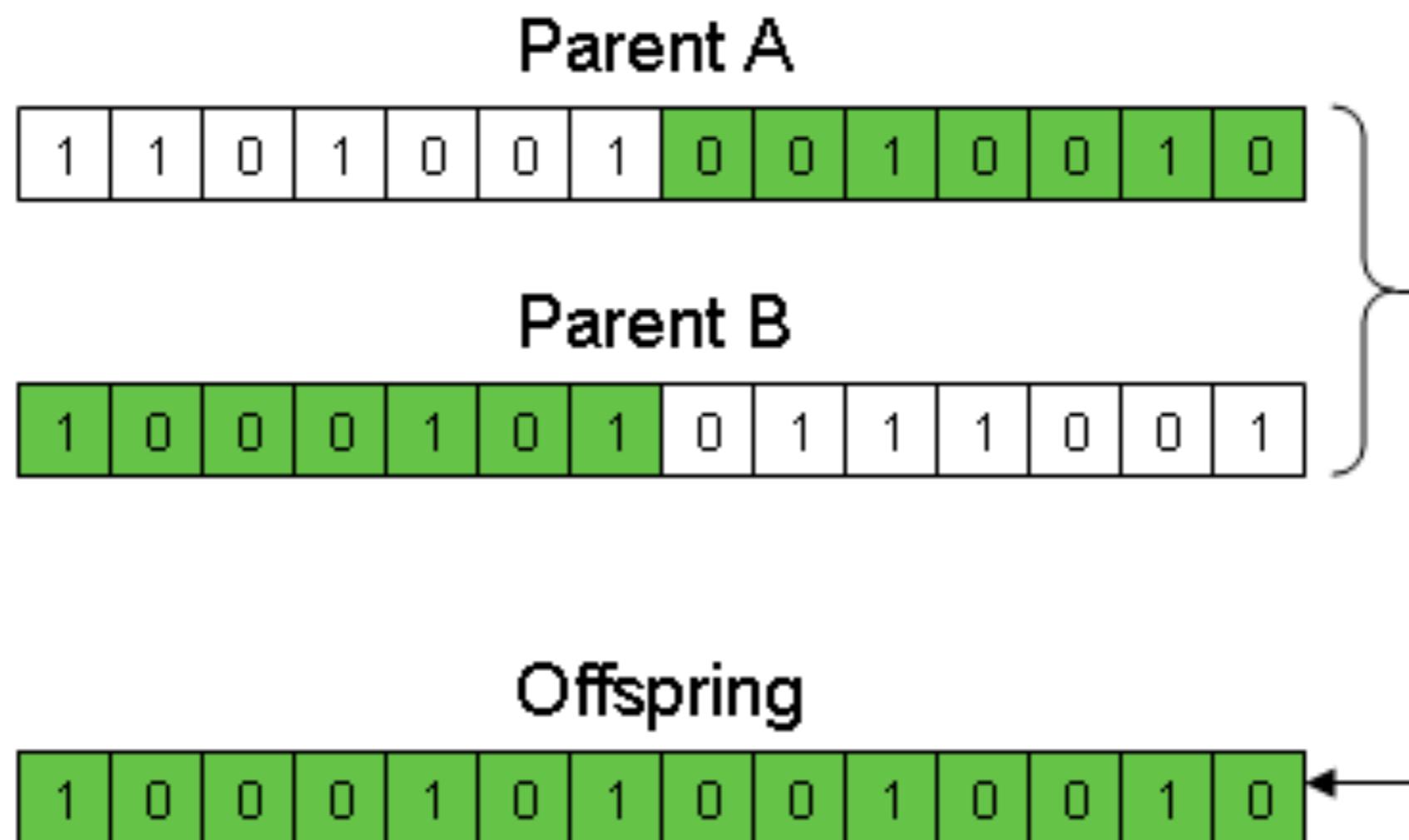


Genetic algorithms



Genetic algorithms

Cross over



Genetic algorithms

- Pros
 - It sorta sometimes works
 - You get to feel like a god!
- Cons
 - Convergence not guaranteed
 - Requires a lot of samples
 - A lot of parameter tuning
 - Hard to scale on large state spaces

Genetic algorithms

Examples:

Genetic Cars

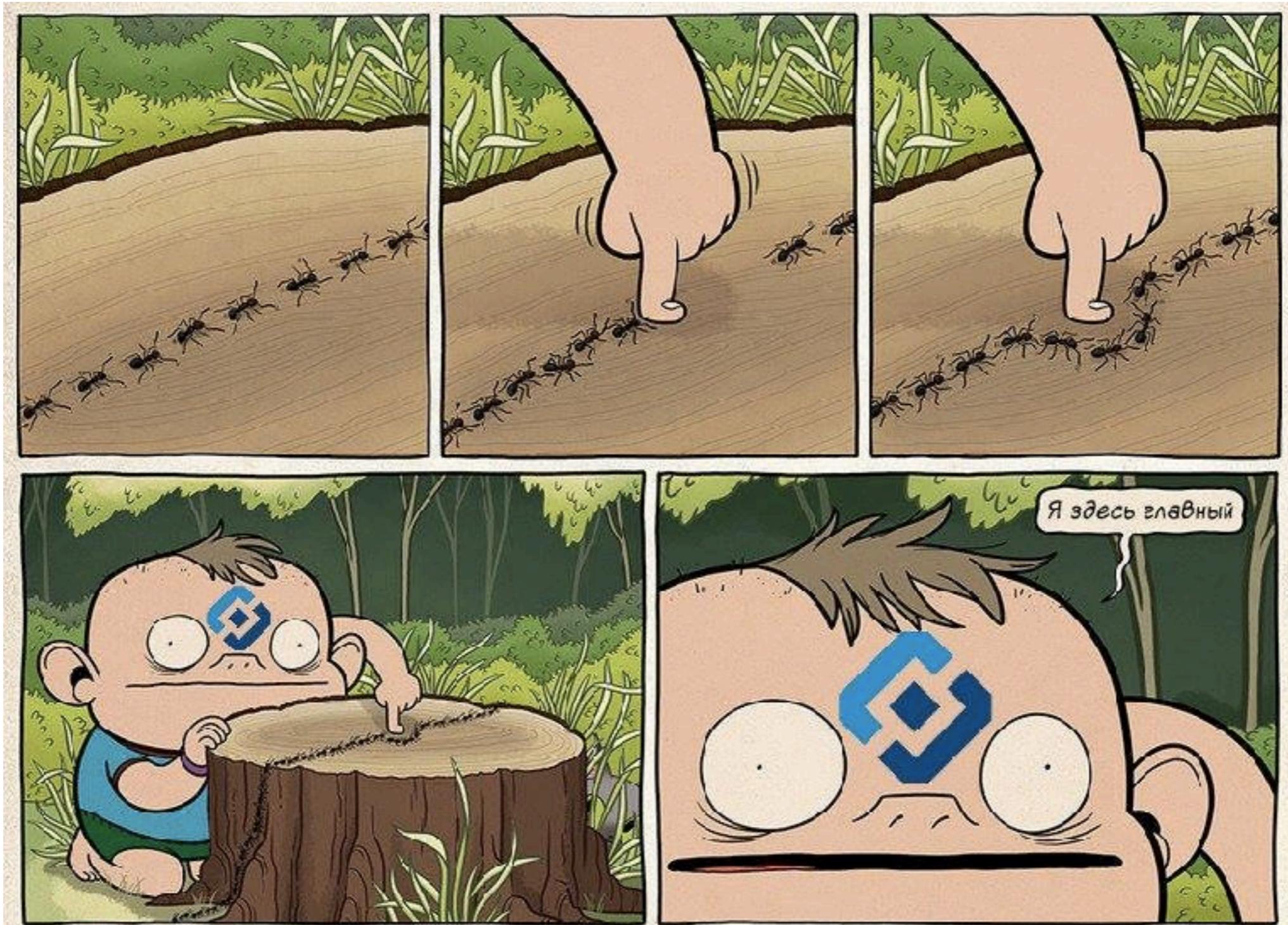
Musical Novelty Search

Stack some layers:

OpenAI paper

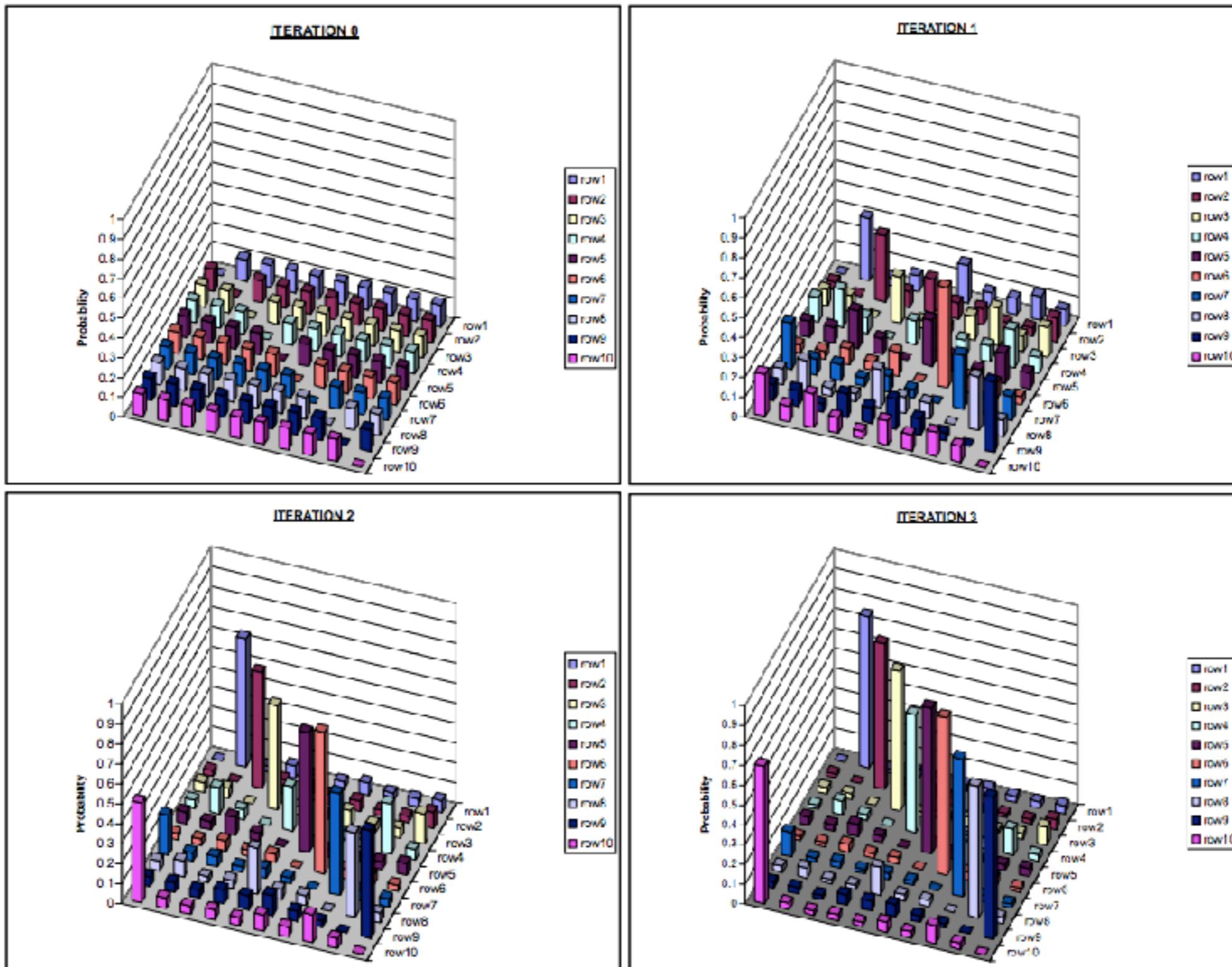
Genetic algorithms

Let's fill like a God!



Cross-entropy method

Cross-entropy method



Cross-entropy method

General idea:

Play a few sessions

Update your policy

Repeat

Cross-entropy method

Initialize policy

Repeat:

- Sample N[100] sessions
- Pick M[25] best sessions, called **elite** sessions
- Change policy so that it prioritizes actions from elite sessions

Cross-entropy method

- Policy is a matrix

$$\pi(a|s) = A_{s,a}$$

- Sample N games with that policy
- Get M best sessions (elites)

$$Elite = [(s_0, a_0), (s_1, a_1), (s_2, a_2), \dots, (s_k, a_k)]$$

Cross-entropy method

- Policy is a matrix

$$\pi(a|s) = A_{s,a}$$

- Sample N games with that policy
- Take M best sessions (elites)
- Aggregate by states

$$\pi(a|s) = \frac{\sum_{\substack{s_t, a_t \in \text{Elite}}} [s_t = s][a_t = a]}{\sum_{\substack{s_t, a_t \in \text{Elite}}} [s_t = s]}$$

Cross-entropy method

- Policy is a matrix

$$\pi(a|s) = A_{s,a}$$

- Sample N games with that policy
- Take M best sessions (elite)
- Aggregate by states

$$\pi(a|s) = \frac{\text{took } a \text{ at } s}{\text{was at } s}$$

In M best games

Cross-entropy method

Practice time!



Cross-entropy method

If your environment has infinite/large state space



Approx cross-entropy method

If your environment has infinite/large state space



Approx cross-entropy method

- Policy is approximated
 - Neural network predicts $\pi_w(a|s)$ given s
 - Linear model / Random Forest / ...

Can't set $\pi(a|s)$ explicitly

All state-action pairs from M best sessions

$$Elite = [(s_0, a_0), (s_1, a_1), (s_2, a_2), \dots, (s_k, a_k)]$$

Approx cross-entropy method

Neural network predicts $\pi_w(a|s)$ given s

All state-action pairs from M best sessions

$$Elite = [(s_0, a_0), (s_1, a_1), (s_2, a_2), \dots, (s_k, a_k)]$$

Maximize likelihood of actions in “best” games

$$\pi = \operatorname{argmax}_{\pi} \sum_{s_i, a_i \in Elite} \log \pi(a_i | s_i)$$

Approx cross-entropy method

- Initialize NN weights $W_0 \leftarrow \text{random}$
- Loop:
 - Sample N sessions
 - $\text{Elite} = [(s_0, a_0), (s_1, a_1), (s_2, a_2), \dots, (s_k, a_k)]$
 - $W_{i+1} = W_i + \alpha \nabla [\sum_{s_i, a_i \in \text{Elite}} \log \pi_{W_i}(a_i | s_i)]$

Approx cross-entropy method

- Initialize NN weights `nn = MLPClassifier(...)`
- Loop:
 - Sample N sessions
 - $Elite = [(s_0, a_0), (s_1, a_1), (s_2, a_2), \dots, (s_k, a_k)]$
 - `nn.fit(elite_states, elite_actions)`

Approx cross-entropy method

- Continuous state space
- Model $\pi_w(a|s) = N(\mu(s), \sigma^2)$
 - Mu(s) is neural network output
 - Sigma is a parameter or yet another network output
- Loop:
 - Sample N sessions
 - elite = take M best sessions and concatenate
 - $W_{i+1} = W_i + \alpha \nabla \left[\sum_{s_i, a_i \in Elite} \log \pi_{w_i}(a_i | s_i) \right]$

Approx cross-entropy method

- Initialize NN weights `nn = MLPRegressor(...)`
- Loop:
 - Sample N sessions
 - $Elite = [(s_0, a_0), (s_1, a_1), (s_2, a_2), \dots, (s_k, a_k)]$
 - `nn.fit(elite_states, elite_actions)`

Approx cross-entropy method

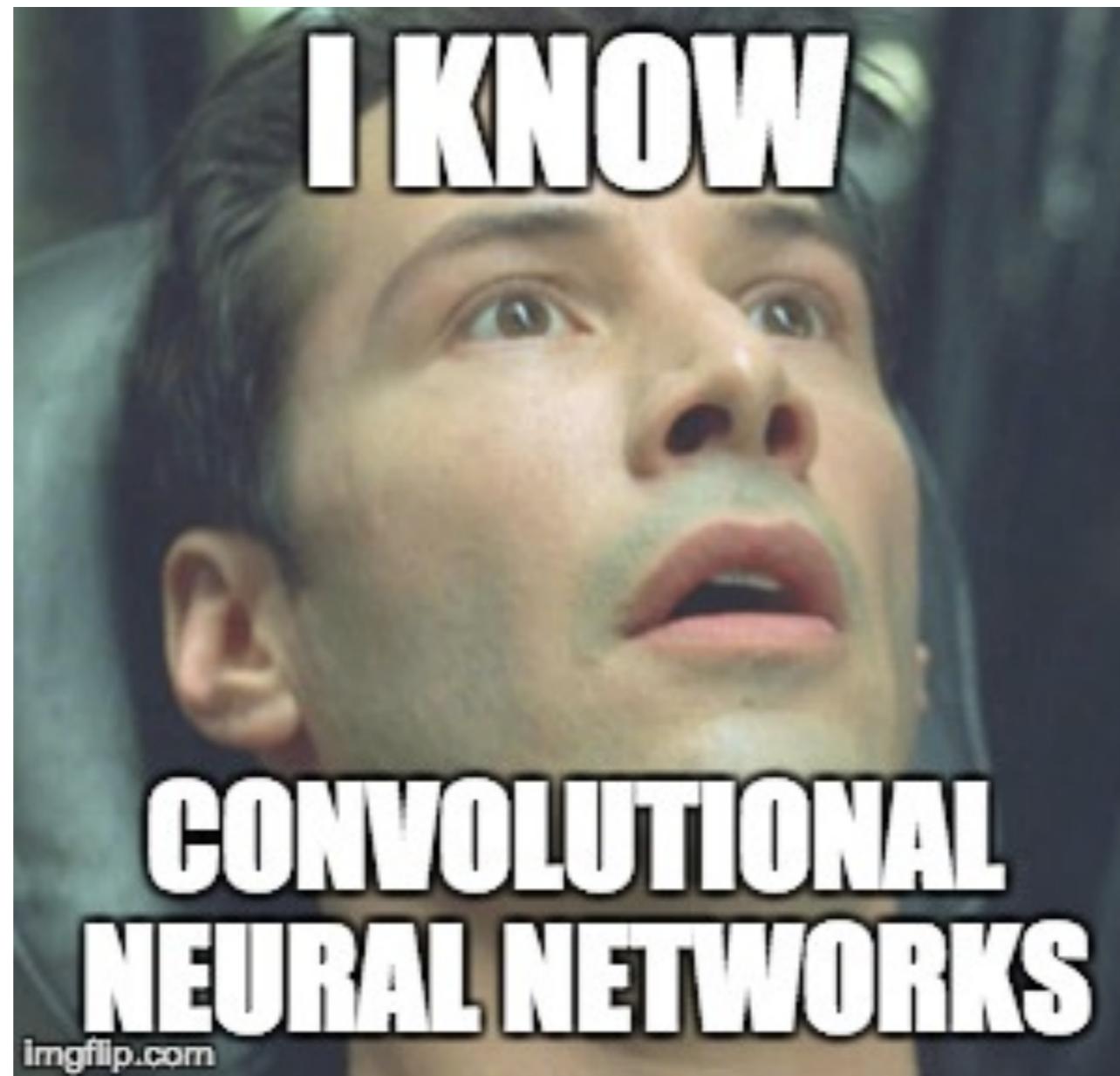
Tricks

- Remember sessions from 3-5 past iterations
 - Threshold and use all of them when training
 - May converge slower if env is easy to solve.
- Regularize with entropy
 - to prevent premature convergence.
- Parallelize sampling
- Use RNNs if partially-observable (later)



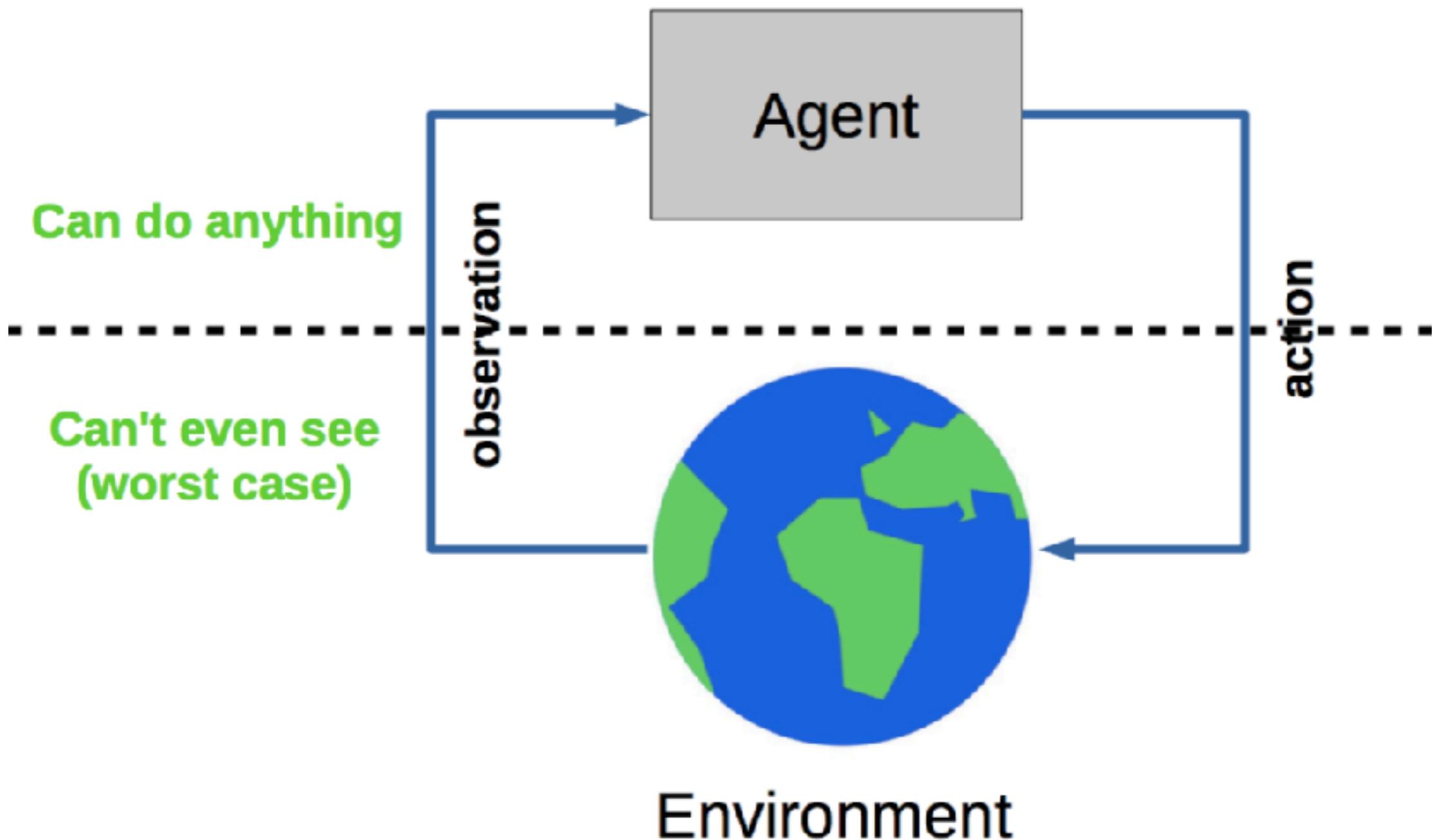
Approx cross-entropy method

Practice time!



Q-learning

Earlier:



Earlier:

- $R(z)$ – evaluated at the very end
- Metaheuristics (genetic algorithm, etc.)
- Stochastic optimization (crossentropy method)

Earlier:

- Both need a full session to start learning
- Requires a lot of interaction
 - A lot of crashed robots / simulations



Now:

MDP formalism: reward on each tick

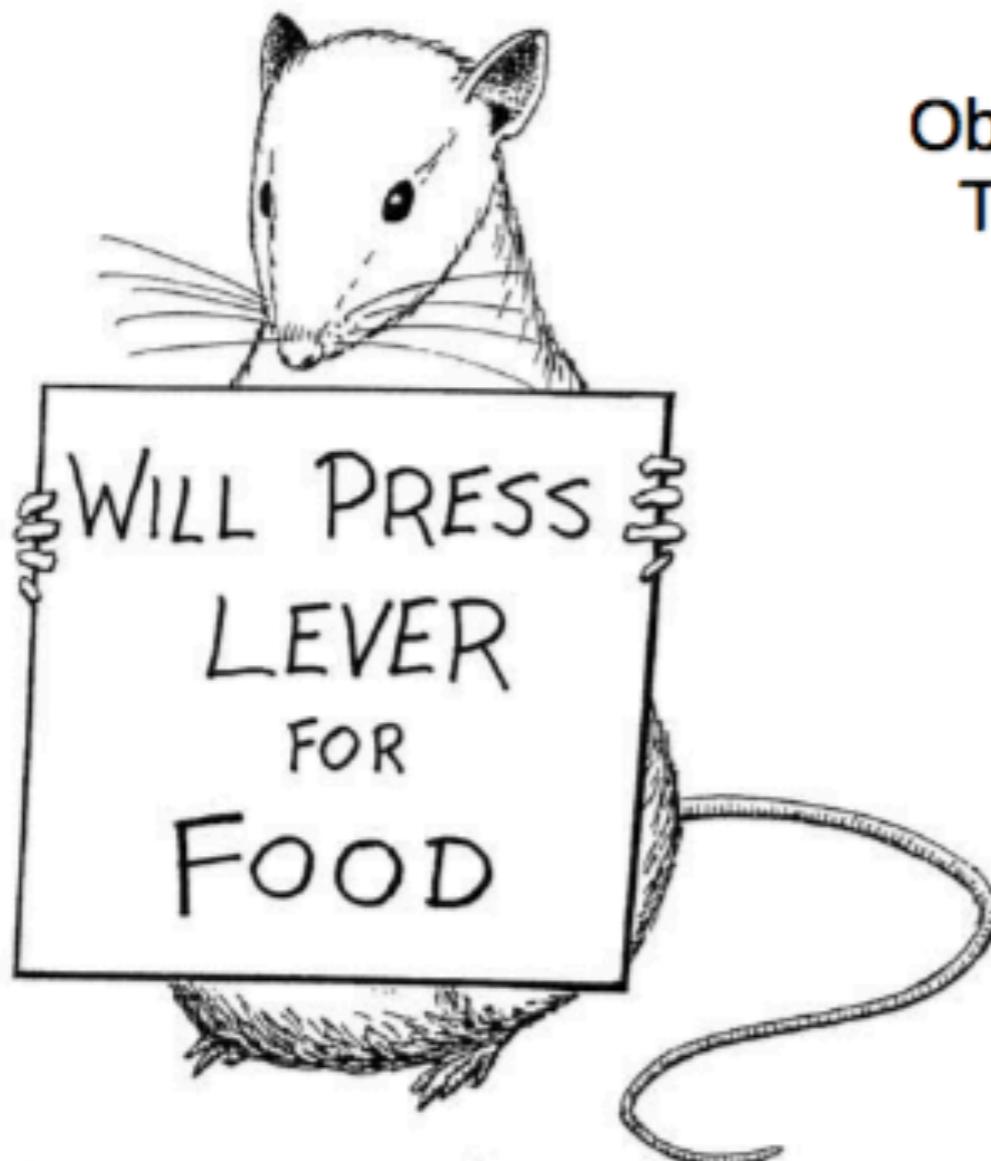


Classic MDP(Markov Decision Process)

Agent interacts with environment

- Environment states: $s \in S$
- Agent actions: $a \in A$
- State transition: $P(s_{t+1}|s_t, a_t)$
- Reward: $r_t = r(s_t, a_t)$

Discounted reward MDP



Objective:
Total action value

$$R_t = r_t + \gamma \cdot r_{t+1} + \gamma^2 \cdot r_{t+2} + \dots + \gamma^n \cdot r_{t+n}$$

$$R_t = \sum_i \gamma^i \cdot r_{t+i} \quad \gamma \in (0,1) \text{ const}$$

$\gamma \sim$ patience

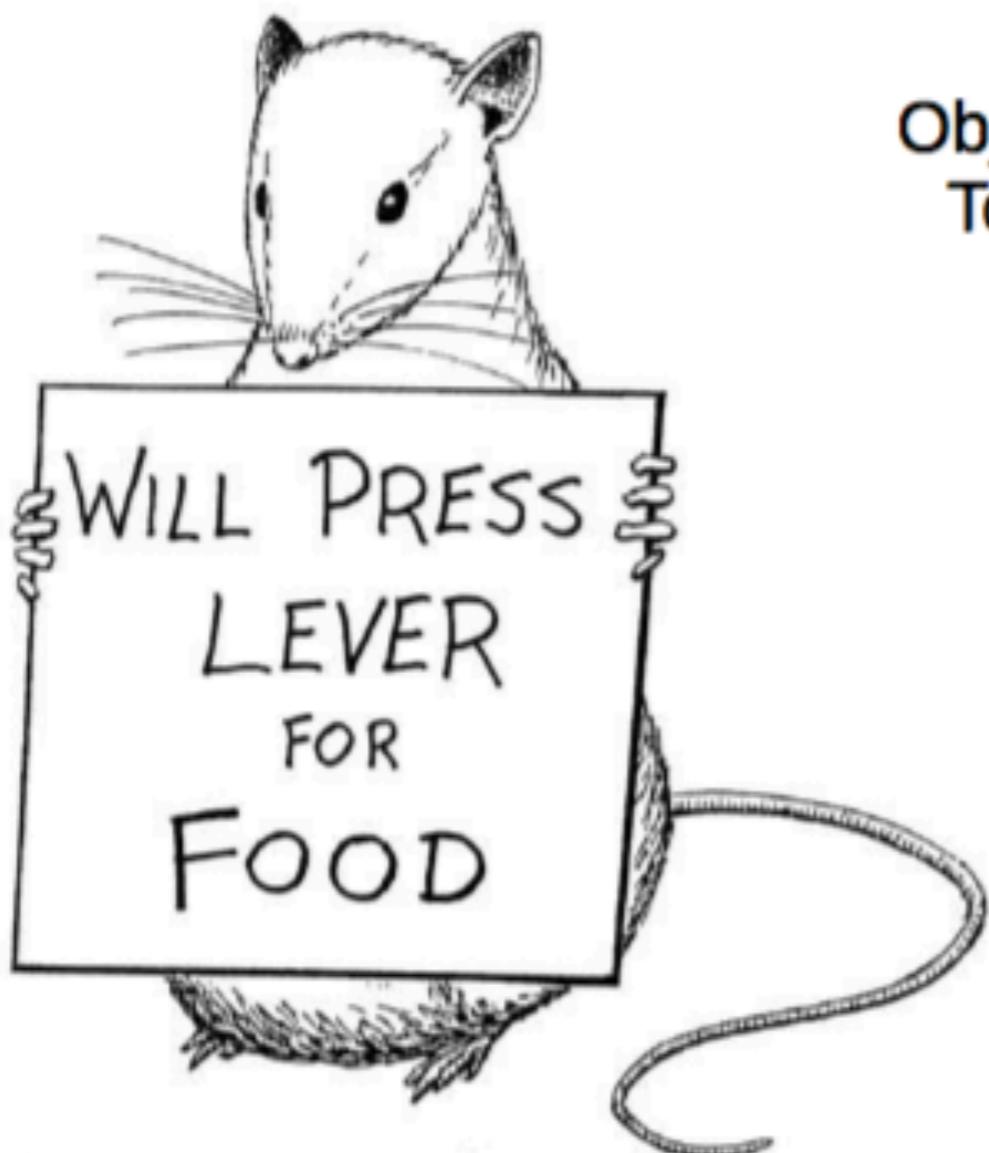
Cake tomorrow is γ as good as now

Reinforcement learning:

- Find policy that maximizes the expected reward

$$\pi = P(a|s): E[R] \rightarrow \max$$

Discounted reward MDP



Objective:
Total action value

$$R_t = r_t + \gamma \cdot r_{t+1} + \gamma^2 \cdot r_{t+2} + \dots + \gamma^n \cdot r_{t+n}$$

$$R_t = \sum_i \gamma^i \cdot r_{t+i} \quad \gamma \in (0,1) \text{ const}$$

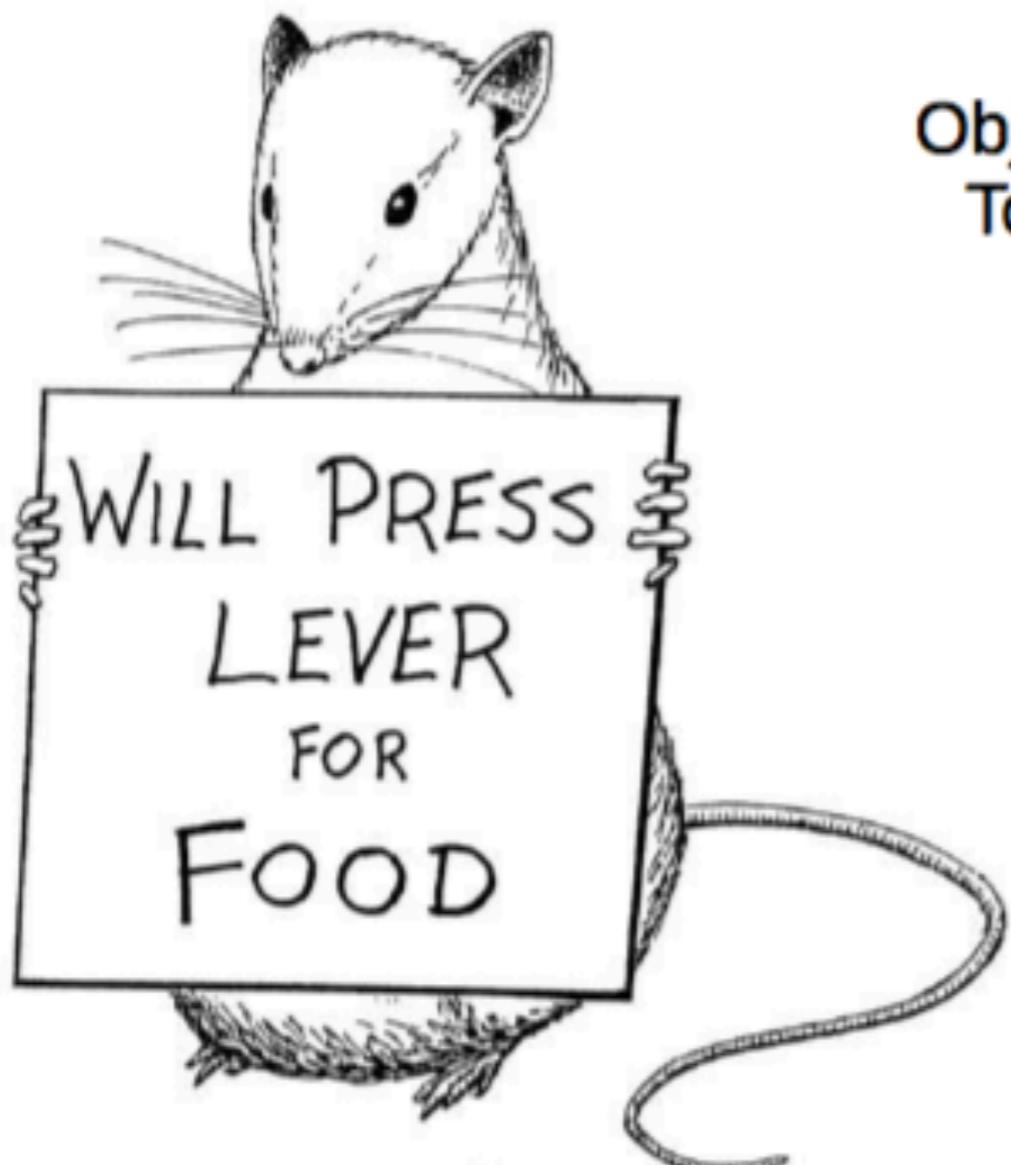
Trivia: which γ corresponds to
“only current reward matters”?

Reinforcement learning:

- Find policy that maximizes the expected reward

$$\pi = P(a|s) : E[R] \rightarrow \max$$

Discounted reward MDP



Objective:
Total action value

$$R_t = r_t + \gamma \cdot r_{t+1} + \gamma^2 \cdot r_{t+2} + \dots + \gamma^n \cdot r_{t+n}$$

$$R_t = \sum_i \gamma^i \cdot r_{t+i} \quad \gamma \in (0,1) \text{ const}$$

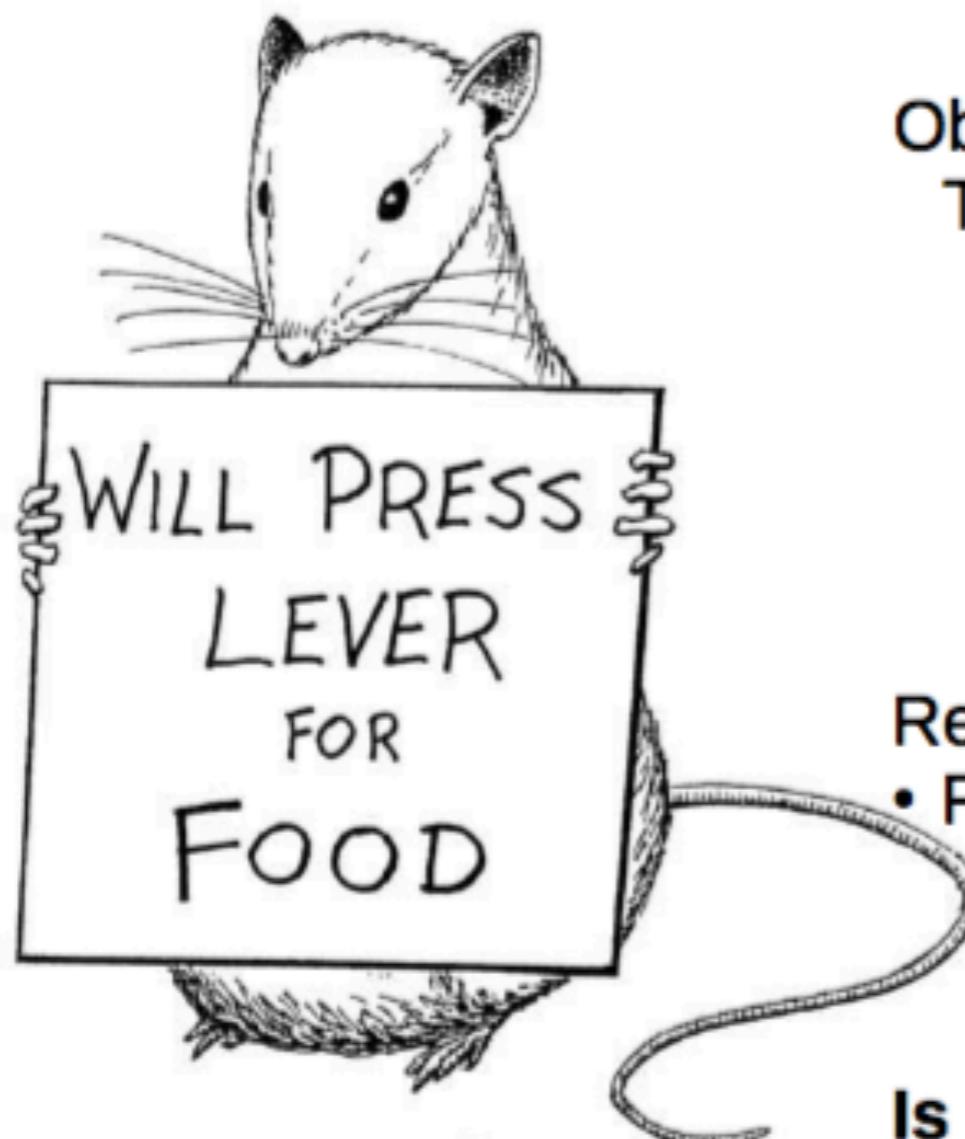
Trivia: which γ corresponds to
“only current reward matters”?

Reinforcement learning:

- Find policy that maximizes the expected reward

$$\pi = P(a|s): E[R] \rightarrow \max$$

Discounted reward MDP



Objective:
Total reward

$$R_t = r_t + \gamma \cdot r_{t+1} + \gamma^2 \cdot r_{t+2} + \dots + \gamma^n \cdot r_{t+n}$$

$$R_t = \sum_i \gamma^i \cdot r_{t+i} \quad \gamma \in (0,1) \text{ const}$$

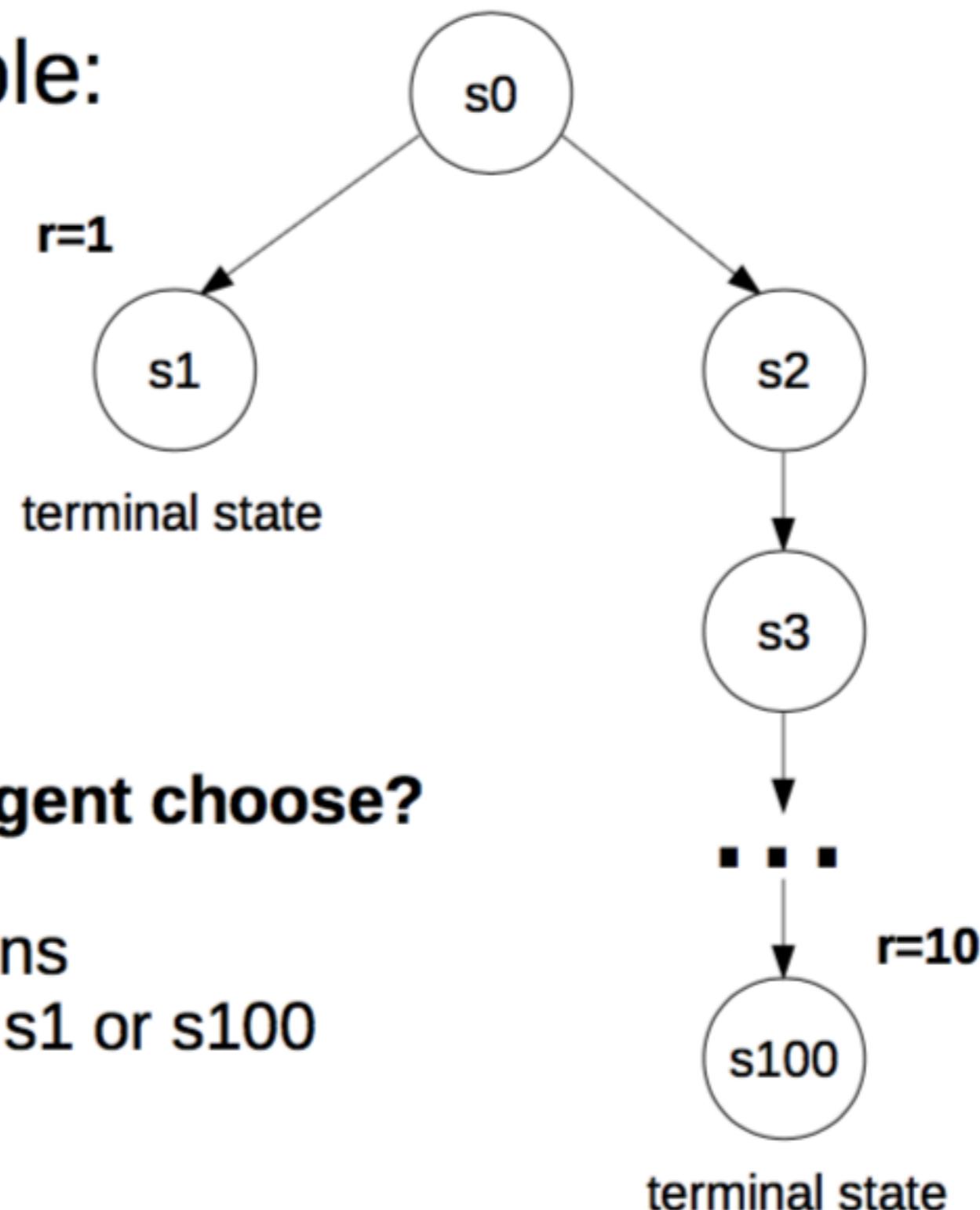
Reinforcement learning:
• Find policy that maximizes
the expected reward

$$\pi = P(a|s) : E[R] \rightarrow \max$$

Is optimal policy same as it would be
in monte-carlo (if we add-up all r_t)?

Discounted reward **fails** #1

Trivial example:

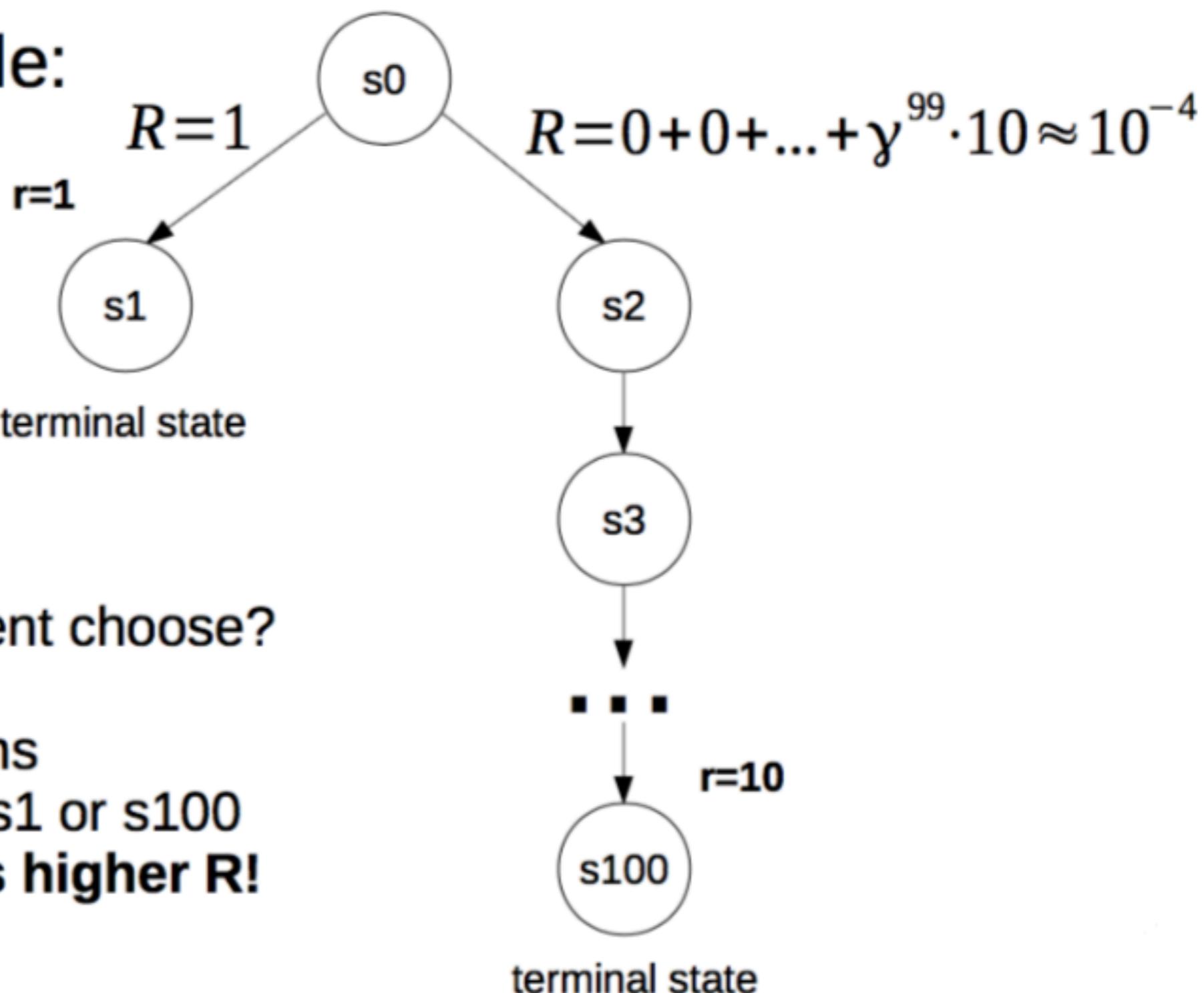


What path will agent choose?

- $\gamma=0.9$
- arrows = actions
- game ends at s_1 or s_{100}

Discounted reward **fails** #1

Trivial example:



What path will agent choose?

- $\gamma=0.9$
- arrows = actions
- game ends at s_1 or s_{100}
- **left action has higher R!**

Discounted reward **fails** #2

Deephack'17 qualification round, Atari Skiing



- You steer the red guy
- Session lasts ~5k steps
- You get -3~-7 reward each tick
(faster game = better score)
- At the end of session, you get up to $r=-30k$
(based on passing gates, etc.)
- Q-learning with $\gamma=0.99$ fails
it doesn't learn to pass gates

What's the problem?

<http://rl.deephack.me/>

Discounted reward **fails** #3

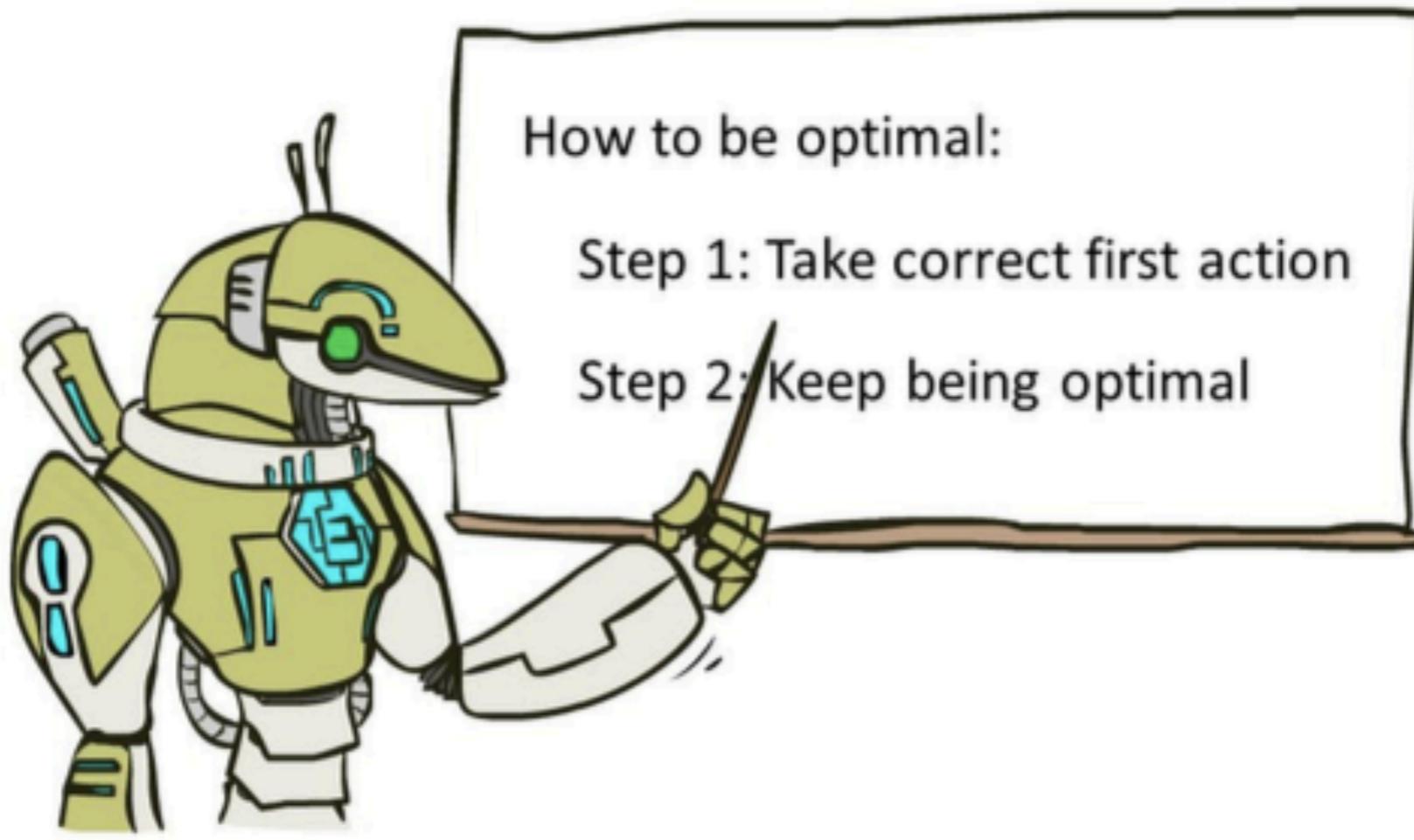
CoastRunner7 experiment (openAI)



- You control the boat
- Rewards for getting to checkpoints
- Rewards for collecting bonuses
- What could possibly go wrong?
- “Optimal” policy video:

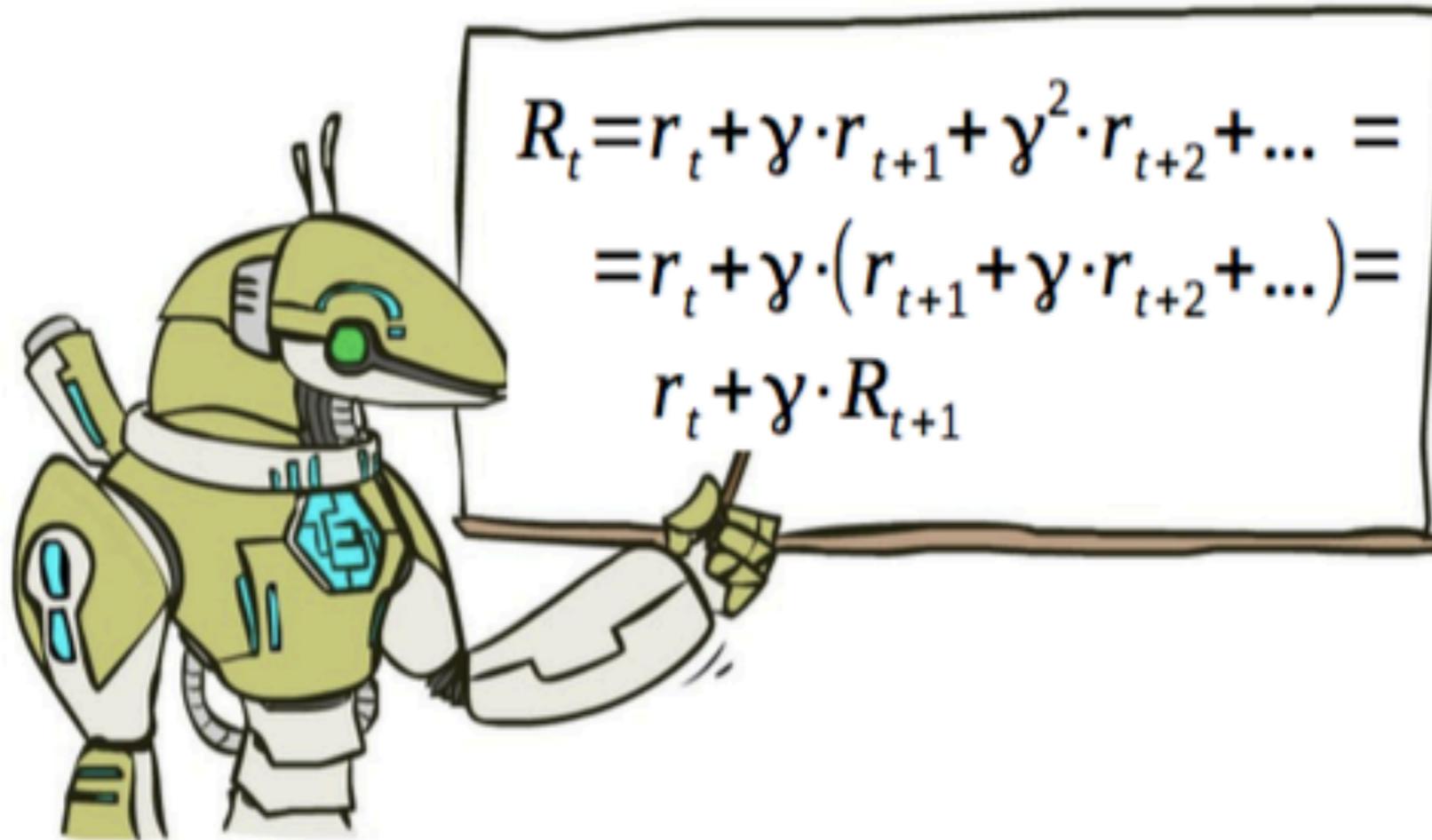
<https://openai.com/blog/faulty-reward-functions/>

Optimal policy



Recurrent optimal strategy definition

Optimal policy



We rewrite R with sheer power of math!

Value iteration (Temporal Difference)

Idea:

- For each state, obtain $V(\text{state})$

$$V(s) = \max_a [r(s, a) + \gamma \cdot V(s'(s, a))]$$

Definition $V(s)$ – expected total reward R that can be obtained starting from state s under optimal policy

Value iteration (Temporal Difference)

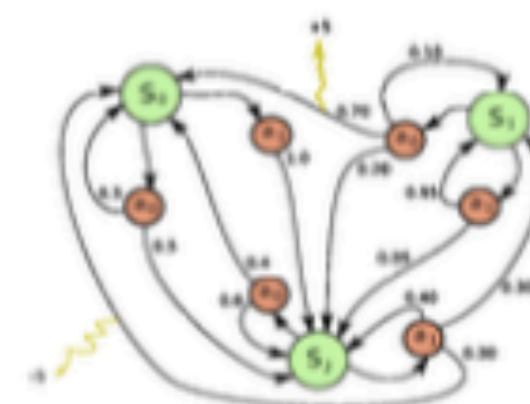
Idea:

- For each state, obtain $V(\text{state})$

$$V(s) = \max_a [r(s, a) + \gamma \cdot E_{s' \sim P(s'|s, a)} V(s')]$$

↑
Stochastic
action outcome

Trivia: if we know the exact $V(s)$ for all states,
how do we determine the best actions?



Value iteration (TD)

Idea:

- Iterative updates

$$\forall s, V_0(s) := 0$$

$$V_{i+1}(s) := \max_a [r(s, a) + \gamma \cdot E_{s' \sim P(s'|s, a)} V_i(s')]$$

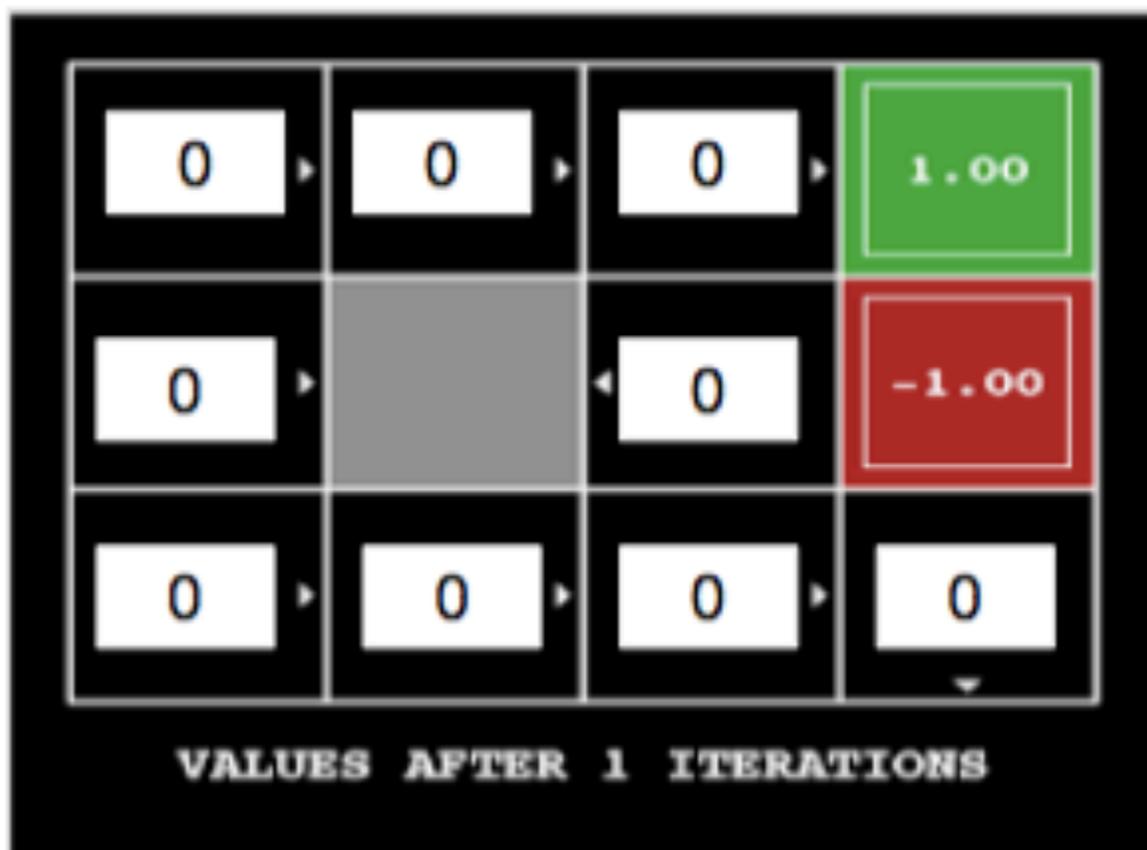
Value iteration (TD)

Idea:

- Iterative updates

$$\forall s, V_0(s) := 0$$

$$V_{i+1}(s) := \max_a [r(s, a) + \gamma \cdot E_{s' \sim P(s'|s, a)} V_i(s')]$$



Value iteration (TD)

Idea:

- Iterative updates

$$\forall s, V_0(s) := 0$$

$$V_{i+1}(s) := \max_a [r(s, a) + \gamma \cdot E_{s' \sim P(s'|s, a)} V_i(s')]$$

0	0	γ	1.00
0		0	-1.00
0	0	0	0

VALUES AFTER 1 ITERATIONS

Value iteration (TD)

Idea:

- Iterative updates

$$\forall s, V_0(s) := 0$$

$$V_{i+1}(s) := \max_a [r(s, a) + \gamma \cdot E_{s' \sim P(s'|s, a)} V_i(s')]$$

0	γ^2	γ	1.00
0		γ^2	-1.00
0	0	0	0

VALUES AFTER 1 ITERATIONS

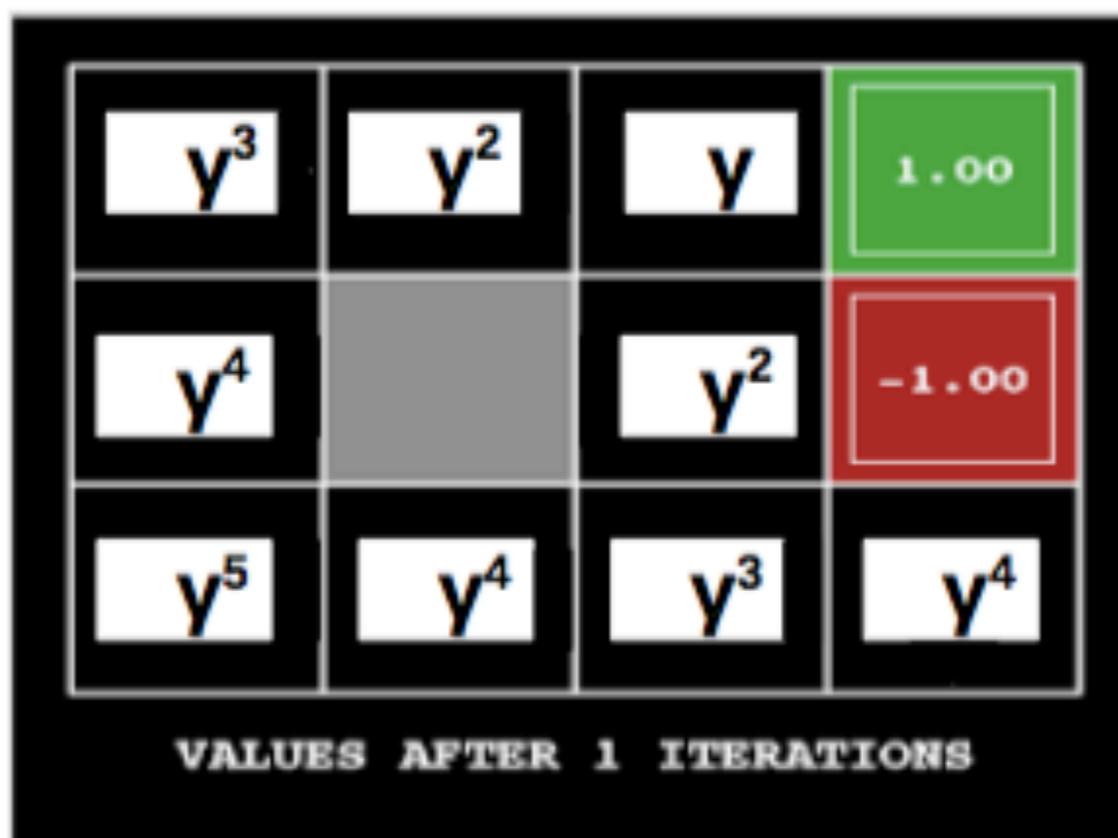
Value iteration (TD)

Idea:

- Iterative updates

$$\forall s, V_0(s) := 0$$

$$V_{i+1}(s) := \max_a [r(s, a) + \gamma \cdot E_{s' \sim P(s'|s, a)} V_i(s')]$$



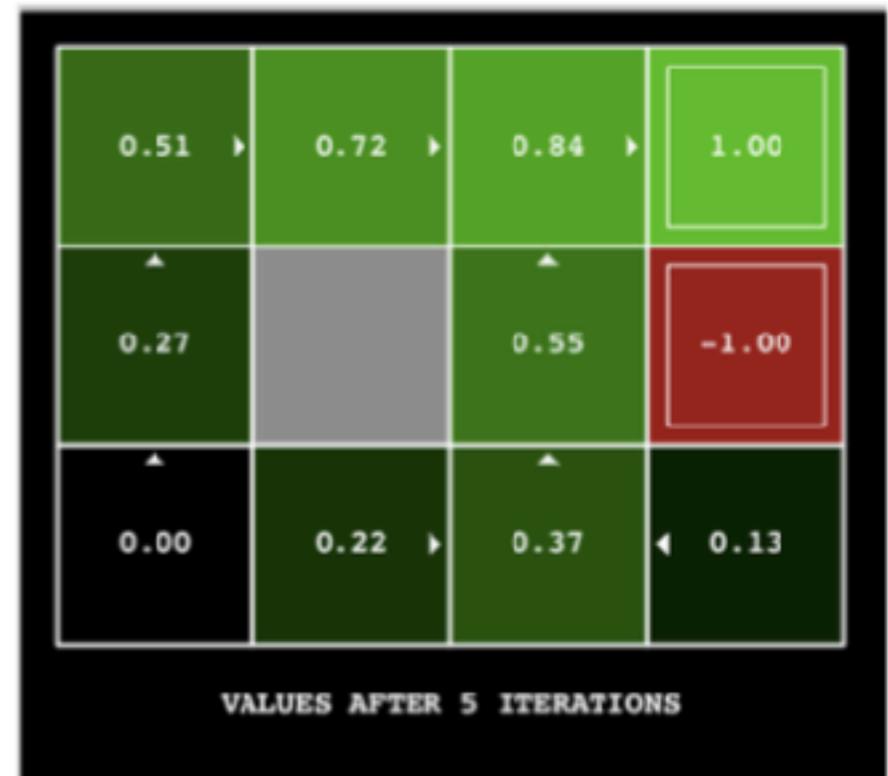
Value iteration (TD)

Idea:

- Iterative updates

$$\forall s, V_0(s) := 0$$

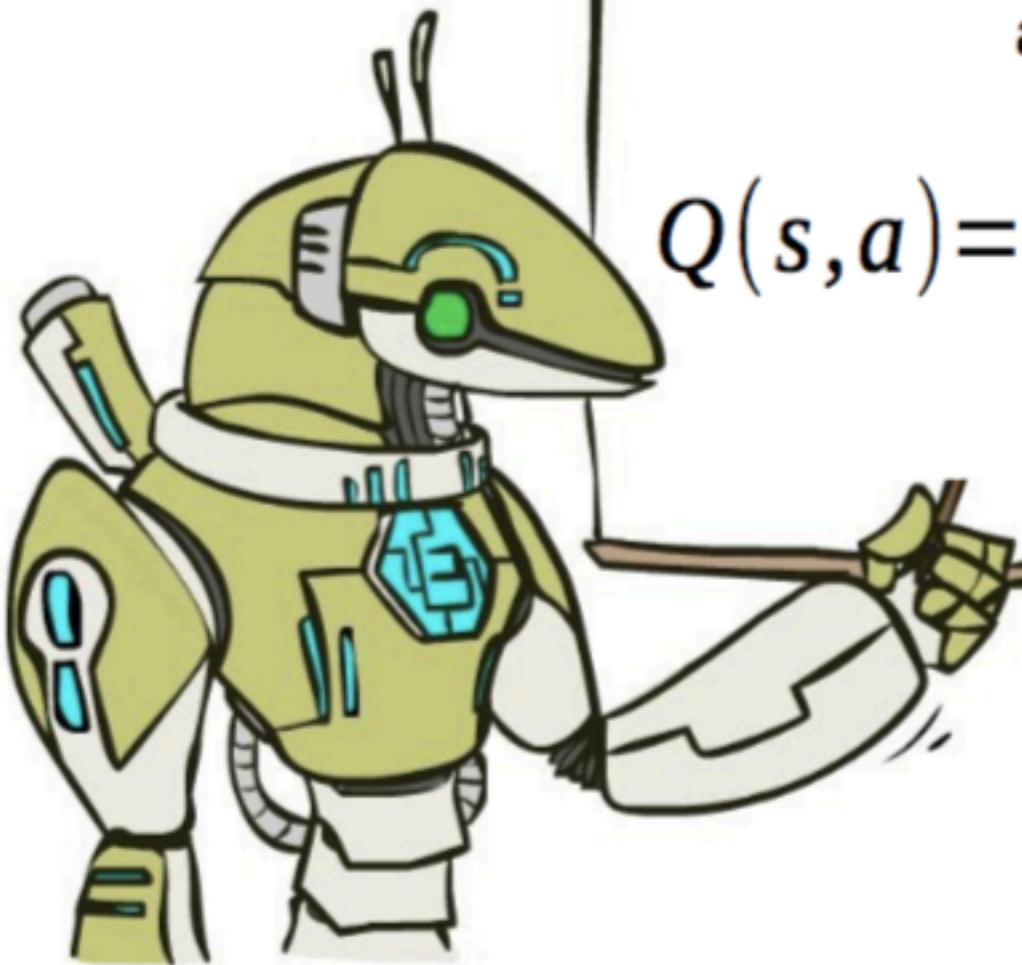
$$V_{i+1}(s) := \max_a [r(s, a) + \gamma \cdot E_{s' \sim P(s'|s, a)} V_i(s')]$$



From V to Q

One approach:
action Q-values

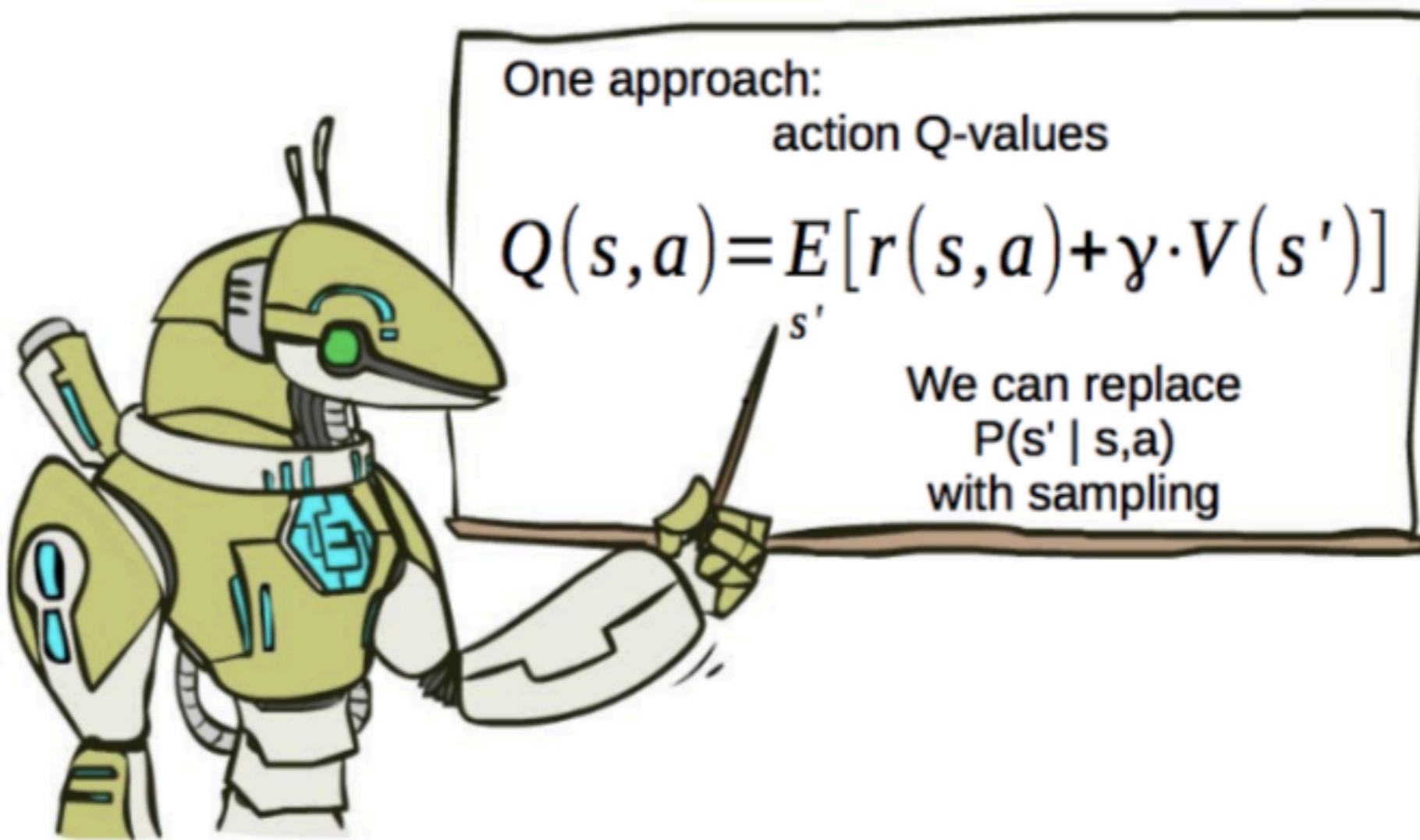
$$Q(s, a) = \underset{s'}{E}[r(s, a) + \gamma \cdot V(s')]$$



Action value $Q(s, a)$ is the expected total reward R agent gets from state s by taking action a and following policy π from next state.

$$\pi(s): \operatorname{argmax}_a Q(s, a)$$

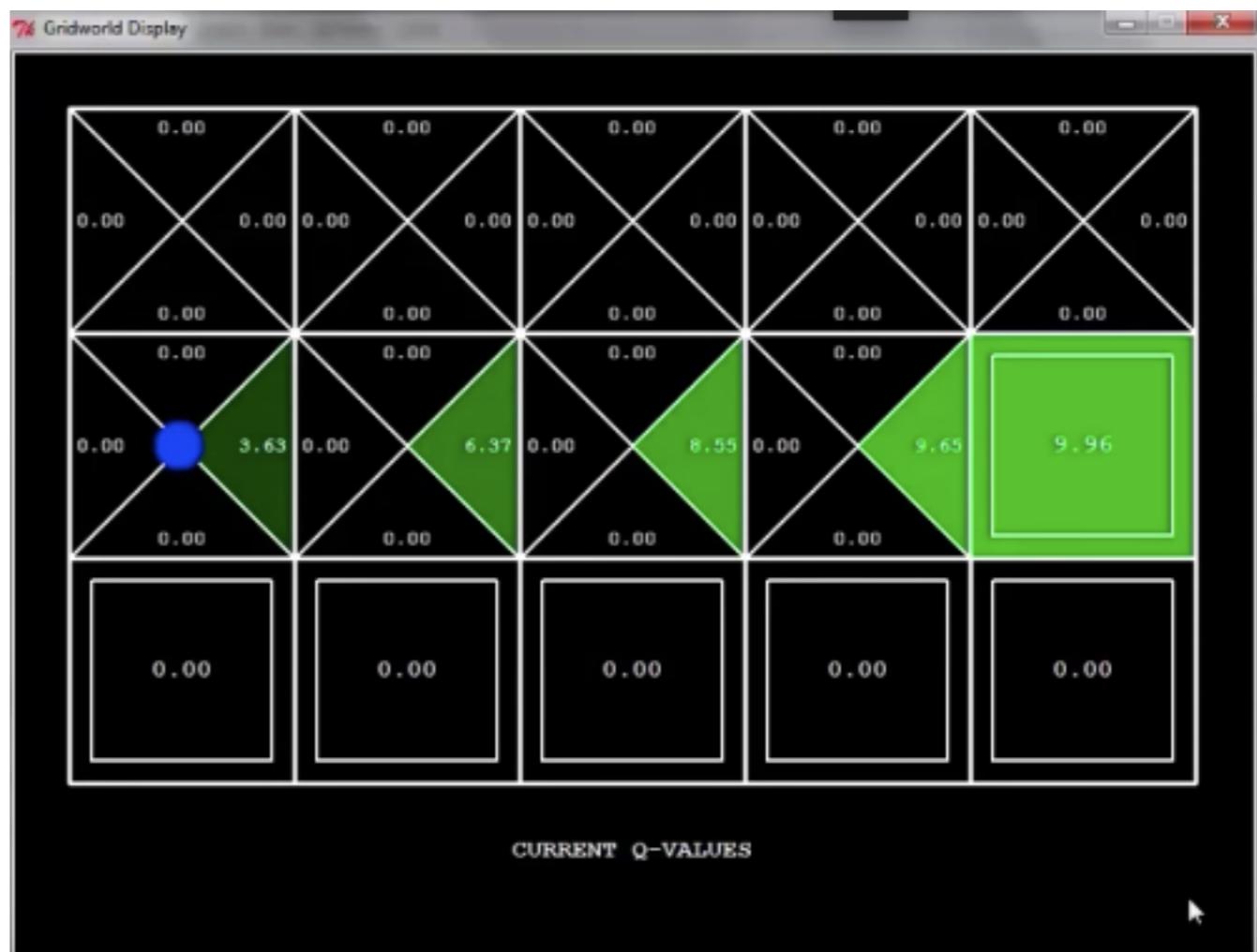
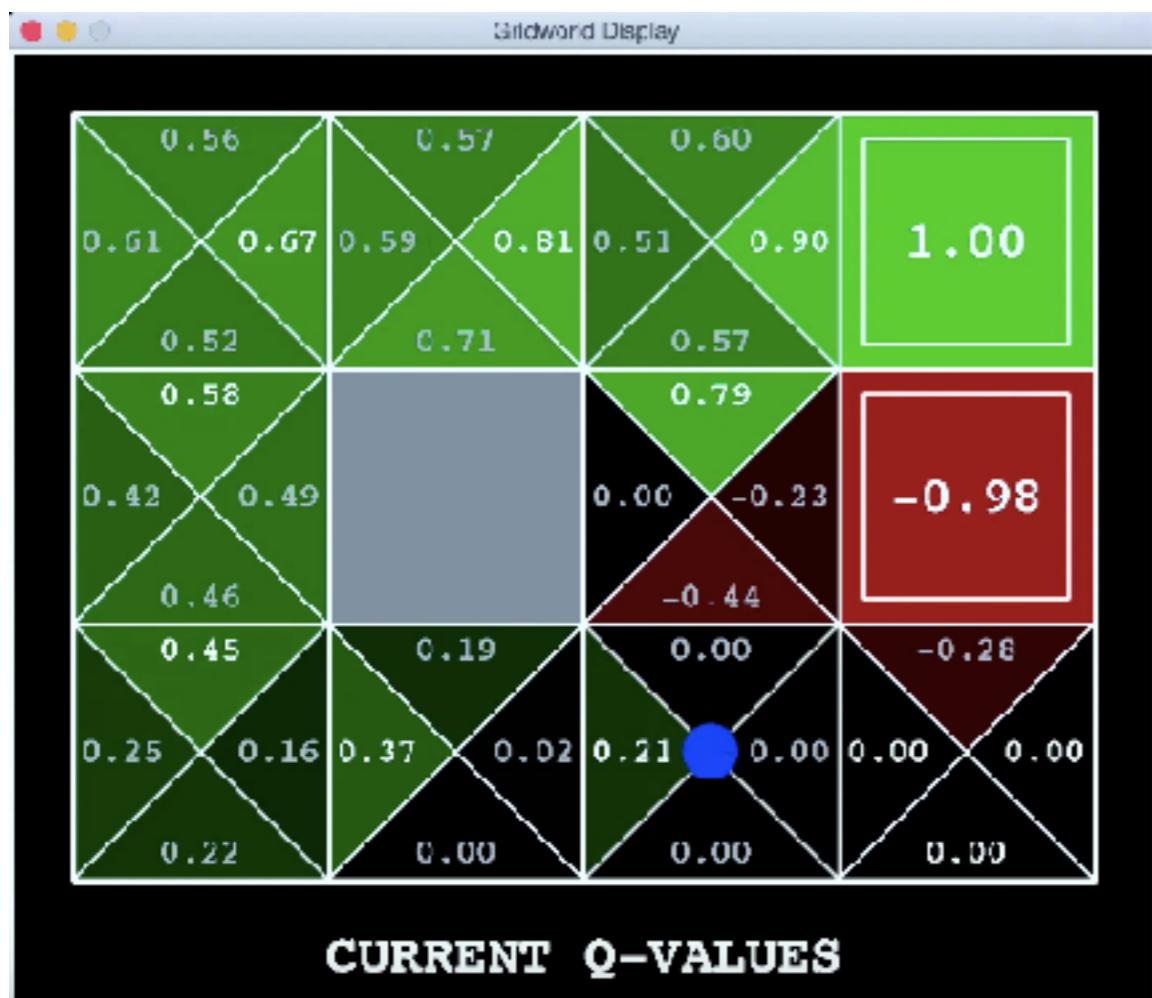
From V to Q



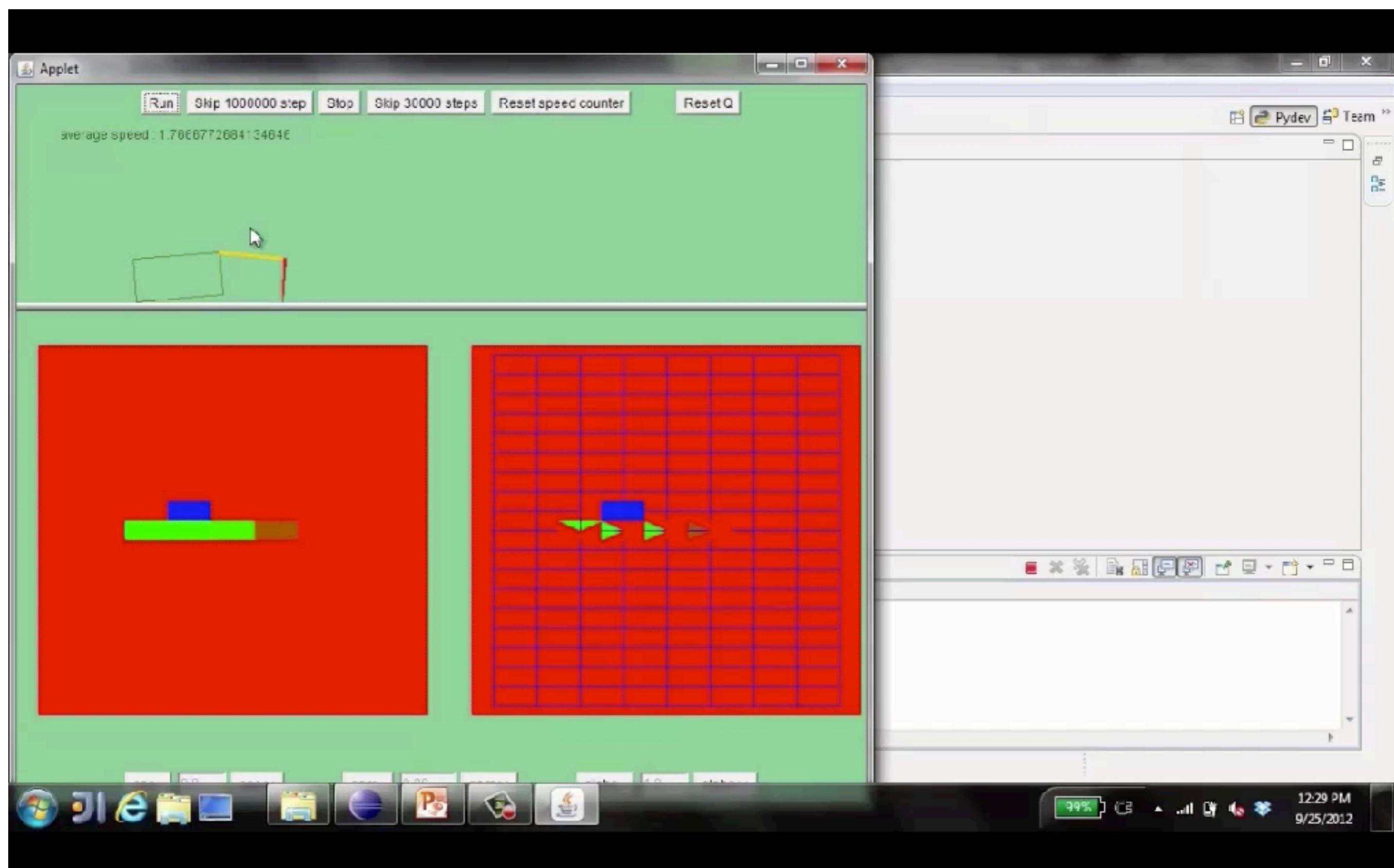
$$Q(s_t, a_t) \leftarrow \alpha \cdot (r_t + \gamma \cdot \max_a Q(s_{t+1}, a')) + (1 - \alpha) Q(s_t, a_t)$$

$$\pi(s): \operatorname{argmax}_a Q(s, a)$$

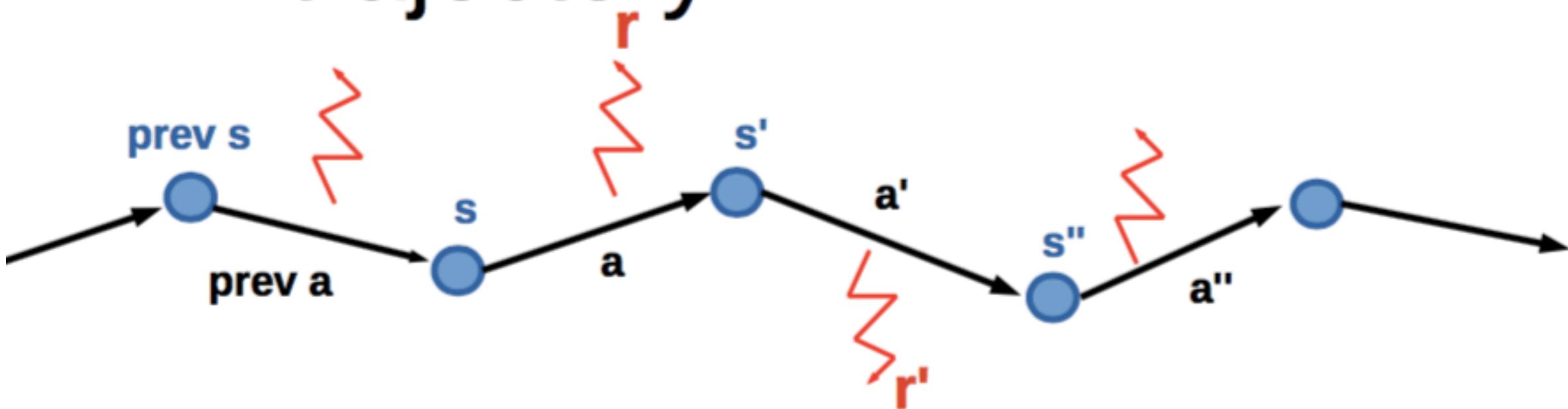
From V to Q



From V to Q

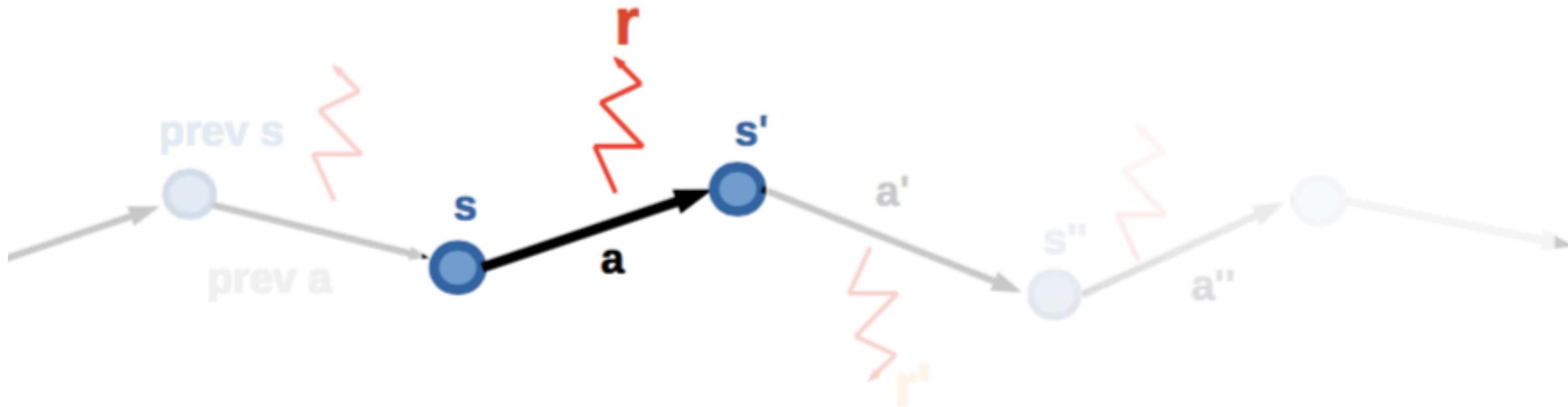


MDP trajectory



- sample sequence of
 - states (s)
 - actions (a)
 - rewards (r)

Q-learning

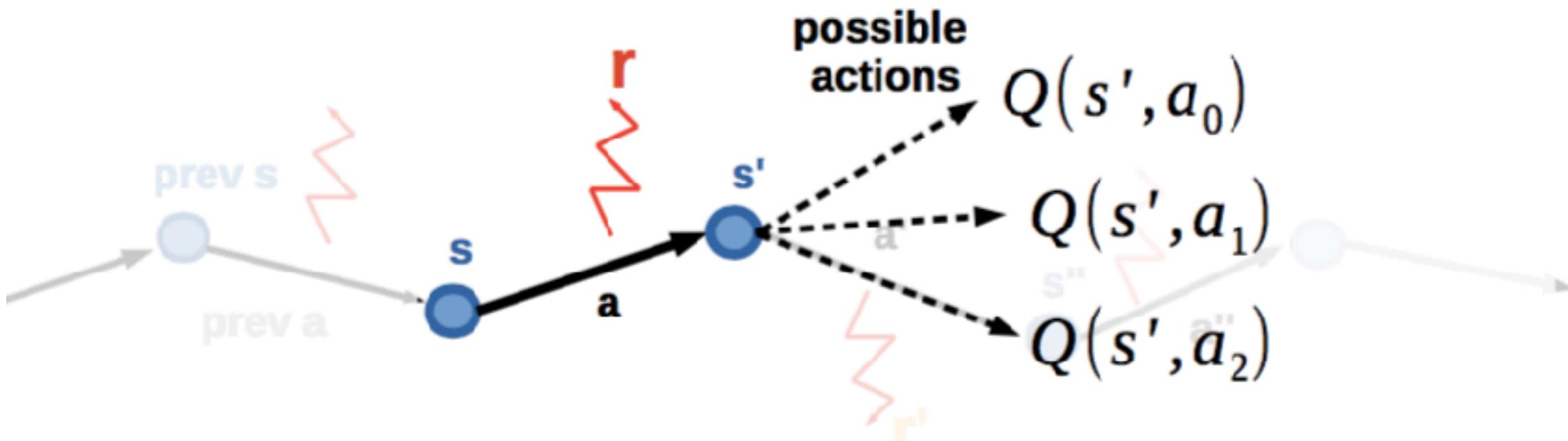


$$\forall s \in S, \forall a \in A, Q(s, a) \leftarrow 0$$

Loop:

- Sample $\langle s, a, r, s' \rangle$ from env

Q-learning

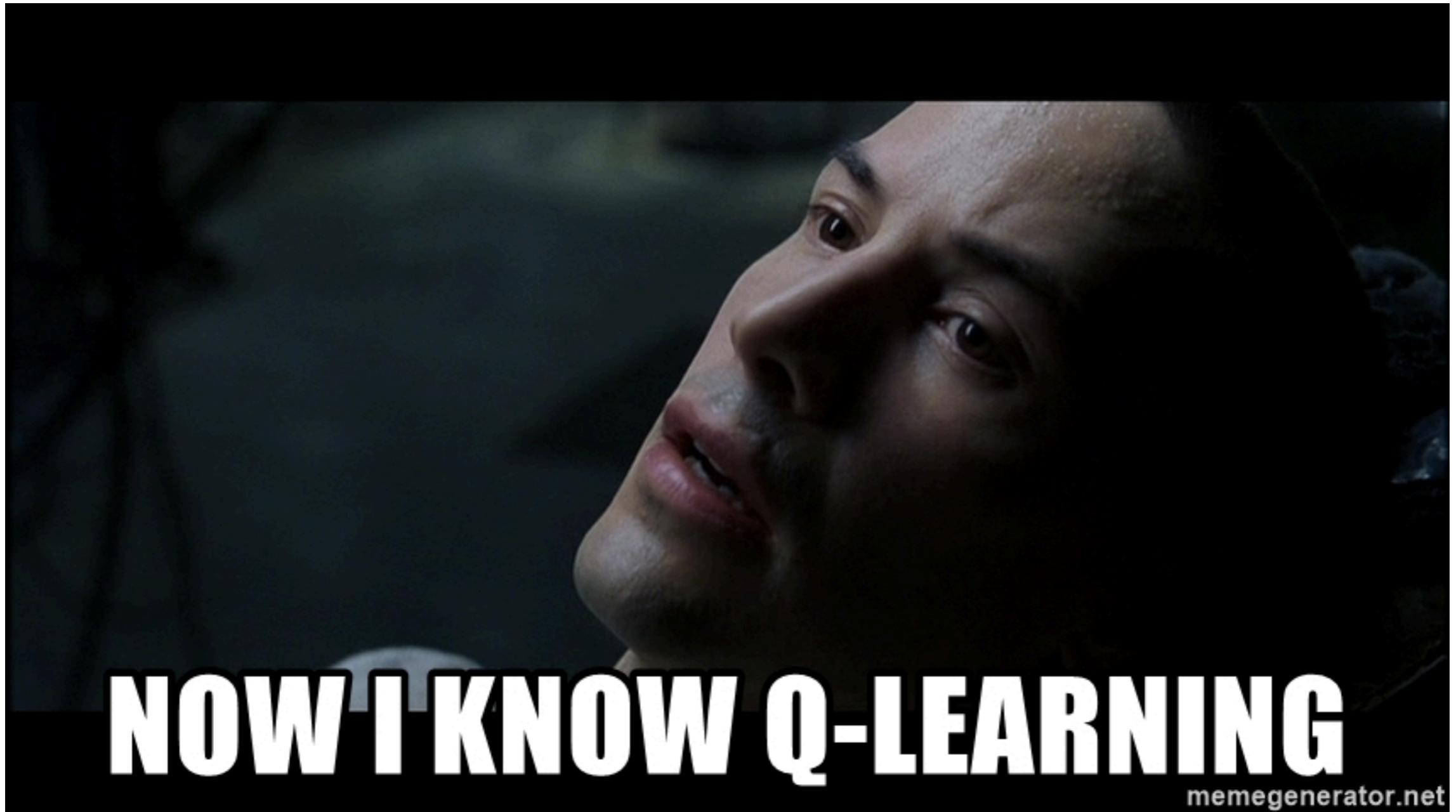


$$\forall s \in S, \forall a \in A, Q(s, a) \leftarrow 0$$

Loop:

- Sample $\langle s, a, r, s' \rangle$ from env
- Compute $\hat{Q}(s, a) = r(s, a) + \gamma \max_{a_i} Q(s', a_i)$
- Update $Q(s, a) \leftarrow \alpha \cdot \hat{Q}(s, a) + (1 - \alpha) Q(s, a)$

Practice time!

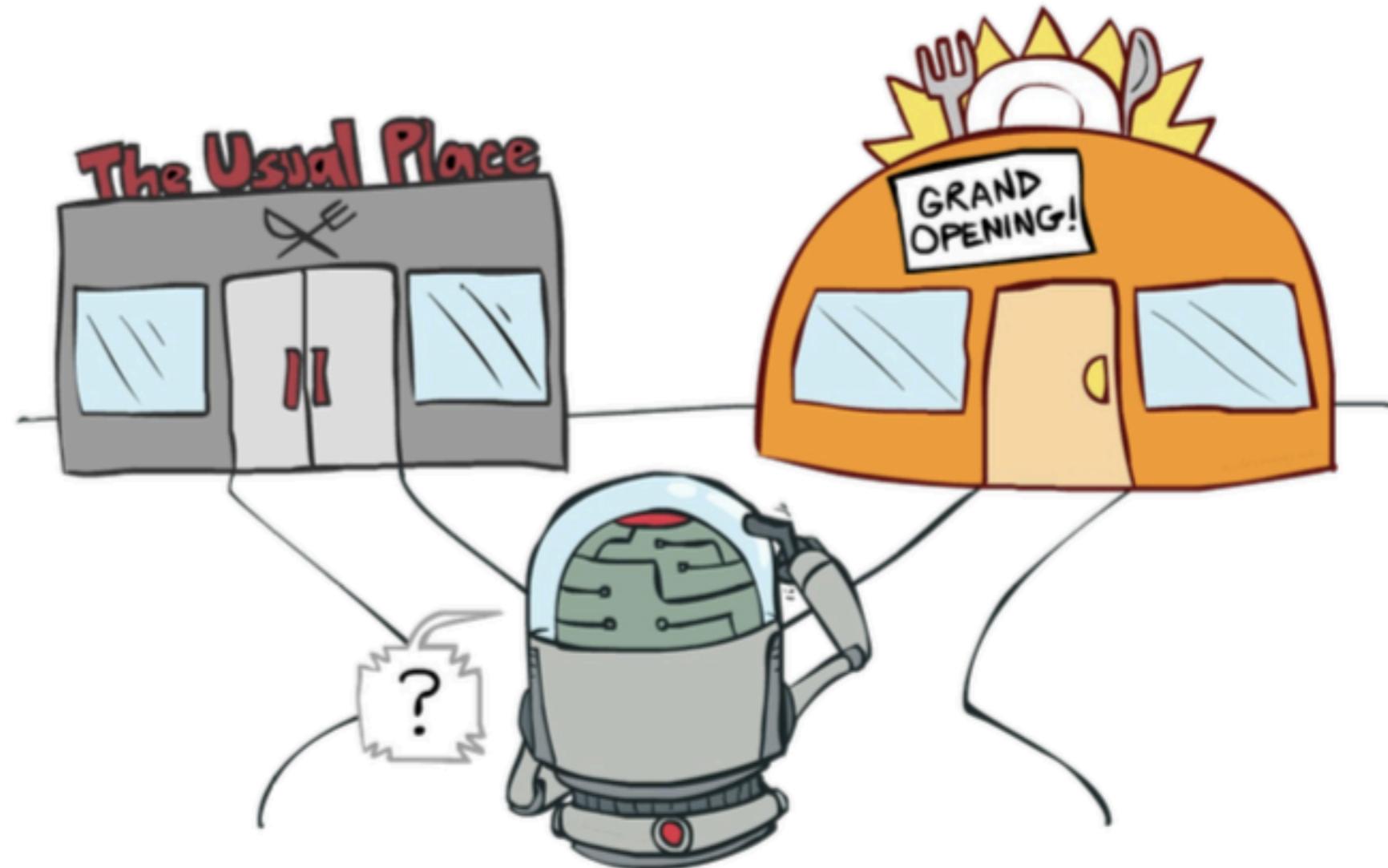


NOW I KNOW Q-LEARNING

memegenerator.net

Exploration Vs Exploitation

Balance between using what you learned and trying to find something even better



Exploration Vs Exploitation

Strategies:

- ϵ -greedy
 - With probability ϵ take a uniformly random action; otherwise take optimal action.
- Softmax
 - Pick action proportional to softmax of shifted normalized Q-values.

$$P(a) = \text{softmax}\left(\frac{Q(a)}{\tau}\right)$$

Nuts and bolts: discount

- Effective horizon $1 + \gamma + \gamma^2 + \dots = \frac{1}{(1 - \gamma)}$

Heuristic: your agent stops giving a damn in *this many* turns.

Typical values:

- $\gamma=0.9$, 10 turns
- $\gamma=0.95$, 20 turns
- $\gamma=0.99$, 100 turns
- $\gamma=1$, infinitely long

Higher γ = less stable algorithm.

$\gamma=1$ only works for episodic MDP (finite amount of turns).

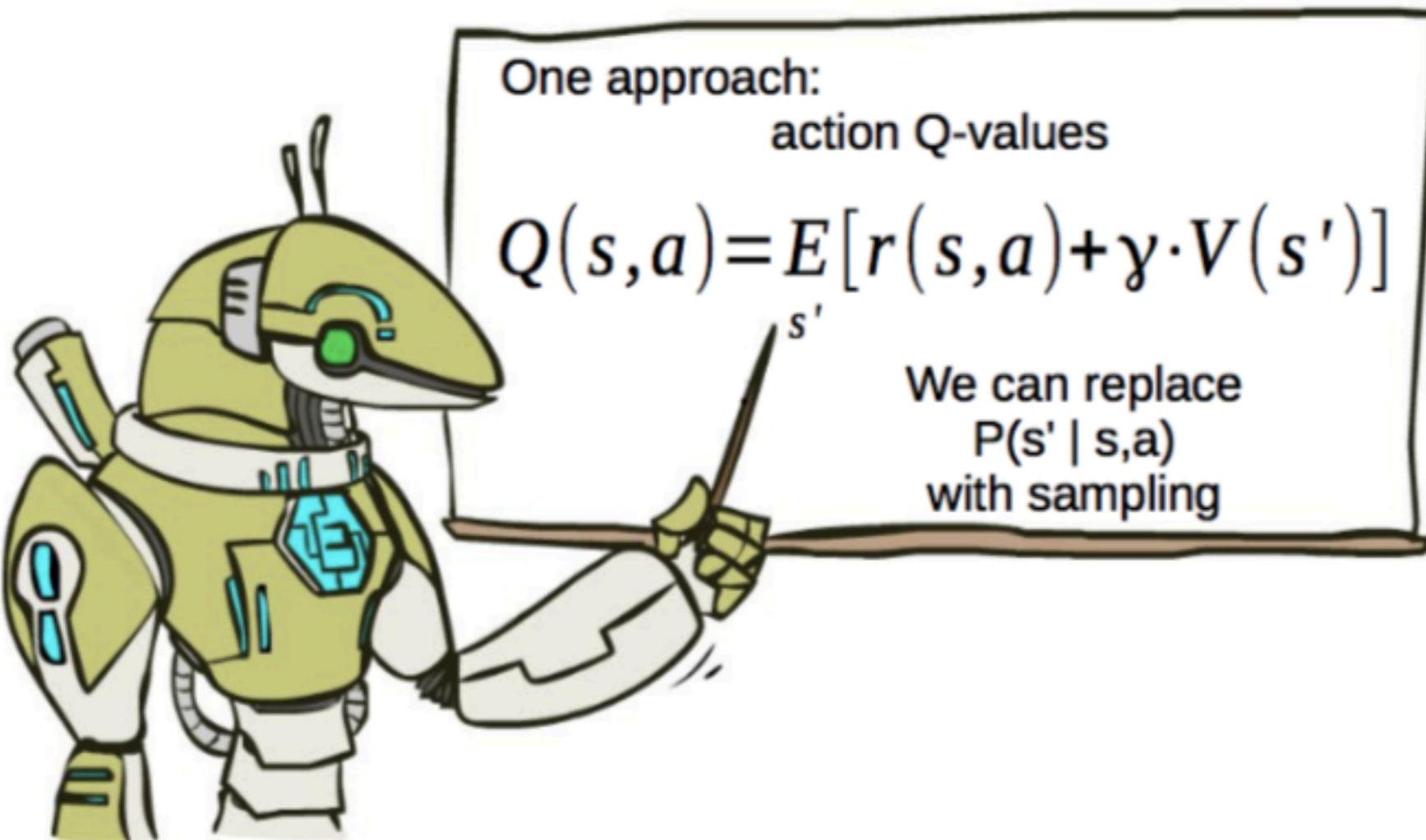
Problem:

**State space is usually large,
sometimes continuous.**

And so is action space;

**However, states do have a structure, similar
states have similar action outcomes.**

Recap: Q-learning



$$Q(s_t, a_t) \leftarrow \alpha \cdot (r_t + \gamma \cdot \max_{a'} Q(s_{t+1}, a')) + (1 - \alpha) Q(s_t, a_t)$$

$$\pi(s): \operatorname{argmax}_a Q(s, a)$$

From tables to approximations

- Before:
 - For all states, for all actions, remember $Q(s,a)$
- Now:
 - Approximate $Q(s,a)$ with some function
 - e.g. linear model over state features

$$\operatorname{argmin}_{w,b} (Q(s_t, a_t) - [r_t + \gamma \cdot \max_{a'} Q(s_{t+1}, a')])^2$$

Q-learning as MSE minimization

Given $\langle s, a, r, s' \rangle$ minimize

$$L = [Q(s_t, a_t) - Q^{\text{true}}(s_t, a_t)]^2$$

$$L \approx [Q(s_t, a_t) - (r_t + \gamma \cdot \max_{a'} Q(s_{t+1}, a'))]^2$$

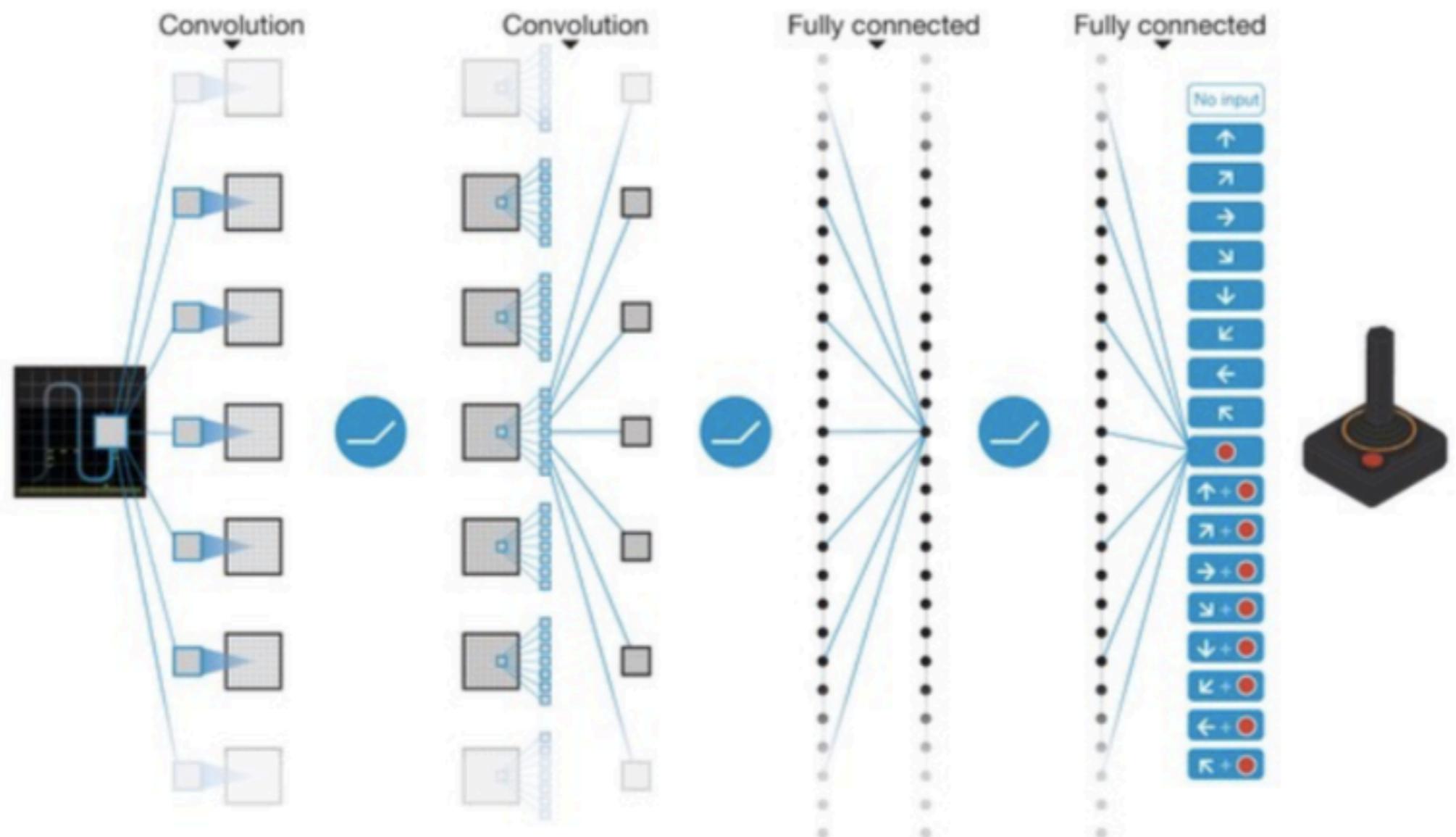
Q-learning as MSE minimization

Given $\langle s, a, r, s' \rangle$ minimize

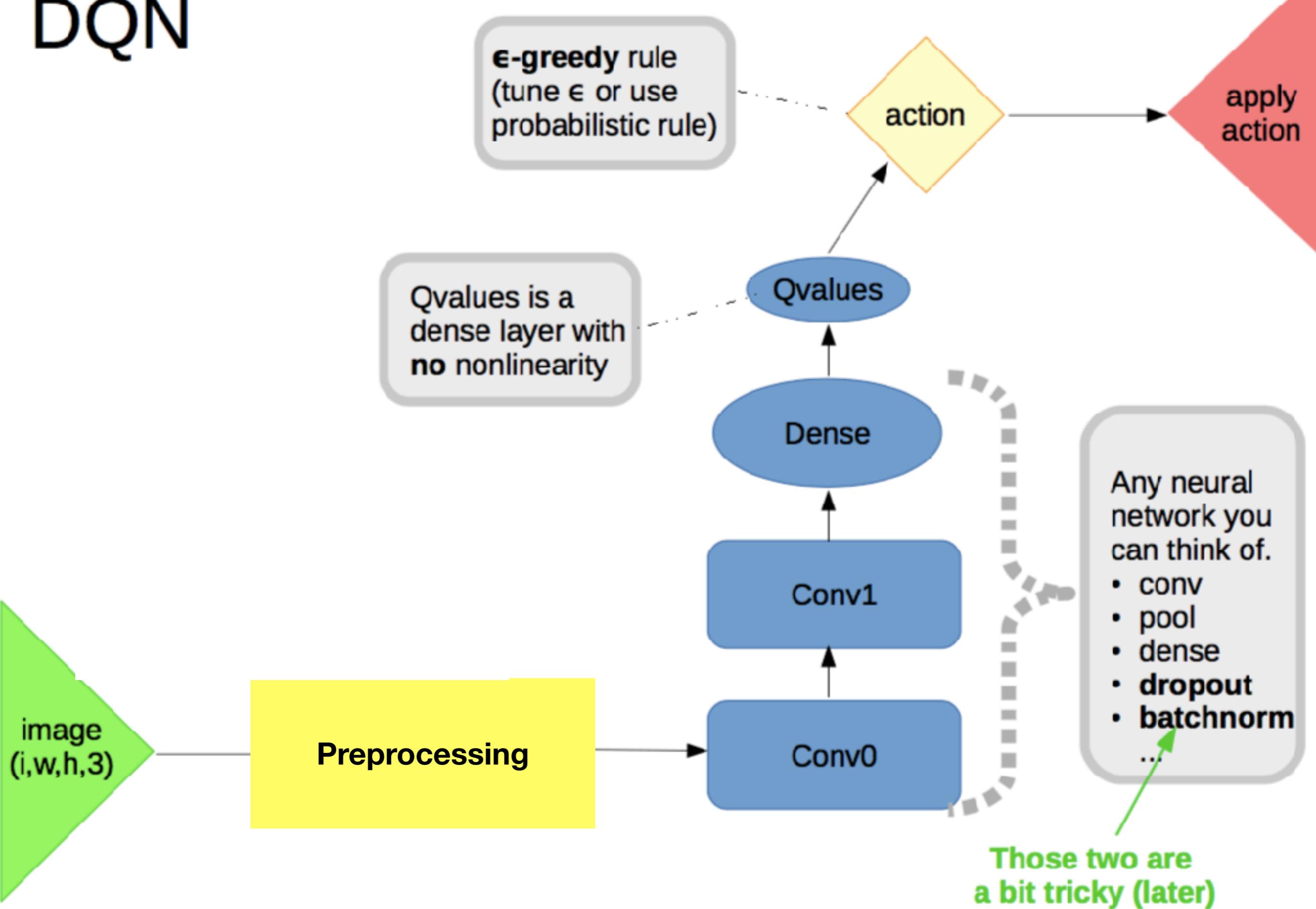
$$L = [Q(s_t, a_t) - \underbrace{Q^{\text{true}}(s_t, a_t)}_{\text{const}}]^2$$

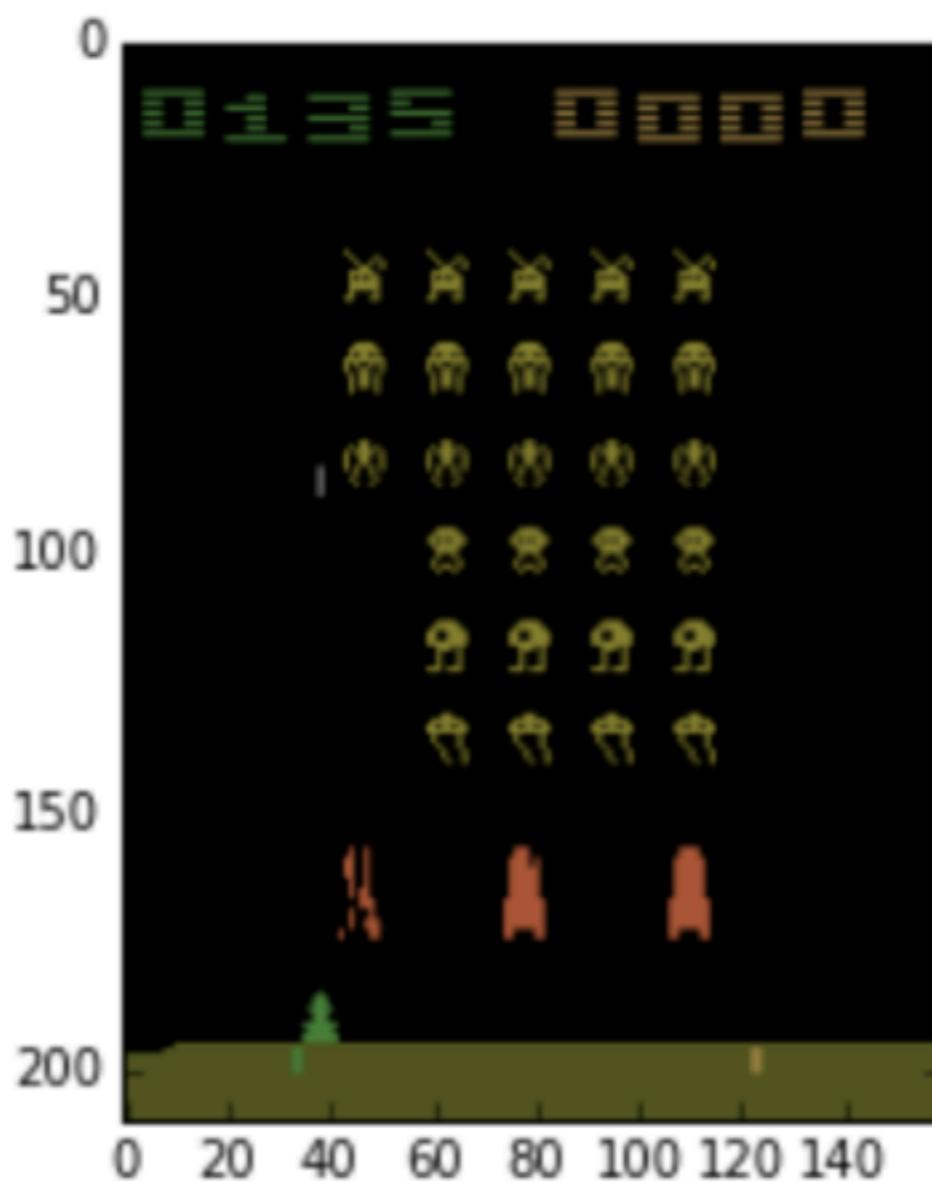
$$L \approx [Q(s_t, a_t) - \underbrace{(r_t + \gamma \cdot \max_{a'} Q(s_{t+1}, a'))}_{\text{const}}]^2$$

Smells like a neural network



DQN

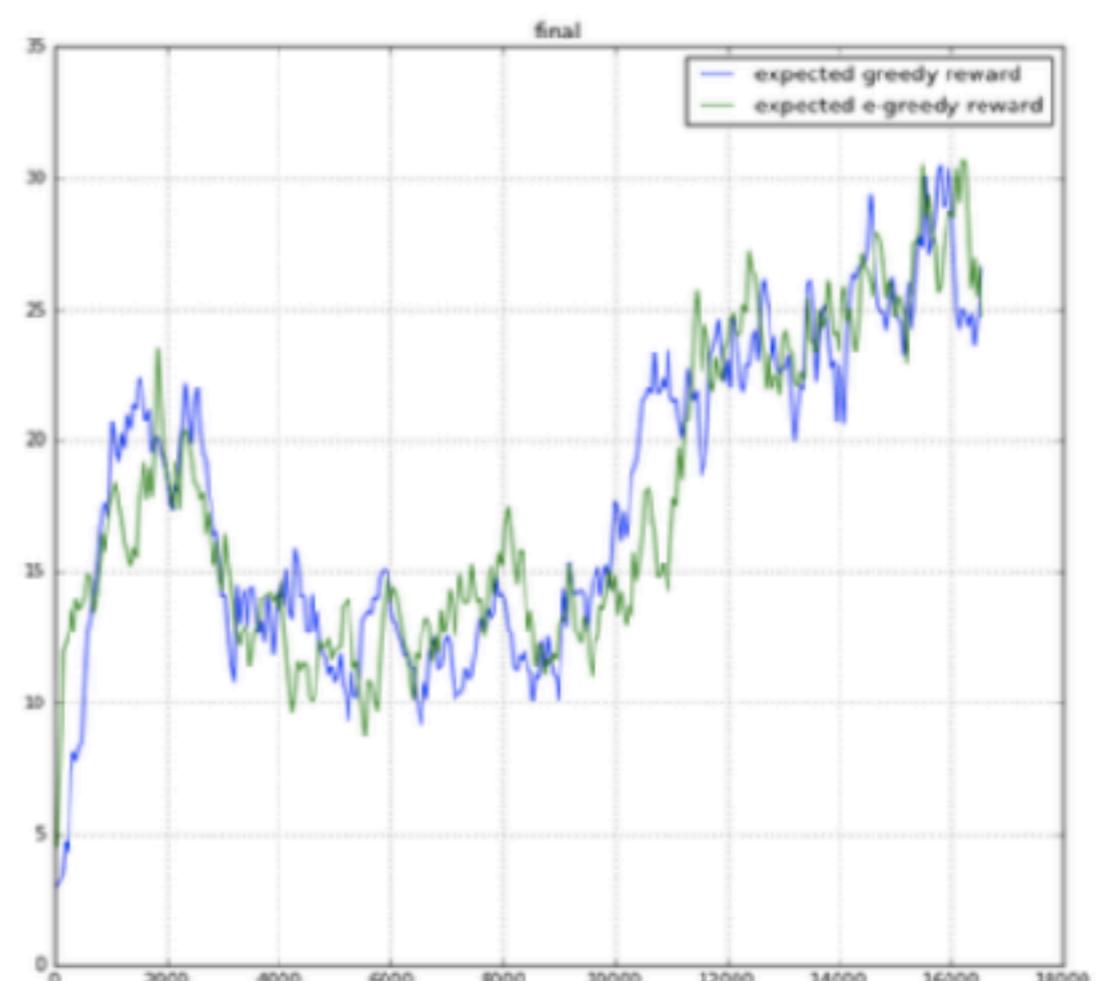




How bad it is if agent spends
next 1000 ticks under the left rock?
(while training)

Problem

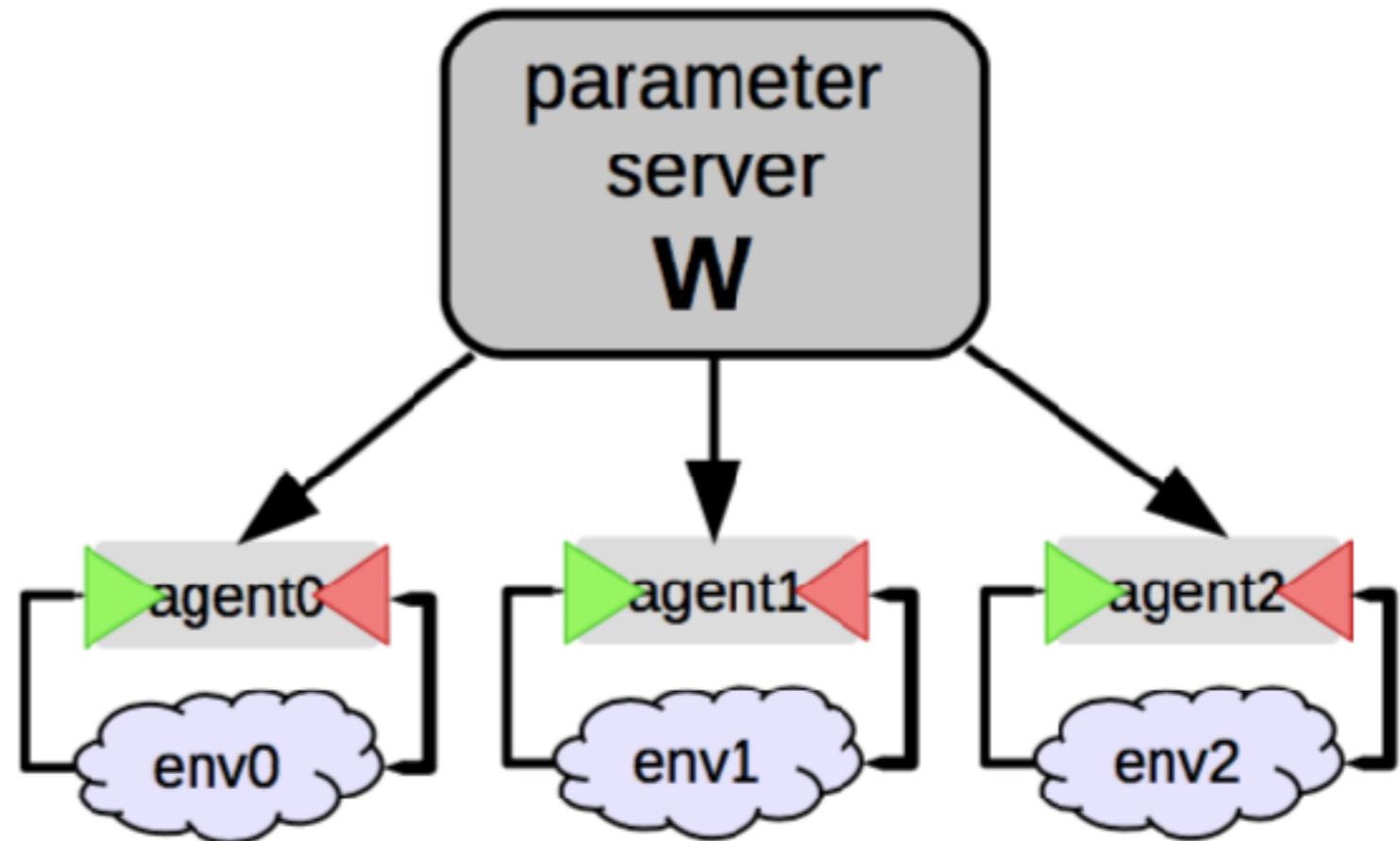
- Training samples are **not** “**i.i.d**”,
- Model forgets parts of environment it hasn't visited for some time
- Drops on learning curve
- **Any ideas?**



Multiple agent trick

Idea: Throw in several agents with shared \mathbf{W} .

- Chances are, they will be exploring different parts of the environment,
- More stable training,
- Requires a lot of interaction



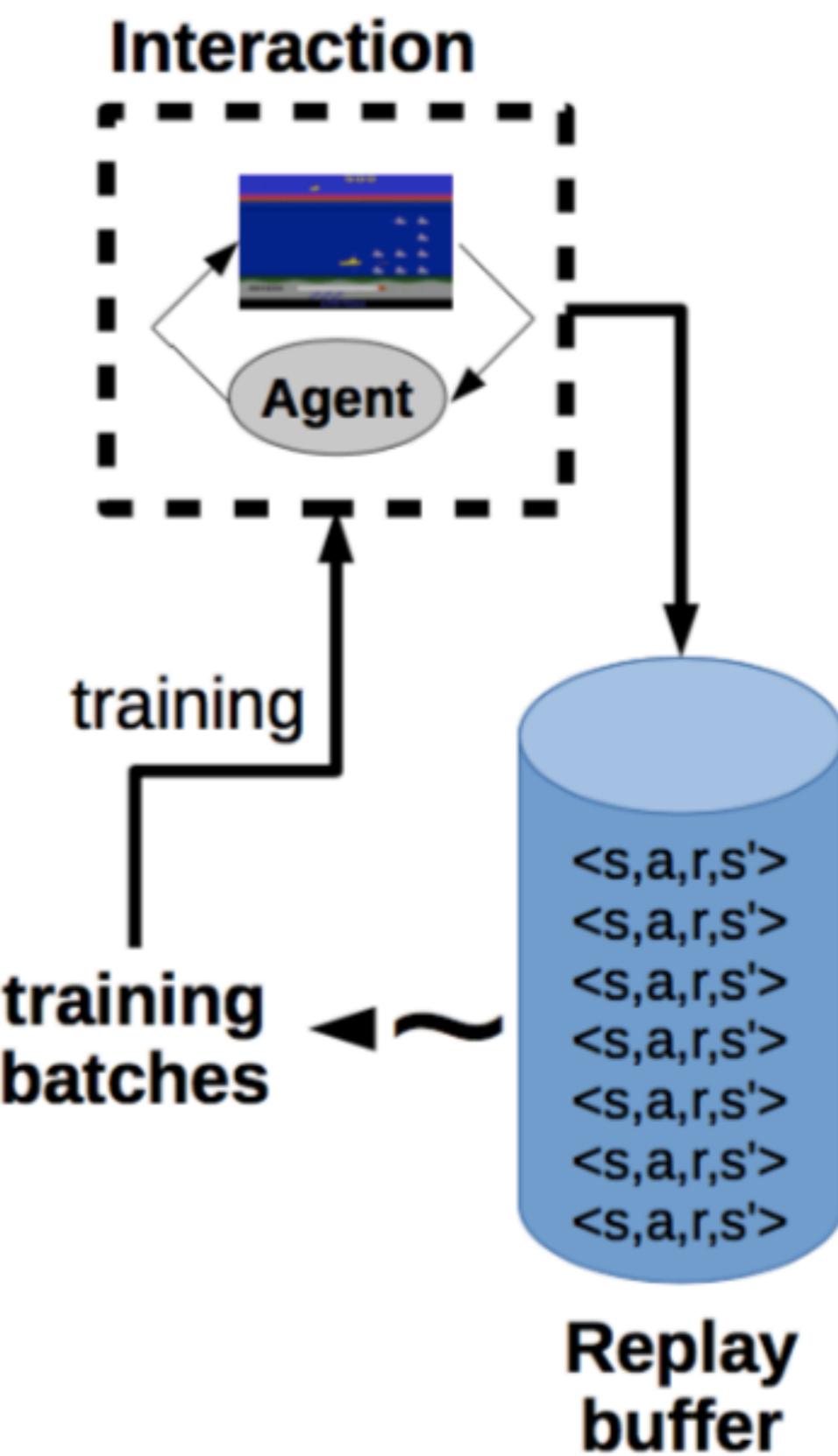
Trivia: your agent is a real robot car. Any problems?

Experience replay

Idea: store several past interactions
 $\langle s, a, r, s' \rangle$

Train on random subsamples

- Atari DQN: $>10^5$ interactions
 - Closer to i.i.d
pool contains several sessions
 - Older interactions were obtained
under weaker policy



Practice time!



made on imgur