# Introduction to Deep Reinforcement Learning

Ludovic Righetti (MPI / NYU)
Nicolas Mansard (LAAS/CNRS)
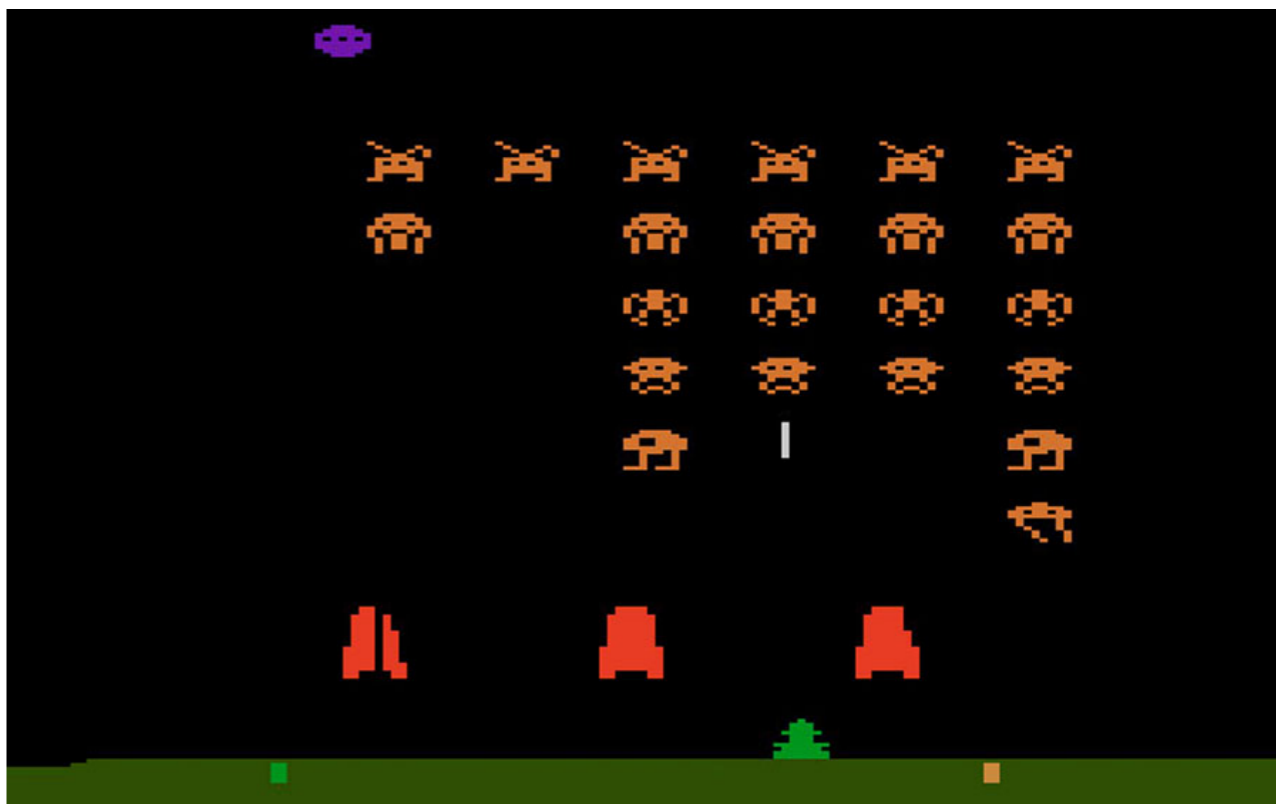
Learning "how to act" from direct interaction with the environment with the goal to maximize a "reward"
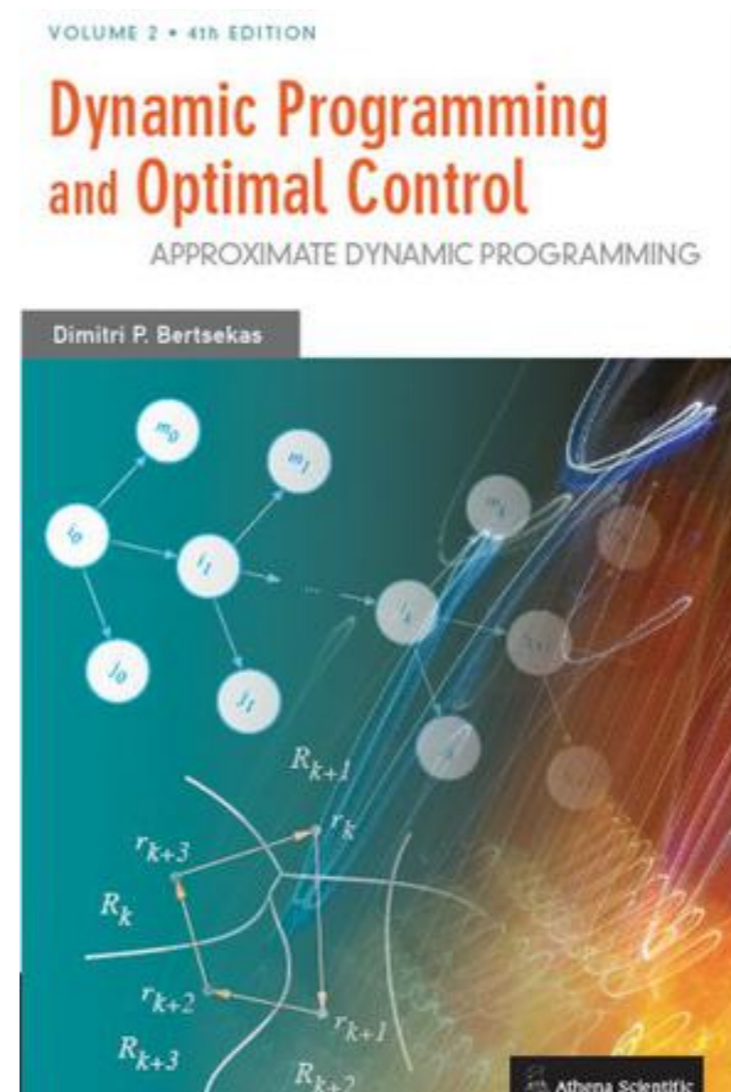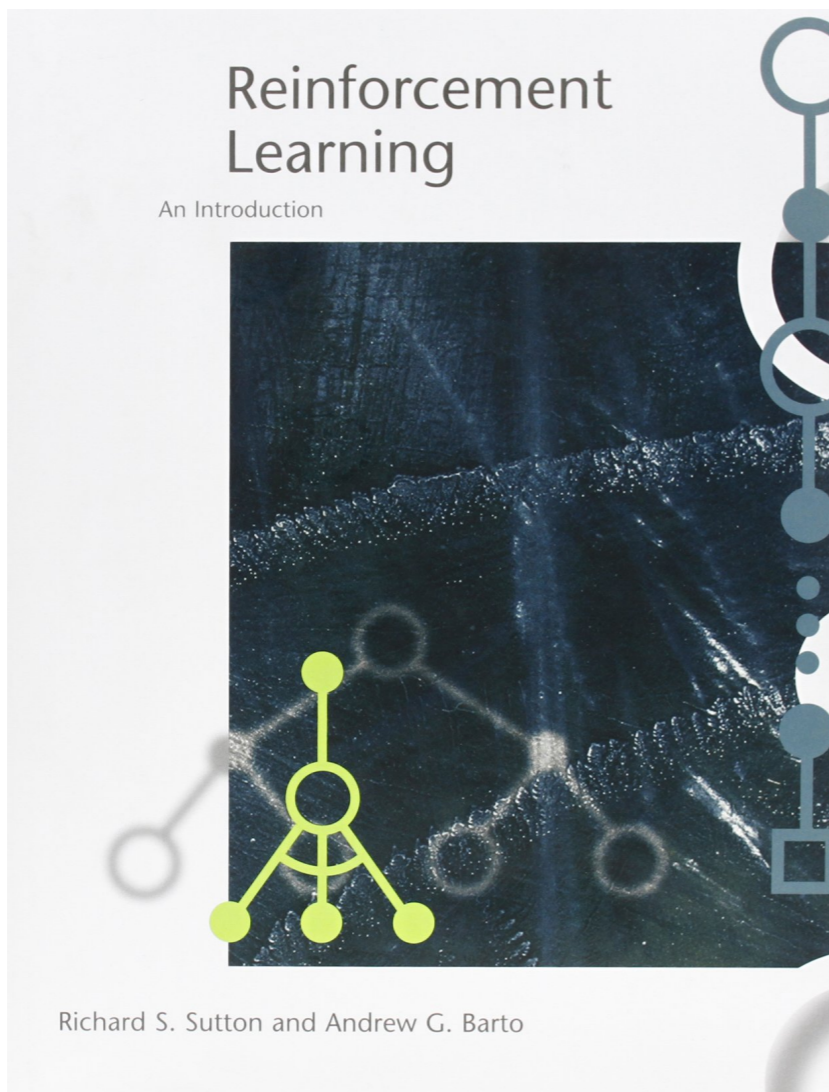


[Silver et al., Nature, 2016]



[Mnih et al., Nature, 2015]

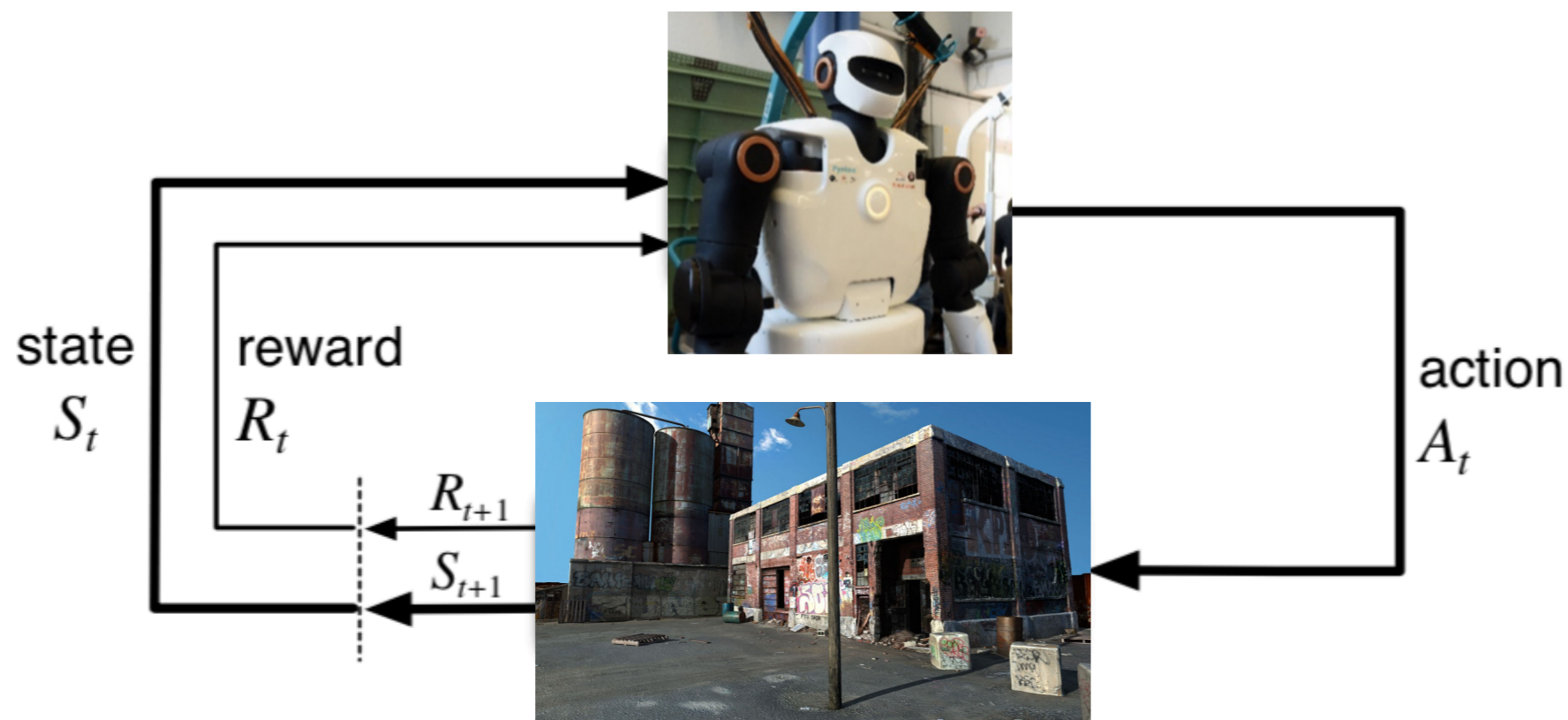Jens Kober et al. "Reinforcement Learning in Robotics: A Survey", IJRR 2013

$$\max \sum_t R_t(S_t, A_t)$$

state
$S_t$

reward
$R_t$

$R_{t+1}$

$S_{t+1}$

action
$A_t$

Policy
$A_t = \pi(S_t)$ **?**

$$S_{t+1} = f(S_t, A_t)$$

## Markov Decision Process

# What is reinforcement learning?

## Typical RL problems

State is discrete (countable)

Set of actions is discrete
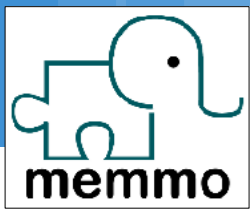
## Robotics RL problems

State is continuous

Action space is continuous

> Most methods designed for discrete state/action models do not carry over to continuous sate/action models

## Reinforcement learning

$$\max_t \sum_t R_t(S_t, A_t)$$

$$S_{t+1} = \ ?$$

In RL we do not know the dynamic model

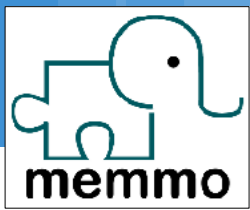## Optimal Control

$$\min \sum_t C_t(S_t, A_t)$$

$$S_{t+1} = f(S_t, A_t)$$

Our typical optimal control setup (cf. DDP tutorial)

Reinforcement learning

$$\min \sum_t C_t(S_t, A_t)$$

$$S_{t+1} = \text{ ? }$$
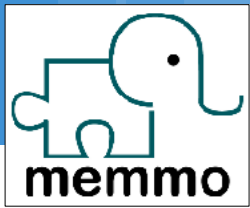
Model-based reinforcement learning => Sample the world, learn a model (i.e. system identification) and then do optimal control

[Schaal and Atkeson, Control Systems Magazine, 1994]

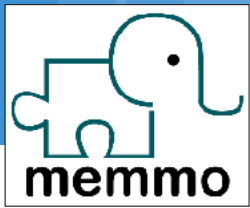[Levine and Koltun, ICML, 2013]

Reinforcement learning

$$\min \sum_t C_t(S_t, A_t)$$

$$S_{t+1} = \text{ ?}$$

Ignore the model and learn directly the policy (or an approximation of the value function to infer the policy) => Q-learning, actor-critic algorithms, policy gradients, etc

$$\min_{A_t} \sum_{t=0}^{\infty} \gamma^t C_t(S_t, A_t)$$
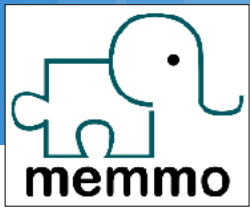
$\gamma$ Discount factor

$$S_{t+1} = f(S_t, A_t)$$

## Bellman's Equations (infinite horizon case)

$$V(S_t)$$

Valu

i.e. optimal

obtained

opti

> $V(S_t)$ is the unique solution of this equation within the class of bounded functions (necessary and sufficient condition for optimality)

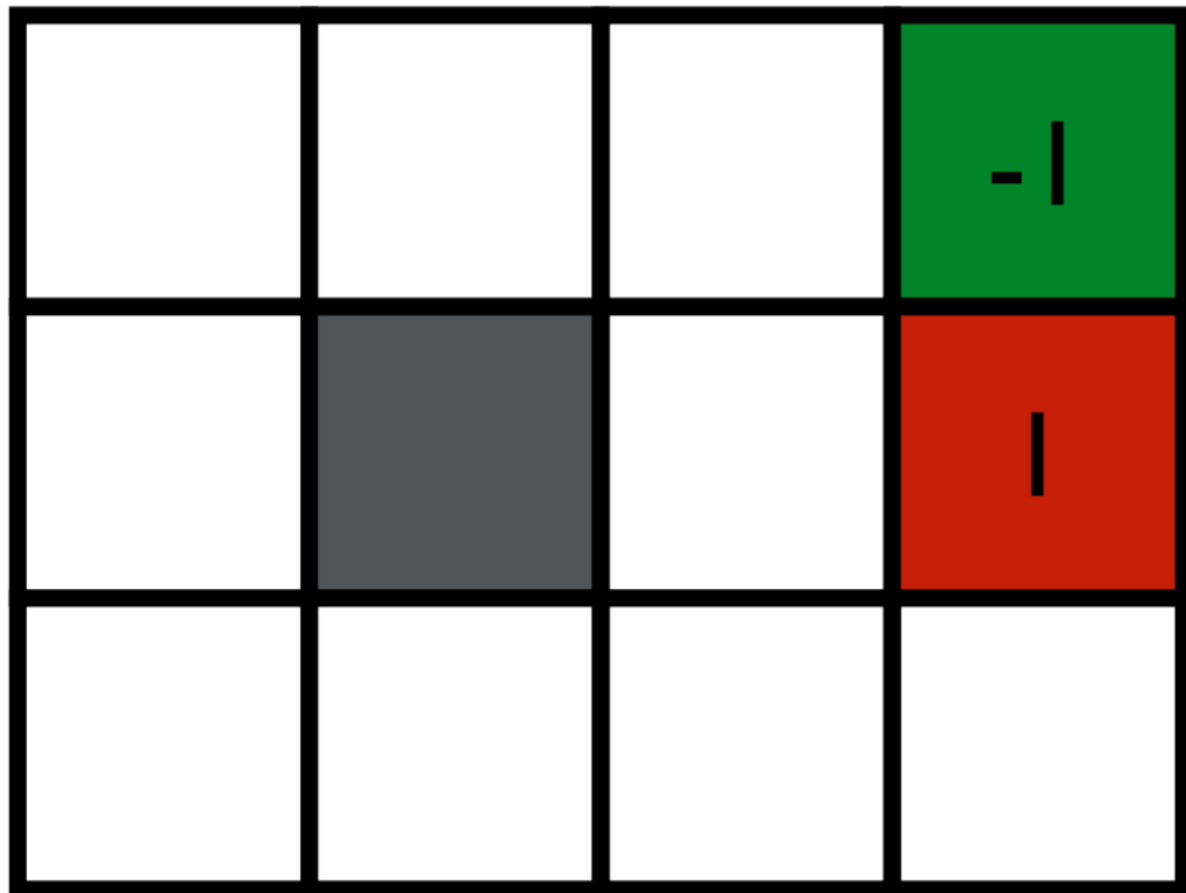$$V(S_t) = \min_{A_t} \left( C_t(S_t, A_t) + \gamma V(S_{t+1}) \right)$$

Start with an initial guess for V $\to V_0$

Iterate the Bellman Equation for each state $S_t$

$$V_{n+1}(S_t) = \min_{A_t} (C_t(S_t, A_t) + \gamma V_n(S_{t+1}))$$

One can show that $V_n \to V$ when $n \to \infty$

Get out of the maze
- Red is bad (+1 cost)
- Green is good (-1 cost)
- Possible actions (N,E,W,S)
- $\gamma = 0.9$

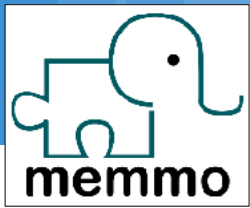| 0 | 0 | 0 | -1 |
|---|---|---|---|
| 0 |  | 0 | 1 |
| 0 | 0 | 0 | 0 |

- Initialize $V_0$

| 0 | 0 | -0.9 | -1.9 |
|---|---|------|------|
| 0 |   | 0 | 1.9 |
| 0 | 0 | 0 | 0 |

- First iteration of Bellman (we update every state)

| 0 | -0.81 | -1.71 | -2.71 |
|---|-------|-------|-------|
| 0 |       | -0.81 | 2.71 |
| 0 | 0 | 0 | 0 |

- 2nd Iteration

| -0.73 | -1.54 | -2.44 | -3.44 |
|-------|-------|-------|-------|
| 0 | | -1.54 | 3.44 |
| 0 | 0 | -0.73 | 0 |

- 3rd Iteration

| -1.39 | -2.2 | -3.1 | -4.1 |
|-------|------|------|------|
| -0.66 |      | -2.2 | 4.1  |
| 0     | -0.66 | -1.39 | -0.66 |

- 4th Iteration

| | | | |
|---|---|---|---|
| -4.15 | -4.96 | -5.86 | -6.86 |
| -3.42 | | -4.96 | 6.86 |
| -2.77 | -3.42 | -4.15 | -3.42 |

- 10th Iteration

| -7.29 | -8.1 | -9 | -10 |
|-------|------|-----|-----|
| -6.56 | | -8.1 | 10 |
| -5.9 | -6.56 | -7.29 | -6.56 |

– 100th Iteration

| -7.29 | -8.1 | -9 | -10 |
|-------|------|----|-----|
| -6.56 | | -8.1 | 10 |
| -5.9 | -6.56 | -7.29 | -6.56 |

- 1000th Iteration

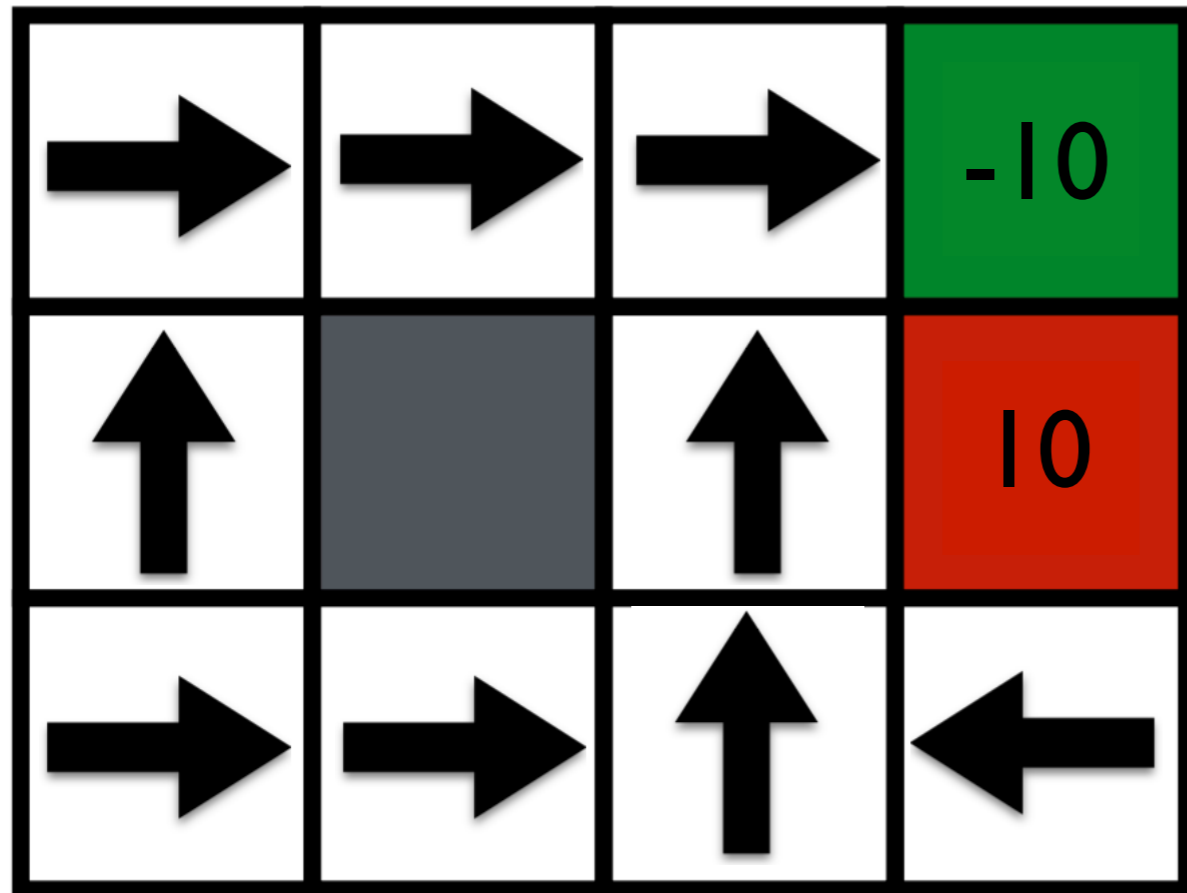| | | | |
|---|---|---|---|
| -7.29 | -8.1 | -9 | -10 |
| -6.56 | | -8.1 | 10 |
| -5.9 | -6.56 | -7.29 | -6.56 |

- 1000th Iteration

We have converged and found the optimal value function

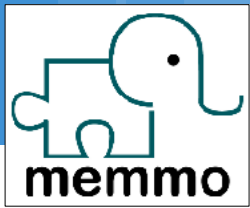The policy is read out by following the action that creates the lowest next value + current cost

- 1000th Iteration

We have converged and found the optimal value function

The policy is read out by following the action that creates the highest next value

Start with an initial guess for the policy $\pi_0(S_t)$
and value function $V(S_t)$
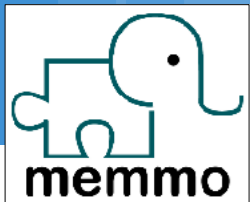
1. Policy evaluation

   Iterate the Bellman Equation for each state $S_t$
   using the current policy for transition

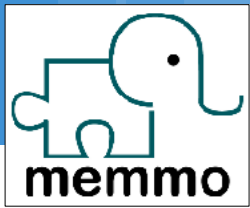   $$V_{n+1}(S_t) = C_t(S_t, \pi_k(S_t)) + \gamma V_n(S_{t+1})$$

2. Policy update

   Given the current guess for V find a new policy
   such that $\quad \pi_{k+1}(S_t) = \arg\min C_t(S_t, A_t) + \gamma V(S_{t+1})$

[Watkins 1989]

Introduce the Q-function $Q(S_t, A_t)$

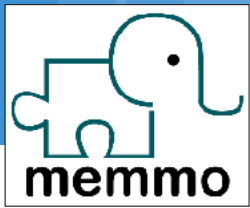which stores for each state/action pairs the cost of performing $A_t$ and then using the optimal policy afterwards

$$Q(S_t, A_t) = C_t + \gamma V(S_{t+1})$$

$$V(S_t) = \min_{A_t} Q(S_t, A_t)$$

Fortunately, the Q-function also satisfies Bellman's equation

$$Q(S_t, A_t) = C_t(S_t, A_t) + \gamma \min_{A_{t+1}} Q(S_{t+1}, A_{t+1})$$

If we choose any random action $A_t$ and observe what reward $C_t$ we get then if our Q-function guess was correct we should have
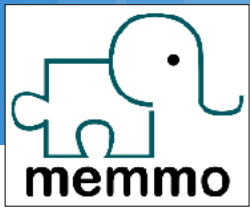
$$Q(S_t, A_t) = C_t(S_t, A_t) + \gamma \min_{A_{t+1}} Q(S_{t+1}, A_{t+1})$$

So if this is not the case, then we can estimate "how wrong we are" and update Q accordingly

$$Q^{new}(S_t, A_t) \leftarrow (1 - \alpha)Q(S_t, A_t) + \alpha(C_t(S_t, A_t) + \gamma \min_A Q(S_{t+1}, A))$$

$\alpha$ is the learning rate between 0 and 1

When dealing with discrete state/action spaces (and moderate size for this space) we can store the Q-function as a table indexed by actions and states

## Q-learning with a table

Choose $\alpha \in (0, 1]$ and small $\epsilon > 0$    Initialize $Q(S, A)$ arbitrarily

For each episode:

Start from an initial state $S_0$
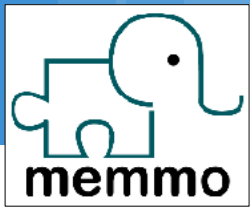
Loop for each step $t$ of the episode:

Choose $A_t$ from $S_t$ using a policy using $Q$ (e.g. $\epsilon$-greedy policy)

Take action $A_t$ and observe cost $C_t(S_t, A_t)$ and next state $S_{t+1}$

Update $Q(S_t, A_t) \leftarrow (1 - \alpha)Q(S_t, A_t) + \alpha(C_t(S_t, A_t) + \gamma \min_A Q(S_{t+1}, A))$
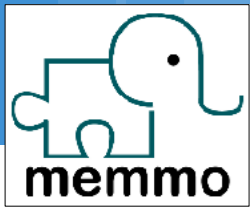
Do until convergence

Break
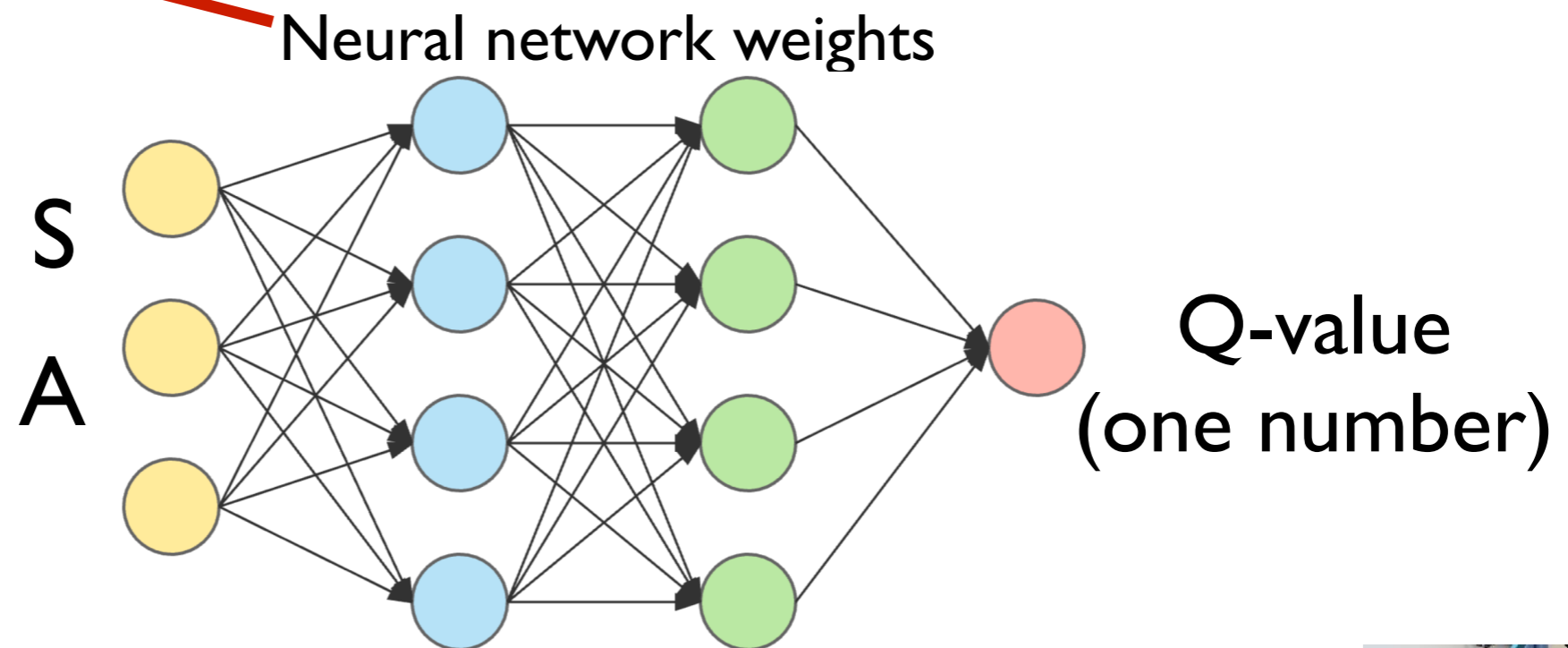Exercise on Q-learning with a table
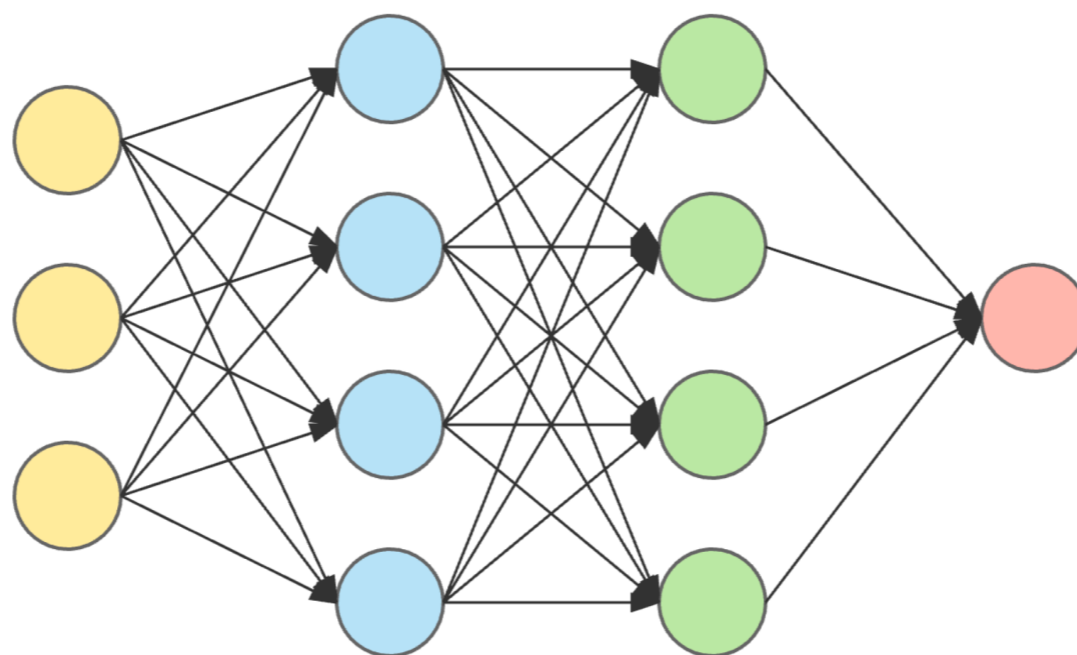Discretized states and actions

Q-learning with a table cannot work for high-dimensional spaces nor for continuous state/action spaces!

<u>Idea</u>: replace the table with a function approximator (e.g. a neural network) - still assume discrete number of actions

$$Q(S, A) \simeq Q(S, A, \theta)$$

Neural network weights



S

A

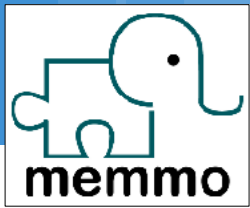Q-value
(one number)

$$Q(S, A) \simeq Q(S, A, \theta)$$



**The problem can be written as a least square problem**

$$\min_{\theta} ||Q(S_t, A_t, \theta) - C_t(S_t, A_t) - \gamma \min_{A_t} Q(S_{t+1}, A, \theta^-)||^2$$

**for set of observed $(S_t, A_t, C_t, S_{t+1})$ samples**

Problem: a direct (naive) approach using solely current episode data tend to be unstable (i.e. it diverges):

- The sequence of observations are correlated
- Small changes in Q can lead to large changes in policy
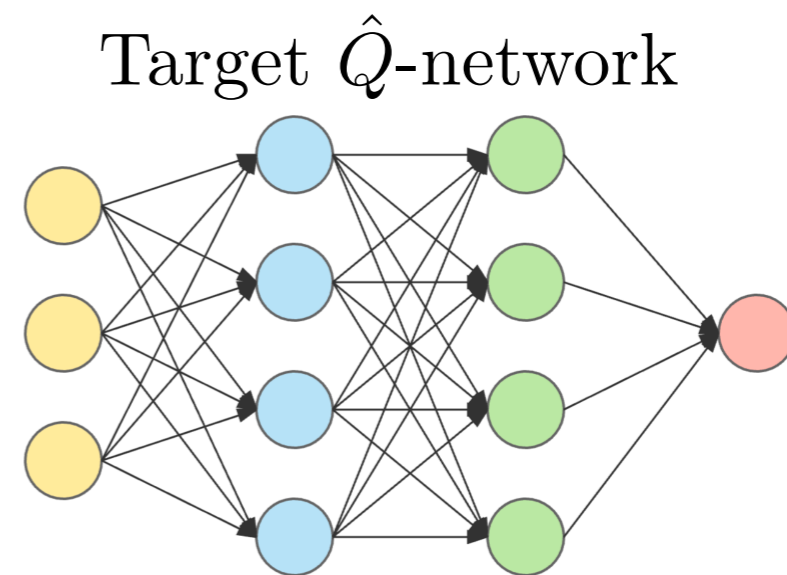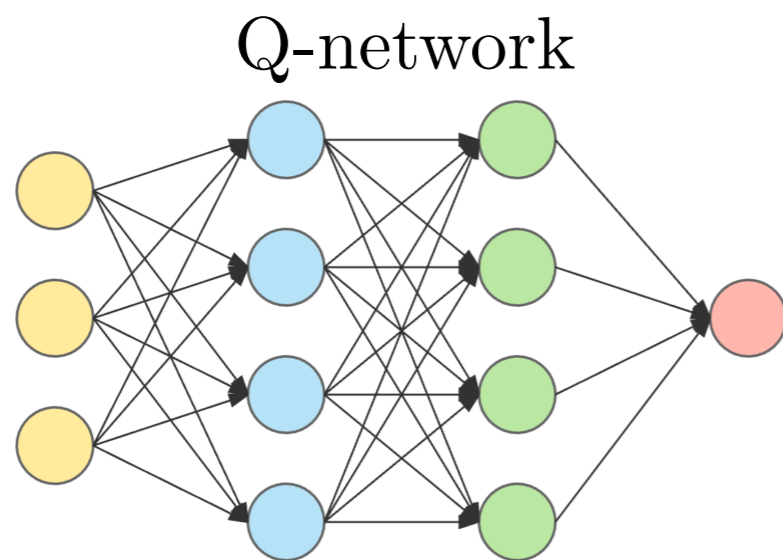
[Mnih et al., Nature, 2015]

Problem: a direct (naive) approach using solely current episode data tend to be unstable (i.e. it diverges):

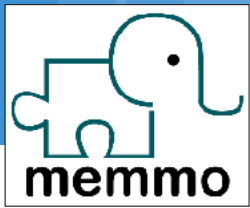- The sequence of observations are correlated
- Small changes in Q can lead to large changes in policy

Solution

1) Use a "replay" memory of a previous samples from which we randomly sample the next training batch (remove correlations)
2) Use 2 Q-network to avoid correlations due to updates

Q-network

Target $\hat{Q}$-network

Initialize replay memory D of size $N$

Initialize Q-network with random weights $\theta$

Initialize target $\hat{Q}$ function with weights $\theta^- = \theta$

For each episode:

  Start from an initial state $S_0$

  Loop for each step $t$ of the episode:

  Choose $A_t$ from $S_t$ using a policy using $Q$ (e.g. $\epsilon$-greedy policy)

  Take action $A_t$ and observe cost $C_t(S_t, A_t)$ and next state $S_{t+1}$

  Store $(S_t, A_t, C_t, S_{t+1})$ in memory D

  Sample minibatch of transitions $(S_j, A_j, C_j, S_{j+1})$ from memory D

  Gradient descent on $\theta$ to minimize $||Q(S_j, A_j, \theta) - C_j - \gamma \min_A \hat{Q}(S_{j+1}, A, \theta^-)||^2$
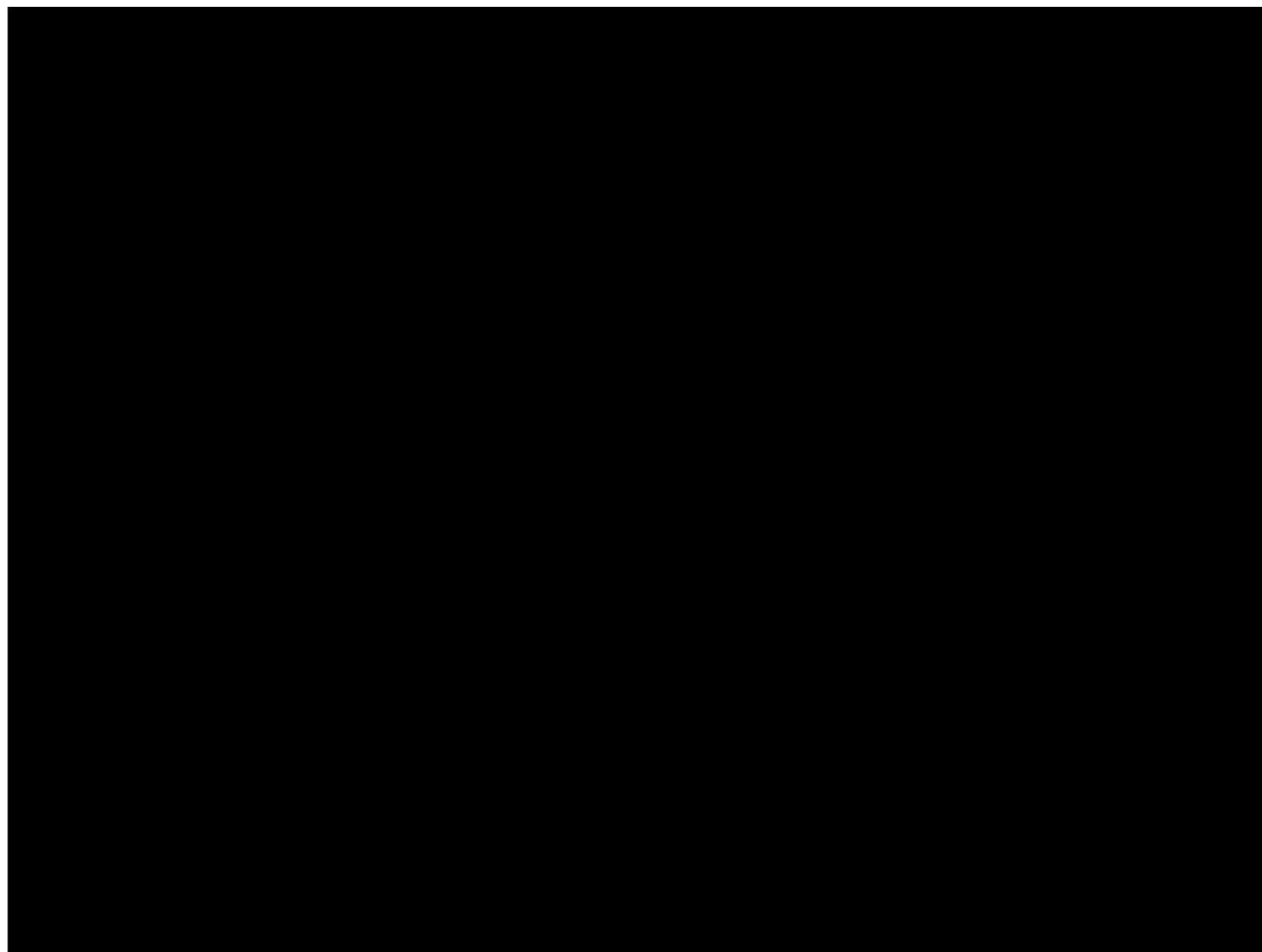
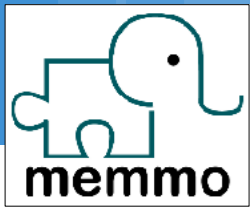  Every $C$ steps reset $\theta^- = \theta$

[Mnih et al., Nature, 2015]

[Mnih et al., Nature, 2015]
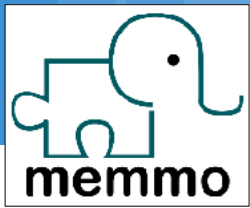
Break

Let's replace the table with a neural network

1) Just replace Q-table with the NN
2) DQN (replay memory, target, etc)

<u>Problem</u>: we need to evaluate the min to be able to do Q-learning with a function approximator

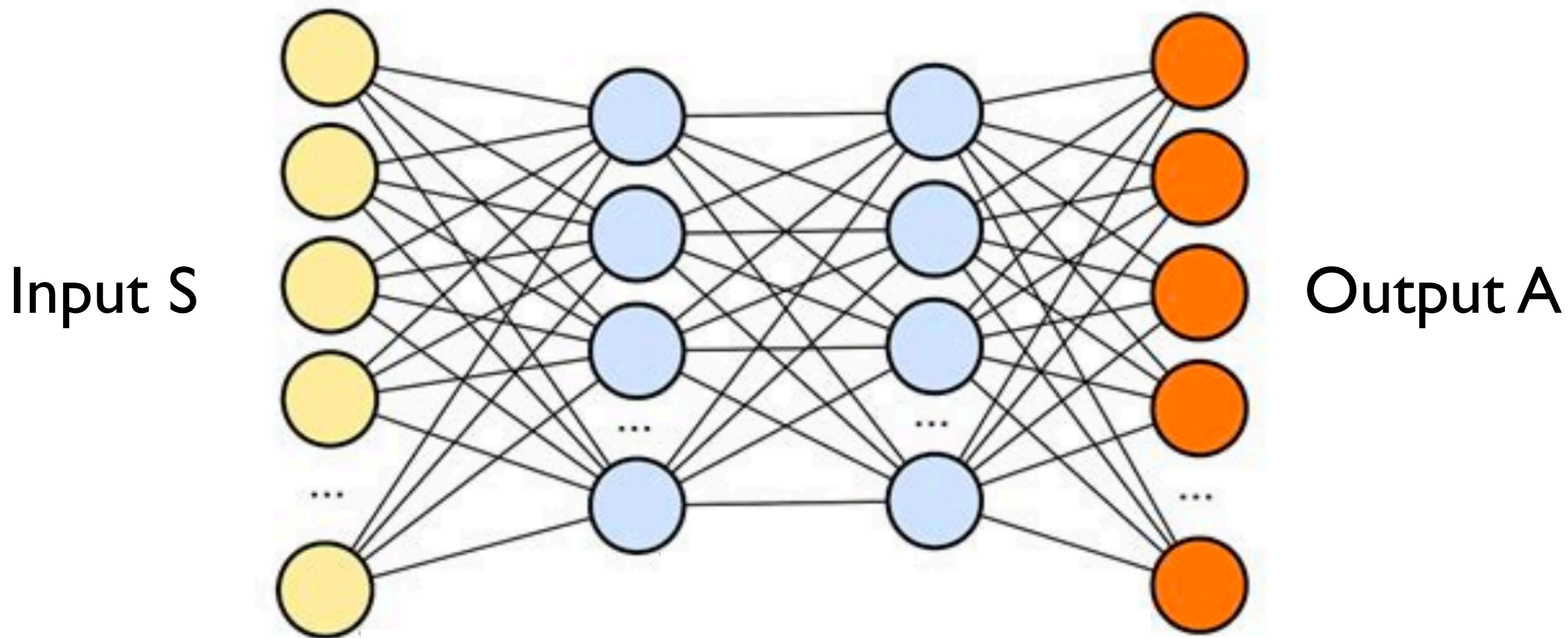$$||Q(S_j, A_j, \theta) - C_j - \gamma \min_A \hat{Q}(S_{j+1}, A, \theta^-)||^2$$

<u>Solution</u>: use another neural network to approximate the min operator (i.e. to approximate the optimal policy) => Actor-critic algorithm
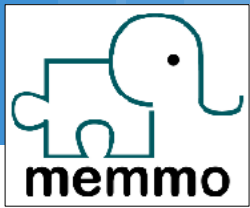
Let $\pi(S_t, \theta^\pi)$ an approximation of a policy with a NN (weights $\theta^\pi$)



Input S
Output A

Let $\pi(S_t, \theta^\pi)$ an approximation of a policy with a NN (weights $\theta^\pi$)

Let $J_\pi(S_t)$ be the value function under policy $\pi$

We would like to find the gradient $\nabla_{\theta^\pi} J_\pi$ to improve $\pi$

The policy gradient theorem (cf. Sutton-Barto book) tells us that

$$\nabla_{\theta^\pi} J \simeq \frac{1}{N} \sum_i \nabla_A Q(S, A, \theta^Q)|_{S=S_i, A=\pi(S_i)} \nabla_{\theta^\pi} \pi(S_i, \theta^\pi)$$

=> we can do gradient descent on the policy
parameters to minimize the value function

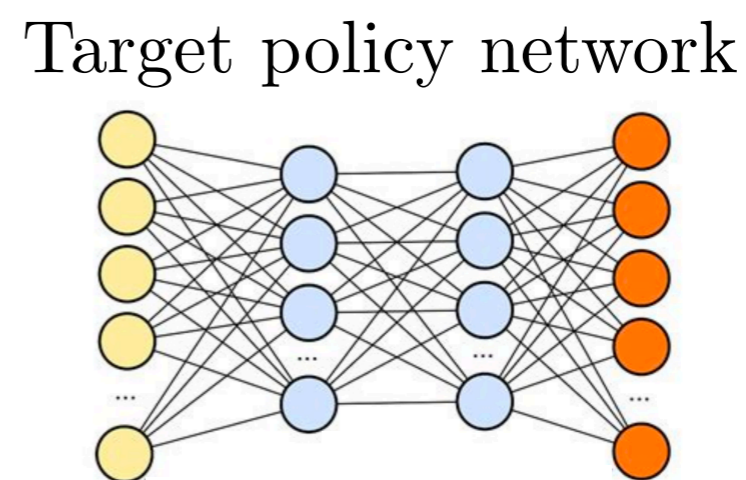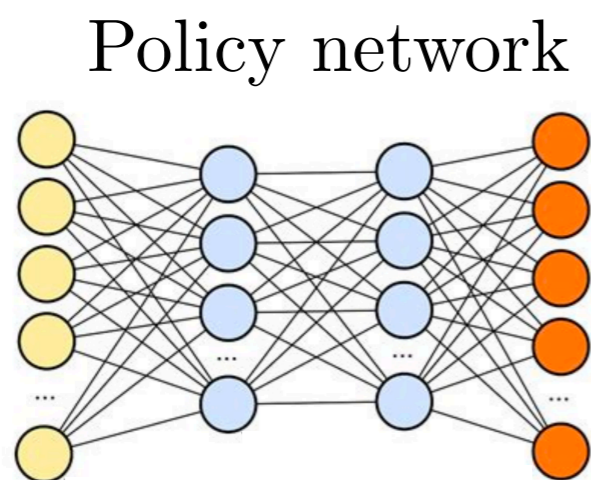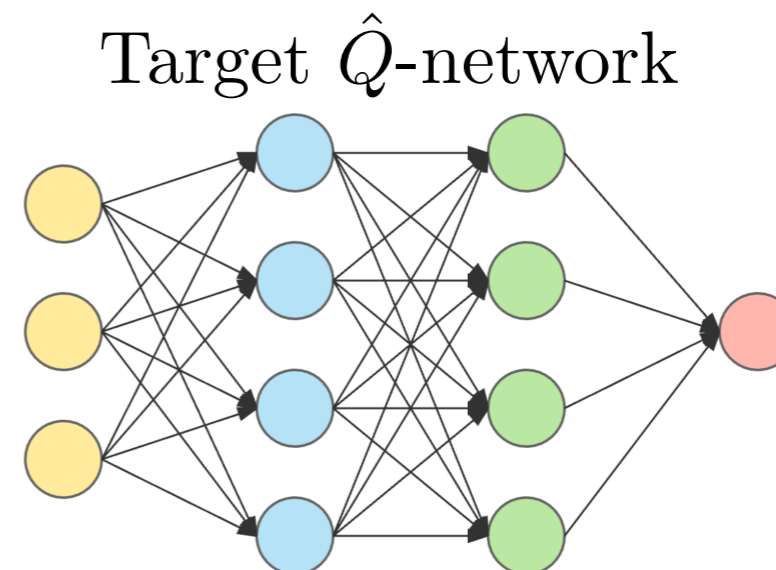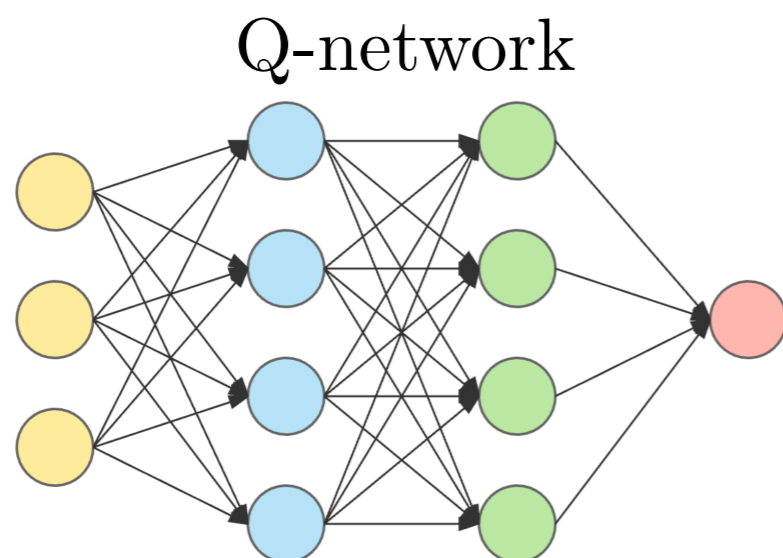[Lillicrap et al., ICML, 2016]

Policy network (actor) - Q-network (critic)
DDPG => Same as DQN + policy network

Q-network



Target $\hat{Q}$-network



Policy network



Target policy network

Initialize replay memory D of size $N$

Initialize Q- and policy networks with random weights $\theta^Q$ and $\theta^\pi$

Set target networks $\theta^{Q'} = \theta^Q$ and $\theta^{\pi'} = \theta^\pi$

For each episode:

Start from an initial state $S_0$

Loop for each step $t$ of the episode:

Choose $A_t = \pi(S_t) + noise$ (to explore a bit)

Take action $A_t$ and observe cost $C_t(S_t, A_t)$ and next state $S_{t+1}$

Store $(S_t, A_t, C_t, S_{t+1})$ in memory D

Sample minibatch of transitions $(S_j, A_j, C_j, S_{j+1})$ from memory D

Gradient descent on $\theta^Q$ to minimize $||Q(S_j, A_j, \theta^Q) - C_j - \gamma Q'(S_{j+1}, \pi'(S_{j+1}))||^2$

Policy update $\nabla_{\theta^\pi} J \simeq \frac{1}{N} \sum_i \nabla_A Q(S, A, \theta^Q)|_{S=S_i, A=\pi(S_i)} \nabla_{\theta^\pi} \pi(S_i, \theta^\pi)$
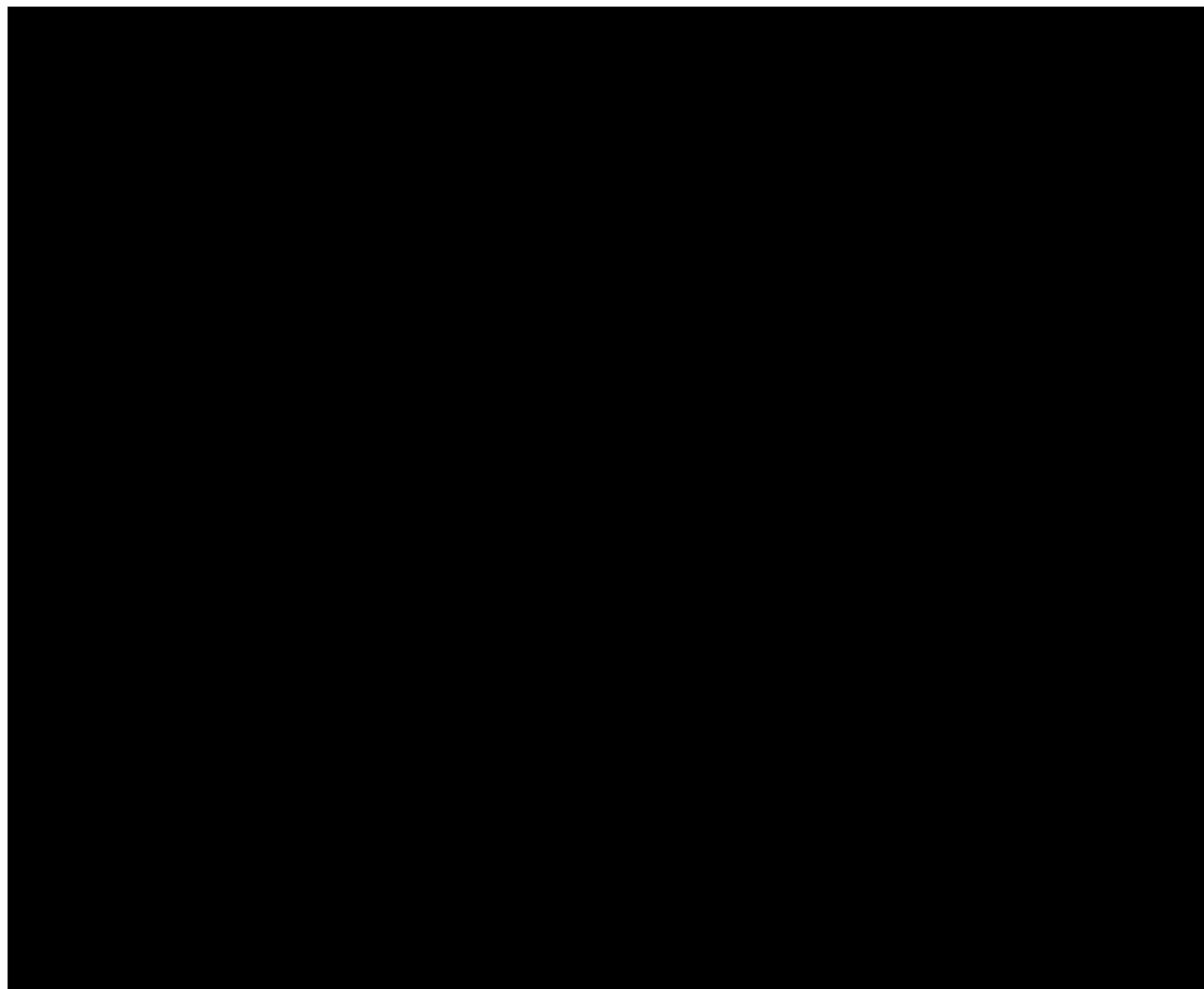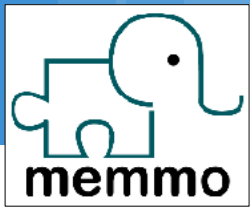
Smooth update of target networks $\begin{aligned} \theta^{Q'} &\leftarrow \tau\theta^Q + (1-\tau)\theta^{Q'} \\ \theta^{\pi'} &\leftarrow \tau\theta^\pi + (1-\tau)\theta^{\pi'} \end{aligned}$
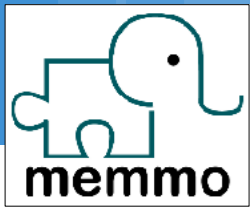
All the presented algorithms are all variations on the same theme: Use <u>Bellman Equations</u> to find iterative algorithms

Severe limitations still exist:
- Algorithms difficult to tune to ensure convergence
- Need lots of samples (not practical on real robots)
- Not clear how to efficiently explore
- Robots can break during learning
- Does not generalize (fixed policy/Q-function)

Learn only the Q-function (Q-learning) => DQN (Atari games)

Learn Q-function + policy function (actor-critic) => DDPG
Very long history of actor-critic algorithms in robotics:
    [Doya, Neural Computation 2000]
    Natural Actor-Critic [Peters et al. 2008]

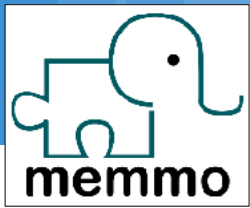Learn directly policy (use policy gradient) => TRPO or PPO
Also long history of policy gradient methods:
    REINFORCE [Williams, 1992]
    [Doya, Neural Computation 2000]
    PI2 [Theodorou et al., JMLR 2010]

Break
Let's add a critic and test DDPG