

Reinforcement Learning Workshop

Day 3

Guido Novati

CSElab

Computational Science & Engineering Laboratory

ETH

Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

What is (Deep) Supervised Learning (DL)

- DL assumes a dataset and/or a data-generating process.
- DL assumes a parameterized model (neural network).
- Typically the model is designed to have easily-computable analytical gradients.
- Objective is:
 - to **learn** to recognize **features of the data** (a.k.a. representation learning)
 - such that **given a datapoint** as input
 - the model's **output minimizes** a prescribed **loss function**.

Parameterized representations of the data

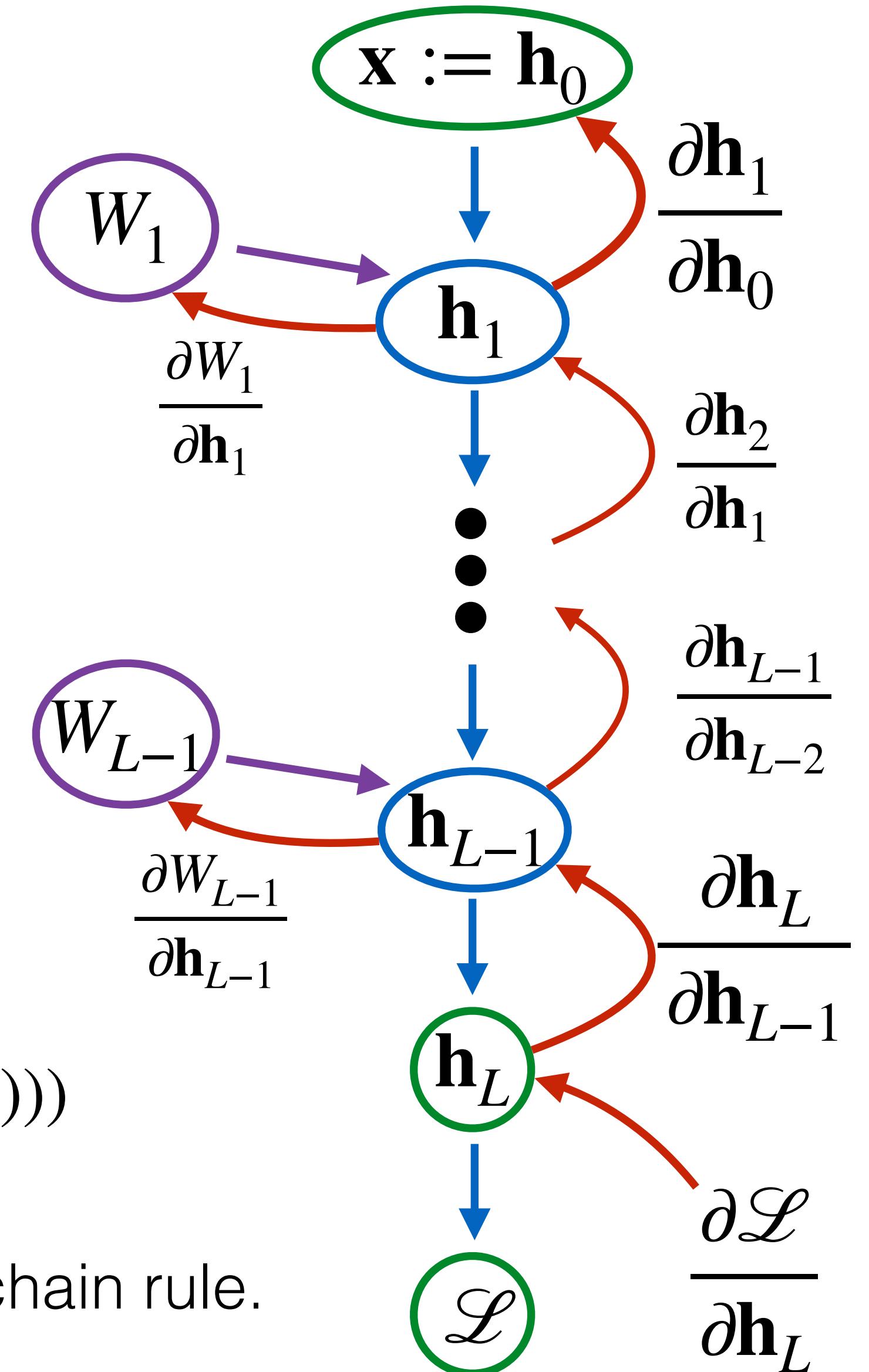
Typically composition of many functions:

- 1) Linear transformations: $\mathbf{h}_{l+1} = W_{l+1} \mathbf{h}_l + \mathbf{b}_{l+1}$
- 2) Non-linear activation functions: $\mathbf{h}_{l+2} = F_{l+2}(\mathbf{h}_{l+1})$
- 3) Loss function on the output

- Mean squared error $\hat{\mathcal{L}}_w^{MSE}(\mathbf{x}, \mathbf{y}) = [\mathbf{y} - F_w(\mathbf{x})]^2$
- Log-likelihood $\hat{\mathcal{L}}_w^{LL}(\mathbf{x}, \mathbf{y}) = -\log \mathbb{P}[F_w(\mathbf{x}) \mid \mathbf{y}]$

Note: $w = \{W_i, \mathbf{b}_i\}_{i \in [1, L]}$ and $F_w(\mathbf{x}) = F_L(W_{L-1}F_{L-2}(\dots F_2(W_1\mathbf{x})))$

back-propagation: compute gradient of loss w.r.t parameters by chain rule.



Deep Learning in practice

Objective: find a local minimum of expected loss on distribution of data:

$$w^* = \arg \min_w \mathcal{L}(w) = \arg \min_w \mathbb{E} \left[\hat{\mathcal{L}}_w(x, y) \middle| \begin{array}{l} x \sim X \\ y \sim Y \end{array} \right]$$

- Estimate loss function through Monte Carlo sampling: $\mathcal{L}(w) \approx \frac{1}{B} \sum_{i=1}^B \hat{\mathcal{L}}_w(x_i, y_i)$
- For each sample, compute gradient of the loss w.r.t. parameters $\frac{\partial \hat{\mathcal{L}}_w(x_i, y_i)}{\partial w}$
- Update the parameters with stochastic gradient descent: $\Delta w = - \frac{\epsilon}{B} \sum_{i=1}^B \frac{\partial \hat{\mathcal{L}}_w(x_i, y_i)}{\partial w}$

Intuition: linear “Neural Network”

One linear layer from input to output:

$$h_w(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b = \sum_{i=1} x_i w_i + b$$

Mean squared error:

$$\mathcal{L}_w^{MSE}(\mathbf{x}, y) = \frac{1}{2} [y - h_w(\mathbf{x})]^2$$

Back-propagation:

$$\boxed{\frac{\partial h_w(\mathbf{x})}{\partial w_i} = x_i}$$

$$\boxed{\frac{\partial \mathcal{L}_w^{MSE}(\mathbf{x}, y)}{\partial h_w} = (h_w(\mathbf{x}) - y)}$$

single datapoint SGD: $\Delta \mathbf{w} = -\epsilon \frac{\partial \mathcal{L}_w^{MSE}(\mathbf{x}, y)}{\partial \mathbf{w}} = -\epsilon \frac{\partial \mathcal{L}_w^{MSE}(\mathbf{x}, y)}{\partial h_w} \frac{\partial h_w(\mathbf{x})}{\partial \mathbf{w}} = -\epsilon [h_w(\mathbf{x}) - y] \mathbf{x}$

in general: SGD: $\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - \epsilon \nabla_{\mathbf{w}} \mathbb{E}_{X, Y} [\hat{\mathcal{L}}_w^{MSE}(\mathbf{x}, \mathbf{y})]$

move network output towards target y

Chain rule: $\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - \epsilon \mathbb{E}_{X, Y} [\nabla_{F_w} \hat{\mathcal{L}}_w^{MSE}(\mathbf{x}, \mathbf{y}) \cdot \nabla_{\mathbf{w}} F_w(\mathbf{x})]$

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + \epsilon \mathbb{E}_{X, Y} [(\mathbf{y} - F_w(\mathbf{x})) \cdot \nabla_{\mathbf{w}} F_w(\mathbf{x})]$$

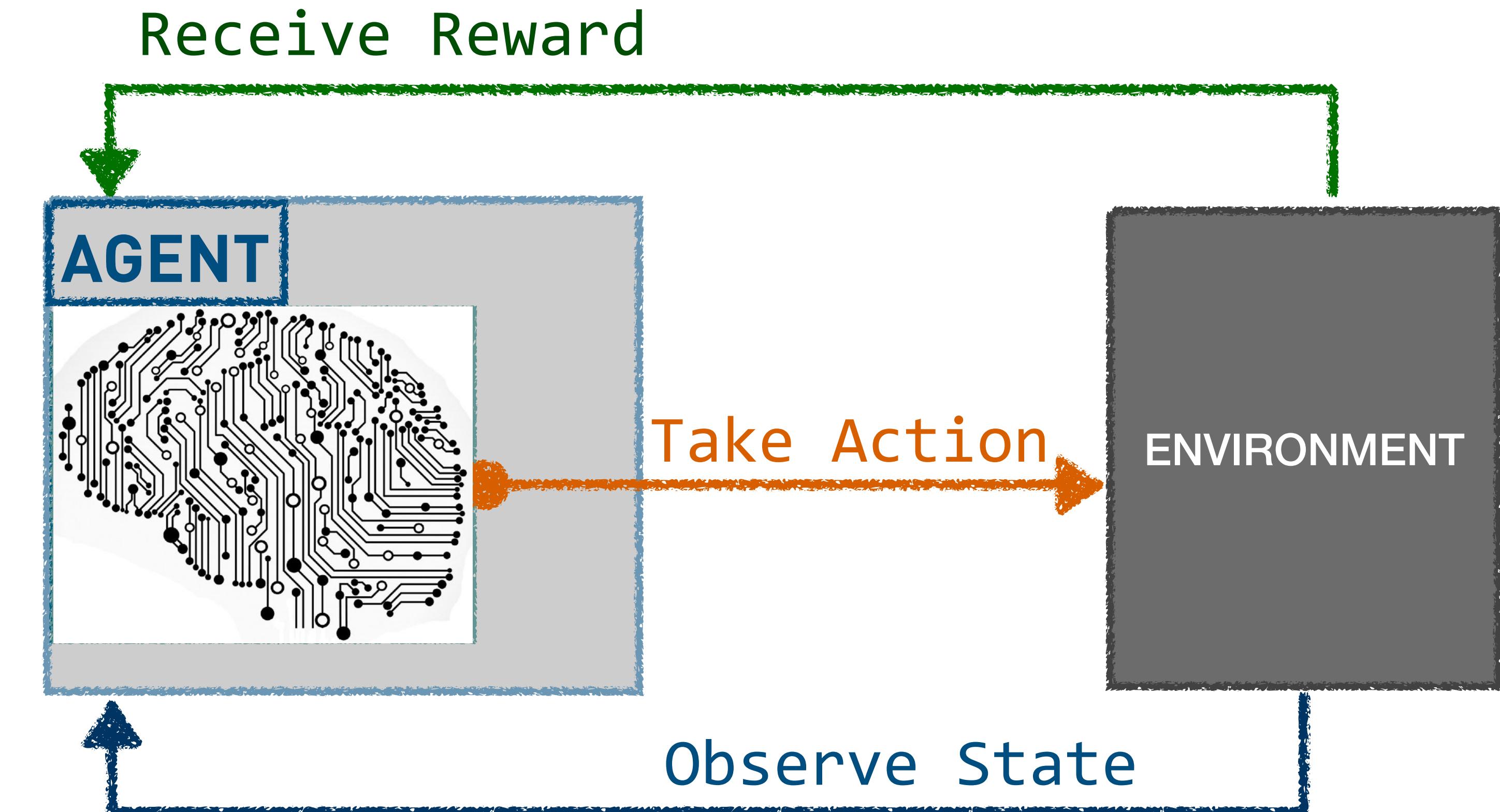
Today's Plan

- What is Deep Supervised Learning and what is Reinforcement Learning
- Introduction to fundamental RL assumptions & tradeoffs
- Core RL concepts:
 - Policy, Value functions, approximate models
- RL via approximating value functions:
 - Monte Carlo, SARSA, Q-learning
- Hands-on RL with pytorch

Reinforcement Learning's loop

Instead of trying to produce a program to simulate the adult mind, why not rather try to produce one which simulates the child's? If this were then subjected to an appropriate course of education one would obtain the adult brain.

Alan Turing,
Computing Machinery and Intelligence



What is Reinforcement Learning

- Mathematical framework to **search optimal strategies** for sequential decision processes.
- RL describes one (or more) **agent** with the capacity to observe its **environment** and act.
- Action influences the state of the env and/or agent's ability to perform future actions.
- Success is measured by total sum of rewards obtained by the agent.
- Reward hypothesis: “Any goal can be described by maximization of expected rewards”

The challenge of RL

- **Exploration/exploitation dilemma**
 - Exploitation : Make the best decision given current information
 - Exploration : Gather more information
 - The best long-term strategy may involve short-term sacrifices.
 - When have we gathered enough information to make the best overall decisions?
- **Temporal credit assignment**
 - Temporal structure of the optimization problem
 - Goodness of action measured over future consequences
 - Which sequence of actions caused the high/low rewards?
- **Representation learning**
 - Deep learning challenge : state may be sensors or sight (CNN).
 - Optimal strategy may require remembering facts from the past (RNN).
 - Learn to approximate probability distributions over arbitrarily large state and action spaces.

What are states / actions / rewards?

The **state** $s_t \in \mathcal{S}$ can be either:

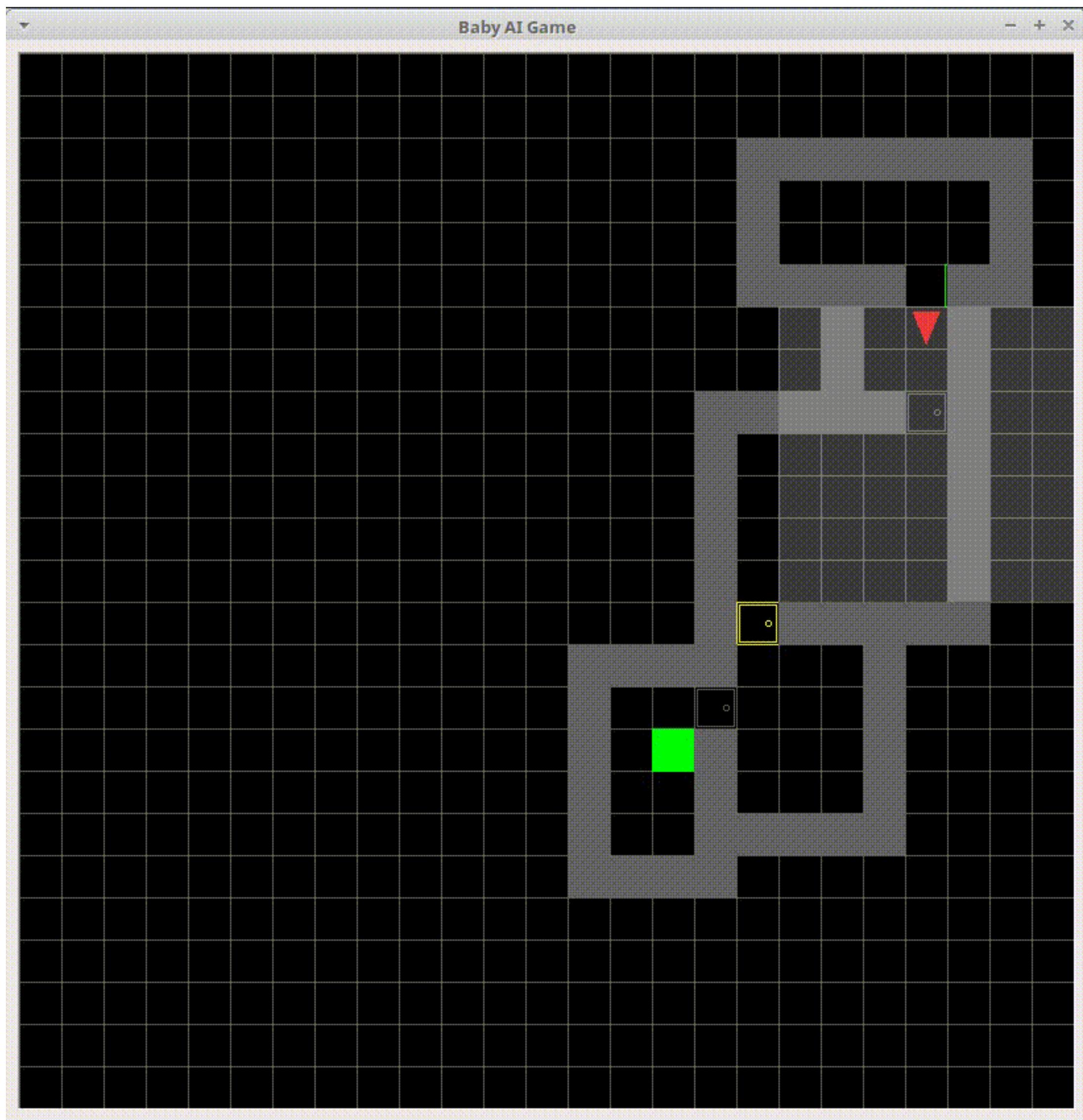
- **discrete** $\mathcal{S} = \{s^{(0)}, s^{(1)}, s^{(2)}, \dots, s^{(N_S)}\}$ i.e. corresponds to one out of N_S alternatives
- **continuous** $\mathcal{S} = \mathbb{R}^{\dim_S}$ i.e. is a vector in a \dim_S space

The **action** $a_t \in \mathcal{A}$ can be either:

- **discrete** $\mathcal{A} = \{a^{(0)}, a^{(1)}, a^{(2)}, \dots, a^{(N_A)}\}$ (policy may be a Bernoulli distribution)
- **continuous** $\mathcal{A} = \mathbb{R}^{\dim_A}$ (policy may be a Gaussian distribution)

The reward $r_t \in \mathbb{R}$ is a scalar.

RL with discrete actions and states

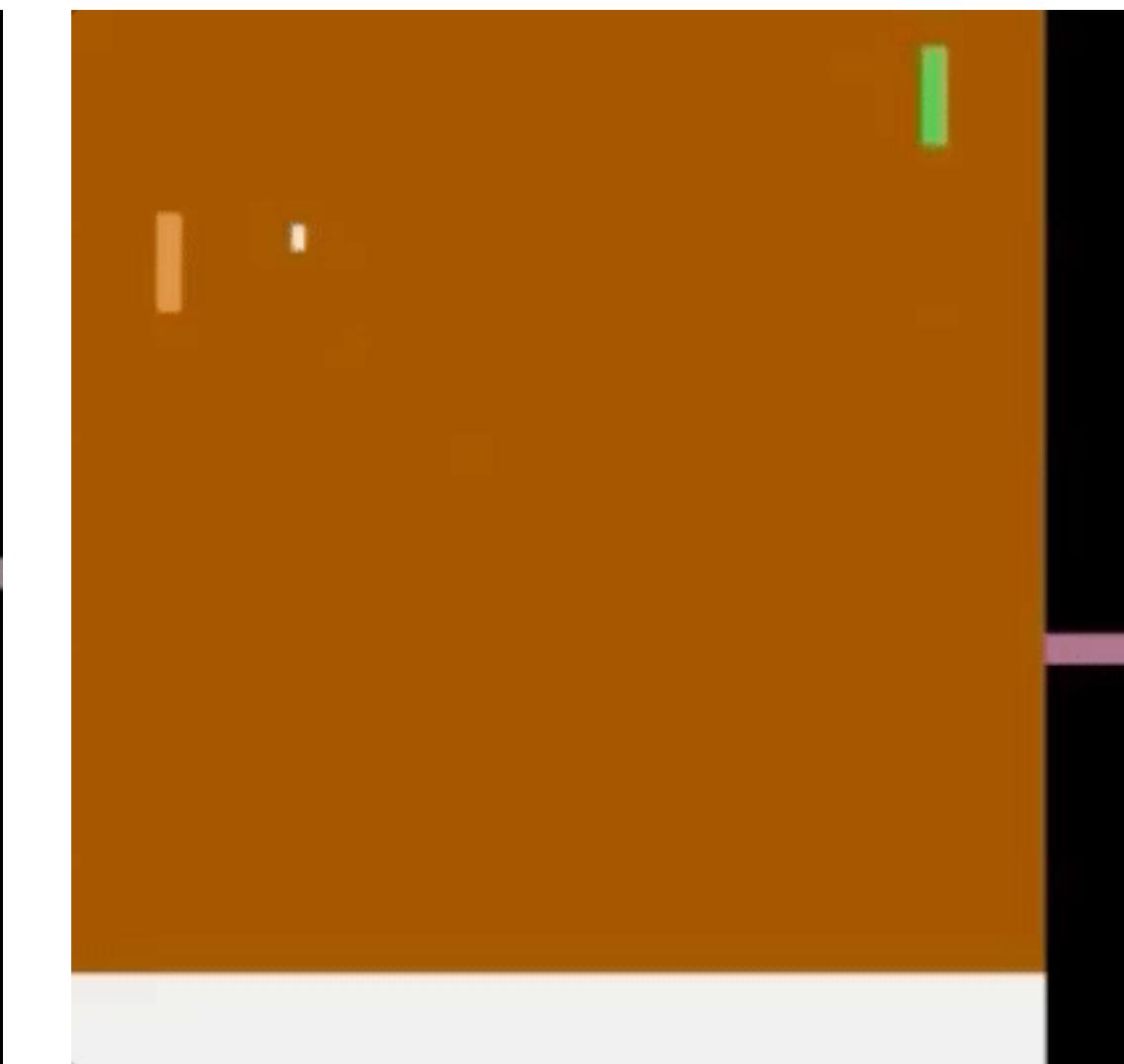


- First applications of RL: grid-world scenarios
- Often simplified by **take Action == Go to State**
- Often amenable to dynamic programming
- RL + deep learning allow mitigating curse of dimensionality and recent breakthroughs in classic games with large state spaces:

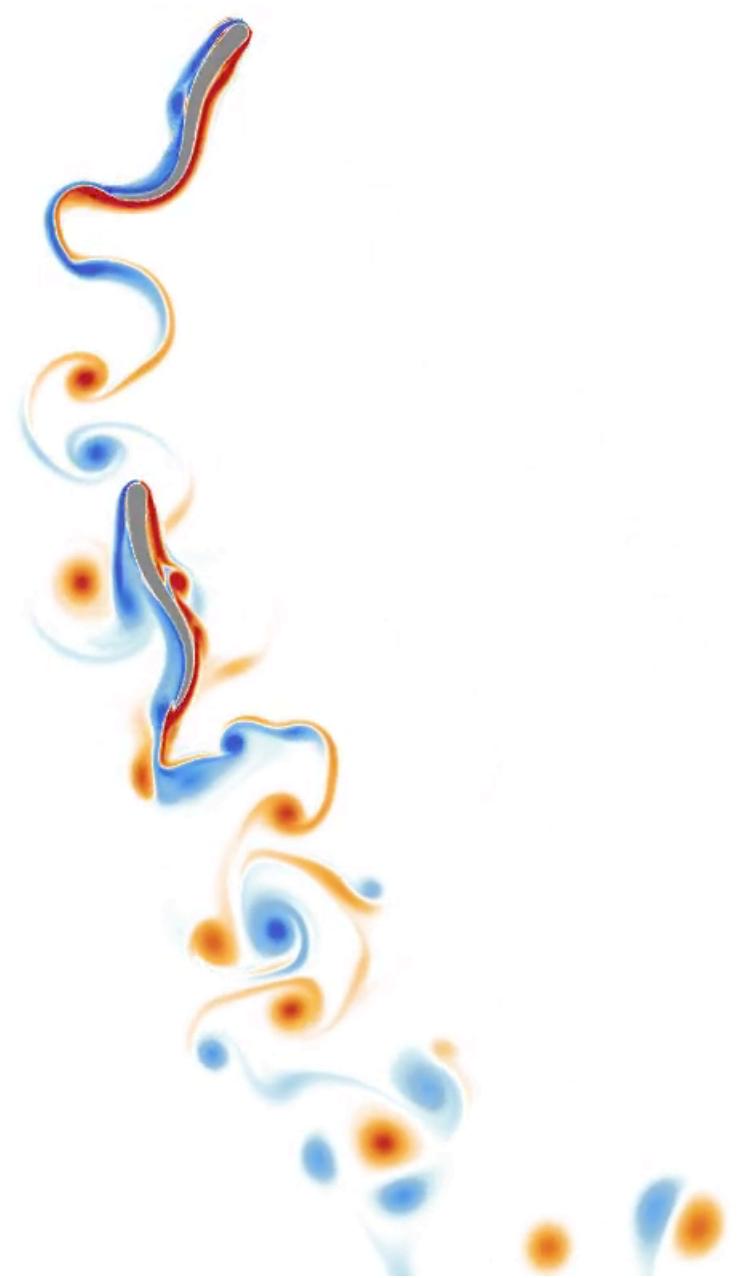
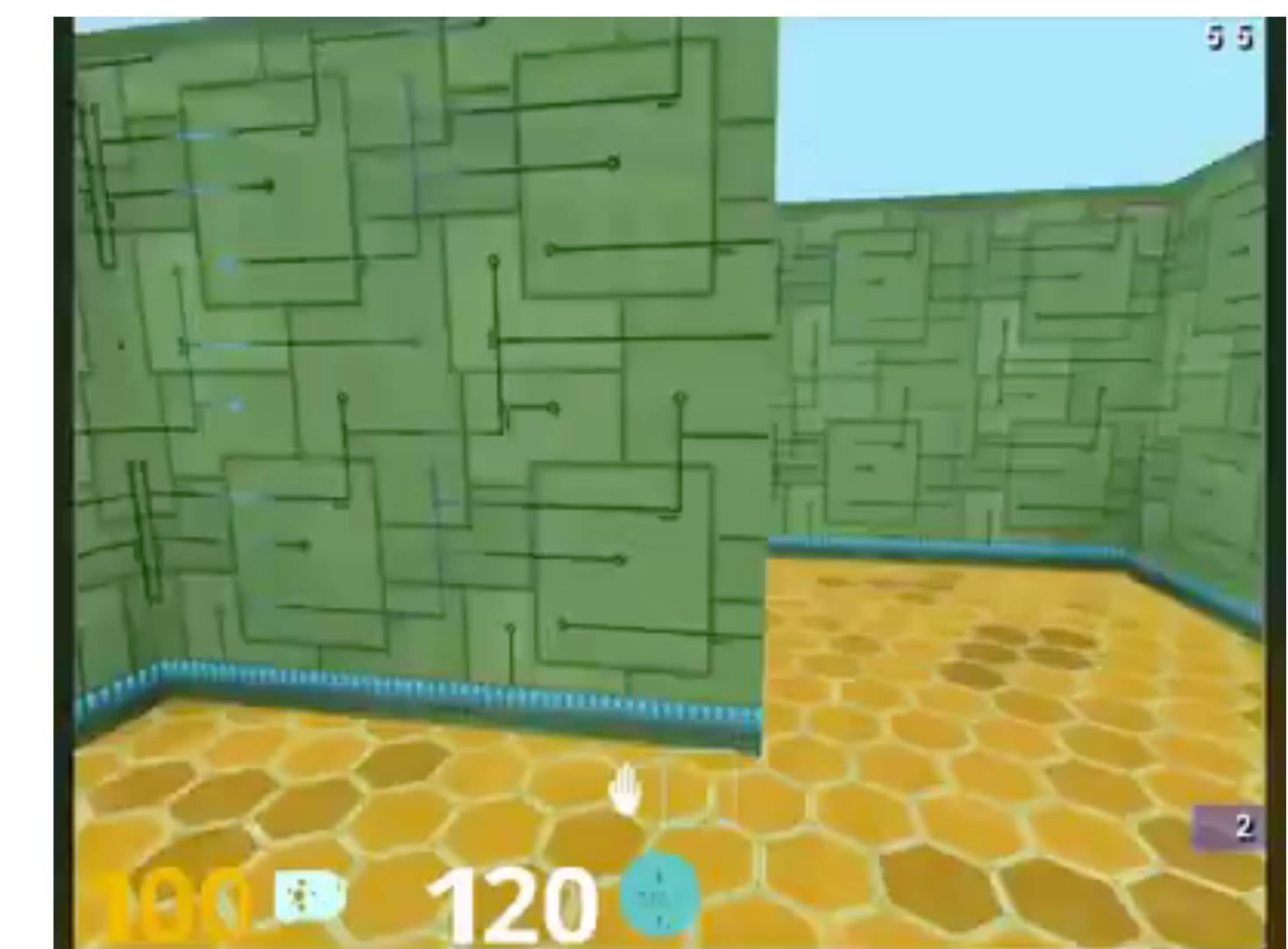


RL with discrete actions and continuous states

Mnih et al. 2015



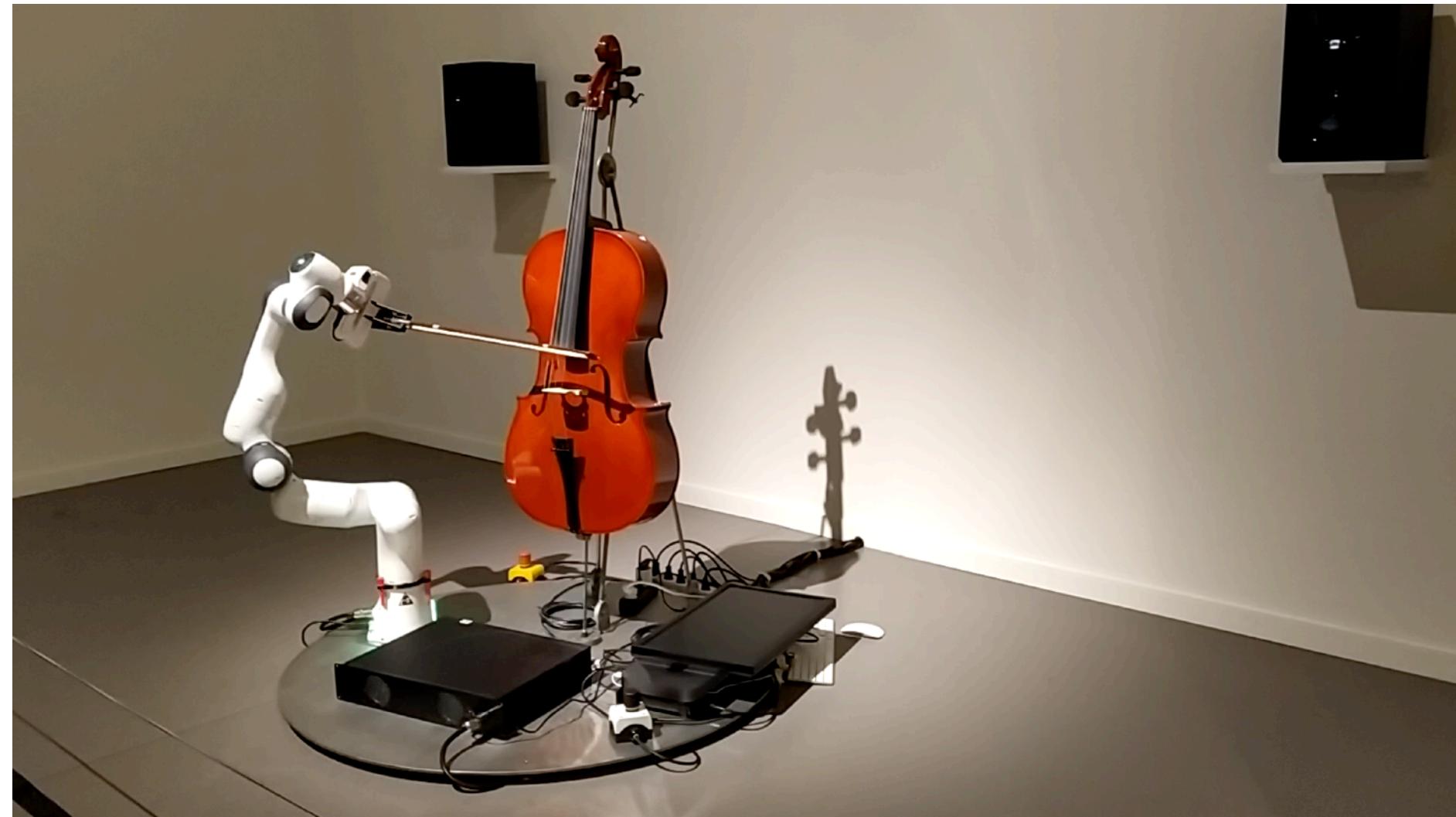
Jaderberg et al. 2016



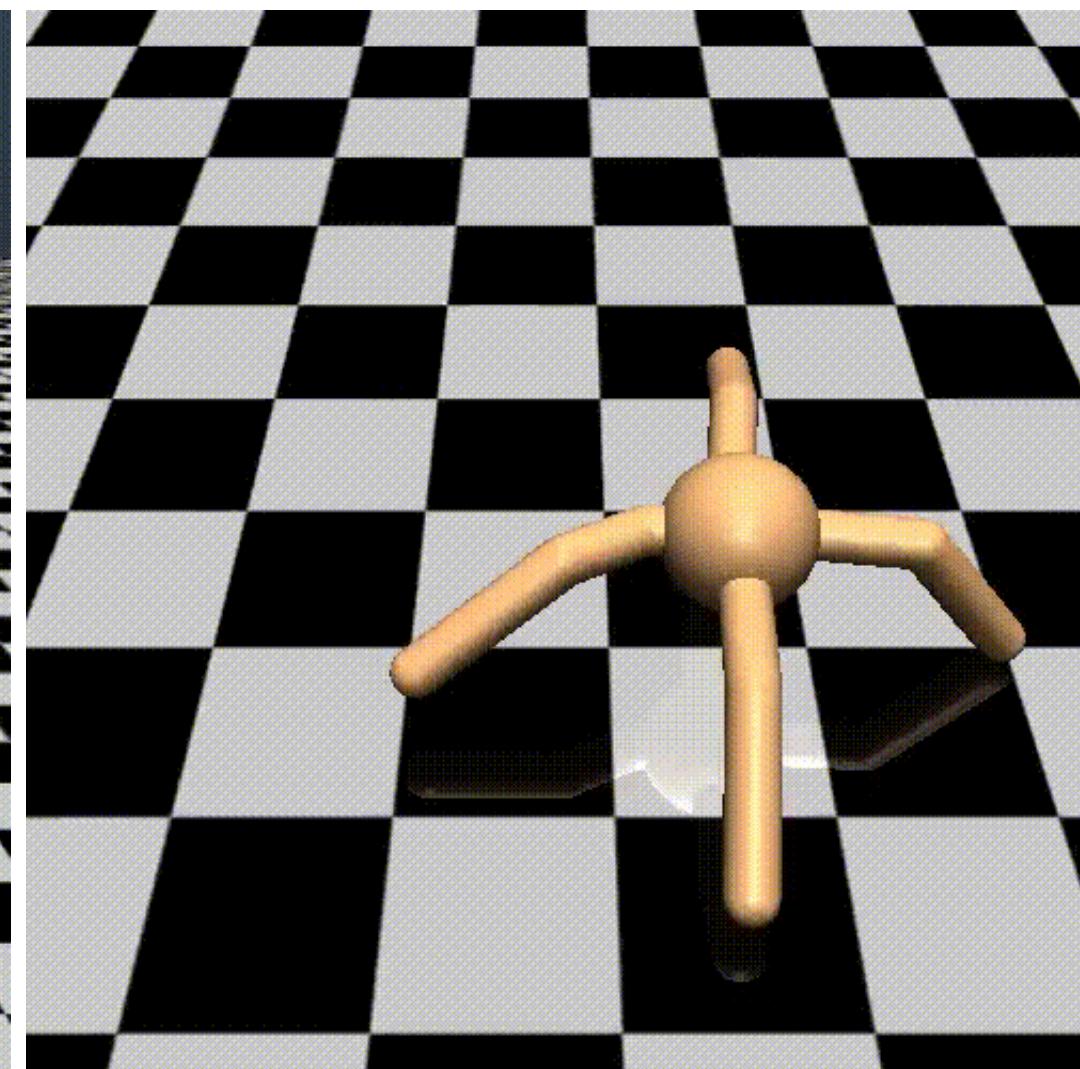
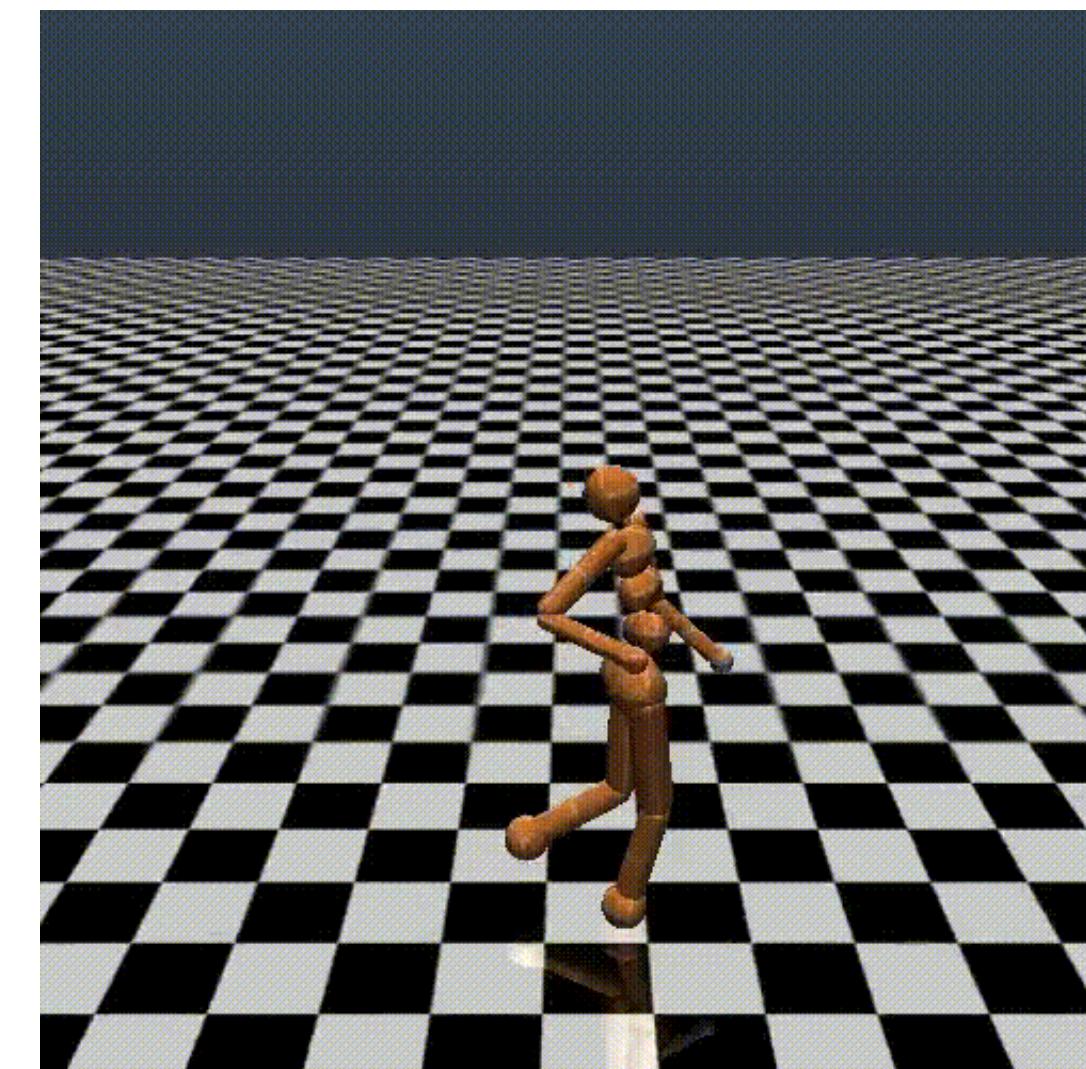
RL with continuous actions and continuous states



Abbeel et al. 2006

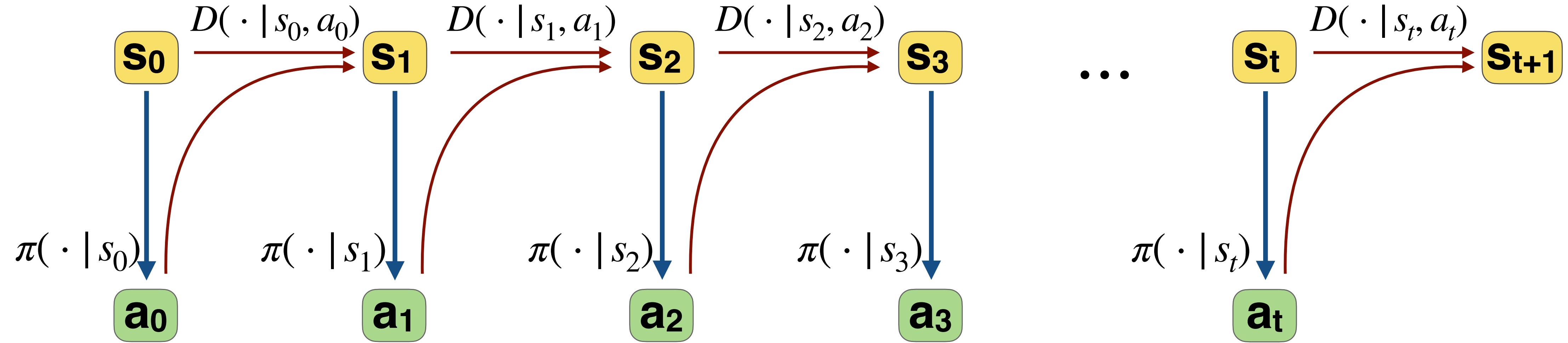


Schulman 2015



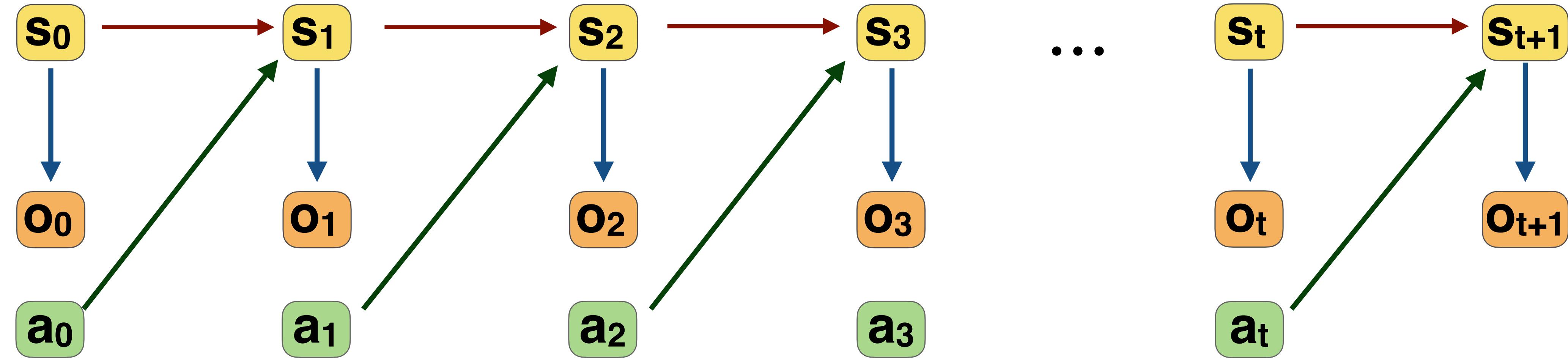
- Most PDE / ODE / robotic systems / real world applications

Markov decision process (MDP)



- Probability of sampling an action a_t depends only on state s_t (may be deterministic)
- Probability of transitioning to state s_{t+1} depends on both s_t and a_t (may be deterministic)
- Reward r_{t+1} is in general function of (s_t, a_t, s_{t+1})
- “The future is independent of the past given the present” : $\mathbb{P}(s_{t+1} | s_t, s_{t-1}, \dots, s_0) = \mathbb{P}(s_{t+1} | s_t)$
- In order to pick optimal action, agent only needs to know current state.
- This is said to be a *fully observable* problem: state fully describes environment.

Partially observable Markov decision process (POMDP)



- Agent receives only partial information about the state of the environment : **observation**
- Probability of current state can only be inferred by looking at all “sensimotor history”:

$$\mathbb{P}(s_t) = \mathbb{F}(o_0, a_0, o_1, a_1, \dots, o_t)$$

- Policy $\pi(\cdot | o_t)$ is probably sub-optimal : does not include enough information to estimate the state
- Optimal policy is function of entire history : $a_t \sim \pi(\cdot | s_t) = \pi(\cdot | o_0, a_0, o_1, a_1, \dots, o_t)$
- Recurrent Neural Networks (RNN) may be used to encode time history in the policy

Variants of the RL problem

- Multi-armed bandit problem



Imagine a gambler at a row of slot machines, who has to decide which machines to play, how many times to play each machine and in which order to play them.

- RL without neither states nor dynamics.
- Purely exploration/exploitation dilemma.

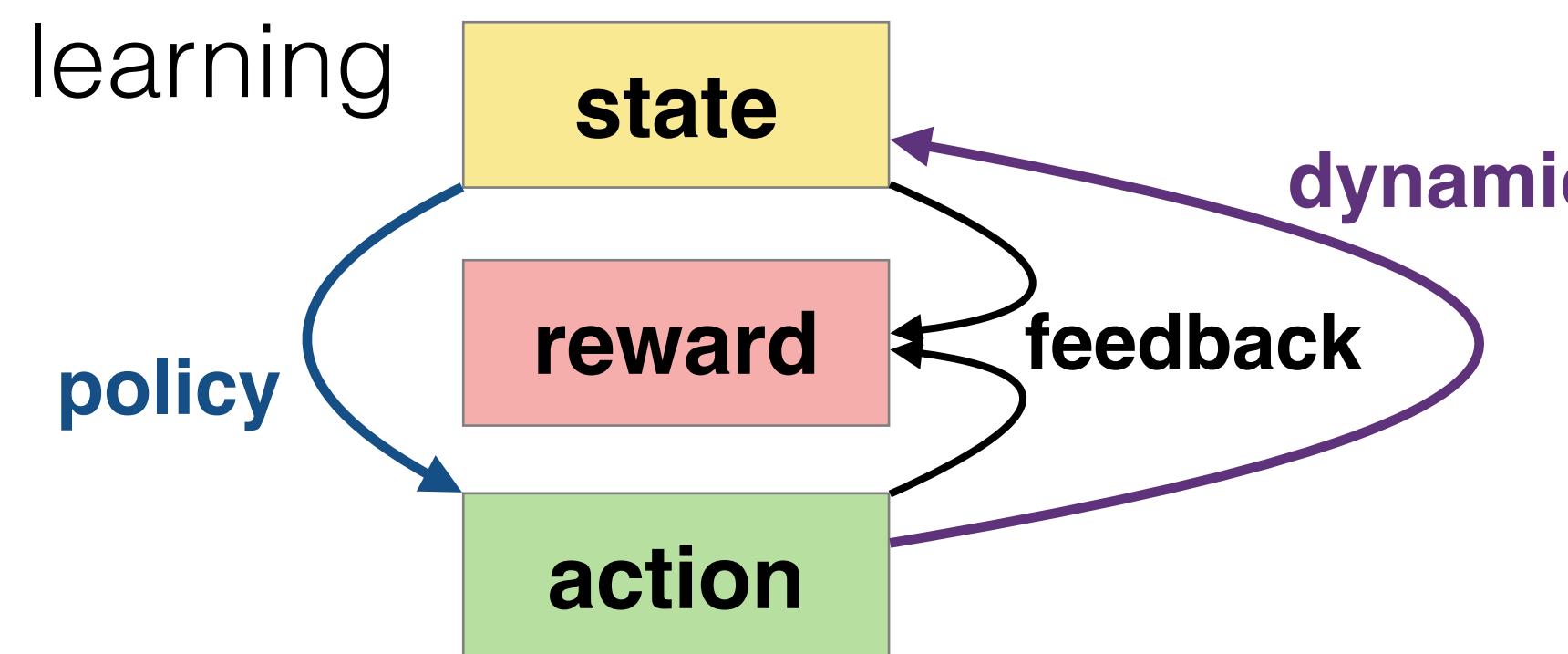
- Contextual bandits problem



Imagine that slot machines have clues (e.g. brand).

- RL without dynamics.
- Exploration/exploitation dilemma
- Feature learning challenge

- Reinforcement learning

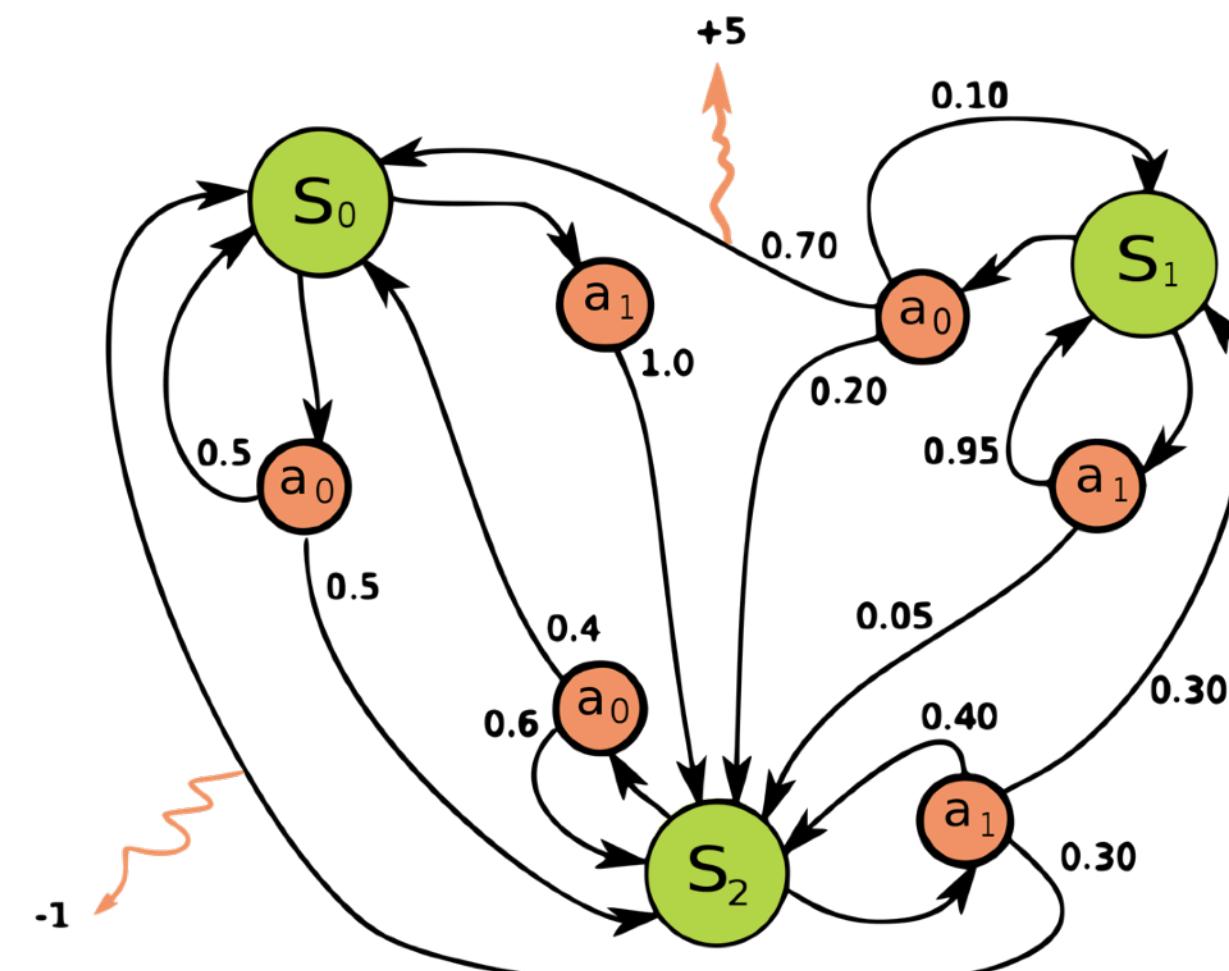


- Environment dynamics.
- Exploration/exploitation dilemma
- Feature learning challenge
- Temporal credit assignment problem

Differences from neighboring fields

Dynamic programming (DP)

- Finite-dimensional problems (i.e. discrete)
- Available model of the environment:
 - $D(s_{t+1} | s_t, a_t)$ is a known transition matrix.
 - “*for each state do ...*”
- e.g. shortest path / resource allocation



Optimal control (OC)

- Sets of ODEs describing the state of the system and a cost functional to minimize.
- May solve Hamilton–Jacobi–Bellman (HJB) equations
- e.g. linear quadratic controller (linear ODE with quadratic actuation cost)

$$\begin{cases} \dot{\mathbf{x}}(t) = A\mathbf{x}(t) + B\mathbf{u}(t) \\ J = \mathbf{x}^T(T)F\mathbf{x}(T) + \int_0^T [\mathbf{x}^T(t)Q\mathbf{x}(t) + \mathbf{u}^T(t)R\mathbf{u}(t)] dt \end{cases}$$

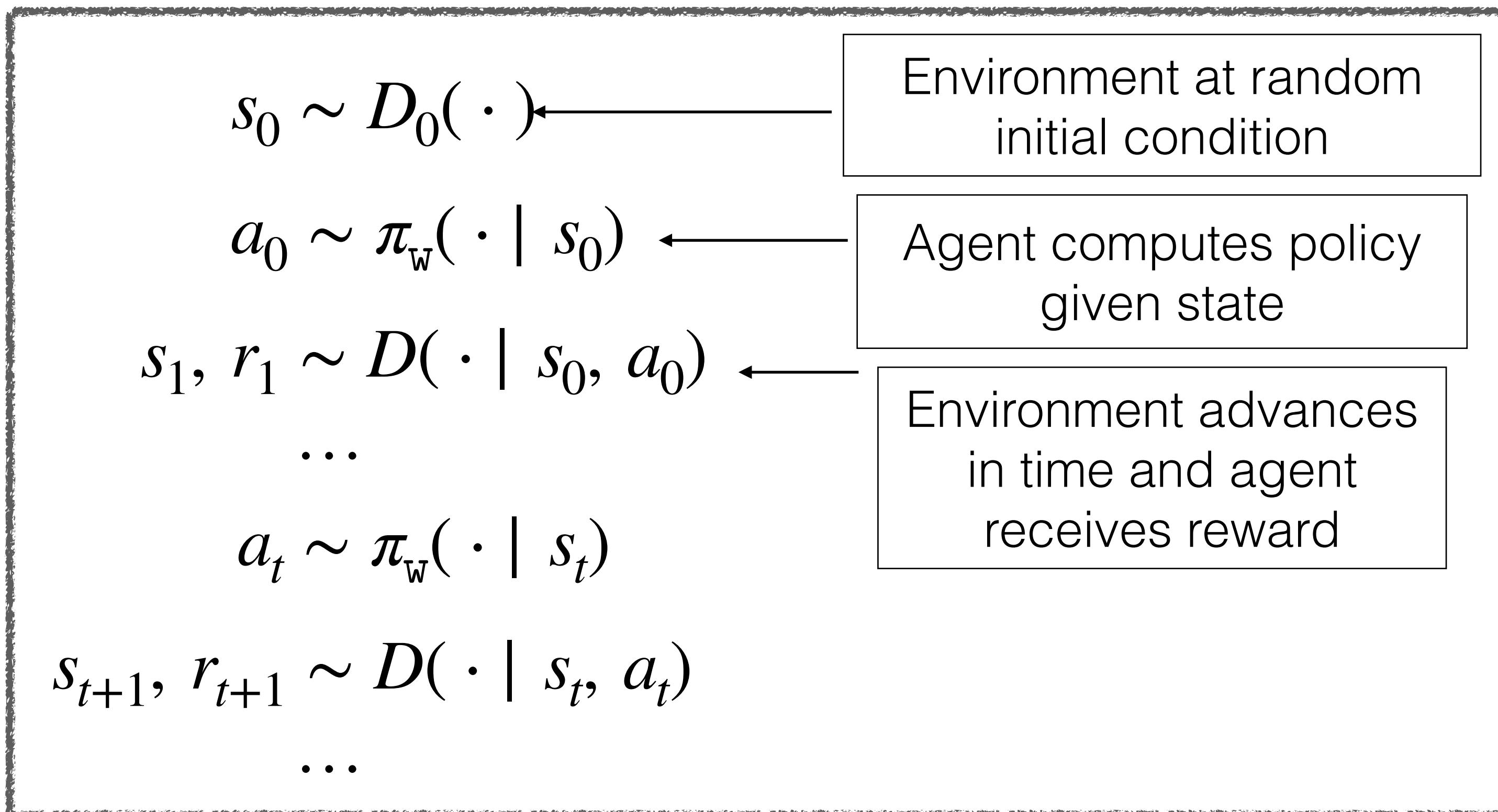
optimal actuation has form:
 $\mathbf{u}(t) = K\mathbf{x}(t)$

RL is characterized by **trial-and-error interaction** with an unknown (**black-box**) environment.

RL: Markov decision process with parameterized policy

RL learns a policy

- Function that computes actions given states.
- **Improvable** i.e. it has parameters w which can be optimized.
- Some stochasticity is required to explore dynamics and possible improvements.



$$J(w) = \mathbb{E} \left[\sum_t^\infty r_t \mid \begin{array}{l} a_t \sim \pi_w(\cdot | s_t) \\ s_{t+1} \sim D(\cdot | a_t, s_t) \end{array} \right]$$

optimal parameters

$$w^* = \arg \max_w J(w)$$

Value functions: estimate of future rewards given a policy

- Q-function “What are expected reward by starting in (s,a) and following policy?”

$$Q^{\pi_w}(s, a) = \mathbb{E} \left[r_1 + r_2 + r_3 + \dots \middle| \begin{array}{ll} s_0 = s & a_t \sim \pi_w(\cdot | s_t), \\ a_0 = a & s_{t+1}, r_{t+1} \sim D(\cdot | s_t, a_t) \end{array} \right]$$

- This estimator is **unattainable!**
- In RL we may never “set s_0 to s and set a_0 to a ” and follow the policy.
- In fact we don’t know if we can ever reproduce/see again a given s .
- We can only keep interacting with the environment.

Value functions: estimate of future rewards given a policy

$$Q^{\pi_w}(s, a) = \mathbb{E} \left[r_1 + r_2 + r_3 + \dots \middle| \begin{array}{ll} s_0 = s & a_t \sim \pi_w(\cdot | s_t), \\ a_0 = a & s_{t+1}, r_{t+1} \sim D(\cdot | s_t, a_t) \end{array} \right]$$

- May include discount $\gamma \in [0, 1]$ to bound expectation:

$$Q^{\pi_w}(s, a) = \mathbb{E} \left[r_1 + \gamma r_2 + \gamma^2 r_3 + \dots \middle| \begin{array}{ll} s_0 = s & a_t \sim \pi_w(\cdot | s_t), \\ a_0 = a & s_{t+1}, r_{t+1} \sim D(\cdot | s_t, a_t) \end{array} \right]$$

- γ affects optimal policy!
 - γ close to 0: optimization is **shortsighted**, but **easier** and more **robust** to uncertainty.
 - γ close to 1: **exact**, non-relaxed, evaluation, but may be unbounded ($Q \rightarrow \infty$).
- Decomposes into Bellman Equation:

$$Q^{\pi_w}(s, a) = \mathbb{E} \left[r + \gamma Q^{\pi_w}(s', a') \middle| \begin{array}{l} s', r \sim D(\cdot | s, a) \\ a' \sim \pi_w(\cdot | s') \end{array} \right]$$

Value functions: estimate of future rewards given a policy

- (γ -)Optimal policy maximizes value function everywhere:

$$\pi^*(s) = \arg \max_a Q^{\pi^*}(s, a)$$

- Similarly, optimal Q-function is maximal for each (s,a) pair:

$$Q^*(s, a) = \max_{\pi} Q^{\pi}(s, a) = Q^{\pi^*}(s, a) \quad \text{for each } (s, a)$$

- Bellman Equation for optimal Q-function:

$$Q^{\pi^*}(s, a) = \mathbb{E} \left[r + \gamma Q^{\pi^*}(s', a') \middle| \begin{array}{l} s', r \sim D(\cdot | s, a) \\ a' \sim \pi^*(\cdot | s') \end{array} \right]$$

$$Q^{\pi^*}(s, a) = \mathbb{E} \left[r + \gamma \max_{a'} Q^{\pi^*}(s', a') \middle| s', r \sim D(\cdot | s, a) \right]$$

Three value functions

- **Q-function:** “What are expected rewards by starting in (s,a) and following policy?”

$$Q^{\pi_w}(s, a) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_{t+1} \mid \begin{array}{l} s_0 = s, \\ a_0 = a, \quad s_{t+1}, r_{t+1} \sim D(\cdot \mid s_t, a_t) \end{array} \right]$$

- **State value:** “What are expected rewards by starting in s and following policy?”

$$V^{\pi_w}(s) = \mathbb{E} \left[Q^{\pi_w}(s, a) \mid a \sim \pi_w(\cdot \mid s) \right] = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_{t+1} \mid \begin{array}{l} s_0 = s, \\ a_0 = \pi_w(\cdot \mid s_0), \quad s_{t+1}, r_{t+1} \sim D(\cdot \mid s_t, a_t) \end{array} \right]$$

- **Action advantage:** “How much better is it to do action a compared to expected rewards from π ?”

$$A^{\pi_w}(s, a) = Q^{\pi_w}(s, a) - V^{\pi_w}(s)$$

- Therefore:

$$\mathbb{E} \left[A^{\pi_w}(s, a) \mid a = \pi_w(\cdot \mid s) \right] = 0$$

Episodes: an extension to the classic RL loop

Instead of infinite interaction loop (continuing tasks) much of RL analysis assumes that loop is split into episodes (episodic tasks).

classic RL loop

$$s_0 \sim D_0(\cdot)$$

$$a_t \sim \pi_w(\cdot | s_0)$$

$$s_1, r_1 \sim D(\cdot | s_0, a_0)$$

...

$$a_t \sim \pi_w(\cdot | s_t)$$

$$s_{t+1}, r_{t+1} \sim D(\cdot | s_t, a_t)$$

...

∞

episode-based loop

∞

...

$$a_{T-1}^{(k)} \sim \pi_w(\cdot | s_{T-1}^{(k)})$$

$$s_T^{(k)}, r_T^{(k)} \sim D(\cdot | s_{T-1}^{(k)}, a_{T-1}^{(k)})$$

$$s_0^{(k+1)} \sim D_0(\cdot)$$

$$a_0^{(k+1)} \sim \pi_w(\cdot | s_0^{(k+1)})$$

...

∞

- Episode may last $T^{(k)}$ or T step.
(i.e. number of steps per episode may vary)
- After s_T no actions, sample initial conditions.
- Agent may reach **terminal state** (& terminal reward) due to success/failure in task or to timeout.
- Or reach a **cutoff state**. This is a convenient and arbitrary time horizon ("test policy for at most T steps").
- Episodic tasks allow $\gamma = 1$!

State visitation frequency

Consider exploring the MDP with given π_w :

for $k = 1, 2, \dots, \infty$:

$$s_0^{(k)} \sim D_0(\cdot)$$
$$a_0^{(k)} \sim \pi_w(\cdot | s_0^{(k)})$$
$$s_1^{(k)}, r_1^{(k)} \sim D(\cdot | s_0^{(k)}, a_0^{(k)})$$
$$\dots$$
$$a_{T-1}^{(k)} \sim \pi_w(\cdot | s_{T-1}^{(k)})$$
$$s_T^{(k)}, r_T^{(k)} \sim D(\cdot | s_{T-1}^{(k)}, a_{T-1}^{(k)})$$

What is the probability of a given state s ?

$$s \sim \eta^{\pi_w}(\cdot) \propto \sum_{k=0}^{\infty} \mathbb{P} \left[s = s_k \middle| \begin{array}{l} s_0 \sim D_0(\cdot) \\ a_k \sim \pi_w(\cdot | s_k) \\ s_{k+1} \sim D(\cdot | s_k, a_k) \end{array} \right]$$

- Even if D and/or π_w are deterministic
- As long as at least one between D_0 or π_w are stochastic
- MDP defines a random process that outputs states

Therefore we could write, for example:

$$J(w) = \mathbb{E} \left[Q^{\pi_w}(s, a) \middle| \begin{array}{l} s \sim \eta^{\pi_w}(\cdot) \\ a \sim \pi_w(\cdot | s) \end{array} \right]$$

prob of a state
given the policy

prob of action
given the state

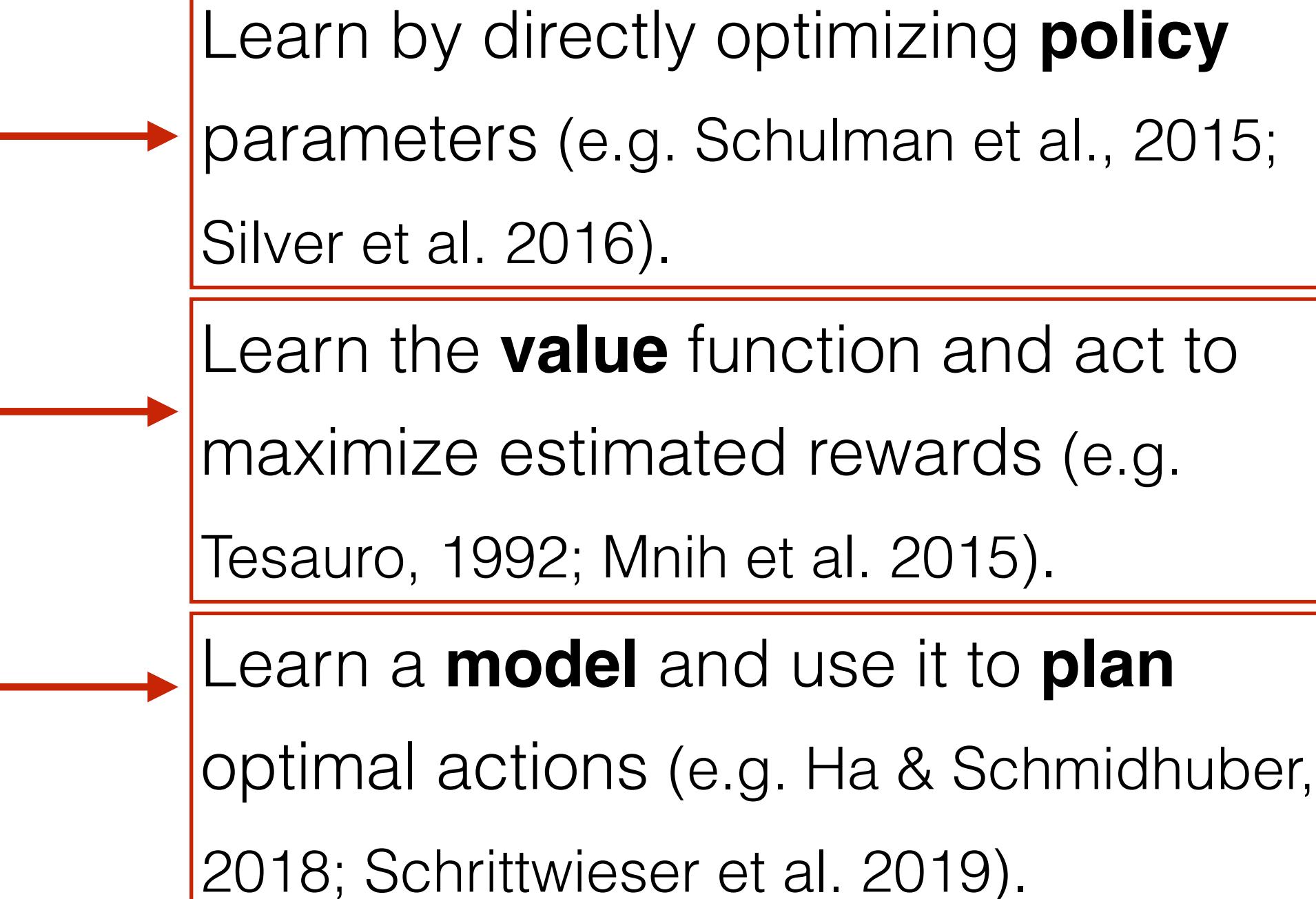
“maximize expected cumulative rewards for any (s, a) reached by the policy”

Components of a RL algorithm

- Common building blocks of RL:

deep learning difficulty ↓

- **Policy**: a function to compute actions given a state
- **Value function**: a function to estimate rewards obtained by a policy
- **Environment model**: a function to estimate environment dynamics.



Why so many RL algorithms?!

- Tradeoffs:

- **Sample efficiency** v.s. **simplicity/ease of use**
- **Stability** v.s. **local optima**
- **Stochastic** v.s. **deterministic**
- Assumptions:
 - **Continuous** v.s. **discrete**
 - **Low dimensional** v.s. **high dimensional**
 - **Fully-observable** v.s. **partially observable**

E.g. **off-policy** vs **on-policy** methods

on-policy: after each update old data is useless

off-policy: re-use data over many iterations

E.g. learn a **policy** vs learn the **value** function

policy: stable, but maybe stuck in local optima

value: with NN is numerically unstable

deterministic (policy/update) algorithms must introduce randomness somewhere to explore

stochastic (policy/update) can naturally learn to explore, but may be less sample-efficient

E.g. **tabular** v.s. **NN** approximators

tabular: more convergence guarantees

NN: reduce curse of dimensionality, generality

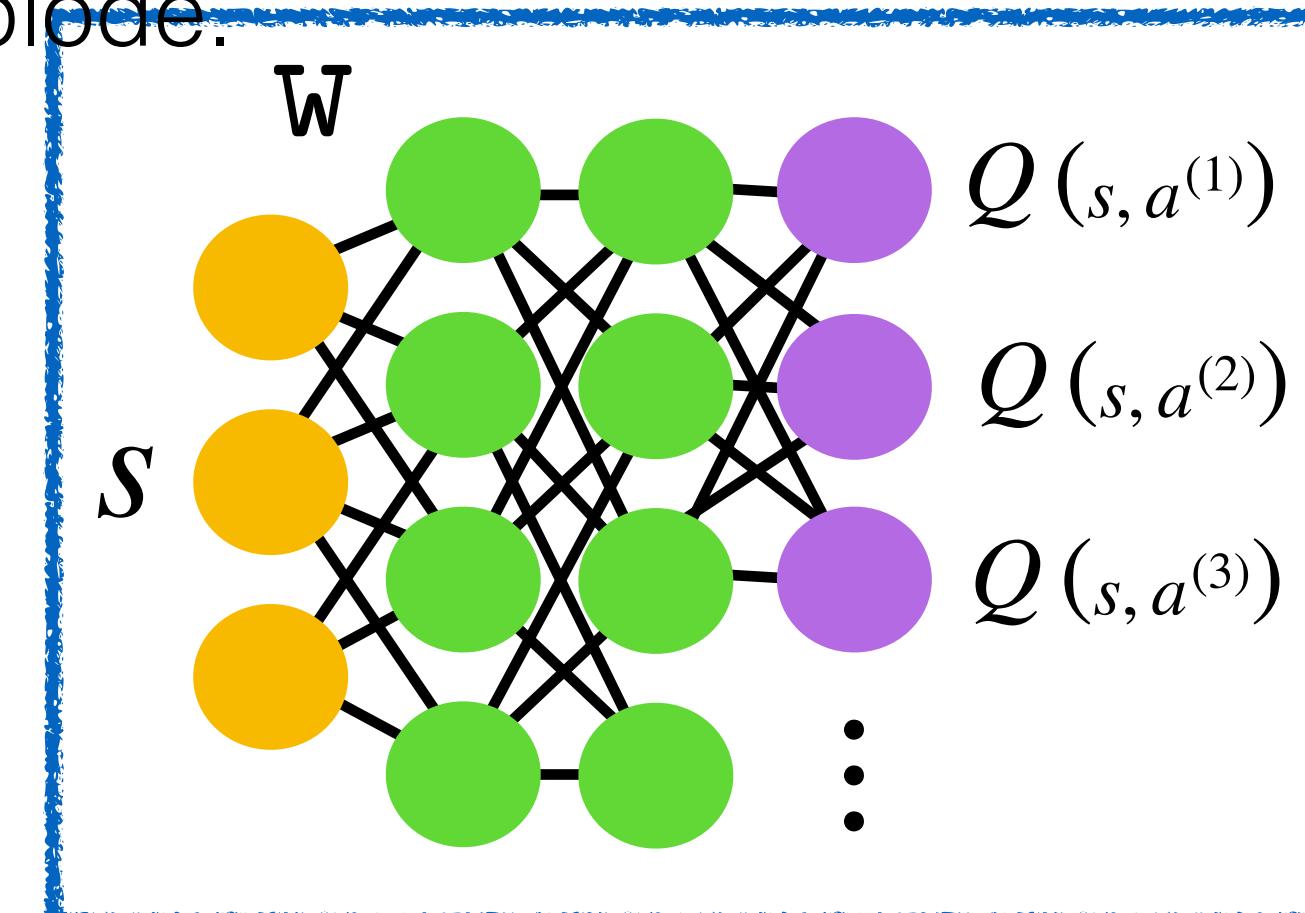
Approximating value functions

- **Goal:** learn $Q^\pi(s, a) = \mathbb{E} \left[\hat{q}_t^\pi \middle| \begin{array}{l} s \sim \eta^{\pi_w(\cdot)} \\ a \sim \pi_w(\cdot | s) \end{array} \right] \quad \forall s \in \mathcal{S}$
 $\forall a \in \mathcal{A}$
- Where $\hat{q}_t^\pi = \sum_{t'=t}^T \gamma^{t'-t} r_{t'+1}$
- With function approximator : $Q_w(s, a) \approx Q^\pi(s, a) \quad \forall s \in \mathcal{S}$
 $\forall a \in \mathcal{A}$
- First assume discrete action space

But **Q-learning is unstable if applied to NN:**

- Updates from temporally correlated data ($s_t, s_{t+1}, s_{t+2}, \dots$),
- max operation leads Q_w to explode.

	$a(1)$	$a(2)$...	$a(N_A)$
$s(1)$	$Q(s^{(1)}, a^{(1)})$	$Q(s^{(1)}, a^{(2)})$...	$Q(s^{(1)}, a^{(N_A)})$
$s(2)$	$Q(s^{(2)}, a^{(1)})$	$Q(s^{(2)}, a^{(2)})$		$Q(s^{(2)}, a^{(N_A)})$
...
$s(N_S)$	$Q(s^{(N_S)}, a^{(1)})$	$Q(s^{(N_S)}, a^{(2)})$...	$Q(s^{(N_S)}, a^{(N_A)})$



$$\hat{q}_t^\pi = \sum_{t'=t}^{\infty} \gamma^{t'-t} r_{t'+1} \quad \begin{array}{l} \text{1 sample estimate of Q:} \\ a_t \sim \pi, \\ s_{t+1}, r_{t+1} \sim D \end{array}$$

analytical estimate of Q:

$$Q^\pi(s, a) = \mathbb{E}_{\pi} [\hat{q}^\pi(s, a)]$$

parametric estimate of Q:

$$Q_w(s, a) \approx Q^\pi(s, a)$$

Estimating value functions

- **Goal:** learn approximate $Q_w^\pi(s, a) \approx Q^\pi(s, a)$ $\forall s \in \mathcal{S}$ $\forall a \in \mathcal{A}$
- **But:** we cannot **ever** compute $Q^\pi(s, a)$ because we do not have a model that we can reset to (s, a) .
 - How can we estimate $Q^\pi(s, a)$ with finite data?
- **Monte Carlo Learning**
 - Learn from full episodes/sequences of actions
 - Update the value towards **empirical mean returns**
- **Temporal Difference Learning**
 - Update the value towards **consistency with next estimated value**
 - Estimates based in part on other estimate: bootstrapping

Monte Carlo Learning

Given a policy π

For $i = 0, 1, 2 \dots$ perform “**roll-outs**”:

$$\begin{aligned}s_0^{(i)} &\sim D_0(\cdot) \\ a_0^{(i)} &\sim \pi(\cdot | s_0^{(i)}) \\ s_1^{(i)}, r_1^{(i)} &\sim D(\cdot | s_0^{(i)}, a_0^{(i)}) \\ &\dots \\ a_{T-1}^{(i)} &\sim \pi(\cdot | s_{T-1}^{(i)}) \\ s_T^{(i)}, r_T^{(i)} &\sim D(\cdot | s_{T-1}^{(i)}, a_{T-1}^{(i)})\end{aligned}$$

For each step, compute: $\hat{q}_t^{\pi(i)} = \sum_{t'=t}^T \gamma^{t'-t} r_{t'+1}^{(i)}$

Minimize approximation error:

$$\mathcal{L}(w)_t^{(i)} = \frac{1}{2} [\hat{q}_t^{\pi(i)} - Q_w(s_t^{(i)}, a_t^{(i)})]^2$$

Example: discrete states & actions

- Each pair (s, a) corresponds to one entry in a multi-dimensional grid
- We can update $Q_w(s_t^{(i)}, a_t^{(i)})$ independently from other entries
- $\frac{d\mathcal{L}(w)_t^{(i)}}{dw} = Q_w(s_t^{(i)}, a_t^{(i)}) - \hat{q}_t^{\pi(i)}$
- Update with SGD:
$$Q_w(s_t^{(i)}, a_t^{(i)}) \leftarrow Q_w(s_t^{(i)}, a_t^{(i)}) + \epsilon [\hat{q}_t^{\pi(i)} - Q_w(s_t^{(i)}, a_t^{(i)})]$$

1 sample estimate of $Q^\pi(s_t^{(i)}, a_t^{(i)})$

Only works for episodic RL!

Monte Carlo Learning

Given a policy π

For $i = 0, 1, 2 \dots$ perform “**roll-outs**”:

$$s_0^{(i)} \sim D_0(\cdot)$$

$$a_0^{(i)} \sim \pi(\cdot | s_0^{(i)})$$

$$s_1^{(i)}, r_1^{(i)} \sim D(\cdot | s_0^{(i)}, a_0^{(i)})$$

...

$$a_{T-1}^{(i)} \sim \pi(\cdot | s_{T-1}^{(i)})$$

$$s_T^{(i)}, r_T^{(i)} \sim D(\cdot | s_{T-1}^{(i)}, a_{T-1}^{(i)})$$

For each step, compute: $\hat{q}_t^{\pi(i)} = \sum_{t'=t}^T \gamma^{t'-t} r_{t'+1}^{(i)}$ **1 sample estimate of $Q^{\pi_w}(s_t^{(i)}, a_t^{(i)})$**

Minimize approximation error:

$$\mathcal{L}(w)_t^{(i)} = \frac{1}{2} [\hat{q}_t^{\pi(i)} - Q_w^{\pi}(s_t^{(i)}, a_t^{(i)})]^2$$

Example: continuous states & actions

- Neural network (maybe linear) with input (s, a)
- Update of $Q_w^{\pi}(s_t^{(i)}, a_t^{(i)})$ affects NN output for other states and actions
- Update with SGD by back-propagating $\frac{d\mathcal{L}(w)_t^{(i)}}{dw}$
- $\Delta w = \epsilon [\hat{q}_t^{\pi(i)} - Q_w^{\pi}(s_t^{(i)}, a_t^{(i)})] \nabla_w Q_w^{\pi}(s_t^{(i)}, a_t^{(i)})$

Only works for episodic RL!

Temporal Difference Learning

Given a policy π

For $i = 0, 1, 2 \dots$ perform “**roll-outs**”:

$$s_0^{(i)} \sim D_0(\cdot)$$

$$a_0^{(i)} \sim \pi(\cdot | s_0^{(i)})$$

$$s_1^{(i)}, r_1^{(i)} \sim D(\cdot | s_0^{(i)}, a_0^{(i)})$$

...

$$a_{T-1}^{(i)} \sim \pi(\cdot | s_{T-1}^{(i)})$$

$$s_T^{(i)}, r_T^{(i)} \sim D(\cdot | s_{T-1}^{(i)}, a_{T-1}^{(i)})$$

- Each step, update Q_w to be consistent with next estimate:
$$\mathcal{L}(w)_t^{(i)} = \frac{1}{2} \left[r_{t+1}^{(i)} + \gamma Q_w(s_{t+1}^{(i)}, a_{t+1}^{(i)}) - Q_w(s_t^{(i)}, a_t^{(i)}) \right]^2$$
- TD target: $r_{t+1} + \gamma Q_w(s_{t+1}, a_{t+1})$
- Can learn before knowing final outcome, with incomplete data or in infinite state/action loops.
- Forces Q_w to be *internally consistent* in its predictions.

One link between MC and TD

- The TD error is:

$$\delta_t = r_{t+1} + \gamma Q_w(s_{t+1}, a_{t+1}) - Q_w(s_t, a_t)$$

- The MC error is:

$$\hat{q}_t^\pi - Q_w(s_t, a_t) = \sum_{t'=t}^T \gamma^{t'-t} r_{t'+1} - Q_w(s_t, a_t)$$

- Unroll the MC error: $\hat{q}_t^\pi - Q_w(s_t, a_t) = r_{t+1} + \gamma \hat{q}_{t+1}^\pi - Q_w(s_t, a_t)$

$$= r_{t+1} + \gamma \hat{q}_{t+1}^\pi - Q_w(s_t, a_t) + \gamma Q_w(s_{t+1}, a_{t+1}) - \gamma Q_w(s_{t+1}, a_{t+1})$$

$$= \delta_t + \gamma (\hat{q}_{t+1}^\pi - Q_w(s_{t+1}, a_{t+1}))$$

$$= \delta_t + \gamma \delta_{t+1} + \gamma^2 (\hat{q}_{t+2}^\pi - Q_w(s_{t+2}, a_{t+2}))$$

$$= \sum_{t'=t}^{T-1} \gamma^{t'-t} \delta_{t'}$$

- Therefore: $\hat{q}_t^\pi = Q_w(s_t, a_t) + \sum_{t'=t}^{T-1} \gamma^{t'-t} \delta_{t'}$

**Monte Carlo target is
a sum of TD errors!**

MC vs TD

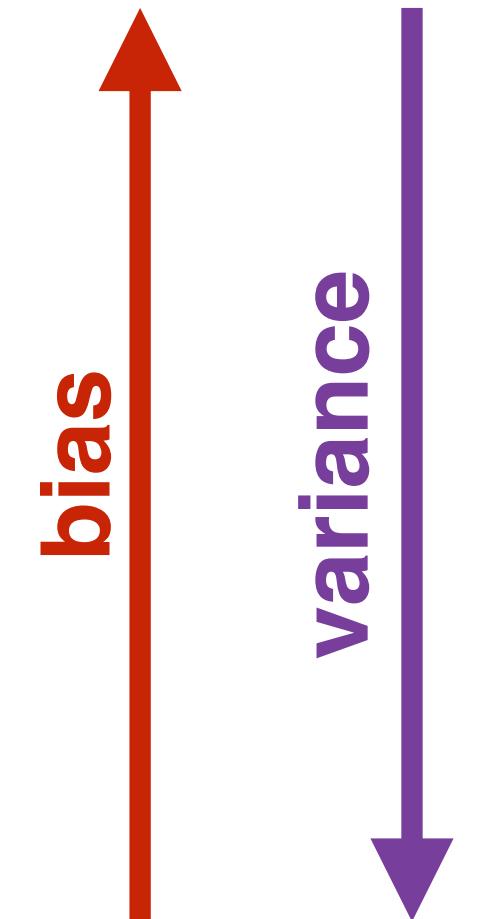
Bias/Variance trade-off

- Monte Carlo target $\sum_{t'=t}^T \gamma^{t'-t} r_{t'+1}$ is:
 - **unbiased** (exact in the limit of ∞ data)
 - **high variance** (requires **many** samples for a good estimate)
 - Why? Because it depends on future $T - t$ states/actions
 - **Good convergence properties**
 - **Not sensitive to initial guess**
 - **Simple to use**
- Temporal Difference target $r_{t+1} + \gamma Q_w(s_{t+1}, a_{t+1})$ is:
 - **low variance** (depends only on next state/action)
 - **biased** (if next guess is inaccurate, update will be wrong)
 - **More data-efficient than MC**
 - **More unstable with NN**
 - **More sensitive to initial guess**

Multi-step targets

- Train approximate Q towards a target value: $\mathcal{L}(w)_t = \frac{1}{2} [\hat{q}_t^\pi - Q_w(s_t, a_t)]^2$
- Options for estimation of target value:

- $\hat{q}_t^\pi = r_{t+1} + \gamma Q_w(s_{t+1}, a_{t+1})$
- $\hat{q}_t^\pi = r_{t+1} + \gamma r_{t+2} + \gamma^2 Q_w(s_{t+2}, a_{t+2})$
- $\hat{q}_t^\pi = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \gamma^3 Q_w(s_{t+3}, a_{t+3})$
- ...
- $\hat{q}_t^\pi = \sum_{t'=t}^{\infty} \gamma^{t'-t} r_{t'+1}$



Temporal Difference learning

Monte Carlo learning

- N-steps estimate aim to balance bias and variance

RL through value learning

- **Policy Evaluation:** learn approximate $Q_w^\pi(s, a)$ by executing policy rollouts
- **Policy Optimization:** optimize $\pi(a | s)$ as to improve $Q^\pi(s, a)$
- Value learning-based RL:
 - Train Q-approximator towards TD/MC target:

$$\mathcal{L}(w) = \mathbb{E} \left[\frac{1}{2} \left[\hat{q}^{\pi_w}(s, a) - Q_w(s, a) \right]^2 \middle| \begin{array}{l} s \sim \eta^{\pi_w(\cdot)}, \\ a \sim \pi_w(\cdot | s), \end{array} \right]$$

- Sample actions corresponding to greater expected rewards:
 - maximum entropy: $\pi_w(a | s) \propto \exp [Q_w(s, a)]$
 - greedy plus noise: $\pi_w(a | s) = \arg \max_a Q_w(s, a) + \xi(a)$
- Learn to estimate returns (**policy evaluation**) and use the estimate to get higher returns (**policy optimization**)

ϵ -greedy policy:

- $\epsilon_t \in (0, 1)$ probability of random action
- else:

$$\pi_w(s_t) = \arg \max_a Q_w(s_t, a)$$

max entropy policy:

- discrete actions:

$$\pi_w(a | s) = \frac{\exp Q_w(s, a)}{\sum_{ai=1}^{d_A} \exp Q_w(s, ai)}$$

SARSA : on-policy value learning

Initialize $w^{(0)}$

for iteration i in 1, 2, ...

- Collect E episodes by executing $\pi_{w^{(i)}}$
 - max entropy: $\pi_{w^{(i)}}(a | s) \propto \exp [Q_{w^{(i)}}(s, a)]$
 - ϵ -greedy
- # of collected steps (i.e. actions) may be $N = TE$
- For each step t , compute TD target:
$$\hat{q}_t^{\pi_{w^{(i)}}} = r_{t+1} + \gamma Q_{w^{(i)}}(s_{t+1}, a_{t+1})$$
- Update $w^{(i)}$ by SGD by minimizing loss:
$$\mathcal{L}^{MSE}(w^{(i)}) = \frac{1}{N} \sum_{t=0}^N \left[\frac{1}{2} [\hat{q}_t^{\pi_{w^{(i)}}} - Q_{w^{(i)}}(s_t, a_t)]^2 \right]$$

SARSA

- Here we describe a batch version of SARSA where multiple steps are collected to perform each update
- May use maximum entropy policy or greedy policy with a small chance of random actions.
- Because we use on-policy actions (a_{t+1}) in Bellman update, this is an **on-policy** method.
- Can generalize to **n-step** SARSA or use **Monte Carlo** rewards (e.g. see Mnih 2016)

Q-learning : off-policy value learning

- Greedy policy is to pick action with $\max Q_w$, but we **need** to explore (add exploratory noise).
- **Bellman optimality equation**: optimal value function is computed by following greedy behavior.
 - **Random action** a_t to explore
 - **Deterministic action** $a' = \arg \max Q_w$ to compute **Bellman update**

Initialize $w^{(0)}$

for iteration i in $1, 2, \dots, N$

- Collect D episodes by executing $\pi_{w^{(i)}}$ (max-entropy or ϵ -greedy)
- For each step, compute target with Bellman Eq:

$$\hat{q}_t^{\pi_{w^{(i)}}} = r_{t+1} + \gamma \max_{a'} Q_{w^{(i)}}(s_{t+1}, a')$$

- Update $w^{(i)}$ by SGD by minimizing loss:

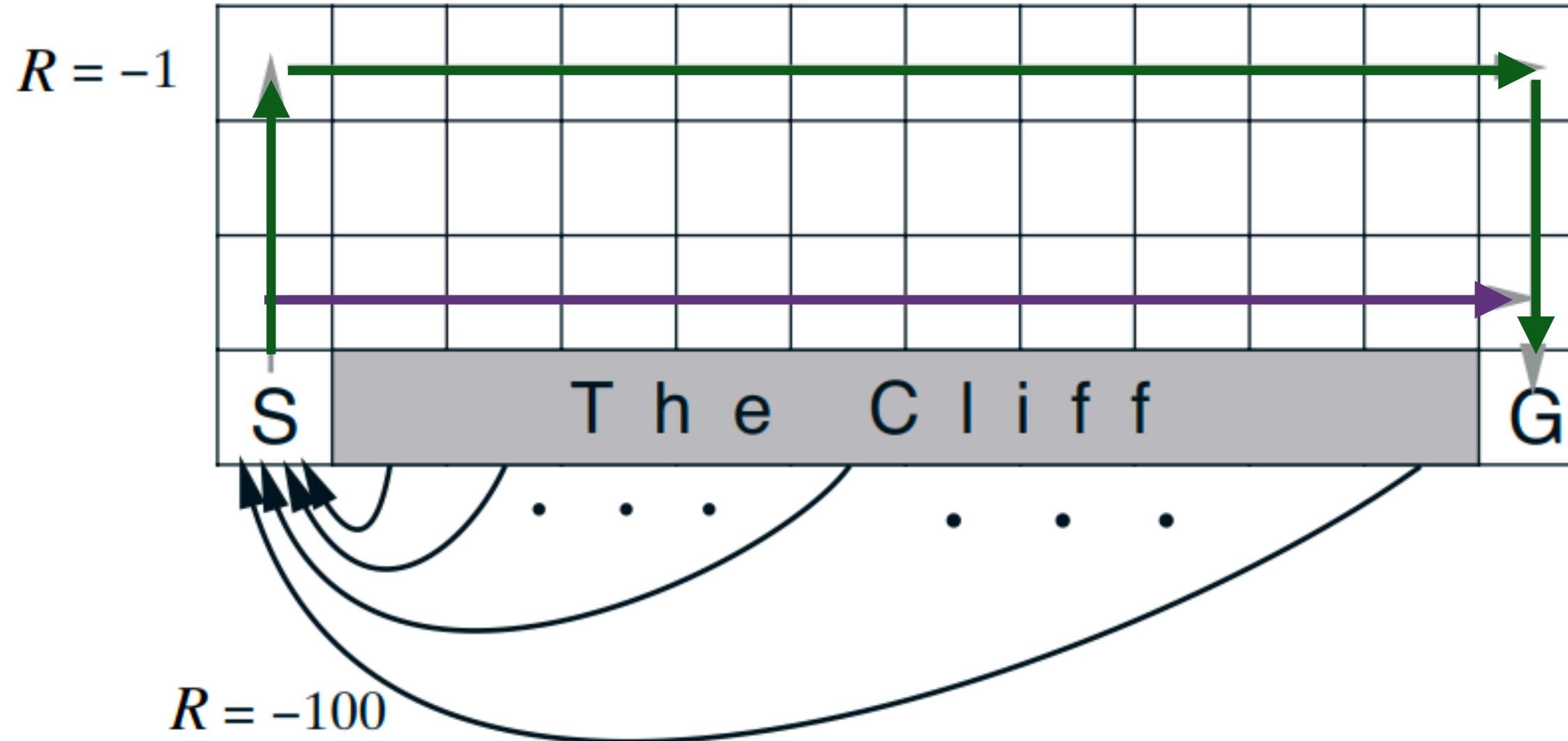
$$\mathcal{L}^{MSE}(w^{(i)}) = \frac{1}{N} \sum_{t=0}^N \left[\frac{1}{2} [\hat{q}_t^{\pi_{w^{(i)}}} - Q_{w^{(i)}}(s_t, a_t)]^2 \right]$$

Q-learning

- **Off-policy** method:
 - learn Q_w for a deterministic optimal **policy**
 - while acquiring data with an exploratory **behavior**
- $\mathbb{E}_{X_1, X_2} \left[\max \{X_1, X_2\} \right] \geq \max \left\{ \mathbb{E}_{X_1} [X_1], \mathbb{E}_{X_2} [X_2] \right\}$

Because Q_w are noisy, Q-learning pathologically over-estimates returns

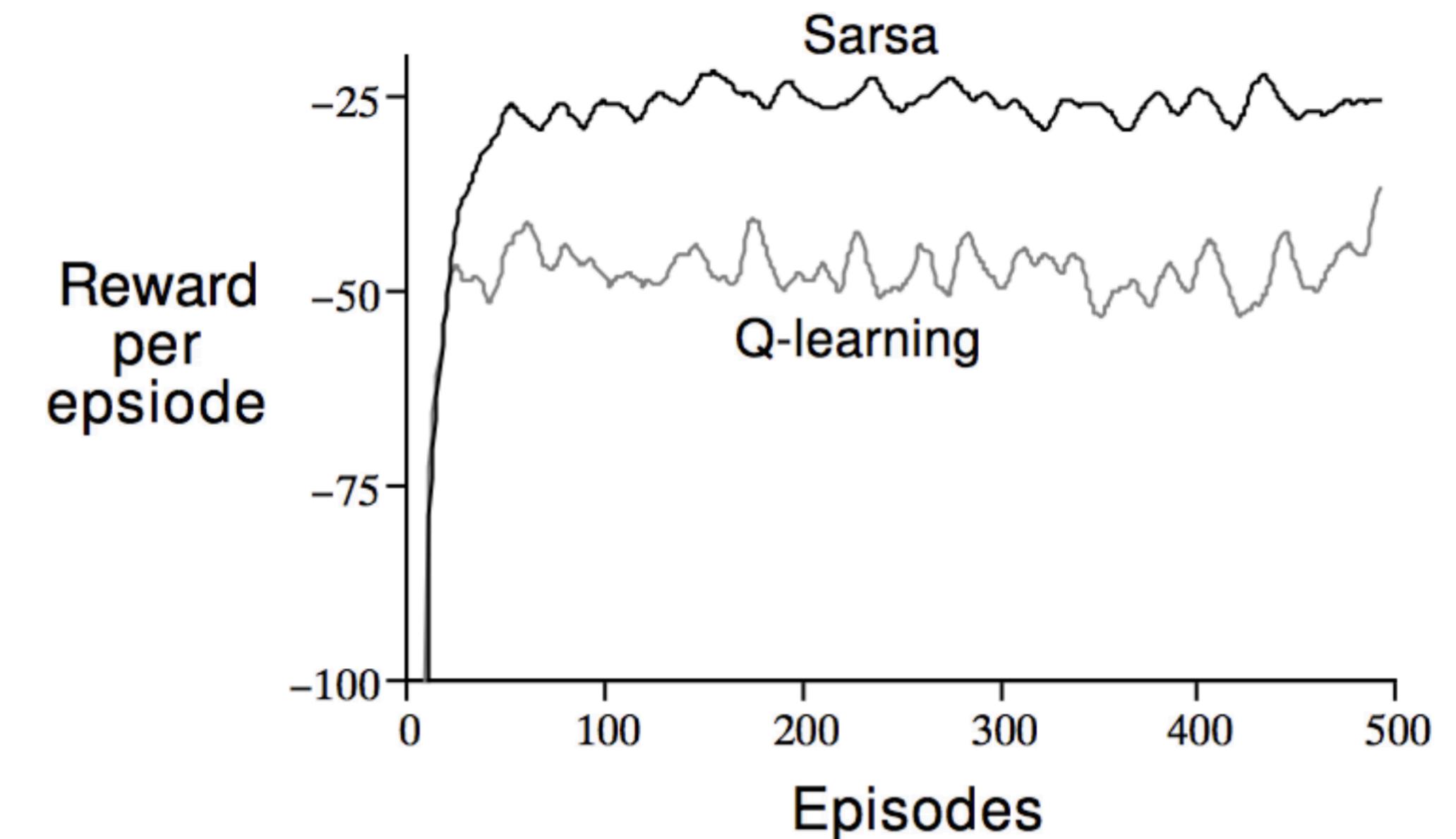
Q-learning vs SARSA



Cliff walking example by Sutton

Suppose we attempt to solve this problem with ϵ -greedy policy with both SARSA and Q-learning

Which method will perform best?

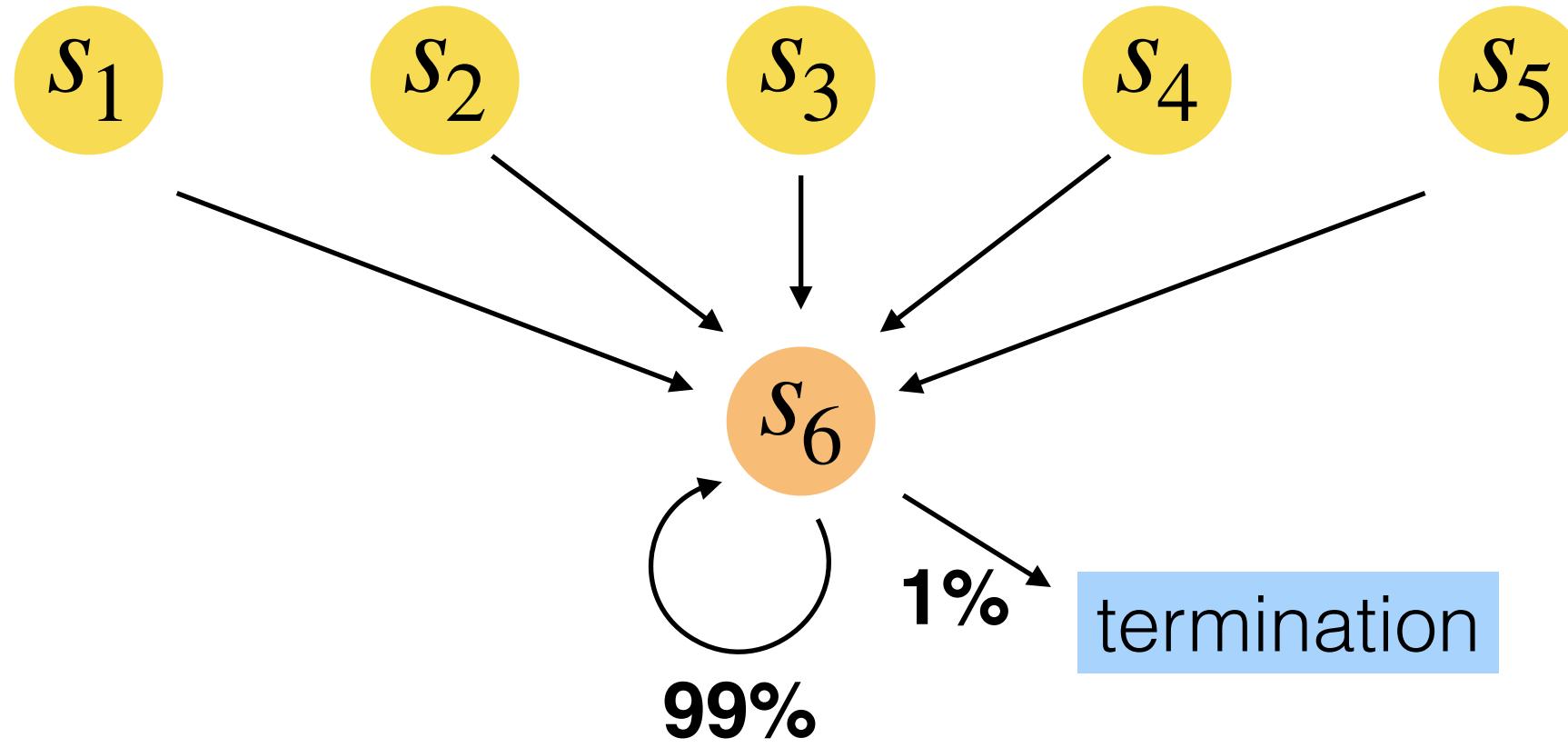


Q-learning approximates Q_w assuming optimal actions:
cannot account for ϵ -chance of random actions

SARSA is on-policy: accounts for ϵ -chance of random actions and takes safe path.

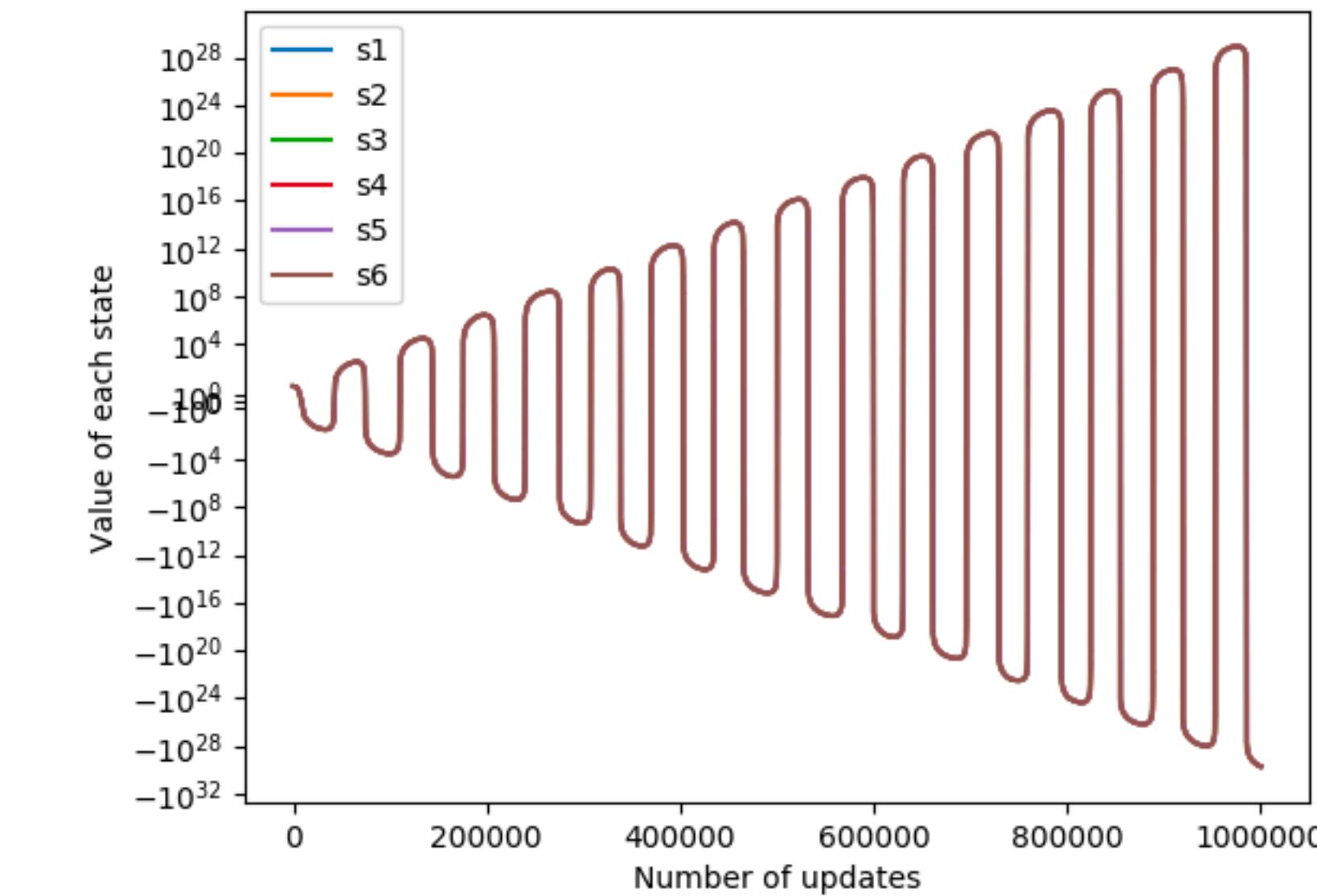
Stability of value-learning based control

Baird's counterexample



- Begin in one of the 5 first states at random
- Always transition to s_6
- Terminate with 1% probability else return to s_6
- All rewards are 0 (including termination)
- Let $V_w(s)$ initialized as:
 - $V_w(s_1) = 2w_1 + w_7$
 - $V_w(s_2) = 2w_2 + w_7$
 - $V_w(s_3) = 2w_3 + w_7$
 - $V_w(s_4) = 2w_4 + w_7$
 - $V_w(s_5) = 2w_5 + w_7$
 - $V_w(s_6) = w_6 + 2w_7$

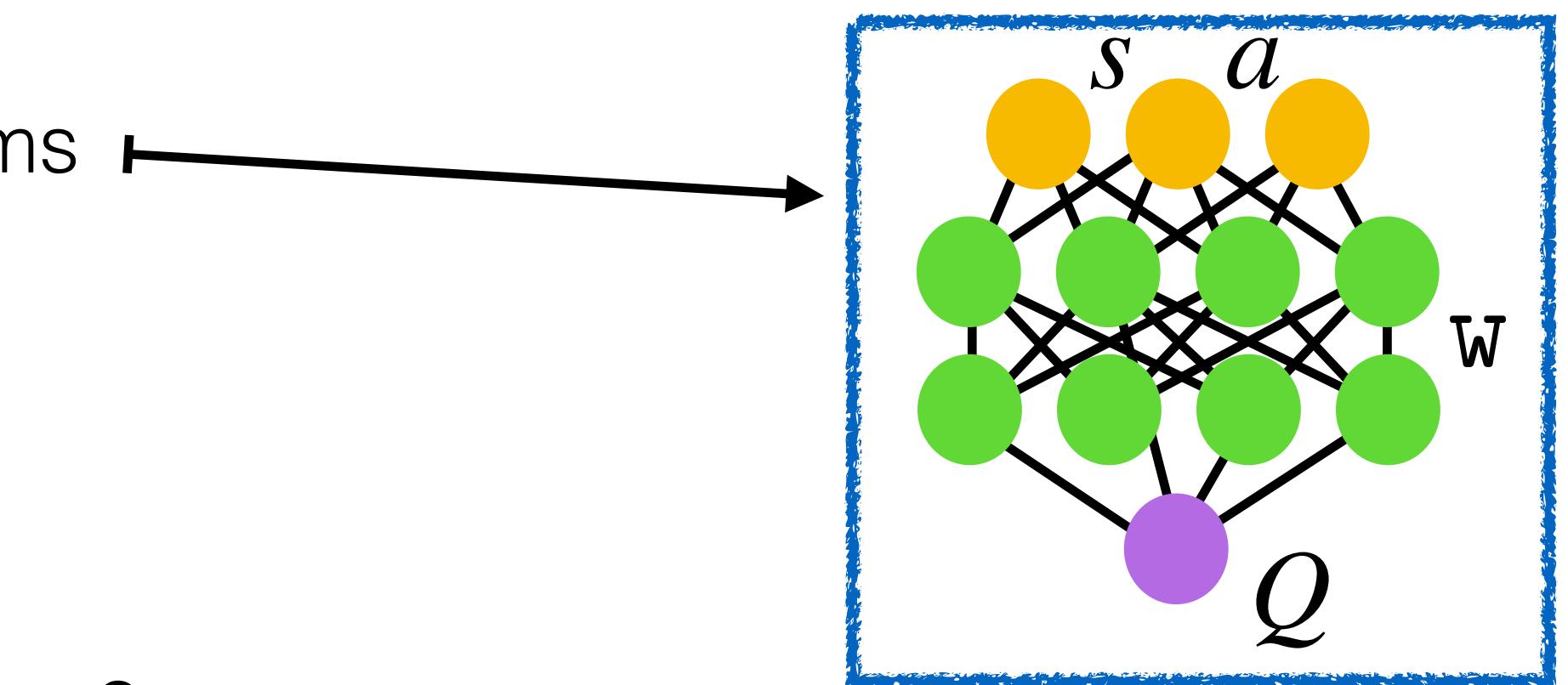
Iteratively update with TD learning:



RL control algorithm:	Tabular	Linear	Non-linear
MC	stable	oscillates around optimum	unstable
SARSA	stable	oscillates around optimum	unstable
Q-learning	stable	unstable	unstable

Q-learning for continuous actions

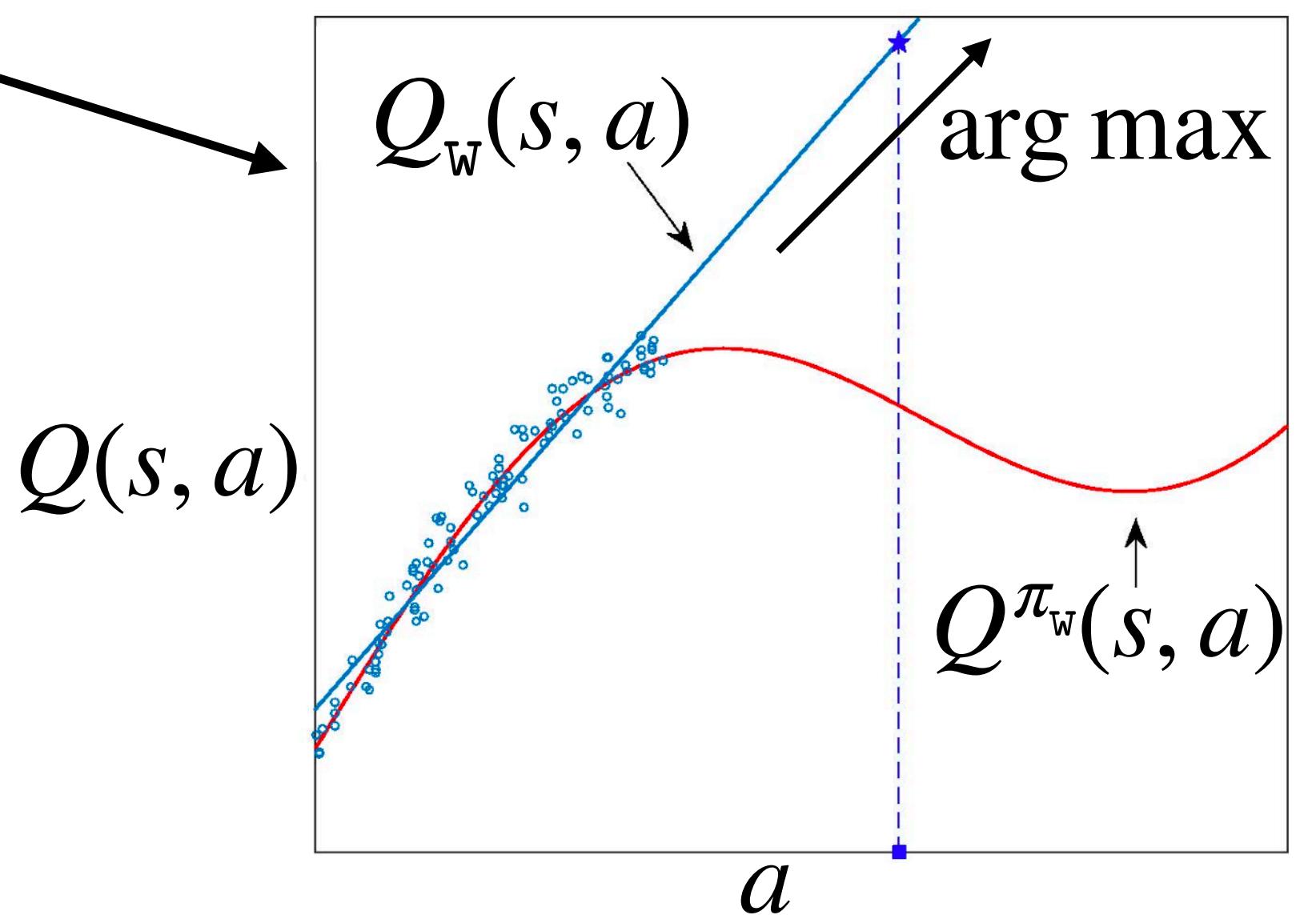
- Continuous-valued action vectors prevalent in robotics / PDEs / ODEs (even modern video-games!)
- Both s and a are vectors: natural choice for Q-network seems



- **But:**

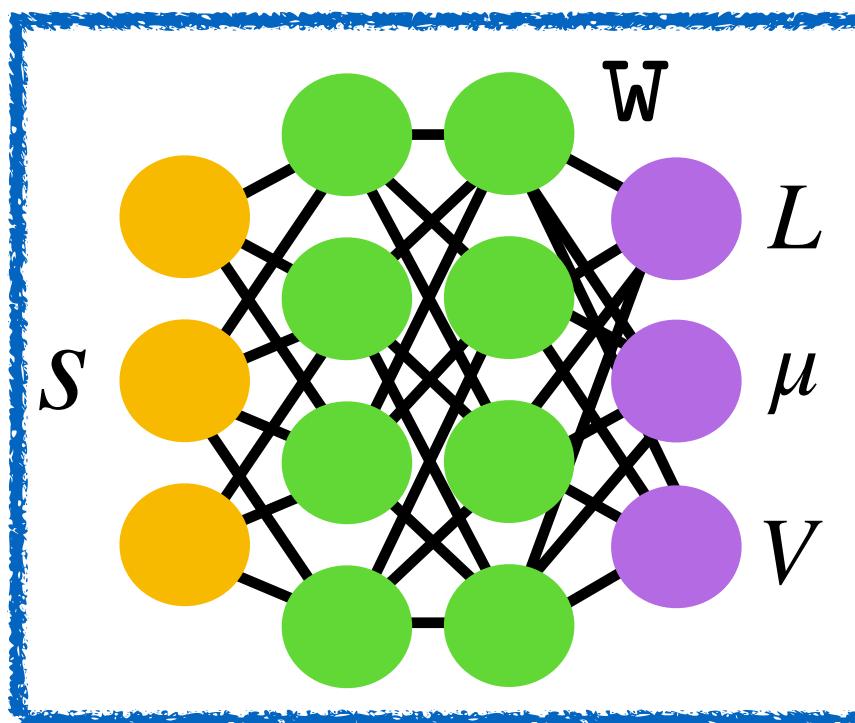
- Expensive to find maxima of a NN (policy)
- How to ensure Q_w is good approximation in all action space?
- If linear approximation: $\arg \max_a Q_w(s, a) \rightarrow \infty$

How can we, for a given s , output a well-behaved approximation of Q_w over the entire action space? (like DQN)



Continuous Q-learning: Normalized Advantage Functions

- NAF network architecture:



Initialize $w^{(0)}$
for iteration t in $1, 2, \dots, N$

- Collect D episodes with $\beta(\cdot | s) = \mathcal{N}(\mu_{w^{(i)}}(s), \sigma^2)$
- For each step, compute target with Bellman Eq:

$$\hat{q}_t^{\pi_{w^{(i)}}} = r_{t+1} + \gamma V_{w^{(i)}}(s_{t+1})$$

- Update $w^{(t)}$ with SGD by minimizing loss:

$$\mathcal{L}^{MSE}(w^{(i)}) = \frac{1}{N} \sum_{t=0}^N \left[\frac{1}{2} [\hat{q}_t^{\pi_{w^{(i)}}} - Q_{w^{(i)}}(s_t, a_t)]^2 \right]$$

- Output $V_w(s)$ approximate state value
- $\mu_w(s)$ vector of size $\dim(A)$
- $L_w(s)$ lower triangular matrix*:

$$L_w(s) = \begin{bmatrix} l_w^{(1,1)} & 0 & 0 & \dots & 0 \\ l_w^{(2,1)} & l_w^{(2,2)} & 0 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots \\ l_w^{(d_A,1)} & l_w^{(d_A,2)} & l_w^{(d_A,3)} & \dots & l_w^{(d_A,d_A)} \end{bmatrix}$$

- $L_w(s)L_w^T(s)$ is positive definite!
- $Q_w(s, a) = V_w - \frac{1}{2}(a - \mu_w)^T L_w L_w^T(a - \mu_w)$
- $\mu_w(s)$ is action with optimal value $V_w(s)$!

* To ease training (bijection) elements on the diagonal should be pos. def. $l_w^{(i,i)} \leftarrow \log [1 + \exp(l_w^{(i,i)})]$

Reinforcement Learning Workshop

Day 4
Policy Optimization

Guido Novati

CSElab

Computational Science & Engineering Laboratory

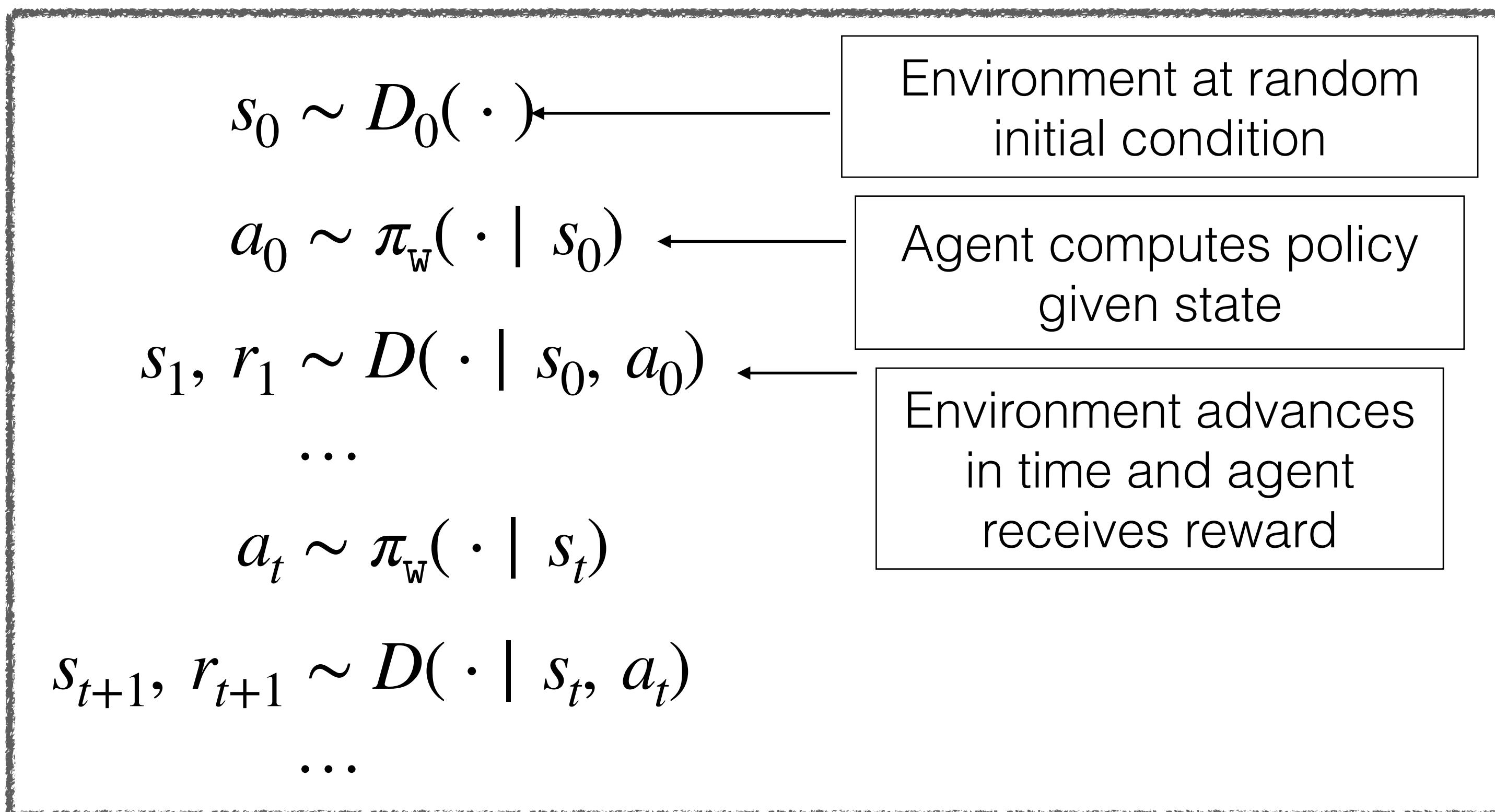
ETH

Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

RL: Markov decision process with parameterized policy

RL learns a policy

- Function that computes actions given states.
- **Improvable** i.e. it has parameters w which can be optimized.
- Some stochasticity is required to explore dynamics and possible improvements.



$$J(w) = \mathbb{E} \left[\sum_t^\infty r_t \mid \begin{array}{l} a_t \sim \pi_w(\cdot | s_t) \\ s_{t+1} \sim D(\cdot | a_t, s_t) \end{array} \right]$$

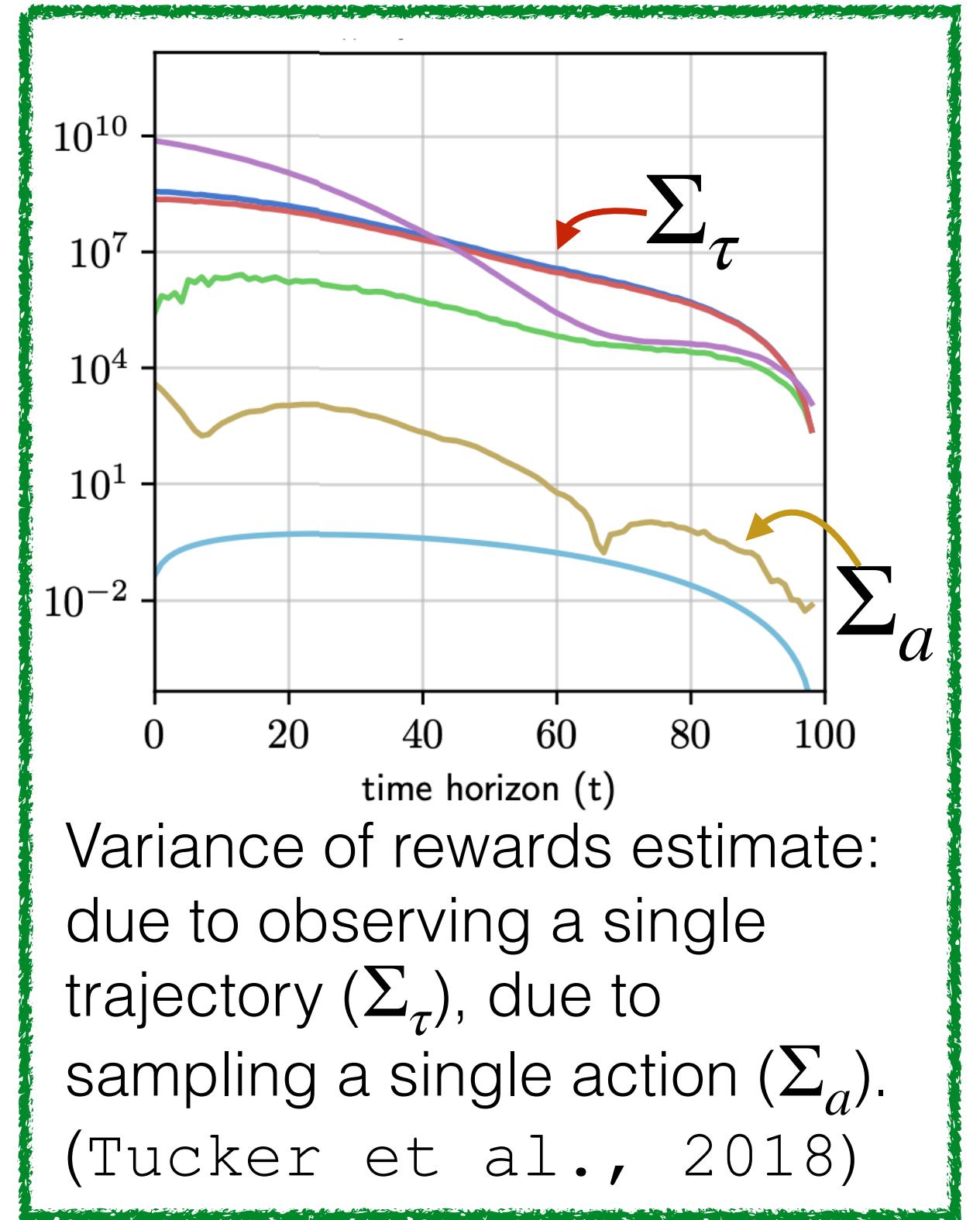
optimal parameters

$$w^* = \arg \max_w J(w)$$

Why policy optimization

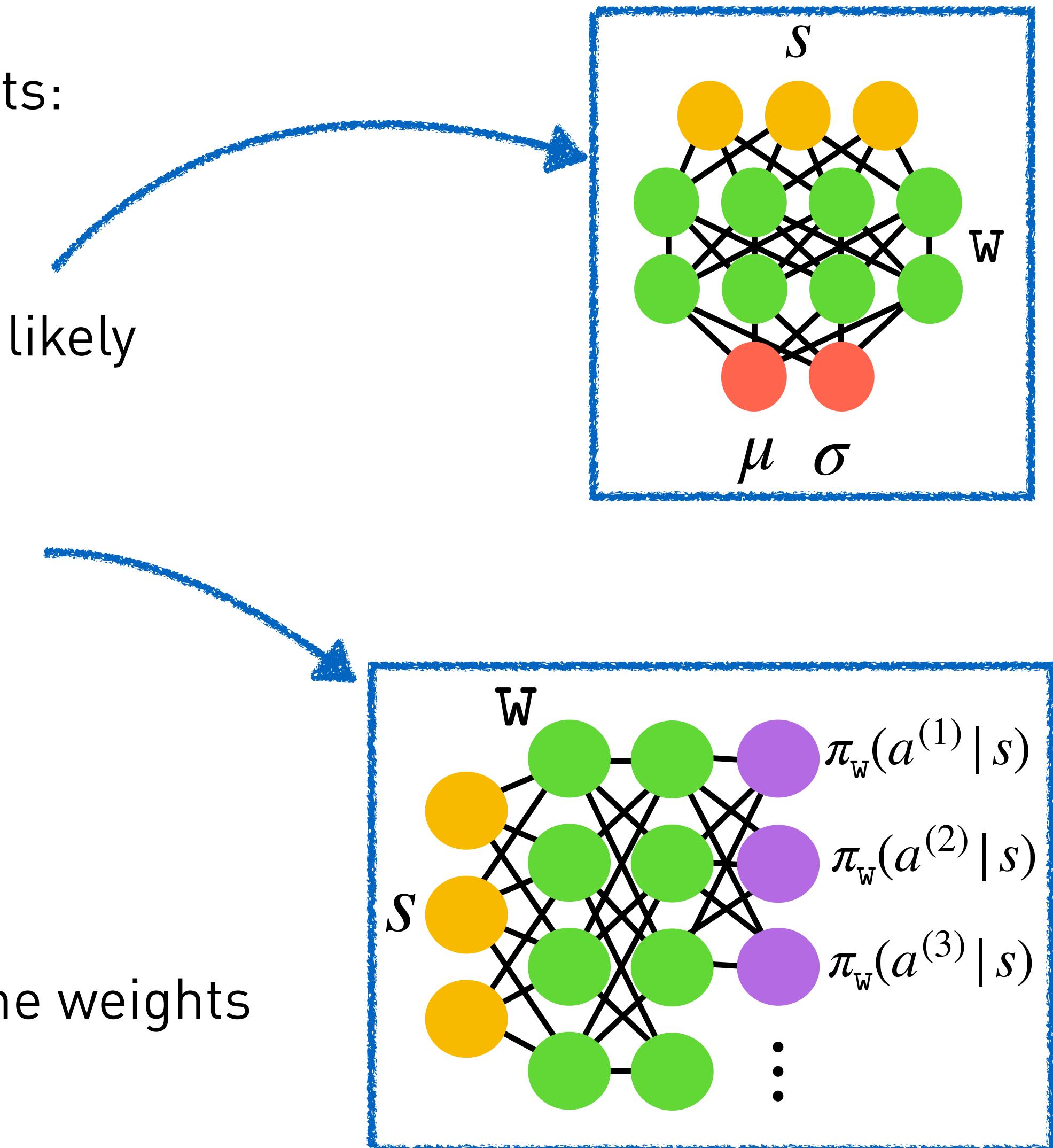
- Generally **policies are simpler** functions **than Q-functions**.
 - Compare *learning entire landscape of rewards* versus *identifying actions correlated with higher rewards*.
 - $Q(s, \mathbf{a})$ value affected by \mathbf{a} and many actions into the future:
 - Need many trials to average out effect of future possible steps (noise) to get at value of individual actions (signal).
 - For most problems, single actions have little effect (can be corrected later) and therefore Q is very flat over \mathbf{a} .
- Even if we know $Q_w(s, a)$, it may be difficult to:
 - sample $\pi_w(a | s) \propto \exp [Q_w(s, a)]$ (rejection sampling)
 - optimize $\pi_w(s) = \arg \max_a Q_w(s, a)$

Not an issue for discrete action problems!



Parametric policies: examples

- Directly parameterize the policy, optimize the weights:
 - Gaussian Policy $a \sim \mathcal{N} [\mu_w(s), \sigma_w^2(s)]$
find Gaussian where optimal actions are most likely
 - (Generalized) Bernoulli Policy
one NN output for each action, must sum to 1,
(softmax output layer)
 - Deterministic policy $a = \pi_w(s)$
output is directly optimal action (e.g. vector)
exploration: must add noise to the action or the weights



Cross-entropy method (CEM) for policy optimization

Initial guess on policy parameter vector $\bar{w}^{(0)}$
for iteration i in $1, 2, \dots, N$
for population member $p = 1, 2, \dots, P$

- Sample $\xi_p \sim \mathcal{N}(0, \mathbf{I})$
- Set $w_p^{(i)} = \bar{w}^{(i)} + \sigma \xi_p$
- Execute E episodes with policy $\pi_{w_p^{(i)}}$
- Measure average utility $J(w_p^{(i)})$

- Sort population according to average utility.
- Select B best candidates with $B < P$

$$\bar{w}^{(i+1)} = \frac{1}{B} \sum_{p=1}^B w_p^{(i)}$$

- Optionally, update sampling noise:

$$\sigma^{2, (i+1)} = \frac{1}{B} \sum_{p=1}^B \left(w_p^{(i)} - \bar{w}^{(i+1)} \right)^T \left(w_p^{(i)} - \bar{w}^{(i+1)} \right)$$

$$\max_w J(w) = \max_w \mathbb{E} \left[\sum_t^T r_t \mid \begin{array}{l} a_t \sim \pi_w(\cdot | s_t) \\ s_{t+1} \sim D(\cdot | a_t, s_t) \end{array} \right]$$

- Evolutionary algorithm.
- Treat $J(w)$ as black box.
- Ignores all other information (states, actions, intermediate rewards).
- Because noise is intrinsic in evolutionary process, policy can be deterministic.

Evolution Strategies (ES) algorithm for policy optimization

Initial guess on policy parameter vector $\bar{w}^{(0)}$
for iteration i in $1, 2, \dots, N$
for population member $p = 1, 2, \dots, P$

- Sample $\xi_p \sim \mathcal{N}(0, \mathbf{I})$
- Set $w_p^{(i)} = \bar{w}^{(i)} + \sigma \xi_p$
- Execute E episodes with policy $\pi_{w_p^{(i)}}$
- Measure average utility $J(w_p^{(i)})$

- Update mean:

$$\bar{w}^{(i+1)} = \bar{w}^{(i)} + \frac{\epsilon}{\sigma P} \sum_{p=1}^P J(w_p^{(i)}) \xi_p$$

- Optionally, update sampling noise:

$$\sigma \leftarrow \sigma + \frac{\epsilon}{\sigma P} \sum_{p=1}^P J(w_p^{(i)}) \frac{\|\xi_p\|^2 - 1}{\dim(w)}$$

$$\max_w J(w) = \max_w \mathbb{E} \left[\sum_t^T r_t \mid \begin{array}{l} a_t \sim \pi_w(\cdot | s_t) \\ s_{t+1} \sim D(\cdot | a_t, s_t) \end{array} \right]$$

- Evolutionary algorithm.
- Treat $J(w)$ as black box.
- Ignores all other information (states, actions, intermediate rewards).
- Because noise is intrinsic in evolutionary process, policy can be deterministic.

Alternative optimization methods

Closely Related Approaches

CEM:

```
for iter i = 1, 2, ...
    for population member e = 1, 2, ...
        sample  $\theta^{(e)} \sim P_{\mu^{(i)}}(\theta)$ 
        execute roll-outs under  $\pi_{\theta^{(e)}}$ 
        store  $(\theta^{(e)}, U(e))$ 
    endfor
     $\mu^{(i+1)} = \arg \max_{\mu} \sum_{\bar{e}} \log P_{\mu}(\theta^{(\bar{e})})$ 
    where  $\bar{e}$  indexes over top p %
endfor
```

- Reward Weighted Regression (RWR)

- Dayan & Hinton, NC 1997; Peters & Schaal, ICML 2007

$$\mu^{(i+1)} = \arg \max_{\mu} \sum_e q(U(e), P_{\mu}(\theta^{(e)})) \log P_{\mu}(\theta^{(e)})$$

- Policy Improvement with Path Integrals (PI²)

- PI2: Theodorou, Buchli, Schaal JMLR2010; Kappen, 2007; (PI2-CMA: Stulp & Sigaud ICML2012)

$$\mu^{(i+1)} = \arg \max_{\mu} \sum_e \exp(\lambda U(e)) \log P_{\mu}(\theta^{(e)})$$

- Covariance Matrix Adaptation Evolutionary Strategy (CMA-ES)

- CMA: Hansen & Ostermeier 1996; (CMA-ES: Hansen, Muller, Koumoutsakos 2003)

$$(\mu^{(i+1)}, \Sigma^{(i+1)}) = \arg \max_{\mu, \Sigma} \sum_{\bar{e}} w(U(\bar{e})) \log \mathcal{N}(\theta^{(\bar{e})}; \mu, \Sigma)$$

- PoWER

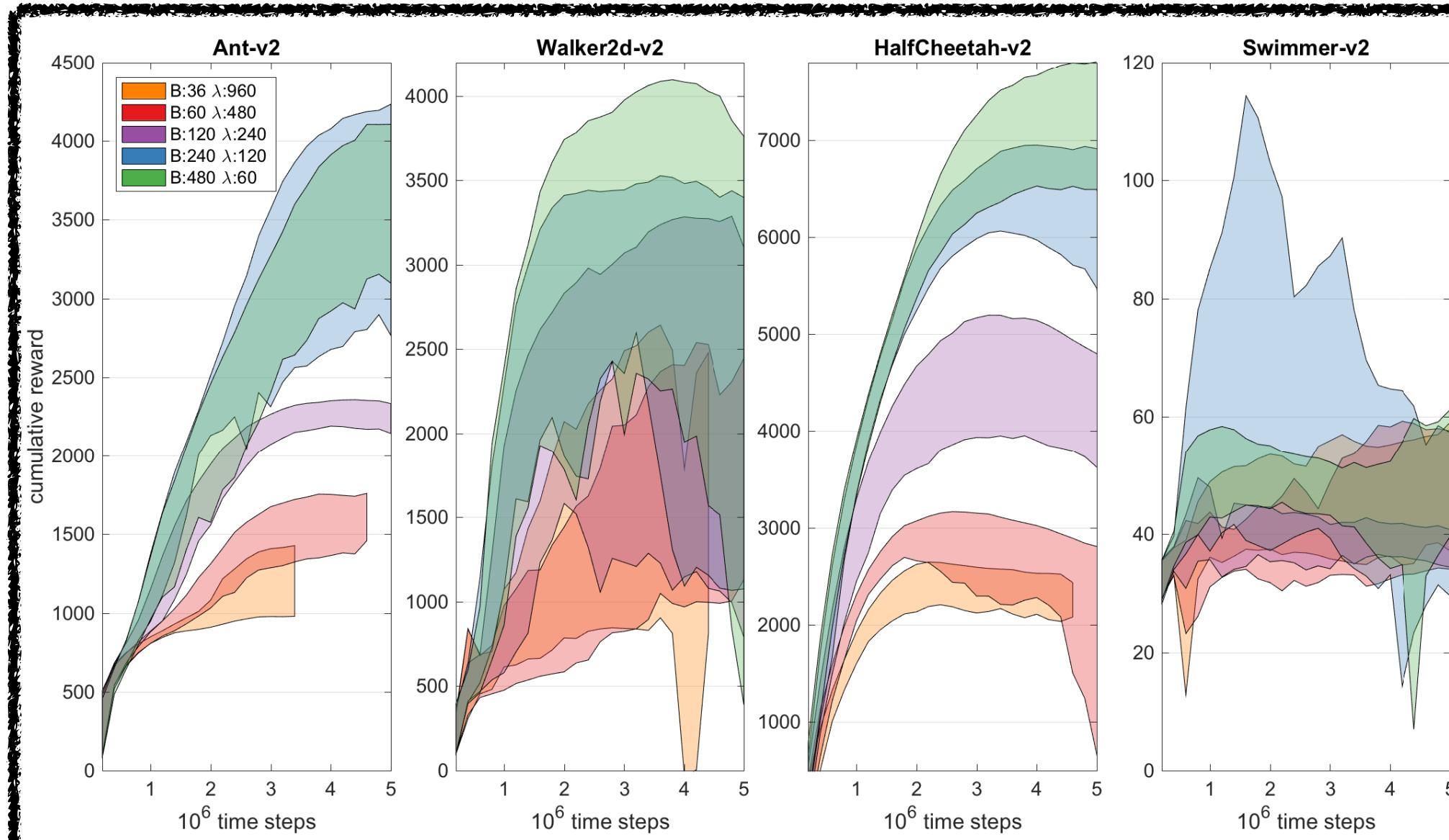
- Kober & Peters, NIPS 2007 (also applies importance sampling for sample re-use)

$$\mu^{(i+1)} = \mu^{(i)} + \left(\sum_e (\theta^{(e)} - \mu^{(i)}) U(e) \right) / \left(\sum_e U(e) \right)$$

From Pieter Abbeel's workshop at ICML'16

Pros&Cons of evolutionary policy optimization methods

- Simple implementation and robust.
- Fundamental limits:
 - **Probability** of finding **good policy update** by random sampling **goes to 0 as $1/\sqrt{\dim(w)}$** (most update directions are orthogonal to the ideal one).
 - There is **no temporal decomposition of the RL task**. For each w we only look at average returns. E.g. we do not exploit $Q^{\pi_w}(s, a)$, we only use $\mathbb{E} [Q^{\pi_w}(s, a)]$.



- **Example** from previous group meeting (not exactly the same, here I optimized off-policy objective)
- Kept **constant** the **workload**: $[\text{pop. size}]^*[\text{N of evaluations per sampled } w]$
- Increasing N evals is better than increasing pop. size.
- Variance of RL objective higher than variance of update estimate
- Temporal decomposition of RL is key for data-efficiency!

Likelihood ratio policy gradient (1/2)

expand notation for conveniency:

- Consider space \mathcal{T} of all possible episodes $\tau = \{s_0, a_0, s_1, \dots, s_{T-1}, a_{t-1}, s_T\}, \tau \in \mathcal{T}$.
- Associated return $R(\tau) = \sum_{i=1}^T r_t$
- Maximize rewards over episodes $J(w) = \mathbb{E} [R(\tau) \mid \tau \sim \mathbb{P}_w(\cdot)]$
- Probability of episode depends on policy parameters w .
- Equivalently write optimization: $\max_w J(w) = \max_w \int_{\mathcal{T}} R(\tau) \mathbb{P}_w(\tau) d\tau$

Likelihood ratio policy gradient (2/2)

Take the gradient of $J(w) = \int_{\mathcal{T}} R(\tau) \mathbb{P}_w(\tau) d\tau$

$$\begin{aligned}\nabla_w J(w) &= \nabla_w \int_{\mathcal{T}} R(\tau) \mathbb{P}_w(\tau) d\tau \\&= \int_{\mathcal{T}} R(\tau) \nabla_w \mathbb{P}_w(\tau) d\tau \\&= \int_{\mathcal{T}} R(\tau) \frac{\mathbb{P}_w(\tau)}{\mathbb{P}_w(\tau)} \nabla_w \mathbb{P}_w(\tau) d\tau \\&= \int_{\mathcal{T}} R(\tau) \mathbb{P}_w(\tau) \nabla_w \log [\mathbb{P}_w(\tau)] d\tau \\&= \mathbb{E}_{\tau \sim \mathbb{P}_w(\cdot)} [R(\tau) \nabla_w \log [\mathbb{P}_w(\tau)]]\end{aligned}$$

Alternative perspective: sample episodes with current w_{old} and estimate J for optimal w :

$$\begin{aligned}J(w) &= \mathbb{E}_{\tau \sim \mathbb{P}_w(\cdot)} [R(\tau)] = \mathbb{E}_{\tau \sim \mathbb{P}_{w_{old}}(\cdot)} \left[R(\tau) \frac{\mathbb{P}_w(\tau)}{\mathbb{P}_{w_{old}}(\tau)} \right] \\ \nabla_w J(w) &= \mathbb{E}_{\tau \sim \mathbb{P}_{w_{old}}(\cdot)} \left[R(\tau) \frac{\nabla_w \mathbb{P}_w(\tau)}{\mathbb{P}_{w_{old}}(\tau)} \right] \\ \nabla_w J(w) \Big|_{w=w_{old}} &= \mathbb{E}_{\tau \sim \mathbb{P}_{w_{old}}(\cdot)} \left[R(\tau) \frac{\nabla_w \mathbb{P}_w(\tau) \Big|_{w=w_{old}}}{\mathbb{P}_{w_{old}}(\tau)} \right] \\ &= \mathbb{E}_{\tau \sim \mathbb{P}_{w_{old}}(\cdot)} \left[R(\tau) \nabla_w \log \mathbb{P}_w(\tau) \Big|_{w=w_{old}} \right]\end{aligned}$$

Intuition: evaluate policy and change parameters to make good trajectory more likely.
However, what is P?

Policy gradient: exploit the sequential structure

Decompose prob. of trajectory in prob of states (given by dynamics D) and prob of actions (policy π_w):

$$\begin{aligned} \nabla_w \log \mathbb{P}_w(\tau) &= \nabla_w \log \left[\prod_{t=1}^T D(s_t | s_{t-1}, a_{t-1}) \cdot \pi_w(a_{t-1} | s_{t-1}) \right] \\ &= \nabla_w \left[\sum_{t=1}^T \log D(s_t | s_{t-1}, a_{t-1}) + \sum_{t=0}^{T-1} \log \pi_w(a_t | s_t) \right] \\ &= \sum_{t=0}^{T-1} \nabla_w \log \pi_w(a_t | s_t) \quad \text{analytically computable!!!} \\ &\quad \text{(backpropagation)} \end{aligned}$$

notation	
$a_t \sim \pi(\cdot s_t)$: policy prob. distribution
$s_{t+1} \sim D(\cdot s_t, a_t)$: dynamics prob. distribution
$\pi(a_t s_t)$: numerical value, prob. of a_t
$D(s_{t+1} s_t, a_t)$: numerical value, prob. of s_{t+1}

$$\nabla_w J(w) = \mathbb{E}_{\tau \sim \mathbb{P}_w} [R(\tau) \nabla_w \log \mathbb{P}_w(\tau)] = \mathbb{E} \left[\left(\sum_{t'=1}^T r_{t'} \right) \left(\sum_{t=0}^{T-1} \nabla_w \log \pi_w(a_t | s_t) \right) \middle| \begin{array}{l} a_t \sim \pi_w(\cdot | s_t), \\ s_{t+1}, r_{t+1} \sim D(\cdot | s_t, a_t) \end{array} \right]$$

Policy gradient: analytical properties (1/2)

$$\nabla_w J(w) = \mathbb{E} \left[\left(\sum_{t'=1}^T r_{t'} \right) \left(\sum_{t=0}^{T-1} \nabla_w \log \pi_w(a_t | s_t) \right) \mid \begin{array}{l} a_t \sim \pi_w(\cdot | s_t), \\ s_{t+1}, r_{t+1} \sim D(\cdot | s_t, a_t) \end{array} \right]$$

- Formula computed performing many simulations with π_w .
- Write more compactly with the state visitation frequency:

$$\eta^{\pi_w}(s) \propto \lim_{T \rightarrow \infty} \sum_{t=0}^T P(s = s_t \mid s_0 = s, a_t \sim \pi_w(\cdot | s_t))$$

- Probability of visiting a state by sampling environment ad infinitum with policy

$$\nabla_w J(w) = \mathbb{E} \left[\left(\sum_{t'=t+1}^T r_{t'} \mid \begin{array}{l} s_t = s \\ a_t = a \end{array} \right) \nabla_w \log \pi_w(a | s) \mid \begin{array}{l} s \sim \eta^{\pi_w}(\cdot), \\ a \sim \pi_w(\cdot | s) \end{array} \right]$$

$$\nabla_w J(w) = \mathbb{E} \left[\hat{q}^{\pi_w}(s, a) \nabla_w \log \pi_w(a | s) \mid \begin{array}{l} s \sim \eta^{\pi_w}(\cdot), \\ a \sim \pi_w(\cdot | s) \end{array} \right]$$

Policy gradient theorem: valid for any differentiable stochastic policy

Policy gradient: analytical properties (2/2)

$$\nabla_w J(w) = \mathbb{E} \left[\hat{q}^{\pi_w}(s, a) \nabla_w \log \pi_w(a | s) \middle| \begin{array}{l} s \sim \eta^{\pi_w(\cdot)} \\ a \sim \pi_w(\cdot | s) \end{array} \right]$$

Consider a baseline function $b(?)$, we can write expectation over actions for a given s :

$$\int_{\mathcal{A}} (\hat{q}^{\pi_w}(s, a) - b(?)) \pi_w(a | s) \nabla_w \log \pi_w(a | s) da$$

What can be the arguments of $b(?)$ such that result is not affected (unbiased operation)?

$$0 = \int_{\mathcal{A}} b(?) \pi_w(a | s) \nabla_w \log \pi_w(a | s) da = \int_{\mathcal{A}} b(?) \nabla_w \pi_w(a | s) da \xrightarrow{\text{Integrate by parts}} = b(?) \int_{\mathcal{A}} \nabla_w \pi_w(a | s) da \\ = b(?) \nabla_w \int_{\mathcal{A}} \pi_w(a | s) da \\ = b(?) \nabla_w 1 = 0$$

As long as $b(?)$ is not a function of a then the result is 0. Including:

- Any reward obtained before performing a
- State-dependent or constant baseline average rewards

E.g.:

$$\nabla_w J(w) = \mathbb{E} \left[(\hat{q}^{\pi_w}(s, a) - \bar{R}) \nabla_w \log \pi_w(a | s) \middle| \begin{array}{l} s \sim \eta^{\pi_w(\cdot)} \\ a \sim \pi_w(\cdot | s) \end{array} \right]$$

Average cumulative reward:

$$\bar{R} = \frac{1}{K} \sum_{k=1}^K \sum_{t=1}^T r_t$$

Policy gradient: recap & intuition

- **Policy-learning methods:** directly find the policy that maximizes the rewards.

- Gaussian policy: $\pi_w(a | s) \propto \exp \left[-\frac{1}{2\sigma^2} (a - \mu_w(s))^2 \right]$

- Update by SGD: $w_{k+1} = w_k + \epsilon \mathbb{E} \left[\hat{A}^{\pi_w}(s, a) \nabla_w \log \pi_w(a | s) \right]$

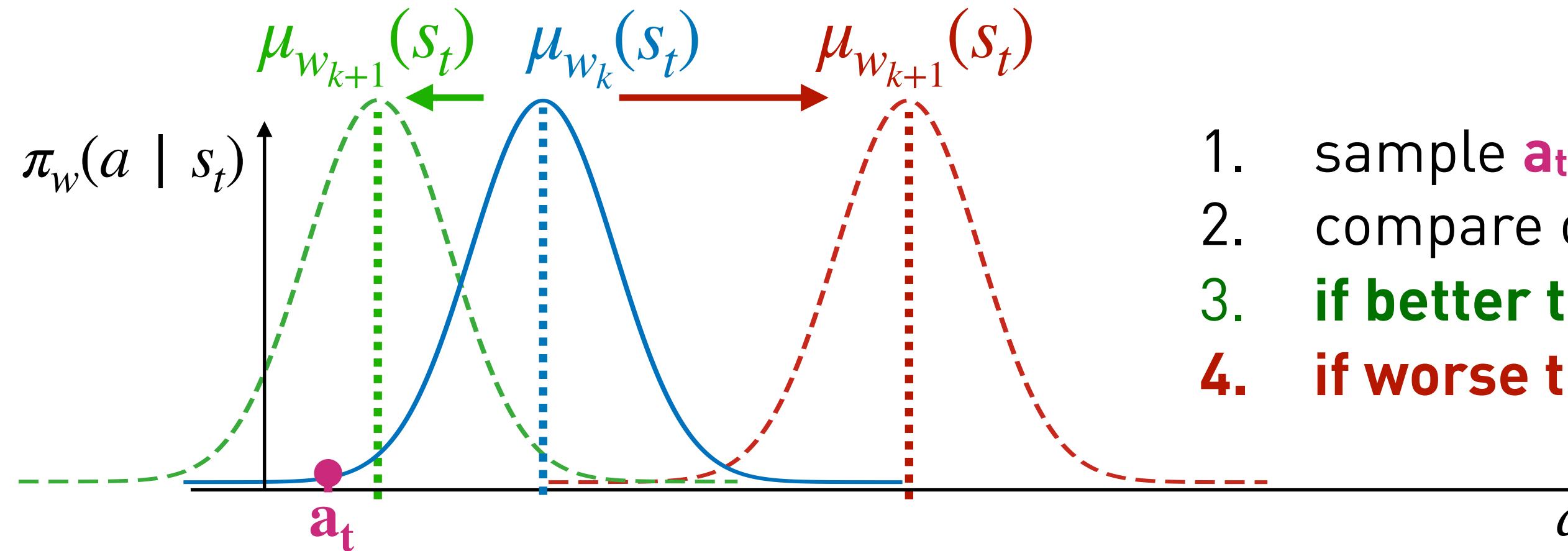
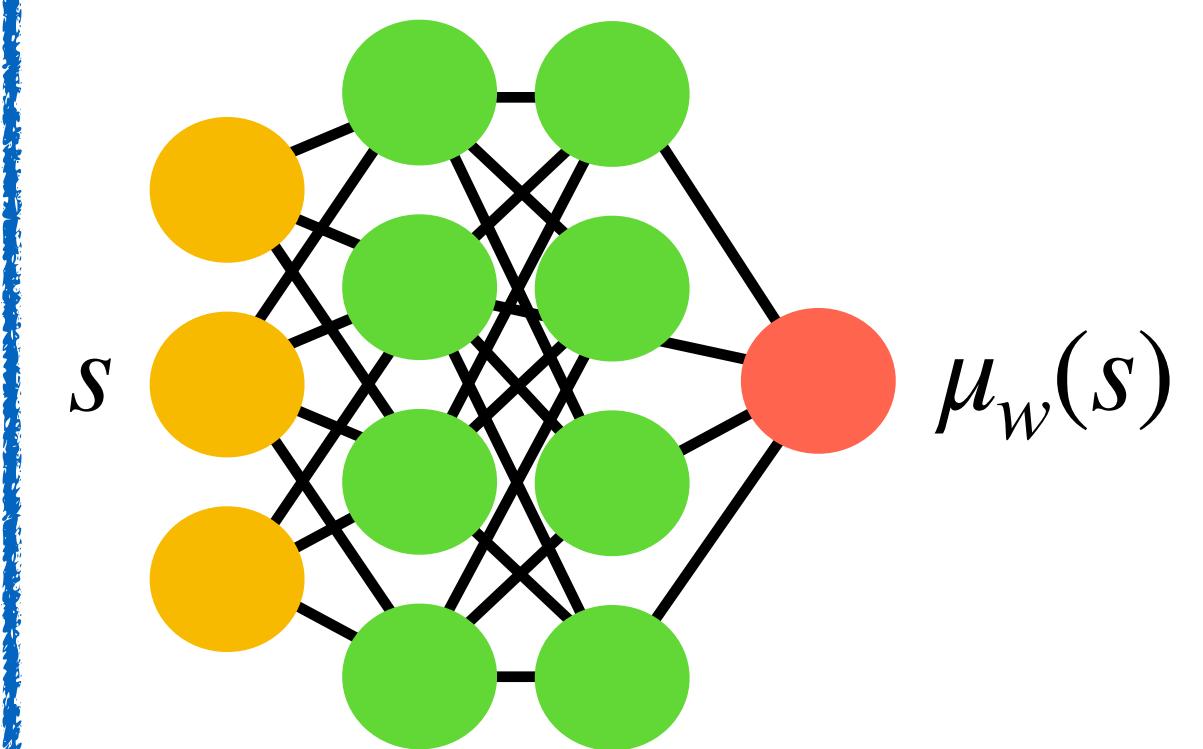
- Chain rule: $w_{k+1} = w_k + \epsilon \mathbb{E} \left[\hat{A}^{\pi_w}(s, a) \frac{a - \mu_w(s)}{\sigma^2} \cdot \nabla_w \mu_w(s) \right]$

prob of a state
given the policy

$s \sim \eta^{\pi_w}(\cdot)$
 $a \sim \pi_w(\cdot | s)$

prob of action
given the state

Assume for simplicity fixed
standard deviation of the policy



1. sample a_t
2. compare obtained rewards relative to expectation
3. **if better than expected ($A_t > 0$) : move μ_w towards a_t**
4. **if worse than expected ($A_t < 0$) : move μ_w away from a_t**

$$\hat{A}_t^{\pi_w} = \sum_{t'=t}^{T-1} \gamma^{t'-t} r_{t'+1} - \bar{R}$$

Pure formalization of “trial and error”!

Vanilla policy gradients: REINFORCE

Initialize $w^{(0)}$ and baseline $b^{(0)} = 0$
for iteration i in $1, 2, \dots, N$

- Collect E episodes by executing $\pi_{w^{(i-1)}}$
- For time step t in each episode compute:

$$\hat{q}_t^{\pi_w} = \sum_{t'=t}^{T-1} \gamma^{t'-t} r_{t'+1}$$

$$\bullet \text{ Update: } b^{(i)} = \frac{1}{E \cdot T} \sum_{t=1}^{E \cdot T} \hat{q}_t^{\pi_w}$$

- Update $w^{(i)}$ by SGD by minimizing loss:

$$\mathcal{L}^{PG}(w^{(i-1)}) = -\frac{1}{E \cdot T} \sum_{t=1}^{E \cdot T} [\hat{q}_t^{\pi_w} - b^{(i-1)}] \log \pi_{w^{(i-1)}}(a_t | s_t)$$

$$w^{(i)} = w^{(i-1)} - \epsilon \nabla_{w^{(i-1)}} \mathcal{L}^{PG}(w^{(i-1)})$$

- As formulated: simple, unbiased, but brittle: policy can diverge, requires small learning rates.
- Baseline can be replaced with NN approximation $V_w(s_t)$ to improve results.
- Algorithm assumes each lasts T steps. Simple extension if each episode has different length.
- Use automatic differentiation library (pytorch/tensorflow) to optimize policy according to \mathcal{L}^{PG} .
- Note: $[\hat{q}_t^{\pi_w} - b^{(i-1)}]$ is a constant in the calculation of the gradients.

Policy Gradient: remarks

- In principle, given enough data, PG guarantees policy improvements.
- 1-sample estimate advantage: $r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots - b(s_t)$
- Baseline reduces variance. Optimally if $b(s_t) \approx V^{\pi_w}(s_t)$.
- But value affected by $a_{t+1}, a_{t+2}, a_{t+3}, \dots$
- For PG, effect of a_t is signal, future actions are noise.
- Discount factor $\gamma \in [0,1]$ reduces importance of long term dependencies.
- Similar to using time horizon (cutoff) of $1/(1 - \gamma)$ time steps.
- Yields a lower variance, but biased (affects optimal policy), gradient.

more variance reduction: TD(λ) (Sutton & Barto, 1998),

Finite Horizon Methods (A2C, Mnih et al., 2016),

Generalized Advantage Function (GAE, Schulman et al., 2016)

Actor-critic : reducing the variance of PG

- Actor-critic methods update two approximators:
 - Actor (policy) perform policy optimization
 - Critic (value function) perform policy evaluation

With an approximate Q_w , or V_w , we can write:

$$\nabla_w J(w) = \mathbb{E} \left[\Psi^{\pi_w}(s, a) \nabla_w \log \pi_w(a | s) \middle| \begin{array}{l} s \sim \eta^{\pi_w}(\cdot) \\ a \sim \pi_w(\cdot | s) \end{array} \right]$$

$$\text{E.g.: } \Psi^{\pi_w}(s, a) = \hat{q}^{\pi_w}(s, a) - \bar{R}$$

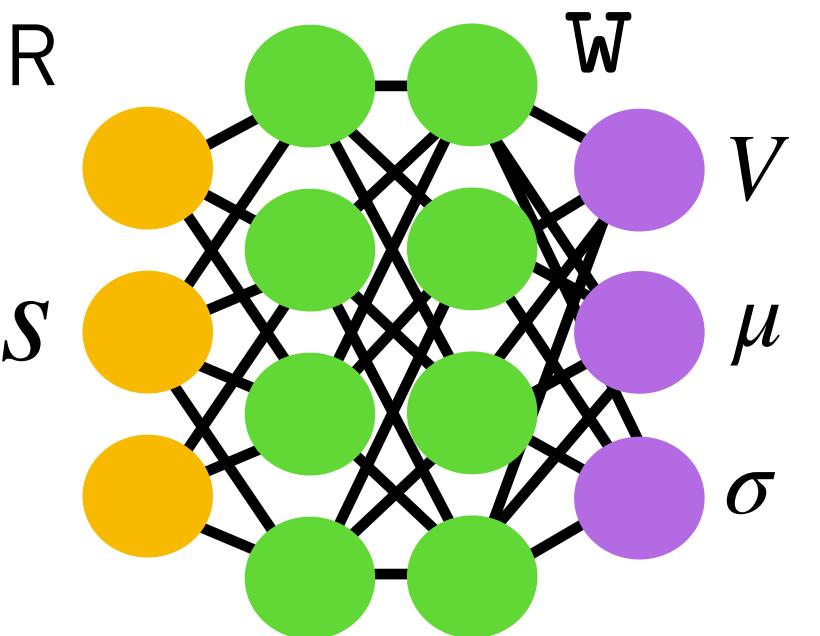
$$\Psi^{\pi_w}(s, a) = r' + \gamma V_{w'}(s') - V_{w'}(s)$$

$$\Psi^{\pi_w}(s, a) = Q_w(s, a)$$

$$\Psi^{\pi_w}(s, a) = \hat{q}^{\pi_w}(s, a) - V_w(s) = \hat{A}^{\pi_w}(s, a)$$

Usual tradeoff: actor-critic methods reduce the variance of the policy gradient at the cost of risking biased estimates (if the critic is wrong)

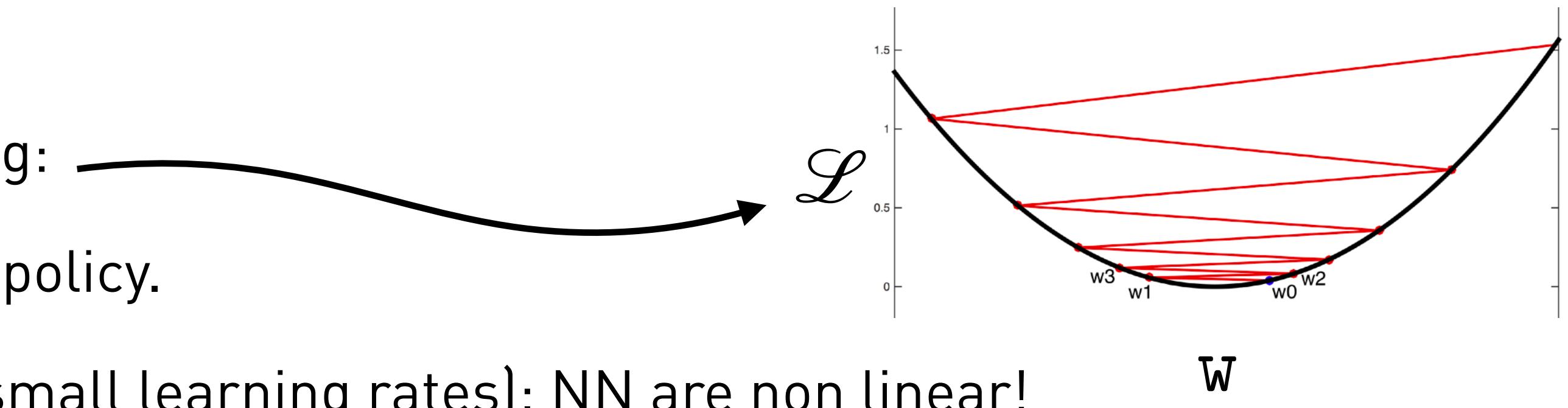
e.g. V-RACER



One NN is both actor and critic: 2 optimization objectives collaborate to learning robust hidden layer structure

Policy step size

- Large steps in supervised learning:
- RL steps may completely change policy.
- Even small parameter changes (small learning rates): NN are non linear!
- On-policy data will be unrelated to past experiences.
- May destroy all previous learning progress (catastrophic forgetting)
- Constrained optimization:



$$\text{maximize: } J(w) = \mathbb{E}_{\pi_{w^{(i)}}} \left[\frac{\pi_w(a | s)}{\pi_{w^{(i)}}(a | s)} A^{\pi_{w^{(i)}}}(s, a) \right]$$

such that: $D_{KL} [\pi_w(\cdot | s) \parallel \pi_{w^{(i)}}(\cdot | s)] < \delta$

$$D_{KL}[\cdot \parallel \cdot] = \int P(x) \log \frac{P(x)}{Q(x)} dx$$

Natural policy gradient (1/2)

$$D_{KL} [\pi_w \parallel \pi_{w^{(i)}}] \approx D_{KL} [\pi_{w^{(i)}} \parallel \pi_{w^{(i)}}] + (w - w^{(i)}) \nabla_w D_{KL} [\pi_w \parallel \pi_{w^{(i)}}]_{w=w^{(i)}} = 0$$

$$J(w) \approx J(w^{(i)}) + (w - w^{(i)}) \nabla_w J(w) \Big|_{w=w^{(i)}}$$

denote $\hat{g} = \nabla_w J(w) \Big|_{w=w^{(i)}}$

policy gradient

Approximate, linearized problem: $w^{(i+1)} = \arg \max_w \hat{g} \cdot (w - w^{(i)})$

$$\text{s.t. } \frac{1}{2} (w - w^{(i)})^T F (w - w^{(i)}) < \delta$$

$$\text{Analytical solution: } w^{(i+1)} = w^{(i)} + \sqrt{\frac{2\delta}{\hat{g}^T F^{-1} \hat{g}}} F^{-1} \hat{g}$$

natural policy gradient

Fisher Information matrix

$$F := \nabla_w^2 D_{KL} [\pi_w \parallel \pi_{w^{(i)}}]_{w=w^{(i)}}$$

Sensitivity of policy to parameter changes (metric tensor).

Natural policy gradient (2/2)

$$w^{(i+1)} = w^{(i)} + \sqrt{\frac{2\delta}{\hat{g}^T F^{-1} \hat{g}}} F^{-1} \hat{g}$$

The diagram illustrates the natural policy gradient update rule. The update step is shown as:

$$w^{(i+1)} = w^{(i)} + \sqrt{\frac{2\delta}{\hat{g}^T F^{-1} \hat{g}}} F^{-1} \hat{g}$$

A green oval encloses the term $\sqrt{\frac{2\delta}{\hat{g}^T F^{-1} \hat{g}}} F^{-1} \hat{g}$. A red arrow labeled "policy gradient" points to the term $F^{-1} \hat{g}$. A green arrow labeled "natural policy gradient" points to the term $\sqrt{\frac{2\delta}{\hat{g}^T F^{-1} \hat{g}}}$.

- Fisher information matrix F : sensitivity of policy to parameter changes.
- Even if policy is non-linear NN: **natural policy gradient produces the same policy change regardless of the model parameterization.**
- **Issue:** quadratic computational complexity (it is a Hessian).
- How to pick δ to guarantee improvements: Trust Region Policy Optimization, Schulman et al. 2015
- Linear complexity approximation (most popular): Proximal Policy Optimization, Schulman et al. 2017

Reinforcement Learning Workshop

Day 5
Experience Replay

Guido Novati

CSElab

Computational Science & Engineering Laboratory

ETH

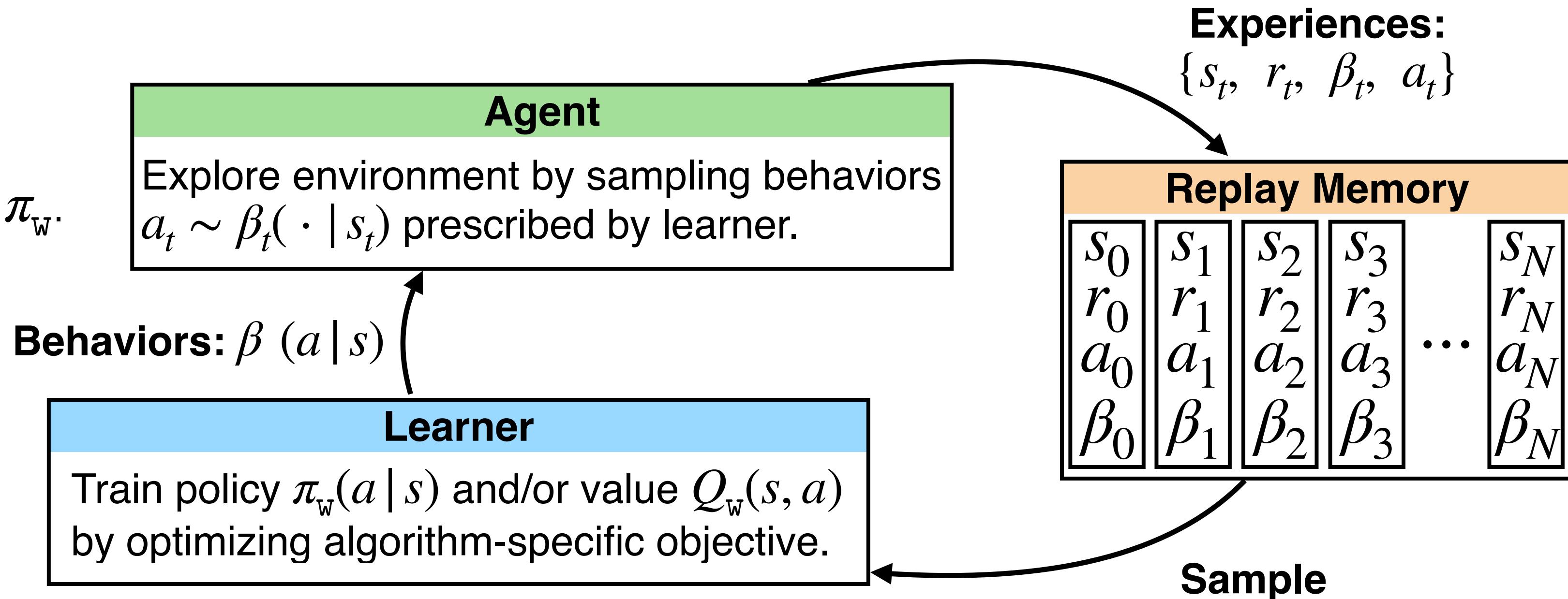
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

Why Experience Replay

- Algorithms seen so far alternate between evaluating π_w and updating π_w :
 - **Simple** implementation.
 - **Sound** theoretical analysis.
 - After update, old experiences are discarded: extremely **data-inefficient**.
 - Experiences gathered sequentially: **temporal correlations** in training data.
 - Temporal correlations: **reduced effective Monte Carlo sample size** (high update variance).
 - High update variance: **more data** is needed **to compute** good **update** estimates.
- **Solution: re-use previously collected experiences** over multiple learning iterations.
 - Allows **high data-efficiency**
 - Allows **sampling uncorrelated data** from a training set
 - Takes RL closer to Supervised Learning (allows employing more SL techniques)
 - More cumbersome implementation and **requires compatible algorithms**.

Experience replay

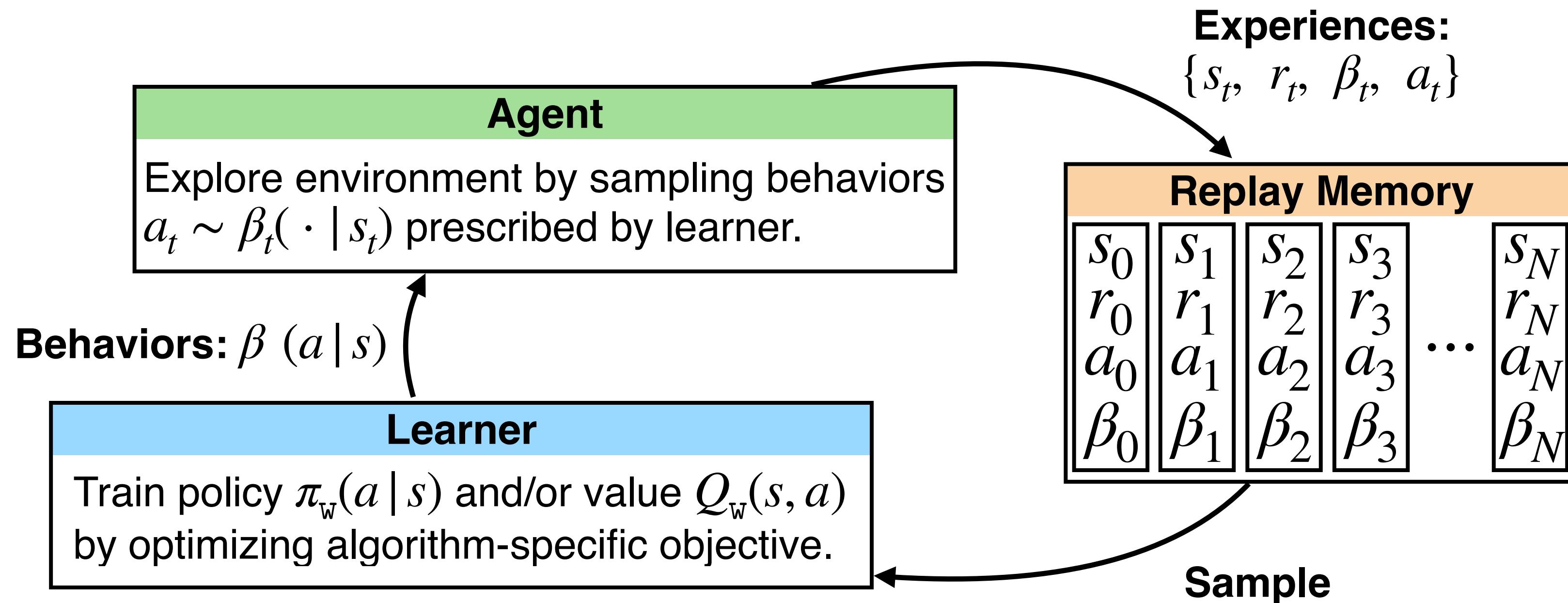
- Off-policy **training data** for policy π_w .
- **Collected by behaviors** β



For example:

- Re-use experiences collected by previous iterations of $\pi_{w^{(i)}}$
- Learn about optimal policy while collecting data with exploration strategies (e.g. Q-learning).
- Learn by observe expert behaviors (expert demonstrations by humans/other agents).
- Learn from dataset of experiences without having access to environment (e.g. safety critical).
- Learn multiple policies at once, each optimal w.r.t. a specific reward function.

Experience replay



- Typically, RM contains N steps with first-in first-out elimination.
- Experiences may be organized in episodes, or individual steps
- Time step counter t tracks total number of actions performed in environment.
- Counter t uniquely identifies the tuple $\{s_t, r_t, \beta_t, a_t\}$
- Let gradient counter k track number of policy updates.
- Many algorithms: “Add F actions to RM per policy update” (typically $F = 1$)

Q-learning for NN: Deep Q-Networks (DQN)

Initialize $w^{(0)}$ and target weights \bar{w} for iteration t in 1, 2, ..., N

DQN

- Advance one environment step with ϵ -greedy policy β_t
- Store experience $\{s_t, r_t, \beta_t, a_t\}$ in RM
- Sample a mini-batch of B experiences:
$$\{s_i, r_i, \beta_i, a_i\}_{i=1:B}$$
- With successor states and reward: $\{s'_i, r'_i\}_{i=1:B}$
- Set Q-learning target for each i :
$$\hat{q}_i = r'_i + \gamma \max_{a'} Q_{\bar{w}}(s'_i, a')$$
- Update $w^{(t)}$ with SGD by minimizing loss:
$$\hat{\mathcal{L}}^{MSE}(w^{(i)}) = \frac{1}{B} \sum_{i=1}^B \left[\frac{1}{2} [\hat{q}_i - Q_{w^{(i)}}(s_i, a_i)]^2 \right]$$
- Every D steps update target weights: $\bar{w} \leftarrow w^{(t)}$

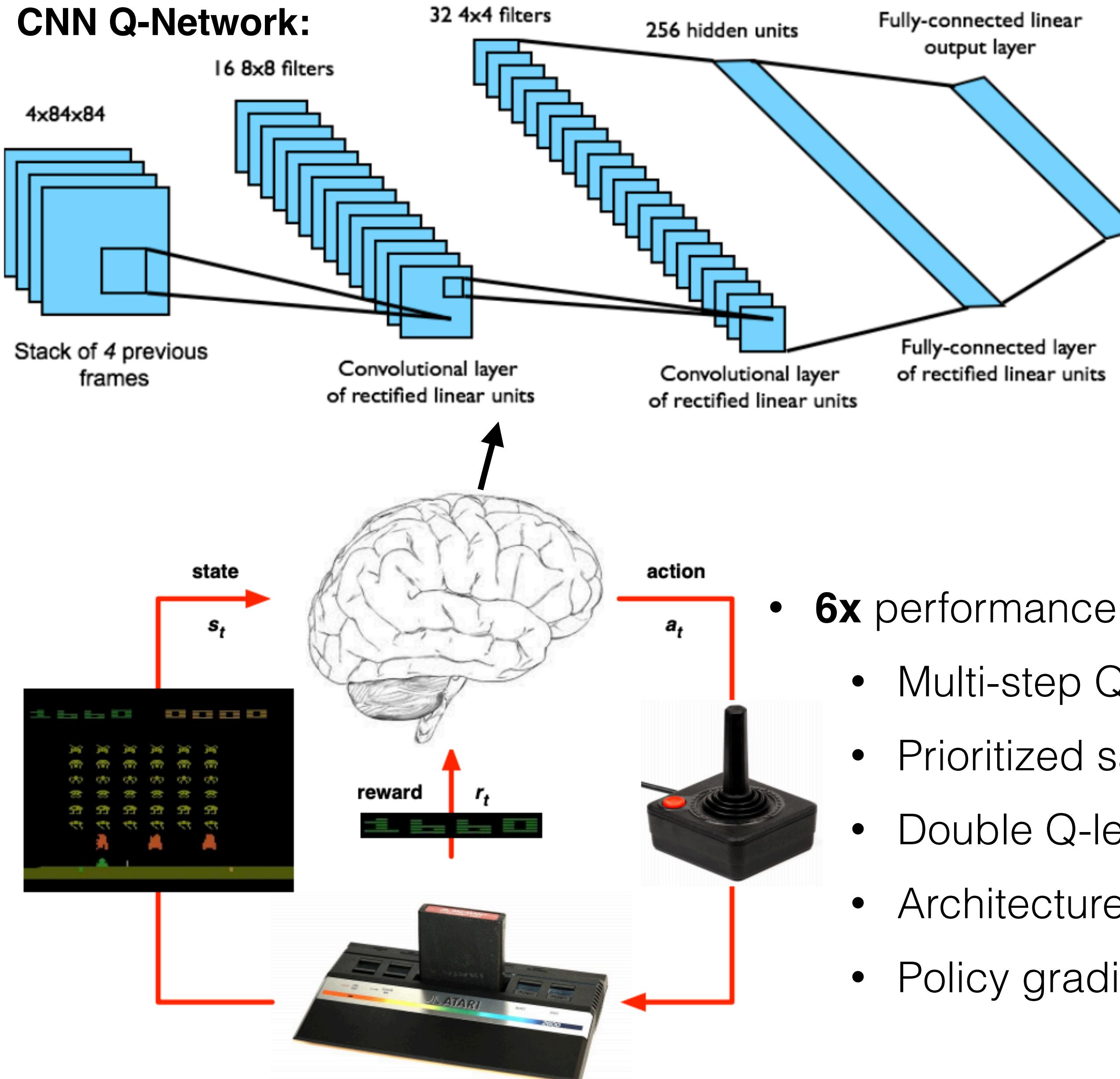
- Fixed-point iteration convergence of Q-learning unstable with deep learning.
- Target-weights are a workaround to stabilize Q-learning. Often unsuccessfully*.
- Target-weights \bar{w} make DQN update similar to loss minimization of supervised learning.
- Behavior β_t is ϵ_t -greedy policy
- Exploration rate ϵ_t annealed during training (from 1 to 0.05).

DQN : Mnih et al. 2013

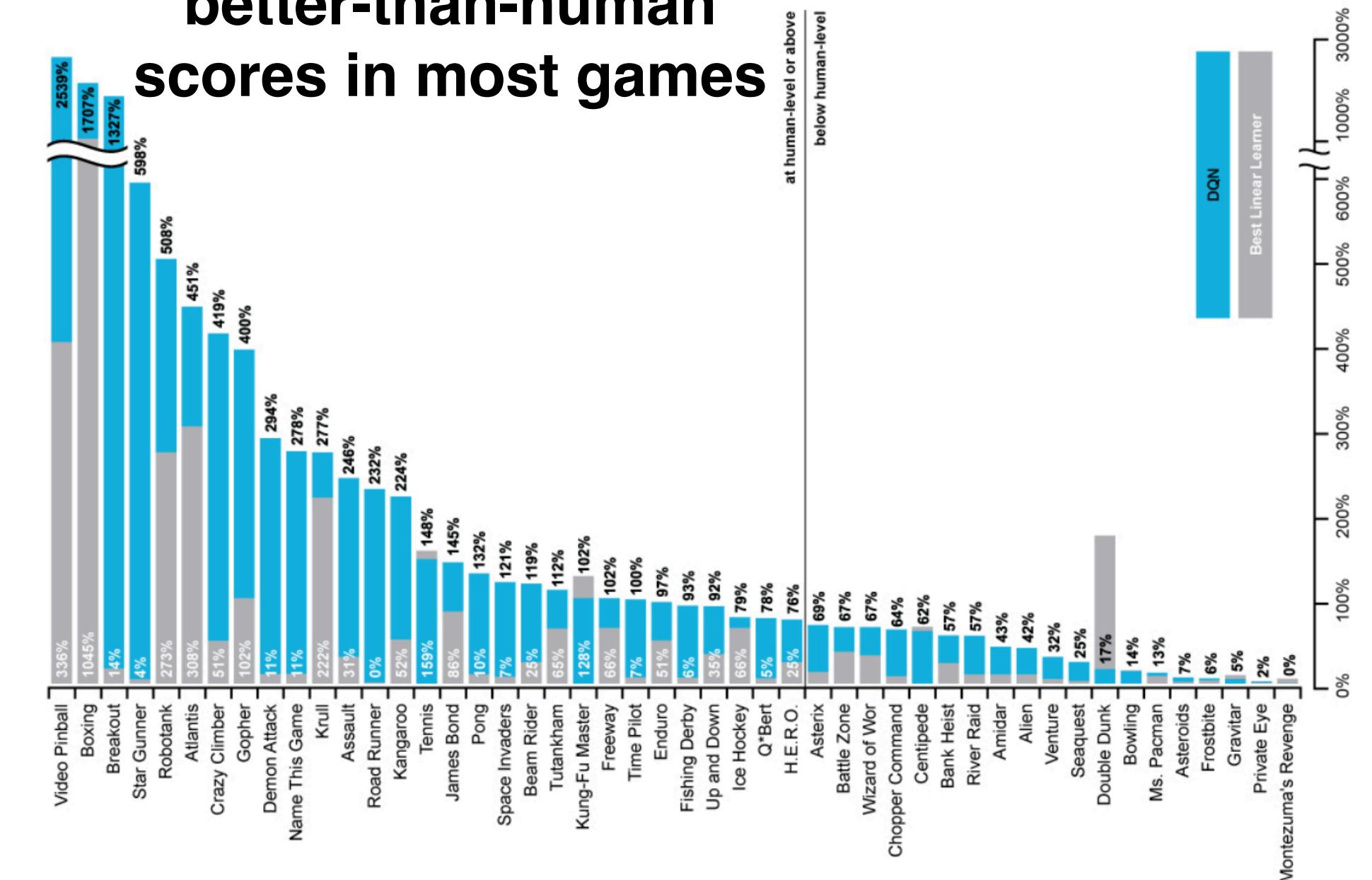
*Double DQN : H van Hasselt et al. 2015

DQN for Atari

CNN Q-Network:



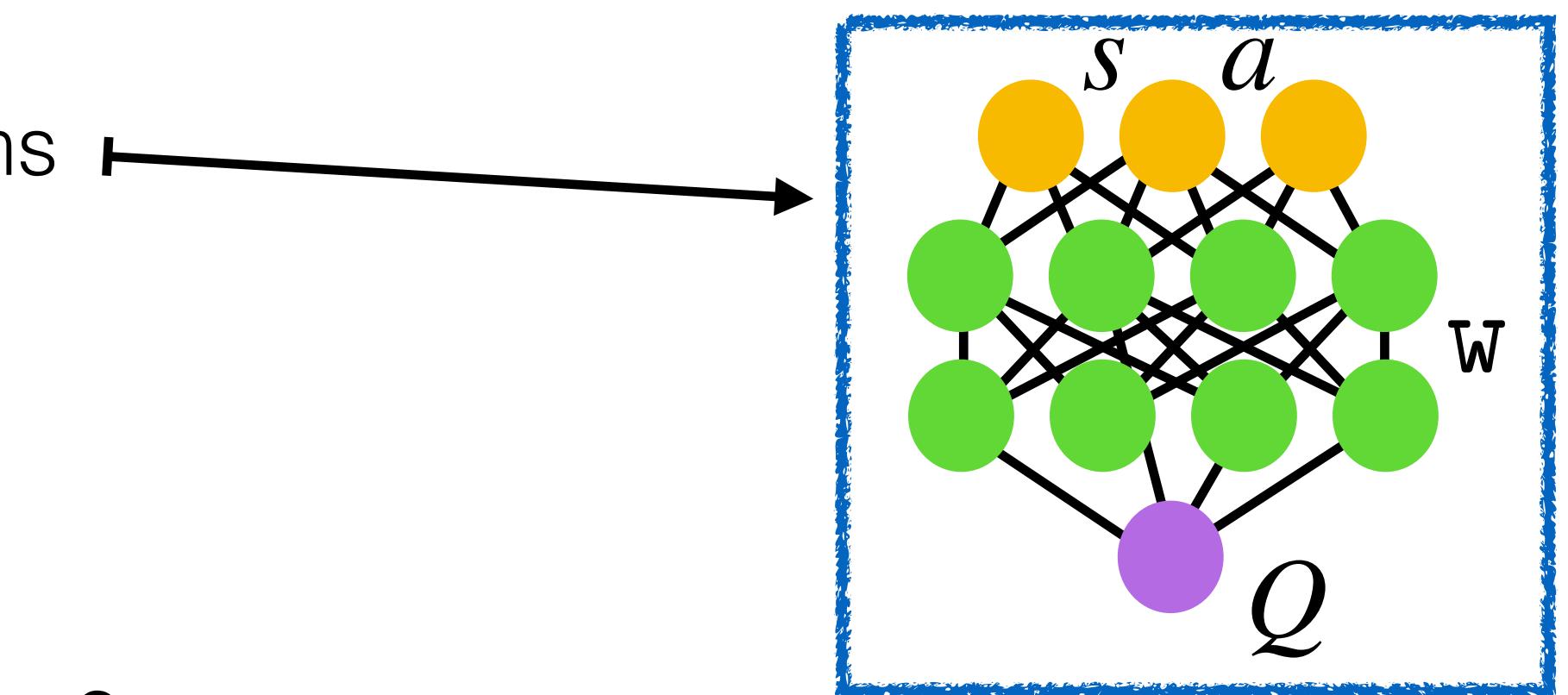
better-than-human
scores in most games



- **6x** performance improvements over the years:
 - Multi-step Q-target (Sutton & Barto, 1998)
 - Prioritized sampling for experience replay (Schaul et al, 2015)
 - Double Q-learning (van Hasselt et al, 2015)
 - Architecture: Dueling (Wang et al., 2016), Noisy (Fortunato et al. 2017)
 - Policy gradients instead of Q-learning (ACER, Wang et al., 2017)

Q-learning for continuous actions

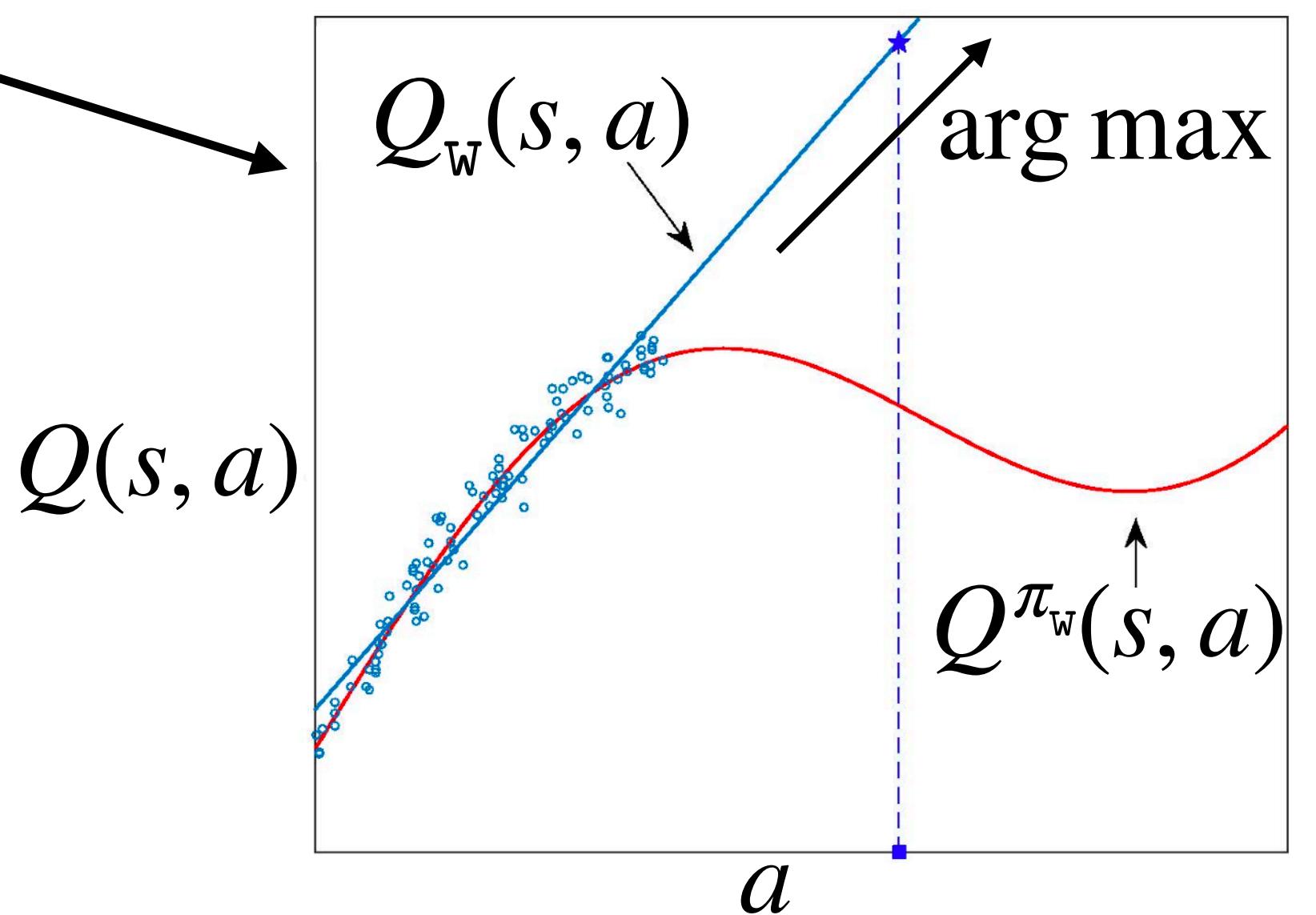
- Continuous-valued action vectors prevalent in robotics / PDEs / ODEs (even modern video-games!)
- Both s and a are vectors: natural choice for Q-network seems



- **But:**

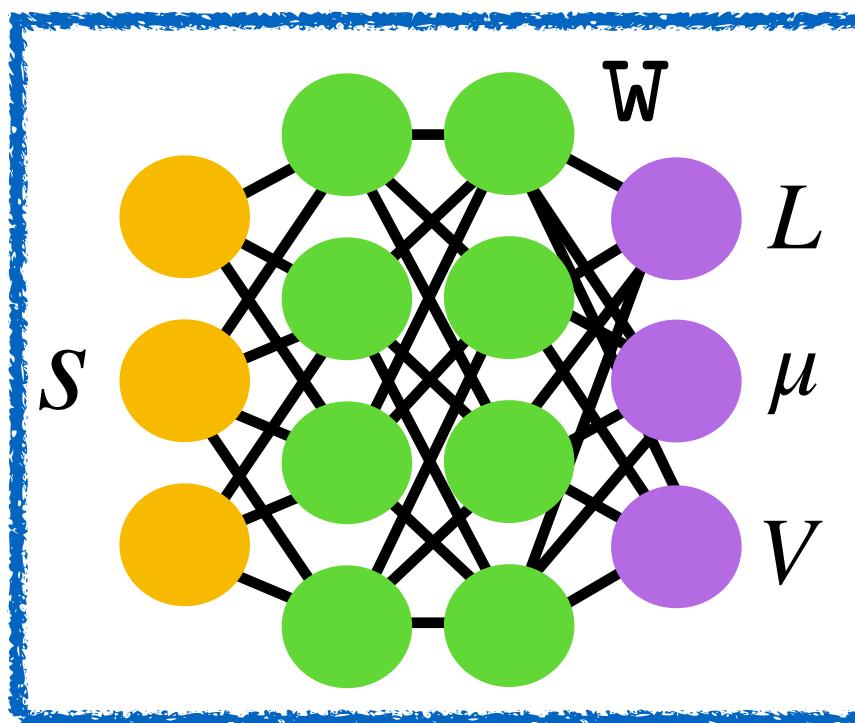
- Expensive to find maxima of a NN (policy)
- How to ensure Q_w is good approximation in all action space?
- If linear approximation: $\arg \max_a Q_w(s, a) \rightarrow \infty$

How can we, for a given s , output a well-behaved approximation of Q_w over the entire action space? (like DQN)



Continuous Q-learning: Normalized Advantage Functions

- NAF network architecture:



Initialize $w^{(0)}$
for iteration t in $1, 2, \dots, N$

- Collect D episodes with $\beta(\cdot | s) = \mathcal{N}(\mu_{w^{(i)}}(s), \sigma^2)$
- For each step, compute target with Bellman Eq:

$$\hat{q}_t^{\pi_{w^{(i)}}} = r_{t+1} + \gamma V_{w^{(i)}}(s_{t+1})$$

- Update $w^{(t)}$ with SGD by minimizing loss:

$$\mathcal{L}^{MSE}(w^{(i)}) = \frac{1}{N} \sum_{t=0}^N \left[\frac{1}{2} [\hat{q}_t^{\pi_{w^{(i)}}} - Q_{w^{(i)}}(s_t, a_t)]^2 \right]$$

- Output $V_w(s)$ approximate state value
- $\mu_w(s)$ vector of size $\dim(A)$
- $L_w(s)$ lower triangular matrix*:

$$L_w(s) = \begin{bmatrix} l_w^{(1,1)} & 0 & 0 & \dots & 0 \\ l_w^{(2,1)} & l_w^{(2,2)} & 0 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots \\ l_w^{(d_A,1)} & l_w^{(d_A,2)} & l_w^{(d_A,3)} & \dots & l_w^{(d_A,d_A)} \end{bmatrix}$$

- $L_w(s)L_w^T(s)$ is positive definite!
- $Q_w(s, a) = V_w - \frac{1}{2}(a - \mu_w)^T L_w L_w^T(a - \mu_w)$
- $\mu_w(s)$ is action with optimal value $V_w(s)$!

* To ease training (bijection) elements on the diagonal should be pos. def. $l_w^{(i,i)} \leftarrow \log [1 + \exp(l_w^{(i,i)})]$

Continuous Q-Learning: Normalized Advantage Functions

Initialize $w^{(0)}$ and target weights \bar{w}
for iteration t in 1, 2, ..., N

- Advance one step with $\beta_t(a | s_t) = \mathcal{N}(\mu_{w^{(i)}}(s_t), \sigma_t^2)$
- Store experience $\{s_t, r_t, \beta_t, a_t\}$ in RM
- Sample a mini-batch of B experiences:

$$\{s_i, r_i, \beta_i, a_i\}_{i=1:B}$$

- With following states and reward: $\{s'_i, r'_i\}_{i=1:B}$
- Set Q-learning target for each i :

$$\hat{q}_i = r'_i + \gamma V_{\bar{w}}(s'_i)$$

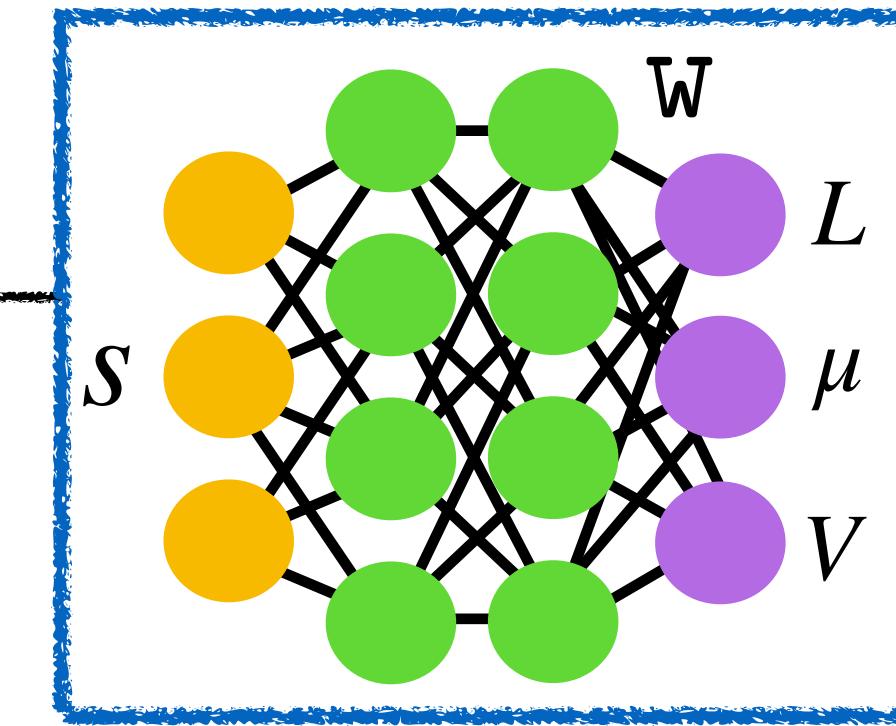
- Update $w^{(t)}$ with SGD by minimizing loss:

$$\hat{\mathcal{L}}^{MSE}(w^{(i)}) = \frac{1}{B} \sum_{i=1}^B \left[\frac{1}{2} [\hat{q}_i - Q_{w^{(i)}}(s_i, a_i)]^2 \right]$$

- Every D steps update target weights: $\bar{w} \leftarrow w^{(t)}$

NAF

- NAF contribution:



- Output $V_w(s)$ approximate state value
- $\mu_w(s)$ vector of size $\dim(A)$
- $L_w(s)$ lower triangular matrix*:

$$L_w(s) = \begin{bmatrix} l_w^{(1,1)} & 0 & 0 & \cdots & 0 \\ l_w^{(2,1)} & l_w^{(2,2)} & 0 & \cdots & 0 \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ l_w^{(d_A,1)} & l_w^{(d_A,2)} & l_w^{(d_A,3)} & \cdots & l_w^{(d_A,d_A)} \end{bmatrix}$$

- $L_w(s)L_w^T(s)$ is positive definite!
- $Q_w(s, a) = V_w - \frac{1}{2}(a - \mu_w)^T L_w L_w^T (a - \mu_w)$
- $\mu_w(s)$ is action with optimal value $V_w(s)$!

Importance sampling for off-policy rewards

Perform “roll-outs” with behavior β :

$$\begin{aligned} s_0^{(i)} &\sim D_0(\cdot) \\ a_0^{(i)} &\sim \beta(\cdot | s_0^{(i)}) \\ s_1^{(i)}, r_1^{(i)} &\sim D(\cdot | s_0^{(i)}, a_0^{(i)}) \\ &\dots \\ a_{T-1}^{(i)} &\sim \beta(\cdot | s_{T-1}^{(i)}) \\ s_T^{(i)}, r_T^{(i)} &\sim D(\cdot | s_{T-1}^{(i)}, a_{T-1}^{(i)}) \end{aligned}$$

Off-policy Episode: $\tau^\beta = \{s_t, \beta_t, a_t, r_t\}_{t=0, \dots, T}$

Off-policy MC returns: $\hat{q}_t^\beta = \sum_{t'=t}^T \gamma^{t'-t} r_{t'+1}^{(i)}$

TD error: $\delta_t^\beta = r_{t+1} + \gamma Q_w(s_{t+1}, a_{t+1}) - Q_w(s_t, a_t)$

MC target: $\hat{q}_t^\pi = Q_w(s_t, a_t) + \sum_{t'=t}^{T-1} \gamma^{t'-t} \delta_{t'}^\pi$

Importance sampling for off-policy SARSA:

$$\hat{q}_t^\pi = r_{t+1} + \gamma \frac{\pi(a_{t+1} | s_{t+1})}{\beta(a_{t+1} | s_{t+1})} Q_w(s_{t+1}, a_{t+1})$$

- Relatively low variance but biased
- Unbounded importance weights may cause numerical instability

Importance sampling for off-policy MC returns:

$$\hat{q}_t^\pi = \frac{\pi(a_{t+1} | s_{t+1})}{\beta(a_{t+1} | s_{t+1})} \cdot \dots \cdot \frac{\pi(a_{T-1} | s_{T-1})}{\beta(a_{T-1} | s_{T-1})} \cdot \hat{q}_t^\beta$$

- Dramatic increase in variance
- Unstable due to exploding product
- Is there a lower variance expression?

Importance sampling for off-policy TD errors

- Given off-policy episode: $\tau^\beta = \{s_t, \beta_t, a_t, r_t\}_{t=0, \dots, T}$
- Introduce:
$$\delta_t^\beta = r_{t+1} + \gamma Q_w(s_{t+1}, a_{t+1}) - Q_w(s_t, a_t)$$

$$\delta_t^\pi = r_{t+1} + \gamma \mathbb{E}_{a' \sim \pi} Q_w(s_{t+1}, a') - Q_w(s_t, a_t)$$

- E.g. $\delta_t^\pi = r_{t+1} + \gamma \max_{a' \sim \pi} Q_w(s_{t+1}, a') - Q_w(s_t, a_t)$
- E.g. $\delta_t^\pi = r_{t+1} + \gamma V_w(s_{t+1}) - Q_w(s_t, a_t)$

General operator to estimate returns from off-policy data:

$$\hat{q}_t^\pi = Q_w(s_t, a_t) + \sum_{t'=t}^{T-1} \gamma^{t'-t} \delta_{t'}^\pi \left(\prod_{s=t+1}^{t'} c_s \right)$$

- Where we define $\prod_{s=t+1}^t c_s = 1$ (i.e. the term for $t' = t$) regardless of c_s
- If $c_s = 0$: we recover Q-learning
- If $c_s = \pi(a_s | s_s) / \beta(a_s | s_s)$: we recover importance sampling MC
- If $c_s = 1$: off-policy N-step methods (i.e. pretend it is not off-policy)
No base in theory, but sometimes performs well and therefore used!
- If $c_s = \min \{1, \pi(a_s | s_s) / \beta(a_s | s_s)\}$ then we use **Retrace** (Munos et al. 2016)

Retrace

$$\hat{q}_t^{retrace} = Q_w(s_t, a_t) + \sum_{t'=t}^{T-1} \gamma^{t'-t} \left(\prod_{s=t+1}^{t'} c_s \right) \left(r_{t+1} + \gamma V_w(s_{t+1}) - Q_w(s_t, a_t) \right)$$

- With $c_s = \min \{1, \pi(a_s | s_s) / \beta(a_s | s_s)\}$
- And $\prod_{s=t+1}^t c_s = 1$

Can we written recursively as:

$$\hat{q}_t^{retrace} = r_{t+1} + \gamma V_w(s_{t+1}) + \gamma \min \left[1, \frac{\pi(a_{t+1} | s_{t+1})}{\beta(a_{t+1} | s_{t+1})} \right] \left[\hat{q}_{t+1}^{retrace} - Q_w(s_{t+1}, a_{t+1}) \right]$$

- Clipping prevents variance explosion of Importance Sampling
- Safely uses all relevant off-policy rewards in the episode
- Provably converges to correct Q^π for **any** behavior policy β (tabular setting)
- Many state-of-the-art results support its efficacy (e.g. Wang 2016, Espeholt 2018)
- If β and π are similar, Retrace converges faster!

The on-policy objectives from off-policy data (1)

On-policy algorithms optimize objectives in the form: $\mathcal{L}(w) = \mathbb{E} \left[L_w(s, a) \mid \begin{array}{l} s \sim \eta^{\pi_w}(\cdot) \\ a \sim \pi_w(\cdot | s) \end{array} \right]$

$$a \sim \pi_w(\cdot | s) \qquad s \sim \eta^{\pi_w}(\cdot)$$

Probability of an action given the current policy

Probability of visiting a state with π_w

Off-policy algorithms compute $L(w)$ from experiences gathered by behaviors β :

$$a \sim \beta(\cdot | s)$$

$$s \sim \eta^\beta(\cdot) \propto \sum_{k=0} P(s = s_k | s_0, a_k \sim \beta(a | s_k))$$

Actions sampled from arbitrary behaviors β .

Probability of observing a state by acting with behaviors β .

The on-policy objectives from off-policy data (2)

On-policy algorithms optimize objectives in the form: $\mathcal{L}(w) = \mathbb{E} \left[L_w(s, a) \mid \begin{array}{l} s \sim \eta^{\pi_w}(\cdot) \\ a \sim \pi_w(\cdot | s) \end{array} \right]$

Off-policy algorithms compute $J(w)$ from expectation over past experiences:

$$a \sim \beta(\cdot | s) \quad s \sim \eta^\beta(\cdot) := \sum_{k=0} P(s = s_k | s_0, a_k \sim \beta(a | s_k))$$

example: first-order approximation of on-policy objectives from off-policy data:

$$\mathcal{L}(w) = \int_S \int_A L_w(s, a) \pi_w(a | s) da \eta^{\pi_w}(s) ds \approx \int_S \int_A L_w(s, a) \pi_w(a | s) da \eta^\beta(s) ds = \tilde{\mathcal{L}}(w)$$

approximation: assume that distribution of states almost constant

$$\tilde{\mathcal{L}}(w) = \int_S \int_A L_w(s, a) \pi_w(a | s) \frac{\beta(a | s)}{\pi_w(a | s)} da \eta^\beta(s) ds = \mathbb{E} \left[L_w(s, a) \frac{\pi_w(a | s)}{\beta(a | s)} \mid \begin{array}{l} s \sim \eta^\beta(\cdot) \\ a \sim \beta(\cdot | s) \end{array} \right]$$

change of distribution:
expectation over old policies

importance sampling

The on-policy objectives from off-policy data (3)

Recall: the on-policy Policy Gradient: $\nabla_w J(w) = \mathbb{E} \left[\hat{q}^{\pi_w}(s, a) \nabla_w \log \pi_w(a | s) \middle| \begin{array}{l} s \sim \eta^{\pi_w}(\cdot) \\ a \sim \pi_w(\cdot | s) \end{array} \right]$

Off-policy approximation of the Policy Gradient:

$$\begin{aligned} \nabla_w J(w) &\approx \mathbb{E} \left[\hat{q}^{\pi_w}(s, a) \nabla_w \log \pi_w(a | s) \middle| \begin{array}{l} s \sim \eta^{\beta}(\cdot) \\ a \sim \pi_w(\cdot | s) \end{array} \right] \\ &= \mathbb{E} \left[\hat{q}^{\pi_w}(s, a) \frac{\beta(a | s)}{\pi_w(a | s)} \nabla_w \log \pi_w(a | s) \middle| \begin{array}{l} s \sim \eta^{\beta}(\cdot) \\ a \sim \pi_w(\cdot | s) \end{array} \right] \\ &= \mathbb{E} \left[\hat{q}^{\pi_w}(s, a) \frac{\pi_w(a | s)}{\beta(a | s)} \nabla_w \log \pi_w(a | s) \middle| \begin{array}{l} s \sim \eta^{\beta}(\cdot) \\ a \sim \beta(\cdot | s) \end{array} \right] \end{aligned}$$

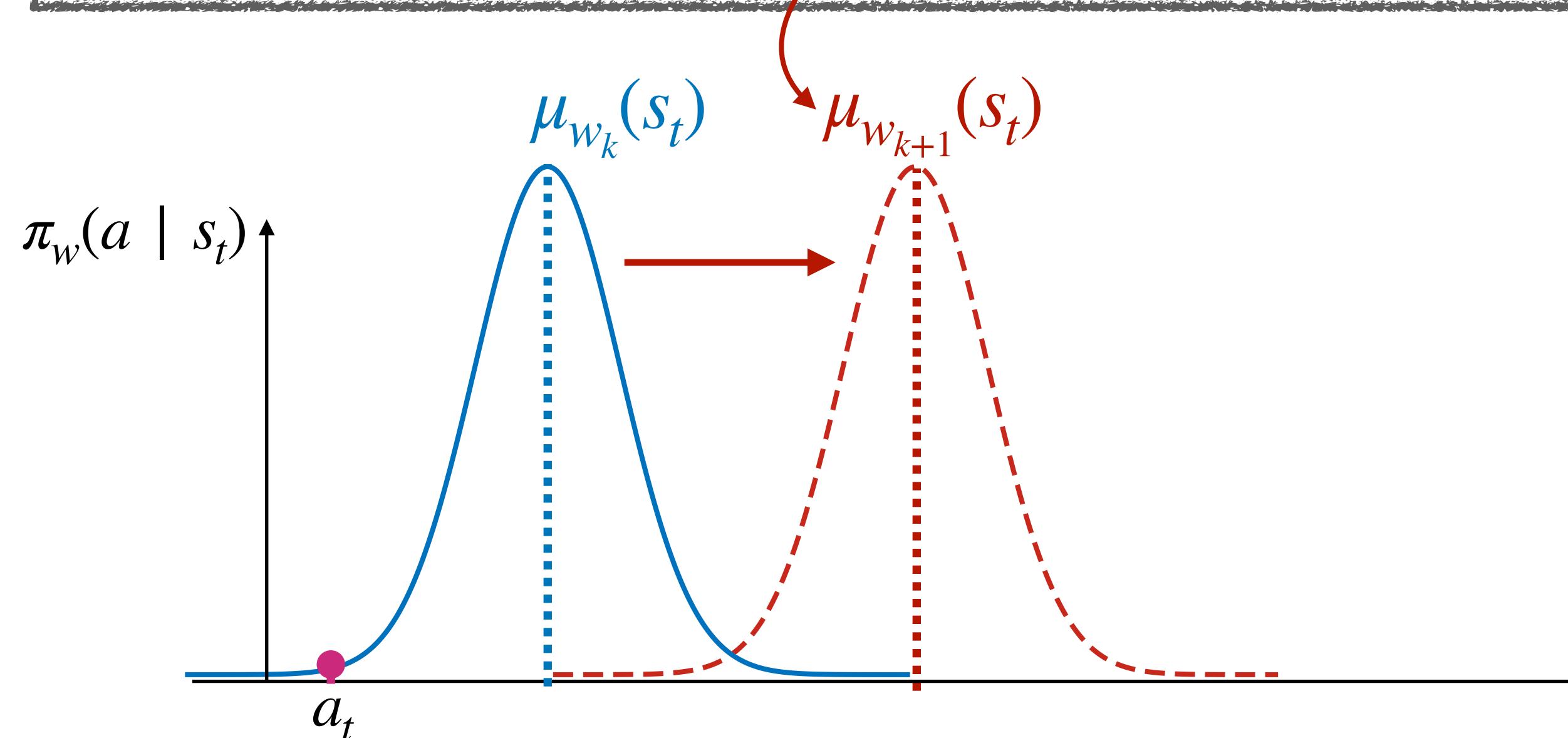
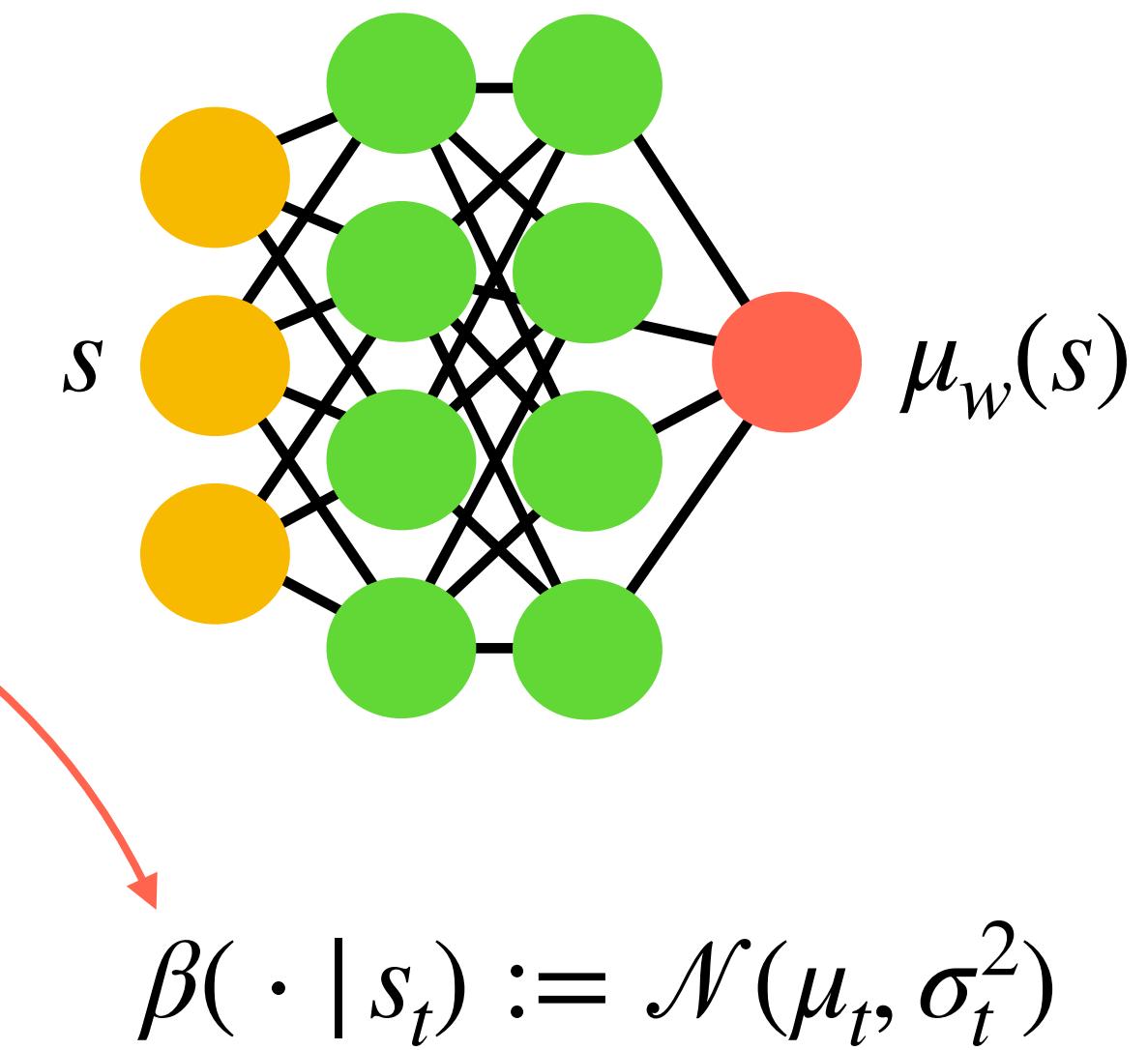
The off-policy Policy Gradient (e.g. V-RACER)

- Approximate the PG by sampling a Replay Memory and iteratively update by SGD:

$$w_{k+1} = w_k + \epsilon \mathbb{E} \left[q^{\pi_w}(s, a) \cdot \frac{\pi_w(a | s)}{\beta(a | s)} \cdot \frac{a - \mu_w(s)}{\sigma^2} \nabla_w \mu_w(s) \right]$$

Monte-Carlo estimate
from past experiences:

- Sample a past experience $\{s_t, \mu_t, \sigma_t, a_t, r_t\}$
- Approximate the returns $q_t^{\pi_w}$ (e.g. with network $Q_w(s_t, a_t)$ or with **Retrace** algorithm)
- if $q^n > 0$: move μ_w (current weights) towards a_t**
- if $q^n < 0$: move μ_w (current weights) away from a_t**



- Policy Gradient does not say **how much** μ_w **should be moved away from a_t** .
- Off-policy updates can easily cause **policy to diverge from training behaviors**.

Issues of off-policy RL

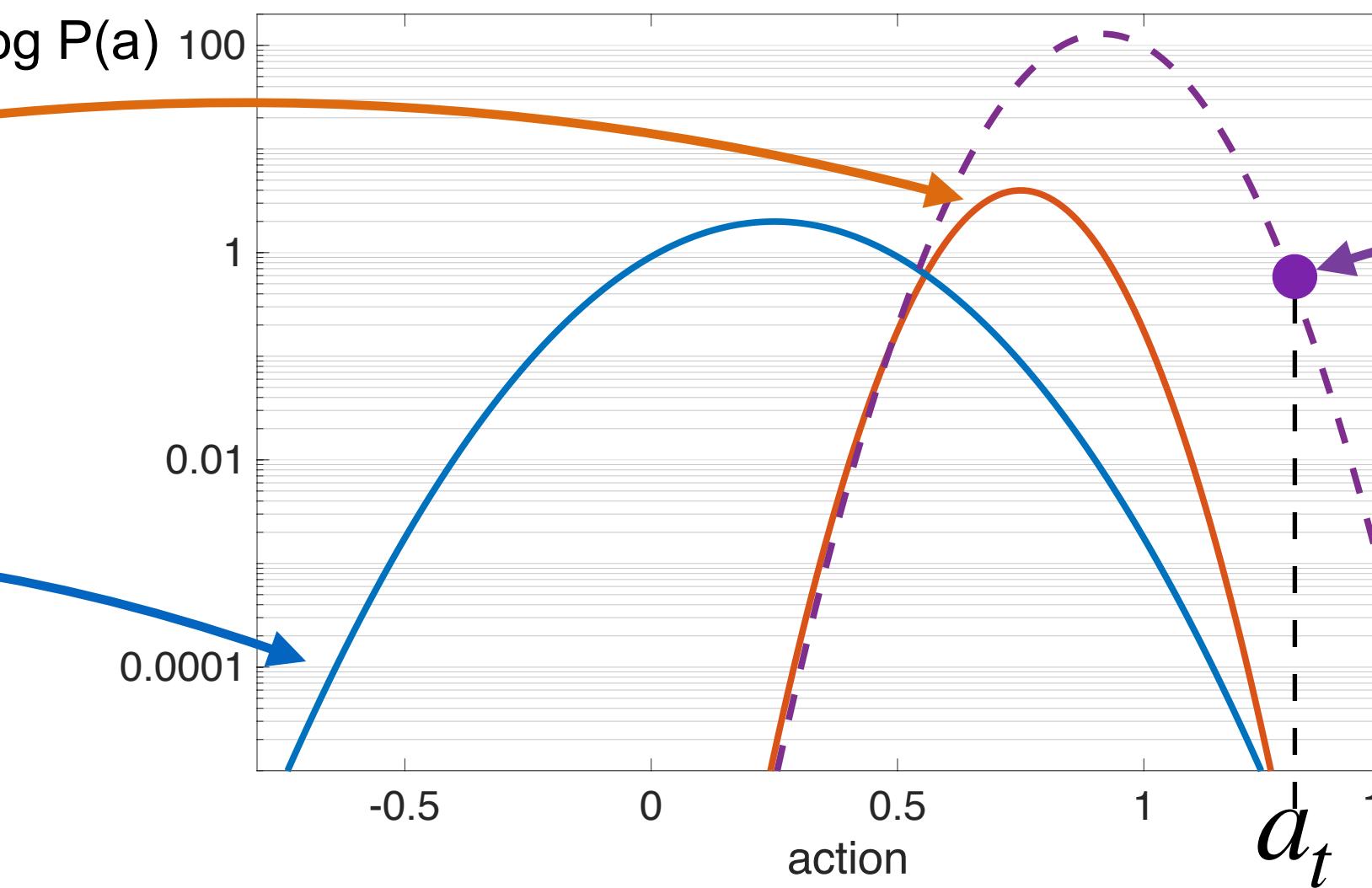
ISSUE #1

$\pi_w(a | s_t)$

Current policy that we aim to update

$\beta(a | s_t)$

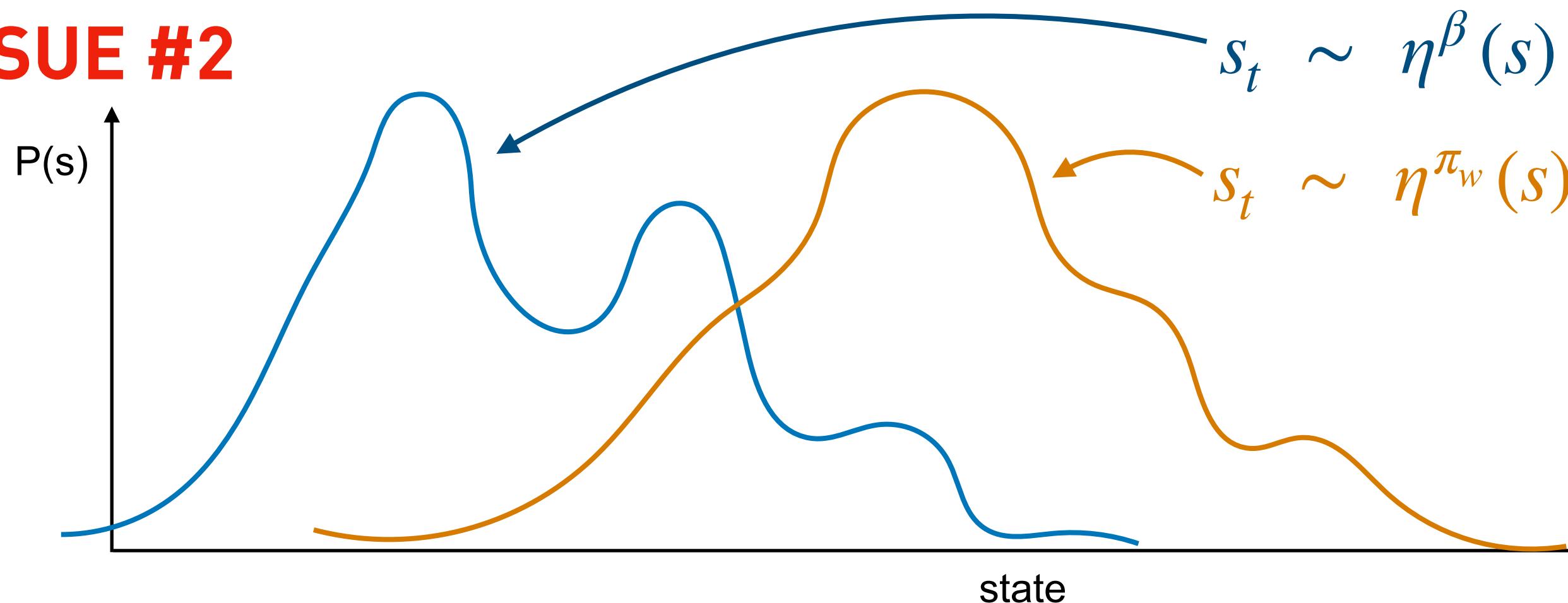
old policy followed by the agent
(training data)



$$\frac{\pi_w(a_t | s_t)}{\beta(a_t | s_t)}$$

Unbounded (0 to ∞) Importance weight, increases the variance and **lowers the accuracy of the cost function $\tilde{J}(w)$.**

ISSUE #2



$s_t \sim \eta^\beta(s)$ Prob of encountering state with old policies (training data)

$s_t \sim \eta^{\pi_w}(s)$ Prob of encountering state with current policy

If distribution of training data is too dissimilar from on-policy outcomes, data may be irrelevant to updating the policy

Proposed solution: constrain policy changes to past policies

Remember and Forget Experience Replay

ReF-ER works with most RL methods that learn a policy by Experience Replay (e.g. Q-learning, off-policy PG, ...)

$$\tilde{\mathcal{L}}(w) = \mathbb{E} \left[\tilde{L}_w(s, a) \mid s \sim \eta^\beta(\cdot), a \sim \beta(\cdot | s) \right]$$

algorithm-specific

Modify objective

of the RL algorithm:

$$\nabla_w \tilde{L}_w^{ReF-ER}(s, a) = \begin{cases} \lambda \nabla_w \tilde{L}_w(s, a) - (1 - \lambda) \nabla_w D_{KL} [\pi_w(\cdot | s) || \beta(\cdot | s)] & \text{if } \frac{1}{c} < \frac{\pi_w(a | s)}{\beta(a | s)} < c \\ 0 - (1 - \lambda) \nabla_w D_{KL} [\pi_w(\cdot | s) || \beta(\cdot | s)] & \text{otherwise} \end{cases}$$

1) Reject samples from objective estimate if importance weight π_w/β lies outside of a trust region.

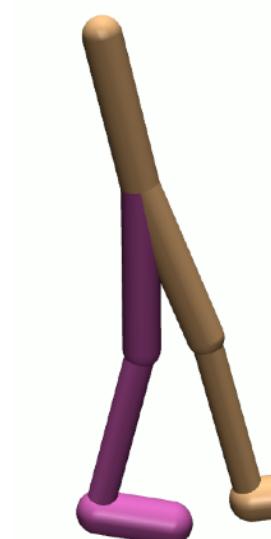
2) Add penalization cost to attract policy back towards training behaviors.

(Coefficient λ is iteratively updated to keep a fixed fraction of the Experience Replay data within the trust region.)

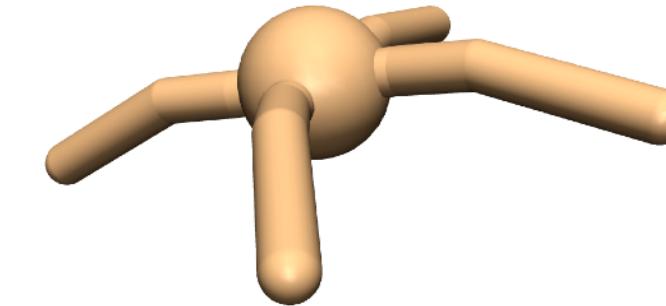
Results on OpenAI gym MuJoCo tasks (1)



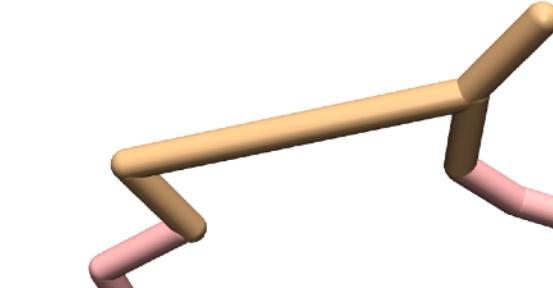
Humanoid-v2



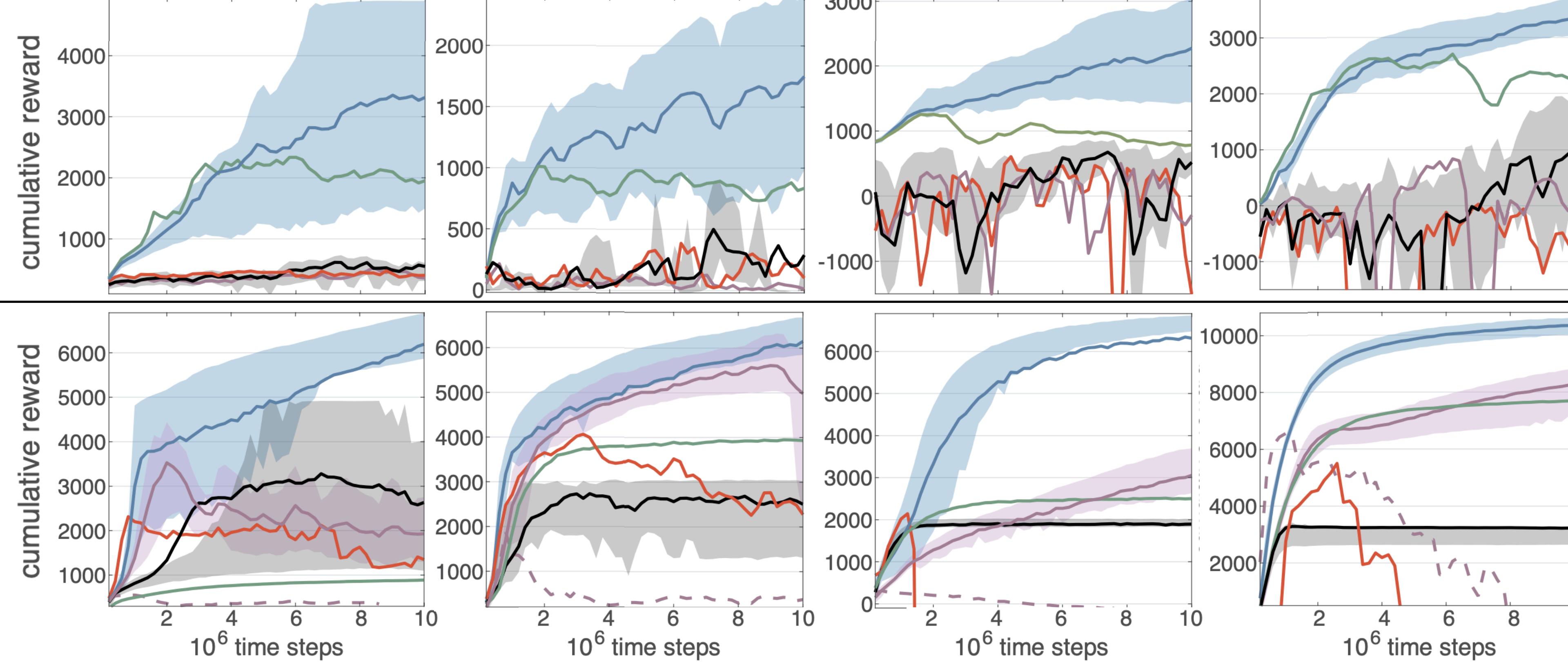
Walker2d-v2



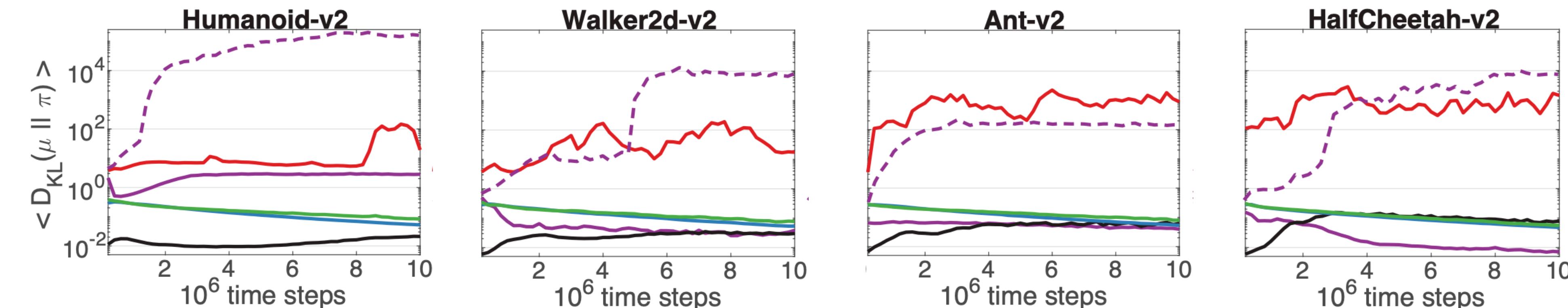
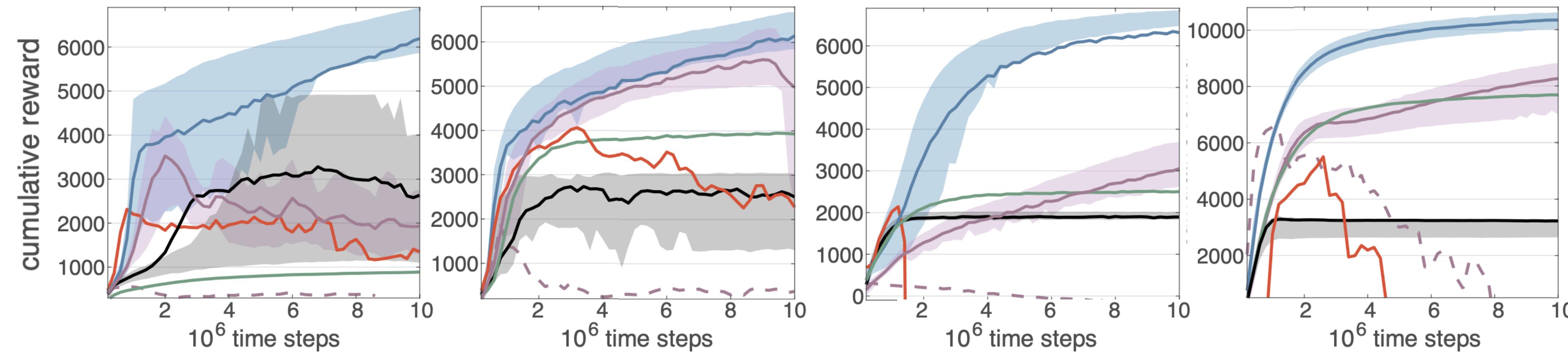
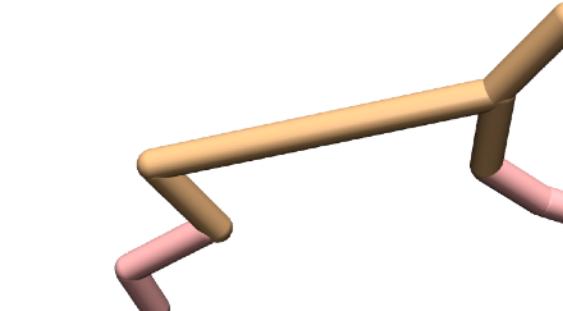
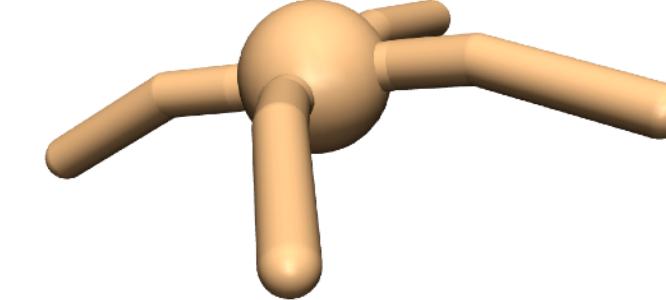
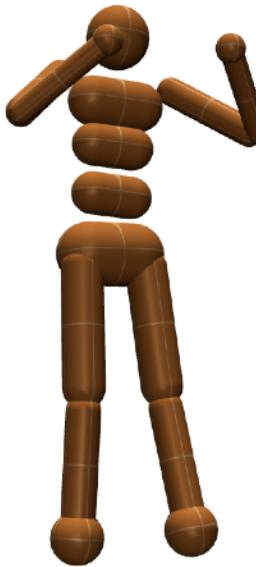
Ant-v2



HalfCheetah-v2



Results on OpenAI gym MuJoCo tasks (2)

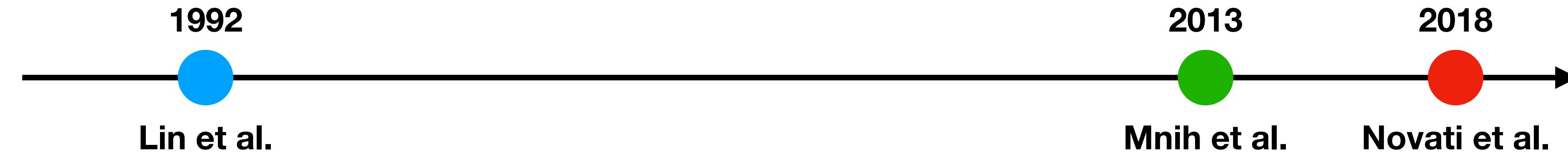


ACER relies on target networks (like DQN) and hyper-parameter optimization.
Increasing learning rate breaks stability measures (**dashed line**).
ReF-ER effectively **constrained** D_{KL} , increased stability & performance.

Off-Pol RL as Supervised Learning

- Years of research in ML gave us reliable techniques to train NN on supervised learning tasks

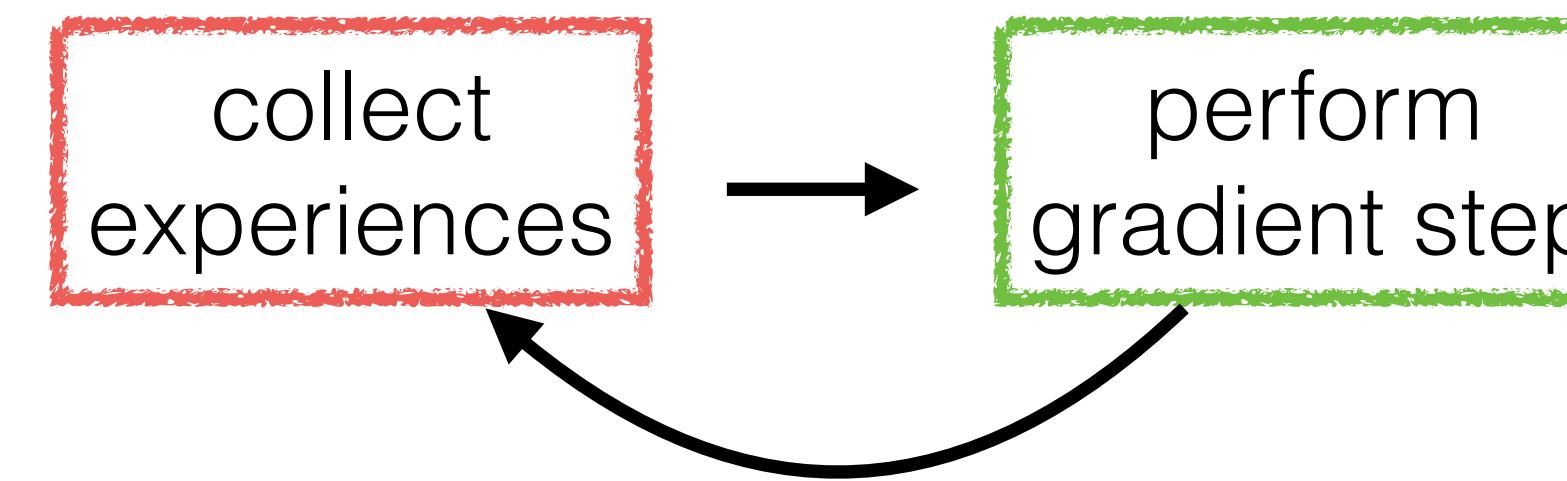
Some techniques increasing affinity of off-policy RL and supervised learning:



- **Experience Replay**: past experiences as training data.
 - Allows RL training with **mini-batches of uncorrelated samples** (closer to i.i.d. assumption)
 - Distribution of **training data is almost constant** during training
- **DQN** (Atari-paper) introduces Target Networks for Q-learning.
 - Target of Q-learning update** is no longer fixed point iteration, but **almost constant**.
- **ReF-ER** controls RL updates and training data.
 - Constrains similarity between distribution of training data and on-policy data.
 - In other words: force **test set and training set to be from almost the same distribution**

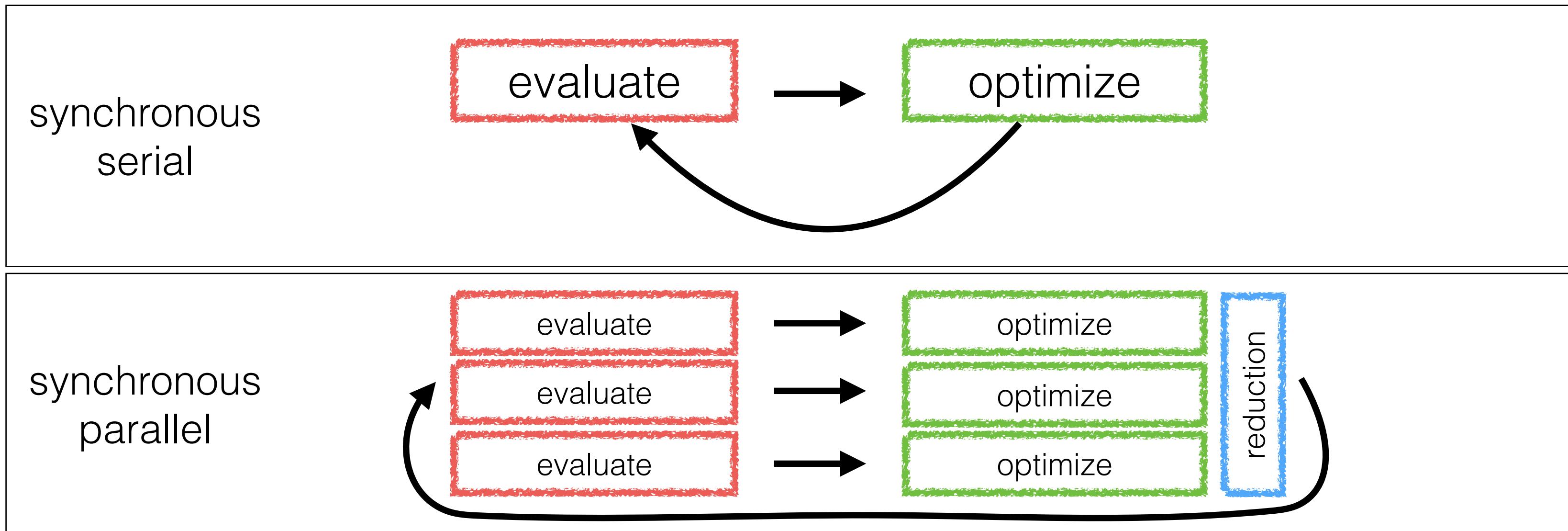
RL software architecture

- Most academic RL software is written in python (using torch or tensorflow for NN)
- Often assumes toy control problems: environment is cheap to simulate



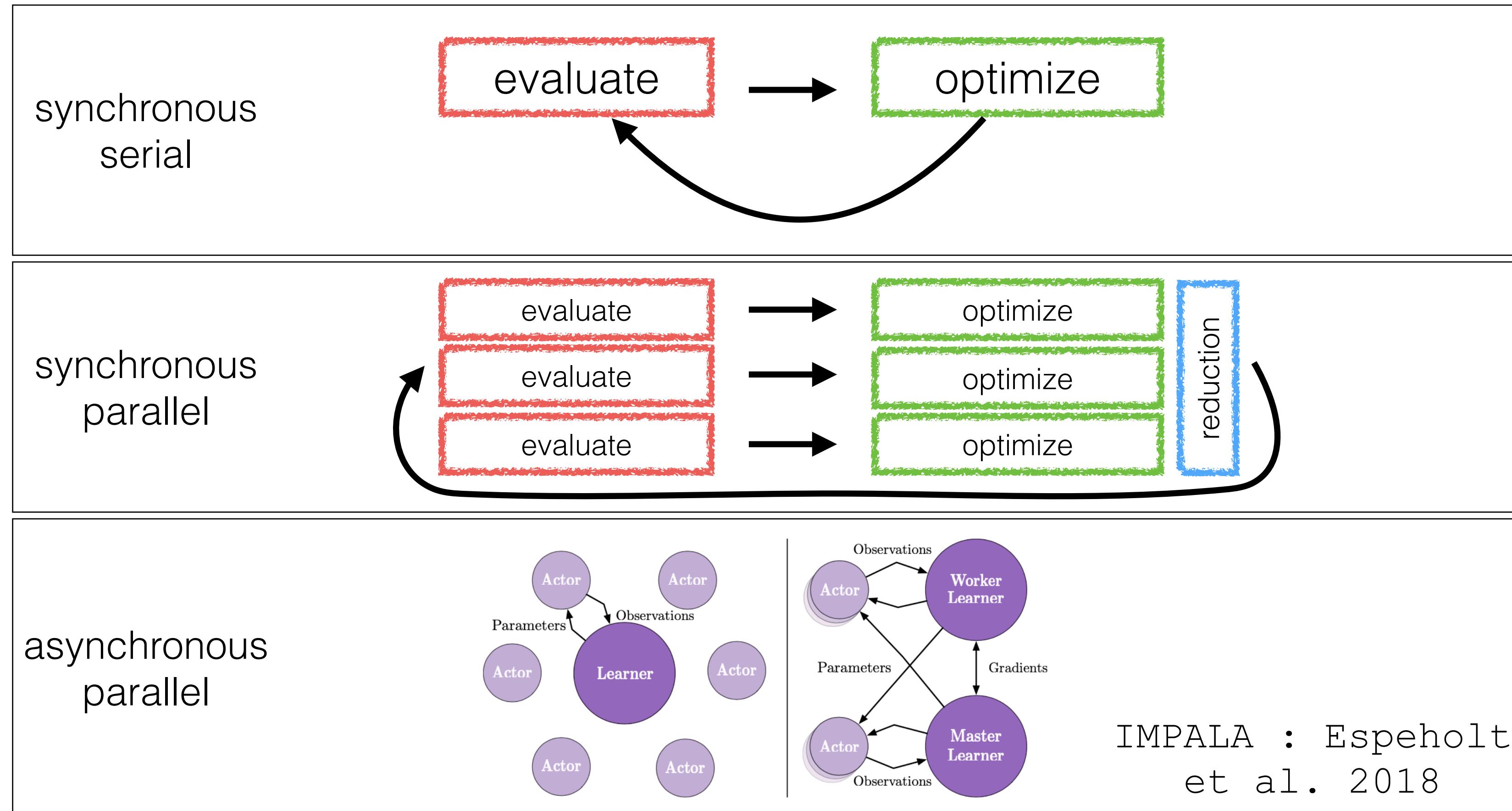
- E.g. Mnih et al. 2015 (Atari Nature paper) followed this approach
- Also most towardsdatascience.com tutorials
- Synchronous software design
- NN on GPU, alternates evaluation and optimization
- Toy environments are typically single threaded and CPU bound
- Environment often ends up being the bottleneck (Amdahl's law)

RL software architecture



- Second round (2016/17) of deep RL papers employed CPU-bound networks
- One “worker” per CPU core, each core runs its own environment
- Number of CPU cores became an hyper parameter in RL papers
- Performance can be bad:
 - Large load imbalances (unevenly sized episodes)
 - Many serial operations are performed

RL software architecture

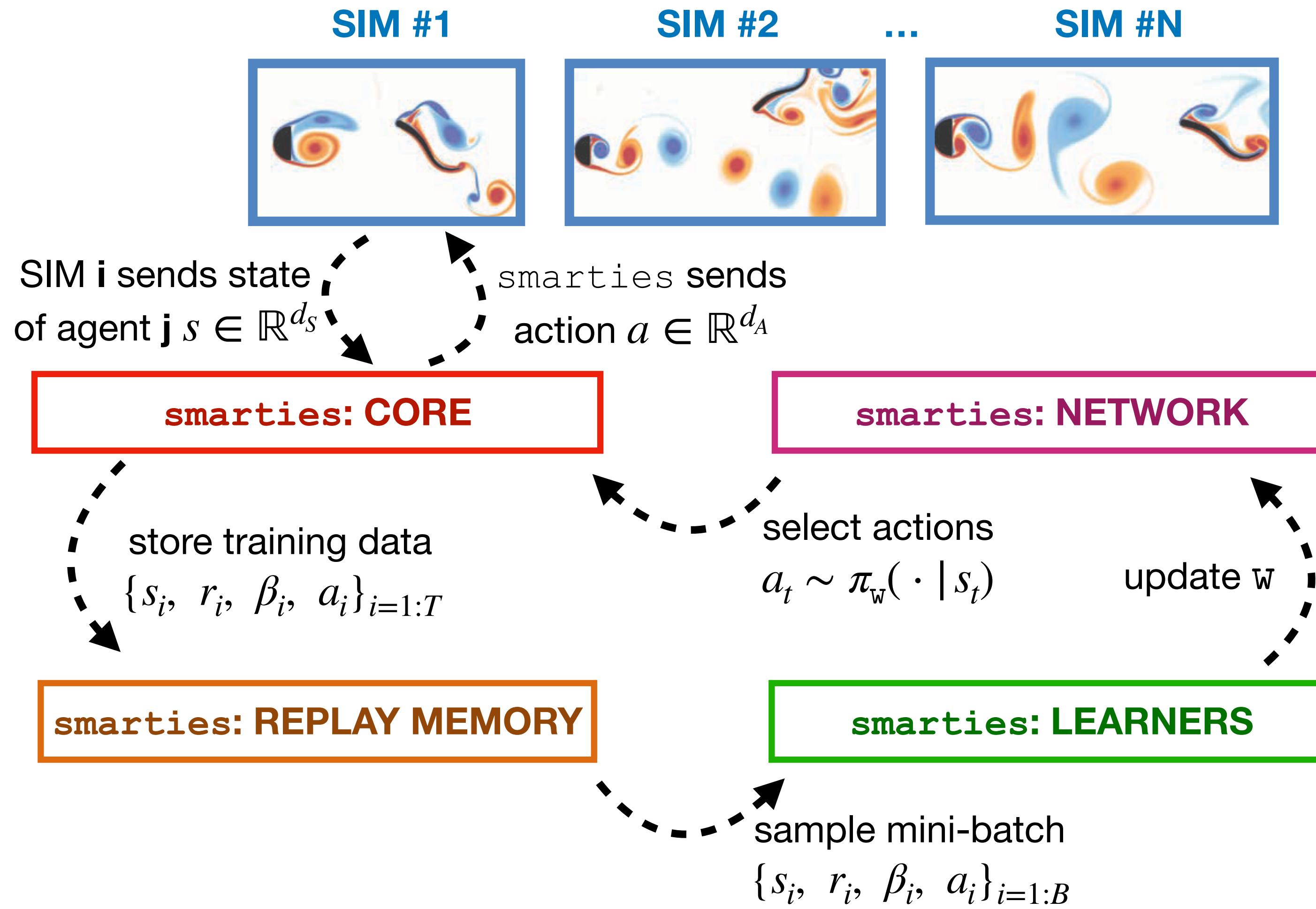


- Third round (2018/19) GPUs are back, price is much more complex software
- ‘Worker’ processes and ‘Learner’ processes, parameter servers...
- Computational resources are hyper parameters (N of workers affects performance)

smarties

- Objectives:
 - Minimally intrusive integration with existing simulation software: python, C++, fortran.
 - OpenAI gym is designed to be used for RL
 - But other environments may not be designed around RL
 - C++, asynchronous and thread-parallel:
 - Latency due to many cheap function calls is reduced
 - Algorithms can simply “*update NN by sampling RM and perform N actions*” without ruining performance: done in parallel
 - Data processing is easily parallelized.
 - Supports multiple learners, multiple workers, simulations that run on multiple CPUs...
 - Supports multiple agents with separate policies, different state/action sizes...
 - Uses RNN for partially observable problems without changing the algorithm
 - Environment communication API is thread-safe
 - One real dependency: MPI (portability?).

smarties : high level view

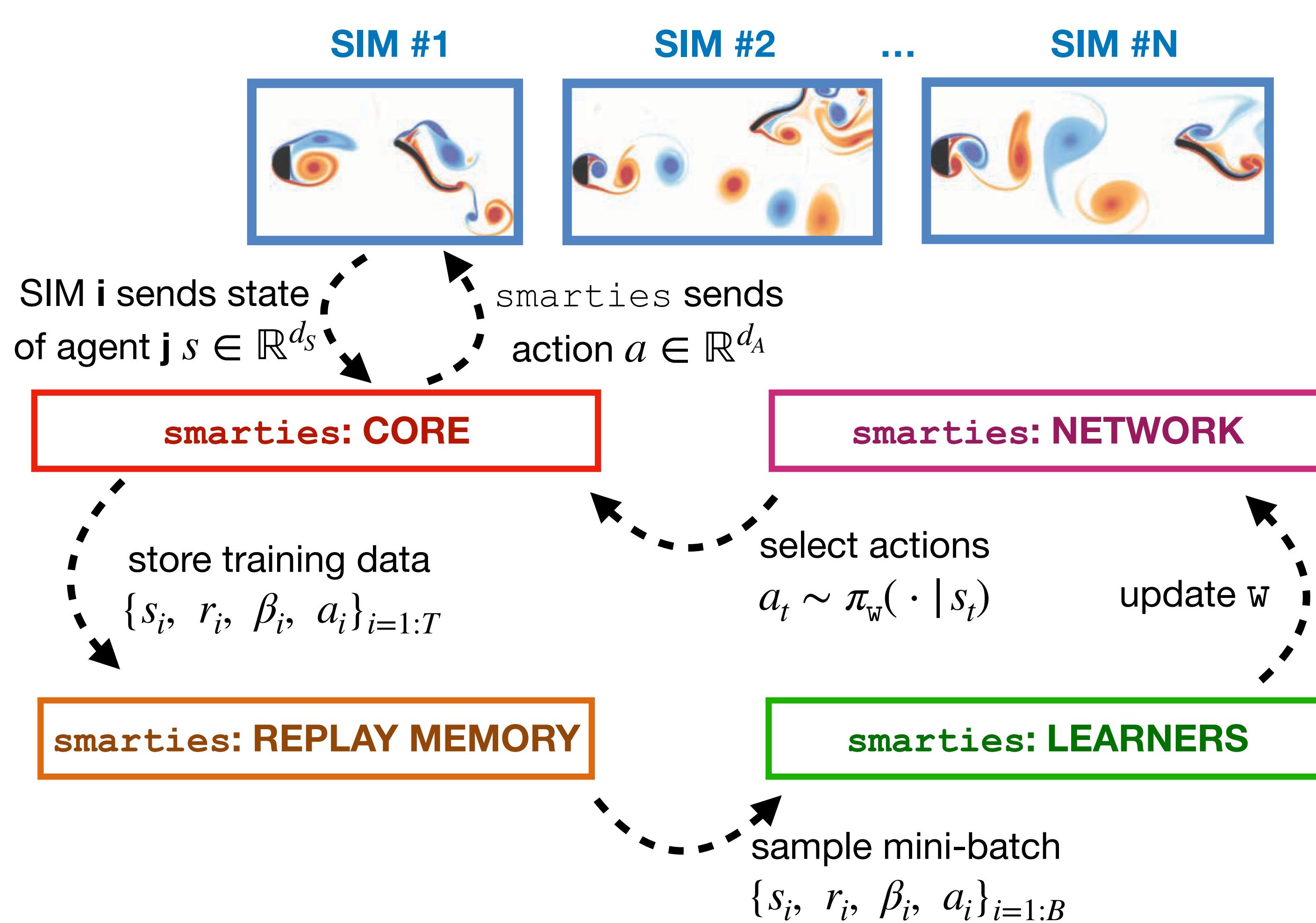


```
import smarties as rl

def app_main(comm):
    comm.set_state_action_dims(state_dim, action_dim)
    env = Environment()
    while 1: #train loop
        env.reset()
        comm.sendInitState(env.getState())
        while 1: #simulation loop
            action = comm.recvAction()
            isTerminal = env.advance(action)
            if isTerminal:
                comm.sendTermState(env.getState(), env.getReward())
                break
            else: # normal state
                comm.sendState(env.getState(), env.getReward())

    if __name__ == '__main__':
        e = rl.Engine(sys.argv)
        if( e.parse() ): exit()
        e.run( app_main )
```

smarties : the sub-components



- **CORE:**

- Common definitions: MDP, agent, environment
- Schedules communication between workers and learners
- Learners that send actions and recv states
- Learners that send parameters and recv episodes

- **REPLAY MEMORY:**

- Store, share, sample, post-process experiences

- **LEARNERS:**

- RACER, DQN, NAF, PPO, CMA, DPG, ACER...
- Each algorithm prescribes how to select actions and how to update network(s)

- **MATH:**

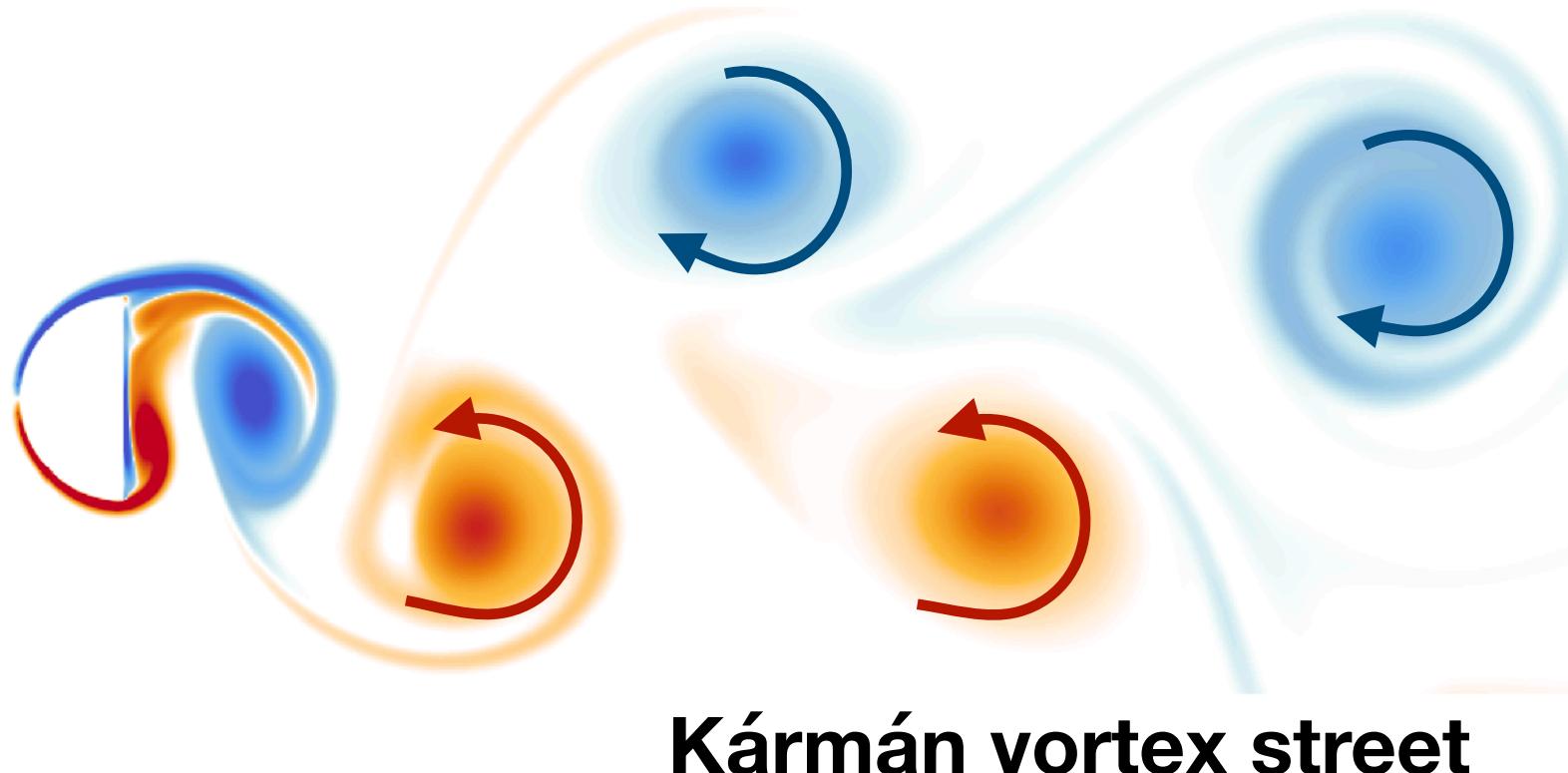
- e.g. Gaussian/Bernoulli/Dirichlet policies and related operations

- **NETWORK:**

- Defines various layer types (linear, LSTM, GRU, Conv2D) and non-linearities
- Optimizers, memory management, ...

ReF-ER and agents in unsteady flows

- Unsteady flows are ubiquitous in nature
- Many animal species learn to capture energy from vortices



- E.g. drafting in flows behind stationary objects

Webb, J. Exp. Bio. 1998

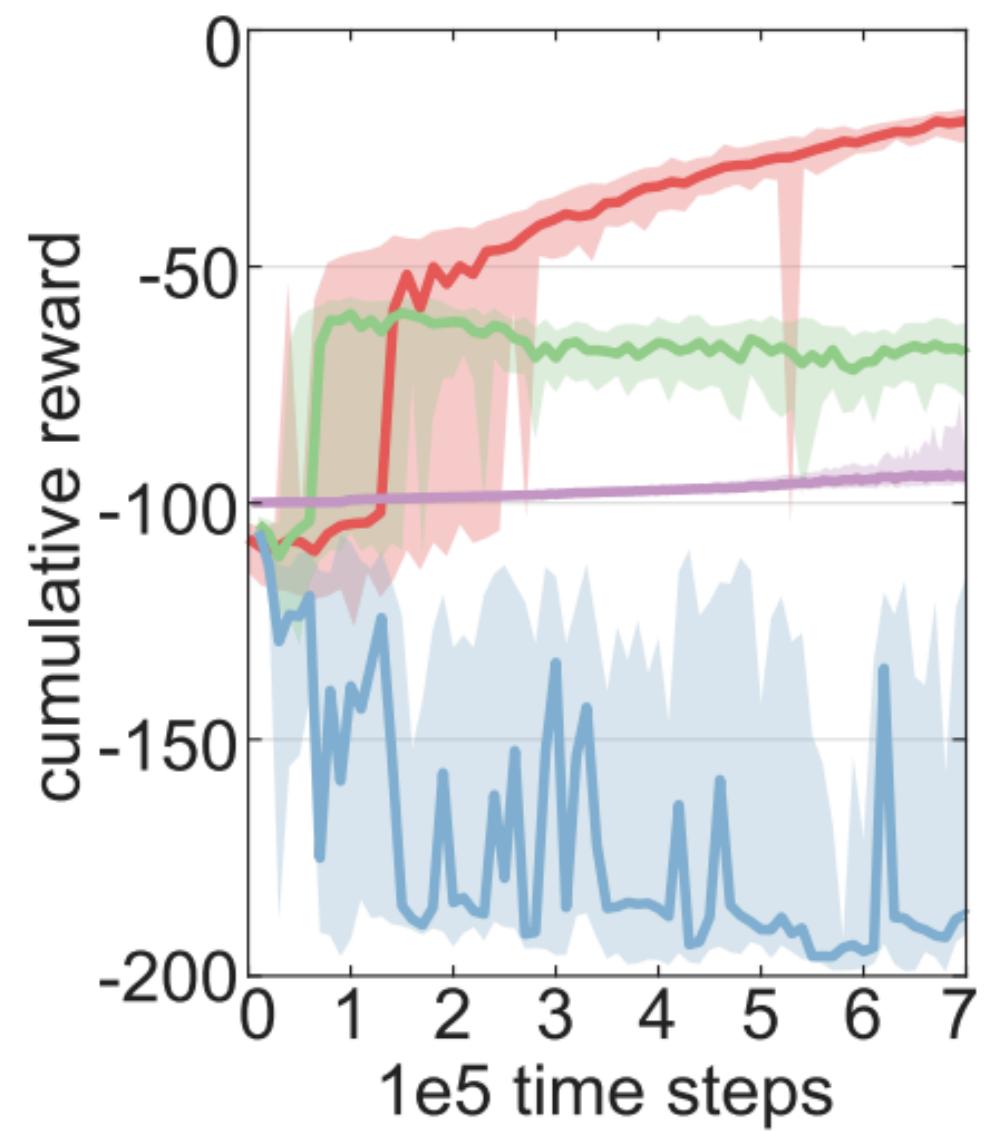
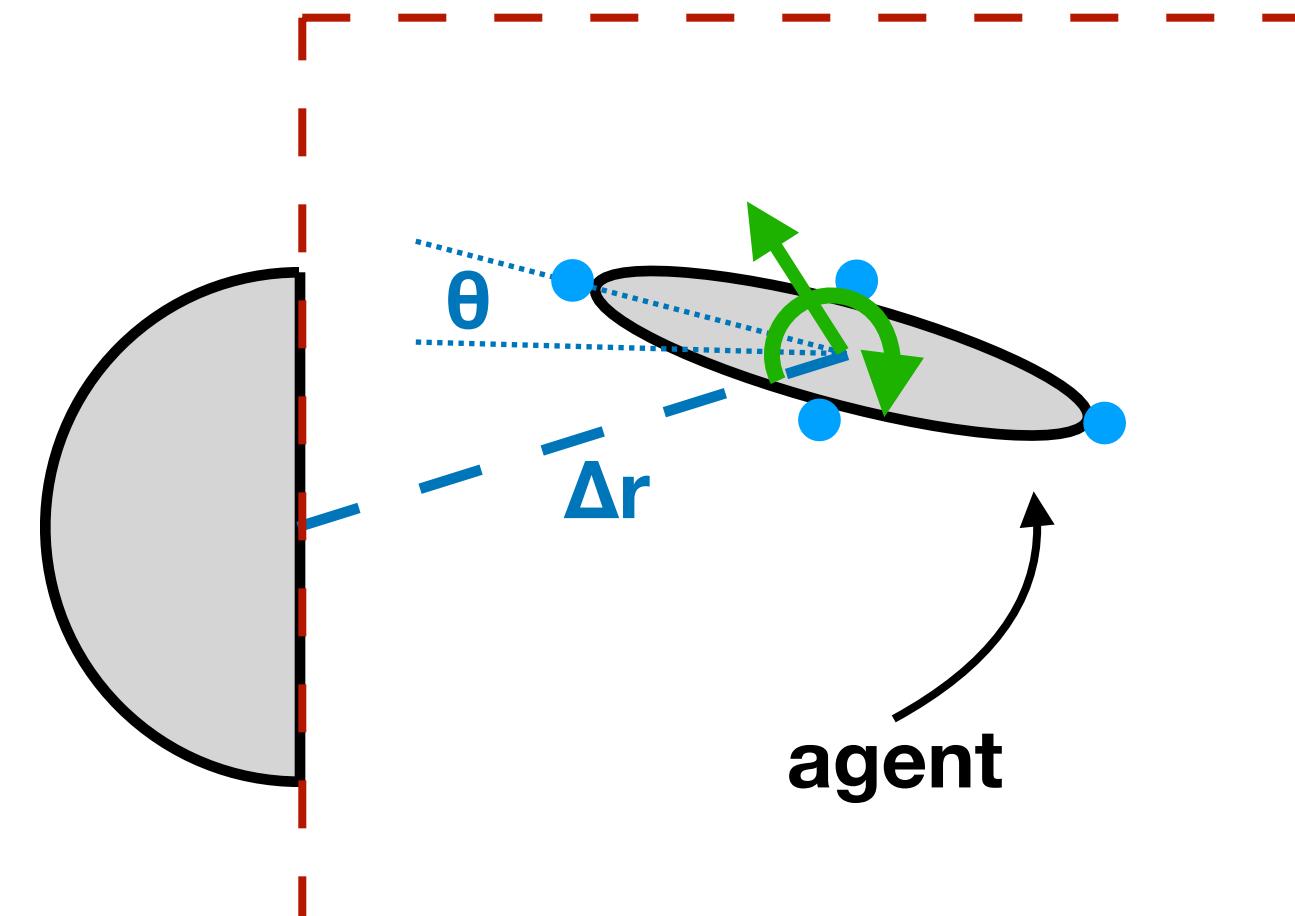
Hinch, Rand, J. Fish. Aquat. Sci. 2000

Liao, Beal, Lauder, Triantafyllou, J. Exp. Bio. 2003

- State:**
- Relative position Δr , angle θ , velocities
 - Flow velocity sensors (at vertices): ●

- Action:**
- Applied force and torque : $\vec{a} := \{f_X, f_Y, \tau\}$

- Reward:**
- Actuation cost: $r = -\|a\|^2$
 - Terminate if reaches border: $r_T = -100$



- Uses RNNs to estimate environment from sensor data
- Legend: **RACER + ReF-ER**, **ACER**, **DPG**, **DPG + ReF-ER**
- RACER at $Re = UD/v = 400$, force to propel ellipse 60% lower than in quiescent flow