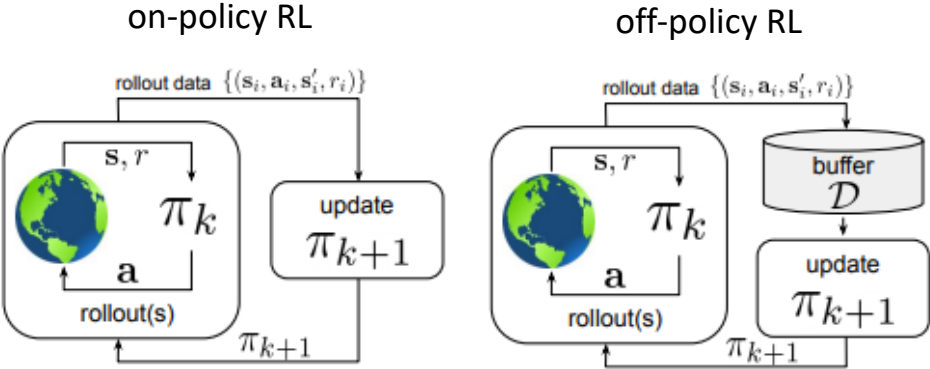


# Batch / Offline Reinforcement Learning

With thanks to Emma Brunskill, Scott Fujimoto, Pieter Abbeel, George Tucker, Sergey Levine, Bilal Piot,  
Yuxin Chen, Yuejie Chi

# On-policy vs off-policy vs offline RL



Formally:

$$\mathcal{D} = \{(s_i, a_i, s'_i, r_i)\}$$

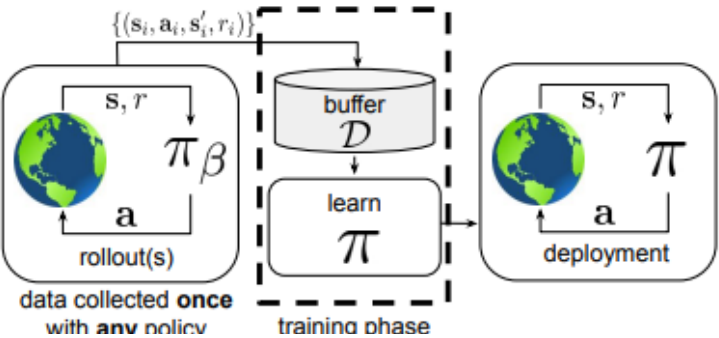
$$s \sim d^{\pi_\beta}(s) \quad \leftarrow \text{generally not known}$$

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

$$s' \sim p(s'|s, a)$$

$$r \leftarrow r(s, a)$$

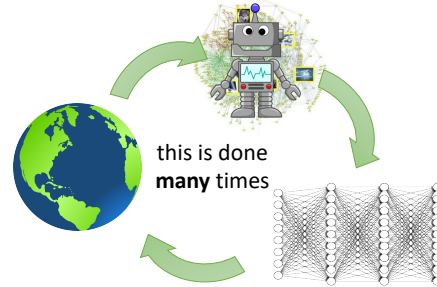
## offline reinforcement learning



$$\text{RL objective: } \max_{\pi} \sum_{t=0}^T E_{s_t \sim d^\pi(s), a_t \sim \pi(a|s)} [\gamma^t r(s_t, a_t)]$$

# Why is this important?

- Collecting new data may be expensive / infeasible
- We may have access to existing/historical data instead



# Problem formulation

**A historical dataset**  $\mathcal{D} = \{(s^{(i)}, a^{(i)}, s'^{(i)})\}$ :  $N$  independent copies of

$$s \sim \rho^b, \quad a \sim \pi^b(\cdot | s), \quad s' \sim P(\cdot | s, a)$$

for some state distribution  $\rho^b$  and behavior policy  $\pi^b$

**Goal:** given some test distribution  $\rho$  and accuracy level  $\varepsilon$ , find an  $\varepsilon$ -optimal policy  $\hat{\pi}$  based on  $\mathcal{D}$  obeying

$$V^*(\rho) - V^{\hat{\pi}}(\rho) = \mathbb{E}_{s \sim \rho} [V^*(s)] - \mathbb{E}_{s \sim \rho} [V^{\hat{\pi}}(s)] \leq \varepsilon$$

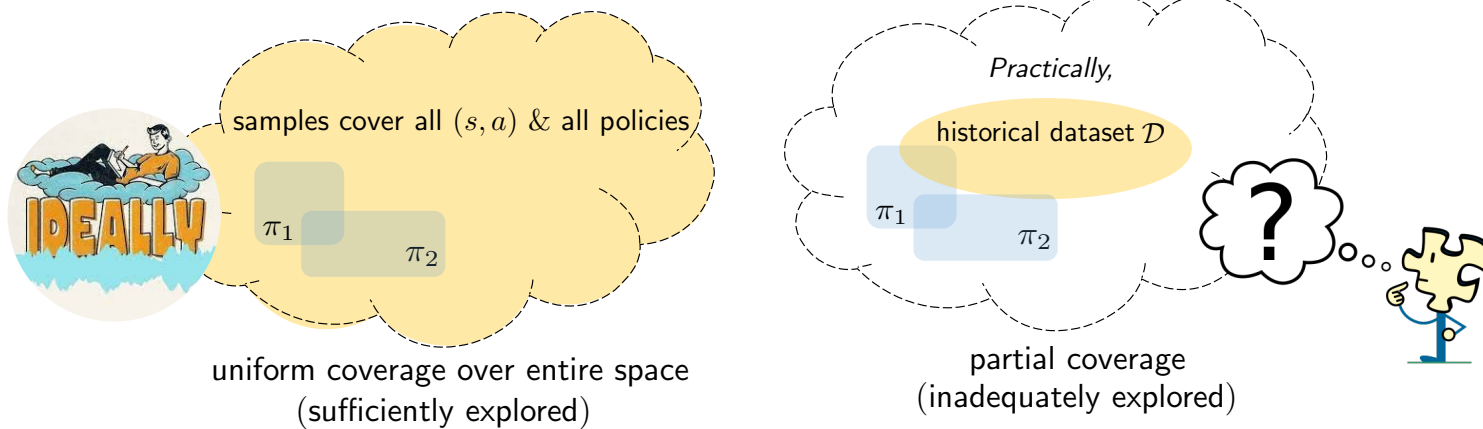
— *in a sample-efficient manner*

# Challenges of offline / batch RL (1)

- **Distribution shift:**

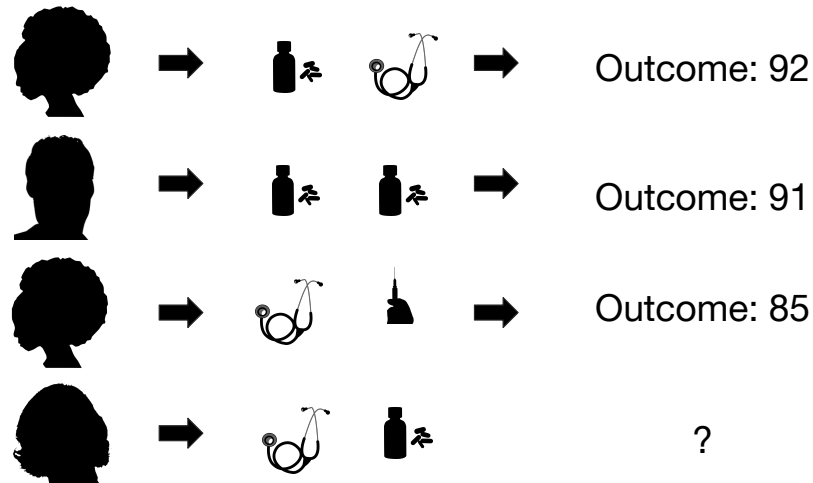
distribution( $\mathcal{D}$ )  $\neq$  target distribution under  $\pi^*$

- **Partial coverage of state-action space:**



## Challenges of offline / batch RL (2)

- Data is *censored*: we only observe outcomes for decisions made (and need to generalize from them)



- Need for *counterfactual inference*: what would happen if one would take a different action?
- Often we do not observe rewards, just states and actions!

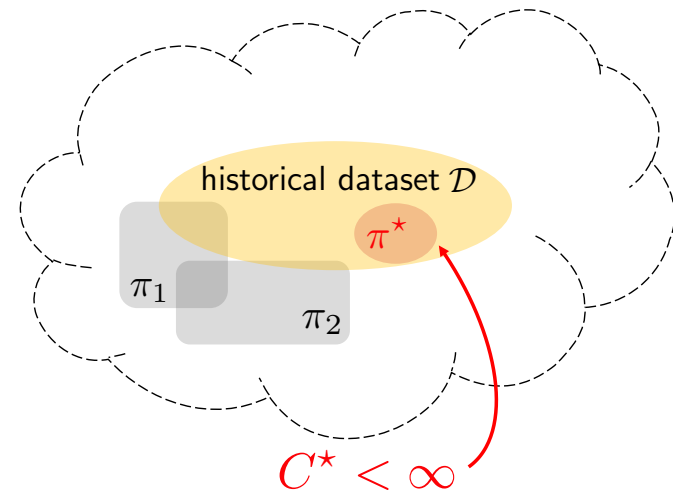
# Dataset quality assessment

## Single-policy concentrability coefficient

$$C^* := \max_{s,a} \frac{d^{\pi^*}(s,a)}{d^{\pi^b}(s,a)} = \left\| \frac{\text{occupancy density of } \pi^*}{\text{occupancy density of } \pi^b} \right\|_{\infty} \geq 1$$

where  $d^{\pi}(s,a) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \mathbb{P}((s^t, a^t) = (s, a) | \pi)$

- captures distributional shift
- allows for partial coverage



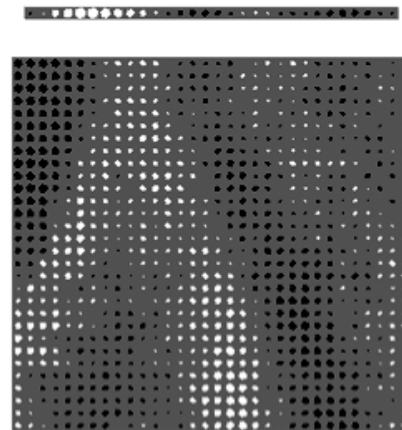
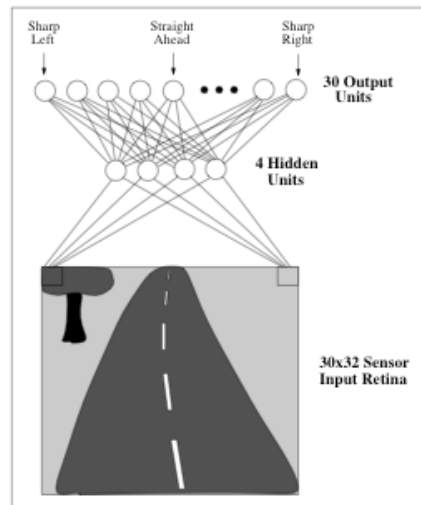
# Classes of algorithms

- Behavior cloning (no rewards required)
- Learn a model, use it for model-based RL (LSTD, LSPI)
- Pessimistic algorithms (require rewards)
- Inverse RL (learn reward function from data, use it for RL agent)

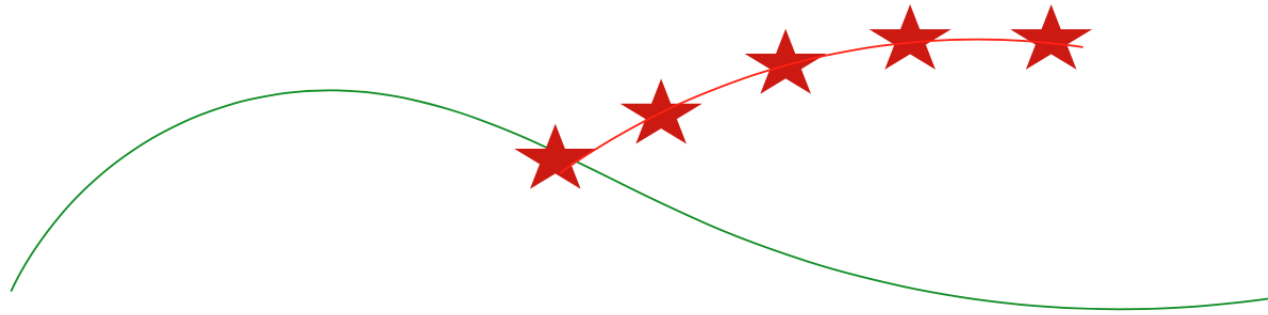


# Behavior cloning

- Take dataset  $\mathcal{D}$ , learn a policy from states to actions
- Often uses a rich policy class (neural net)



## Problem: compounding errors



- Error at time  $t$  with probability  $\epsilon$
- Approximate intuition:  $\mathbb{E}[\text{Total errors}] \leq \epsilon(T + (T - 1) + (T - 2) \dots + 1) \propto \epsilon T^2$

# One solution: dataset aggregation

Initialize  $\mathcal{D} \leftarrow \emptyset$ .  
Initialize  $\hat{\pi}_1$  to any policy in  $\Pi$ .  
**for**  $i = 1$  **to**  $N$  **do**  
    Let  $\pi_i = \beta_i \pi^* + (1 - \beta_i) \hat{\pi}_i$ .  
    Sample  $T$ -step trajectories using  $\pi_i$ .  
    Get dataset  $\mathcal{D}_i = \{(s, \pi^*(s))\}$  of visited states by  $\pi_i$   
    and actions given by expert.  
    Aggregate datasets:  $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_i$ .  
    Train classifier  $\hat{\pi}_{i+1}$  on  $\mathcal{D}$ .  
**end for**  
**Return** best  $\hat{\pi}_i$  on validation.

- Idea: Get more labels of the expert action along the path taken by the policy computed by behavior cloning
- Obtains a stationary deterministic policy with good performance under its induced state distribution

## Pessimism in the face of uncertainty

- *Conservative* approach
- Assume that states or state-action pairs not visited are bad
- Use a penalty to avoid the new policy visiting them

# Value iteration with lower confidence bounds

**Pessimism in the face of uncertainty:** penalize value estimate of those  $(s, a)$  pairs that were poorly visited [Jin et al., 2021, Rashidinejad et al., 2021]

**Algorithm:** value iteration w/ lower confidence bounds

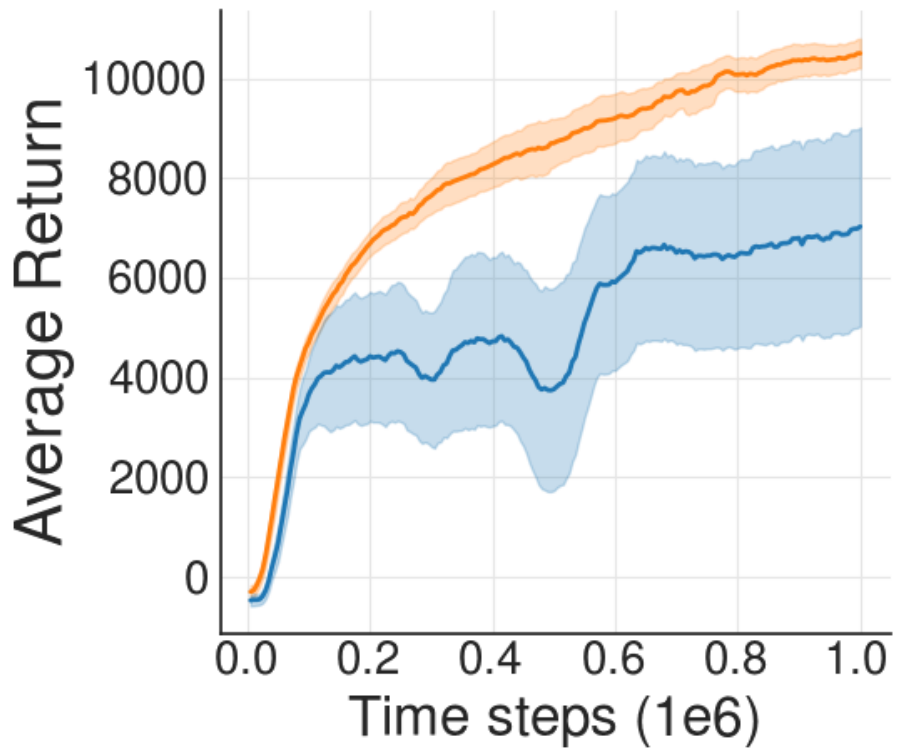
- compute empirical estimate  $\hat{P}$  of  $P$
- initialize  $\hat{Q} = 0$ , and repeat

$$\hat{Q}(s, a) \leftarrow \max \left\{ r(s, a) + \gamma \langle \hat{P}(\cdot | s, a), \hat{V} \rangle - \underbrace{b(s, a; \hat{V})}_{\text{Bernstein-style confidence bound}}, 0 \right\}$$

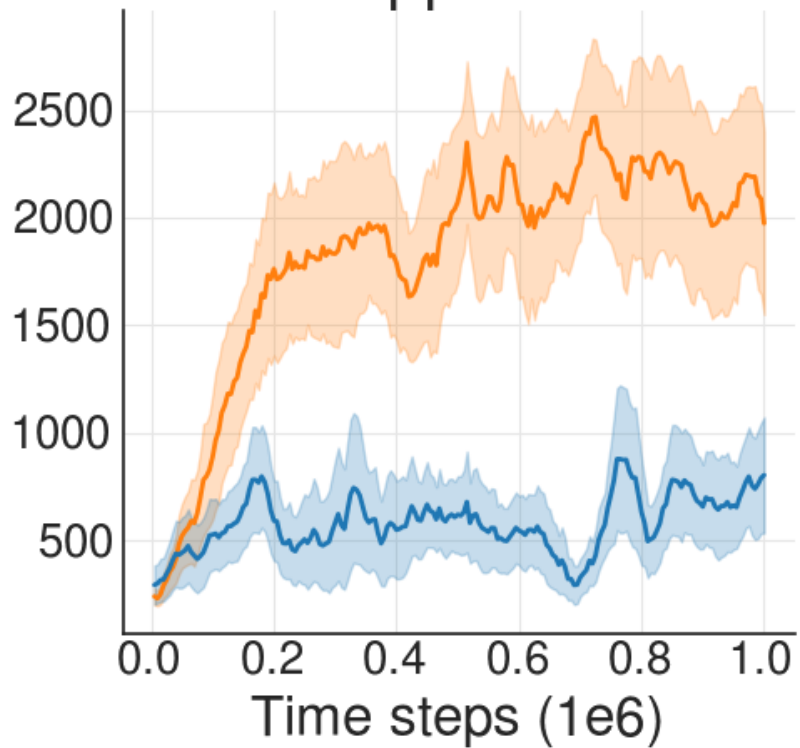
for all  $(s, a)$ , where  $\hat{V}(s) = \max_a \hat{Q}(s, a)$

Q-learning version exists as well

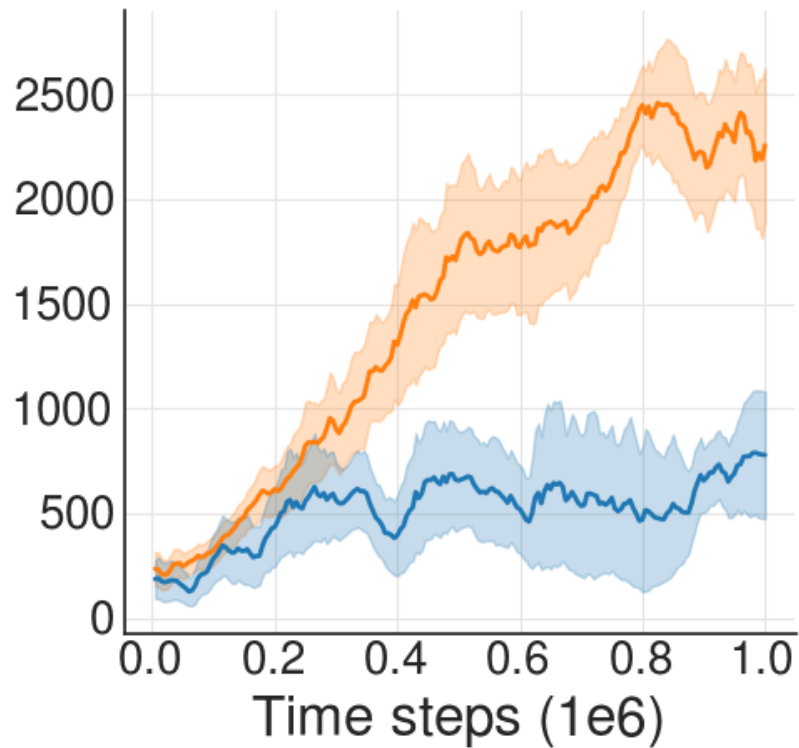
### HalfCheetah-v1



### Hopper-v1



### Walker2d-v1



Surprise!

Agent orange and agent blue are trained with...

1. The **same off-policy algorithm (DDPG)**.
2. The **same dataset**.

# The Difference?

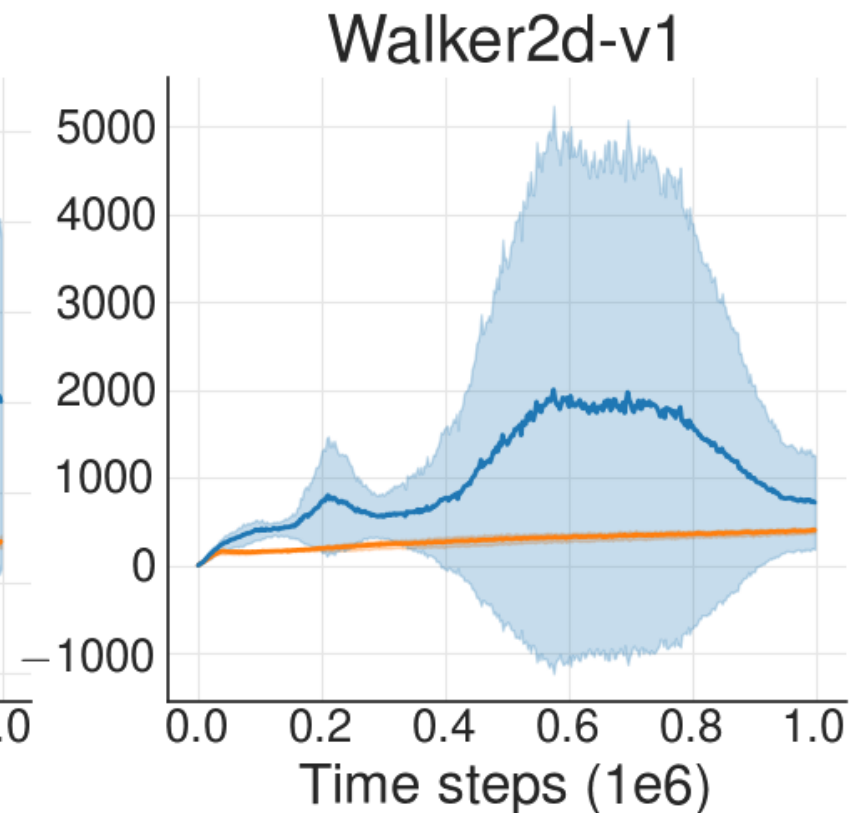
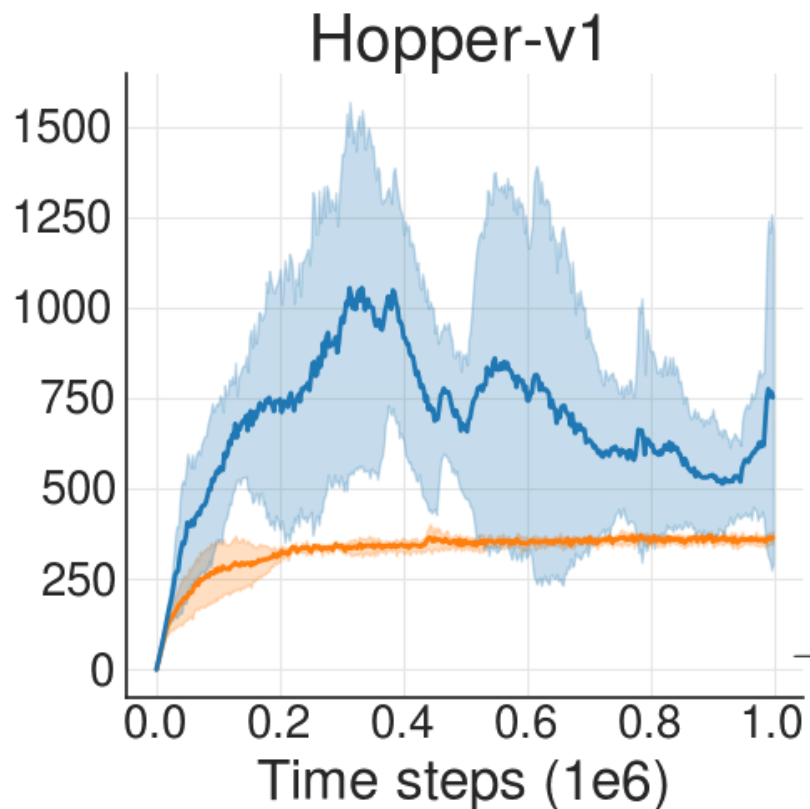
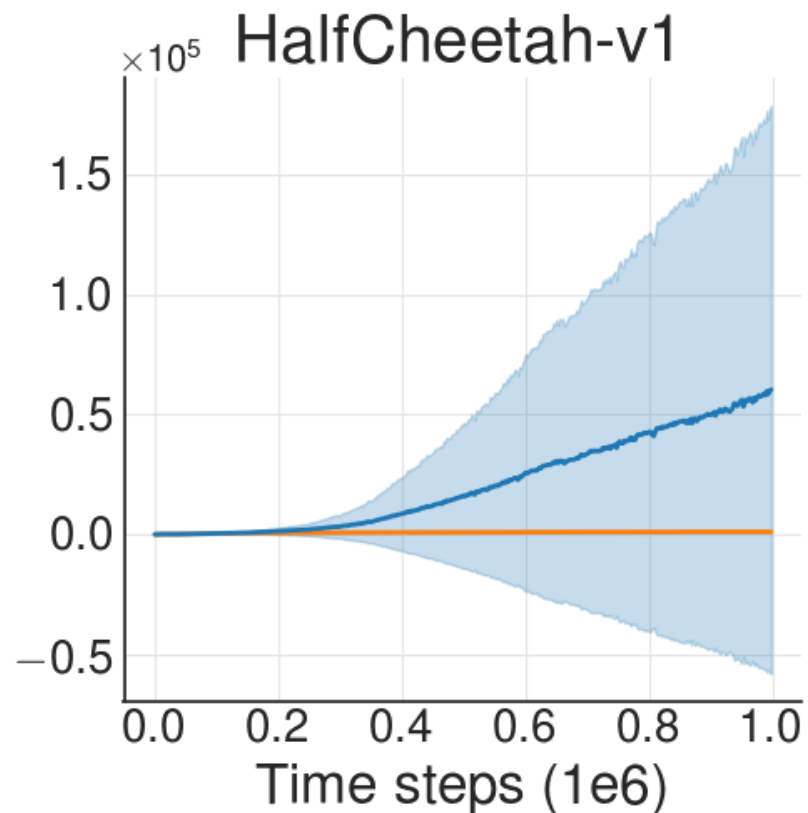
1. **Agent orange:** Interacted with the environment.
  - Standard RL loop.
  - Collect data, store data in buffer, train, repeat.
2. **Agent blue:** Never interacted with the environment.
  - Trained with data collected by agent orange concurrently.



1. Trained with the same off-policy algorithm.
2. Trained with the same dataset.
3. One interacts with the environment. One doesn't.

**Off-policy** deep RL fails when **truly off-policy**.

# Value Predictions



Extrapolation Error

$$Q(s, a) \leftarrow r + \gamma Q(s', a')$$

# Extrapolation Error

$$Q(s, a) \leftarrow r + \gamma Q(s', a')$$

The diagram illustrates the Bellman optimality equation  $Q(s, a) \leftarrow r + \gamma Q(s', a')$  with annotations. The word "GIVEN" is written in red below the state-action pair  $(s, a)$ . Two red arrows point upwards from "GIVEN" to  $s$  and  $a$ . Another red arrow points from the reward  $r$  to the left-hand side of the equation. A second red arrow points from the term  $Q(s', a')$  to the left-hand side. The word "GENERATED" is written in blue below the next state-action pair  $(s', a')$ . A blue arrow points upwards from "GENERATED" to  $a'$ .

# Extrapolation Error

$$Q(s, a) \leftarrow r + \gamma Q(s', a')$$

1.  $(s, a, r, s') \sim \text{Dataset}$
2.  $a' \sim \pi(s')$

# Extrapolation Error

$$Q(s, a) \leftarrow r + \gamma Q(s', a')$$

$(s', a') \notin \text{Dataset} \rightarrow Q(s', a') = \mathbf{bad}$   
 $\rightarrow Q(s, a) = \mathbf{bad}$

# Extrapolation Error

$$Q(s, a) \leftarrow r + \gamma Q(s', a')$$

$(s', a') \notin \text{Dataset} \rightarrow Q(s', a') = \mathbf{bad}$

$\rightarrow Q(s, a) = \mathbf{bad}$



# Extrapolation Error

$$Q(s, a) \leftarrow r + \gamma Q(s', a')$$

$(s', a') \notin \text{Dataset} \rightarrow Q(s', a') = \mathbf{bad}$

$\rightarrow Q(s, a) = \mathbf{bad}$

# Extrapolation Error

Attempting to evaluate  $\pi$  without (sufficient) access to the  $(s, a)$  pairs  $\pi$  visits.

# Batch-Constrained Reinforcement Learning

Only choose  $\pi$  such that we have access to the  $(s, a)$  pairs  $\pi$  visits.

# Batch-Constrained Reinforcement Learning

1.  $a \sim \pi(s)$  such that  $(s, a) \in Dataset$ .
2.  $a \sim \pi(s)$  such that  $(s', \pi(s')) \in Dataset$ .
3.  $a \sim \pi(s)$  such that  $Q(s, a)$  is maxed.

# Batch-Constrained Deep Q-Learning (BCQ)

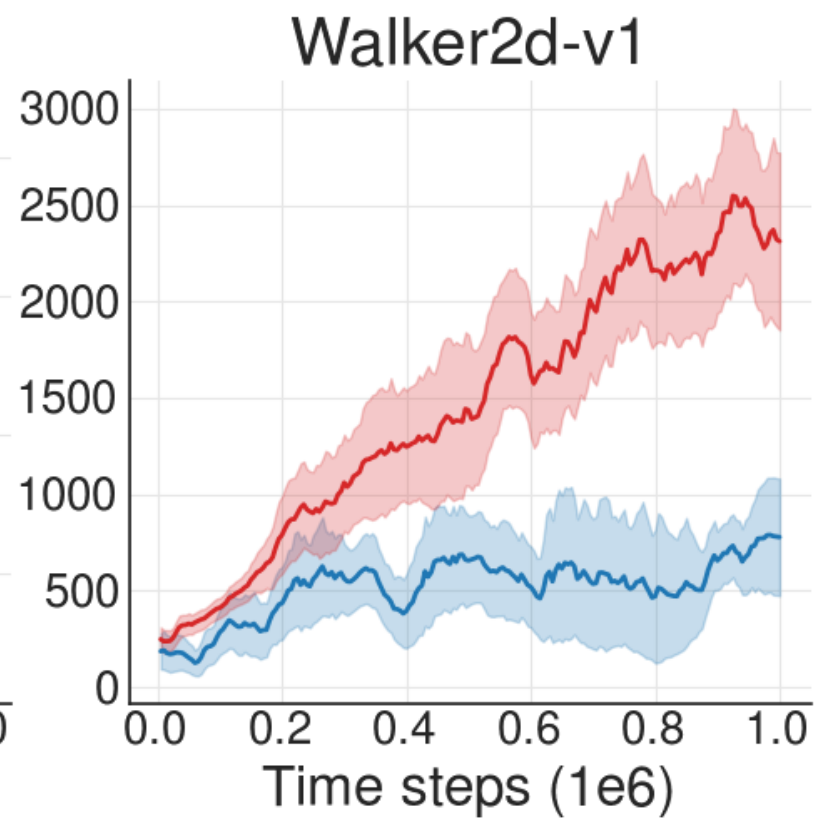
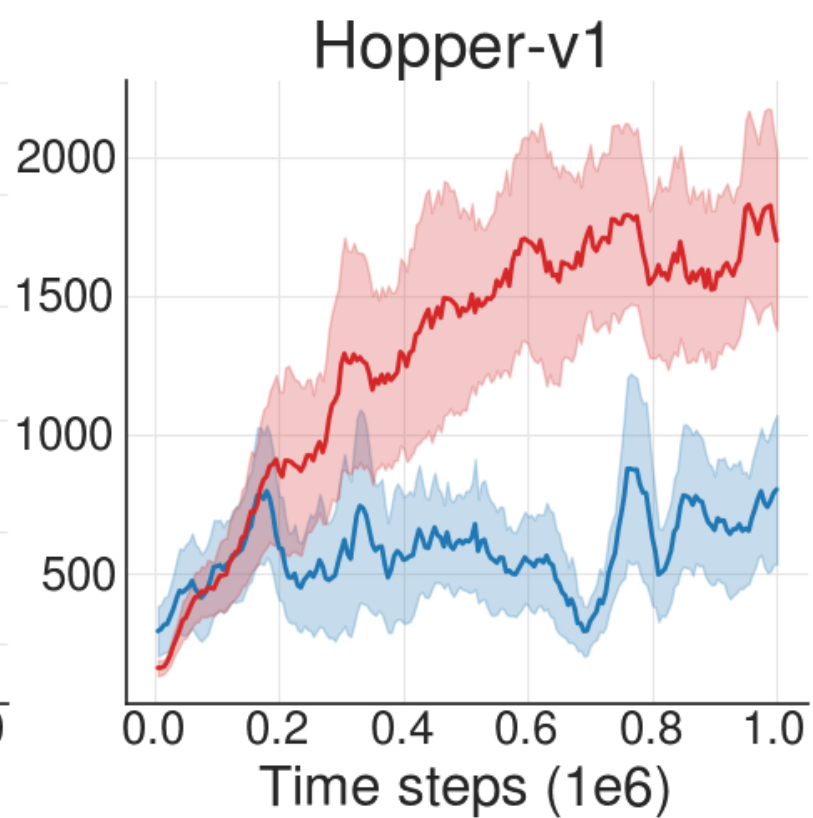
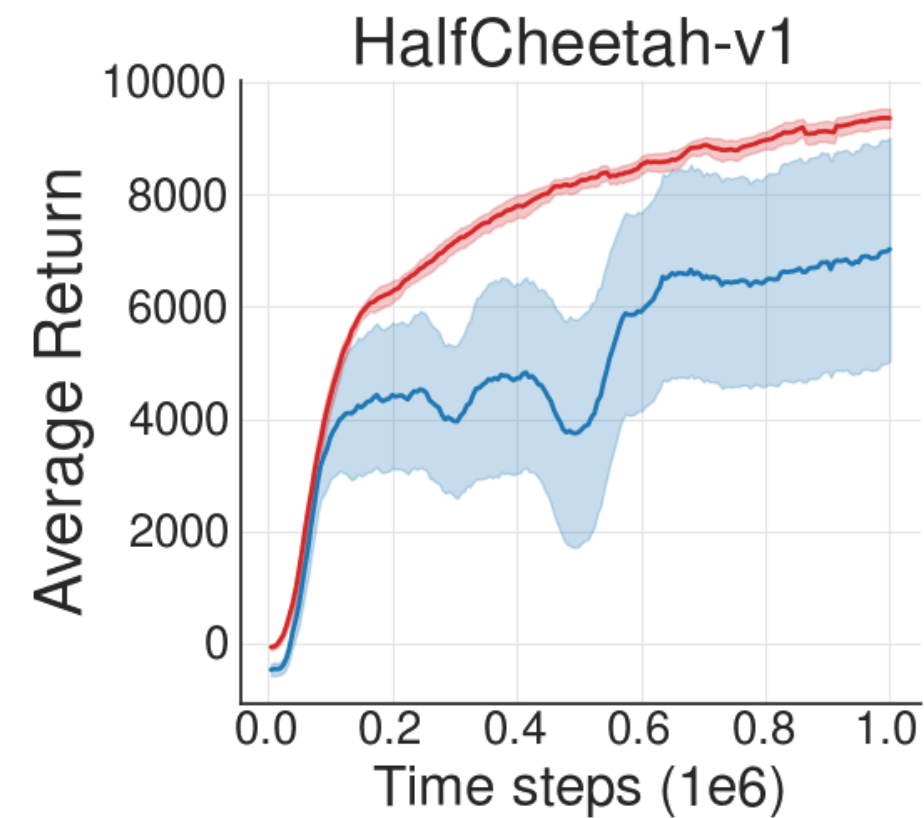
First imitate dataset via generative model:

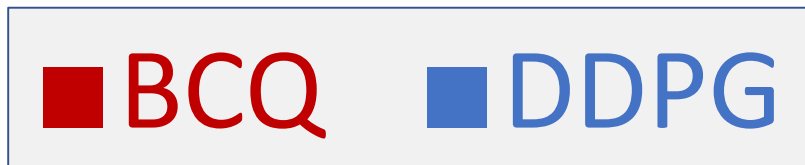
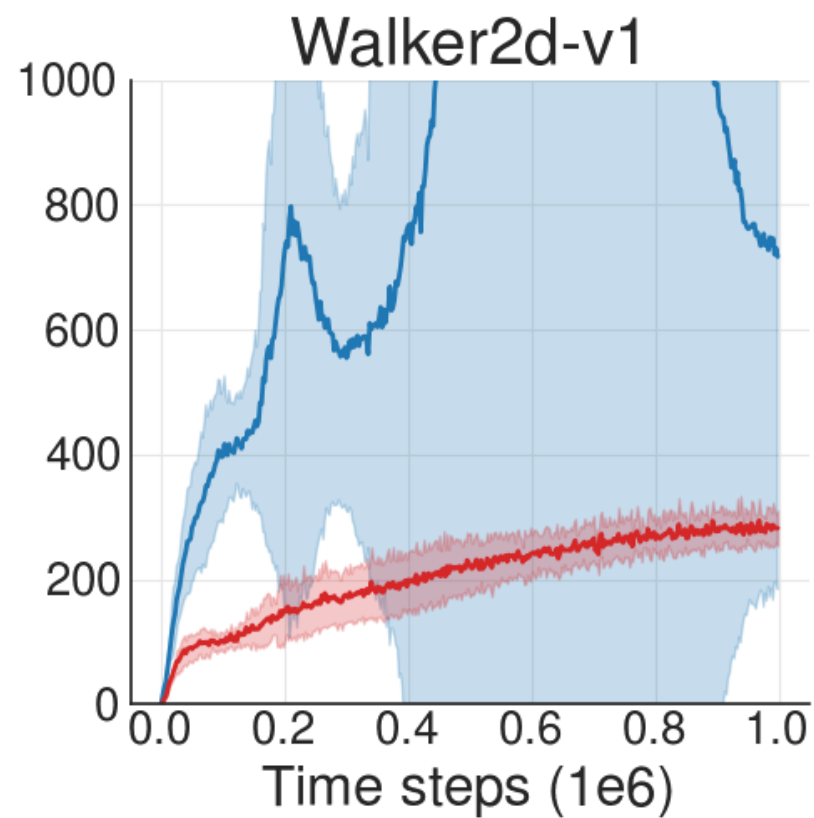
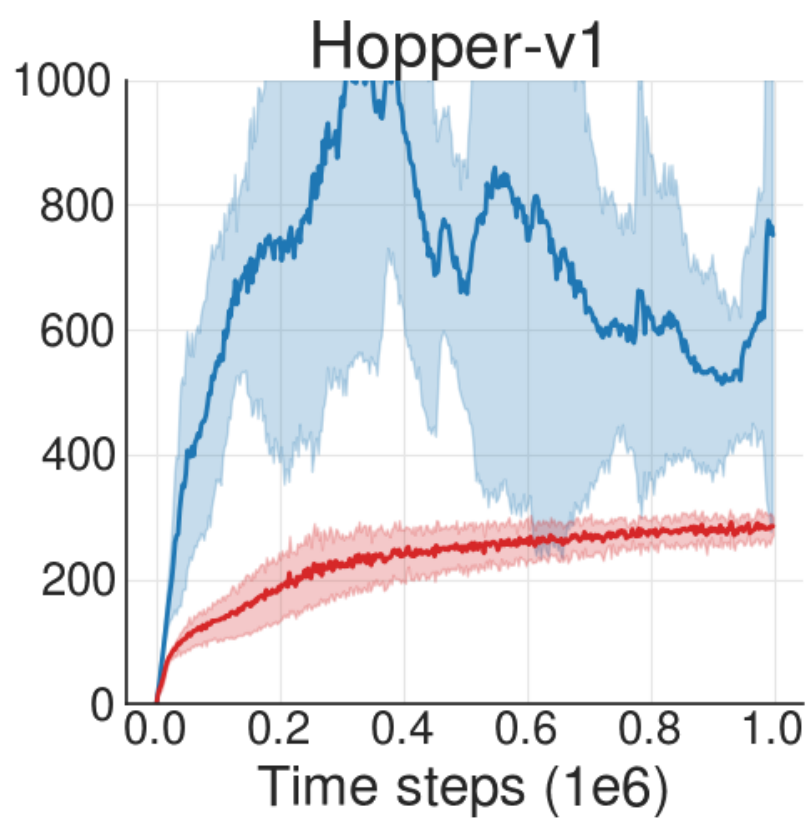
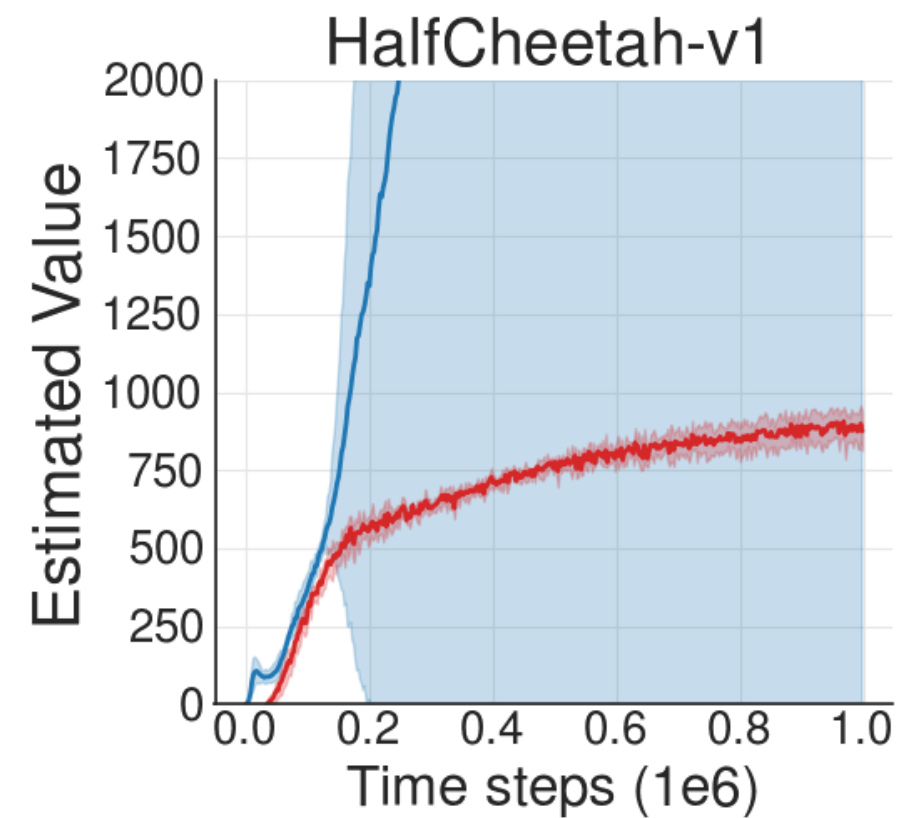
$$G(a|s) \approx P_{Dataset}(a|s).$$

$$\pi(s) = \operatorname{argmax}_{a_i} Q(s, a_i), \text{ where } a_i \sim G$$

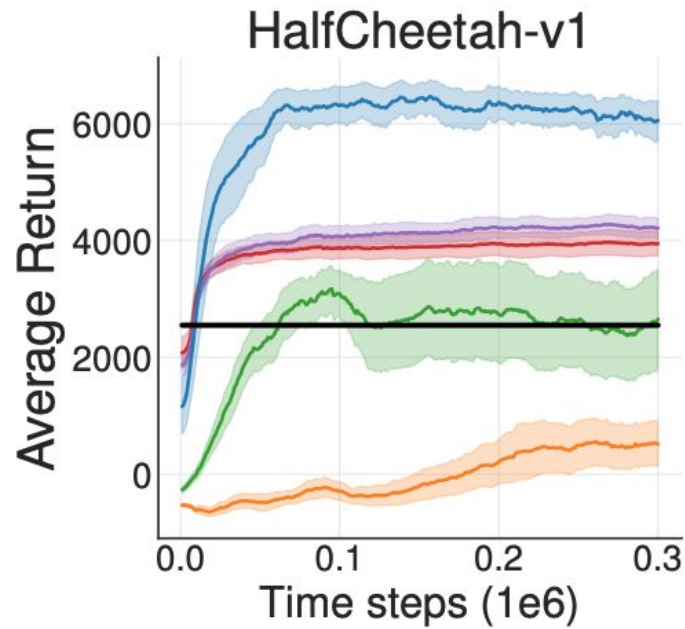
(i.e. select the best action that is likely under the dataset)

(+ some additional deep RL **magic**)





# BCQ comparison



BCQ figure from Fujimoto, Meger, Precup ICML 2019

BCQ DDPG DQN BC VAE-BC Behavioral<sub>17</sub>