



Computer Engineering Department

# Offline Reinforcement Learning

**Mohammad Hossein Rohban, Ph.D.**

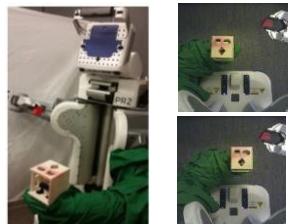
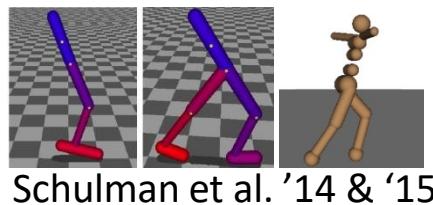
Spring 2024

Slides are adopted from CS 285, UC Berkeley.

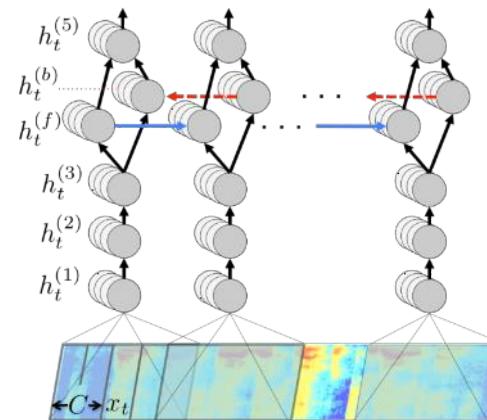
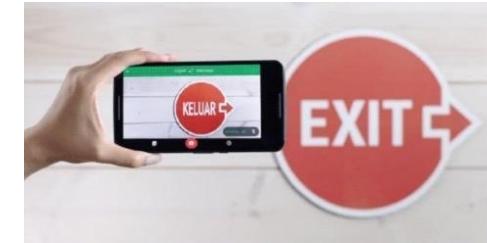
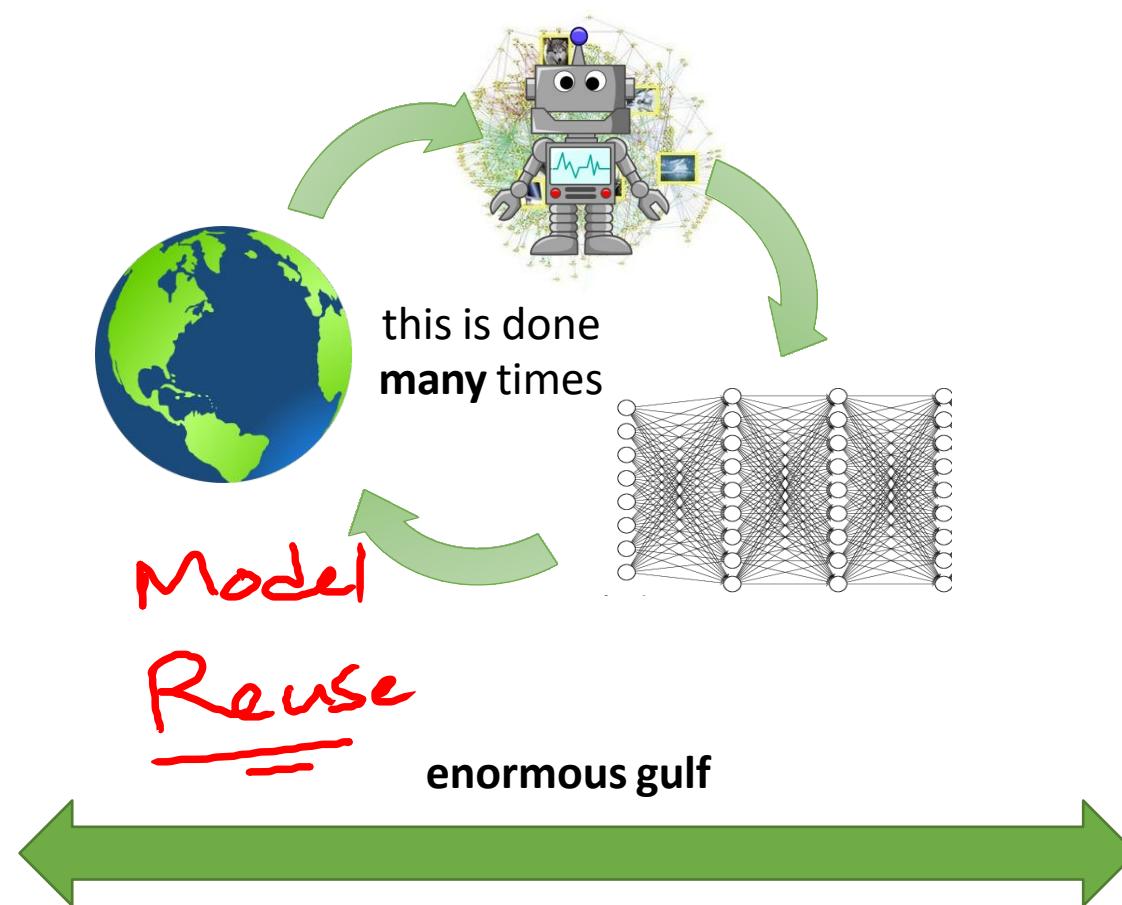
# The generalization gap



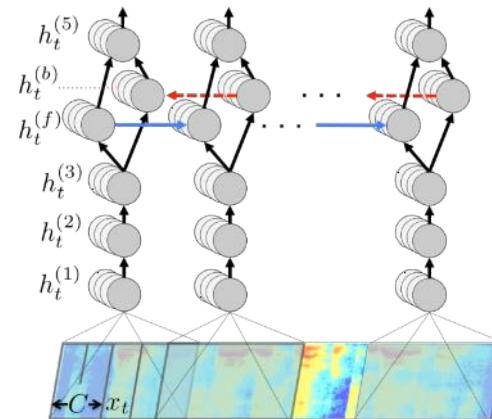
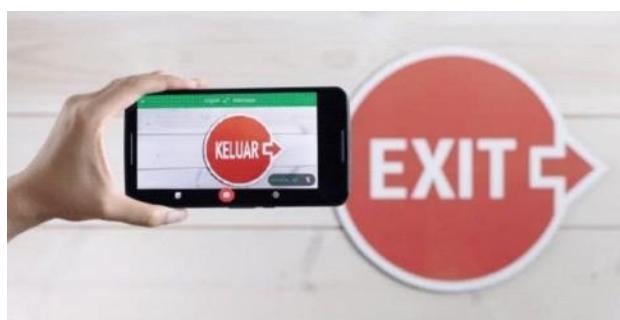
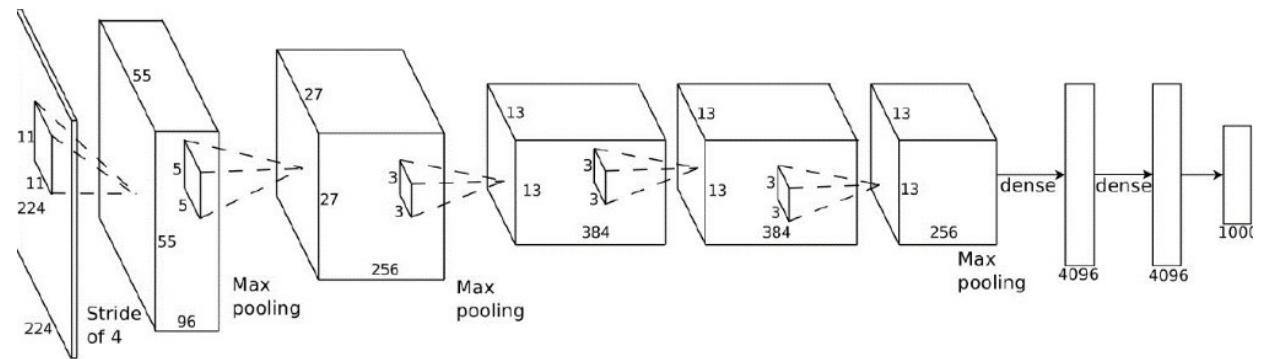
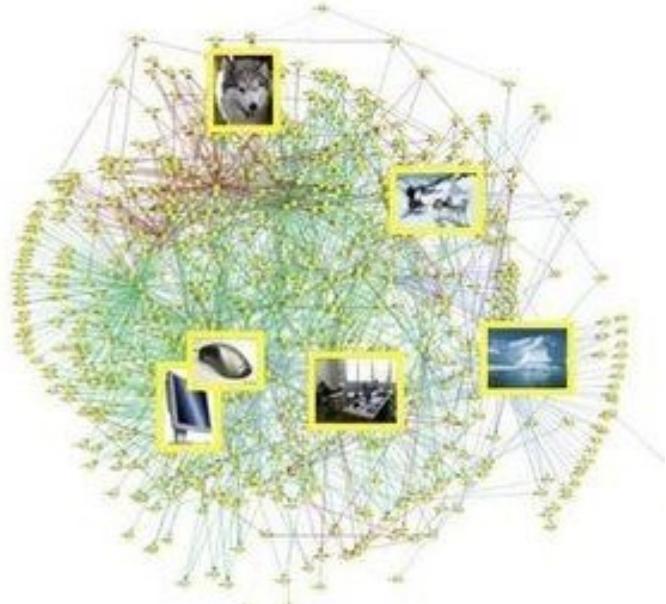
Mnih et al. '13



Levine\*, Finn\*, et al. '16



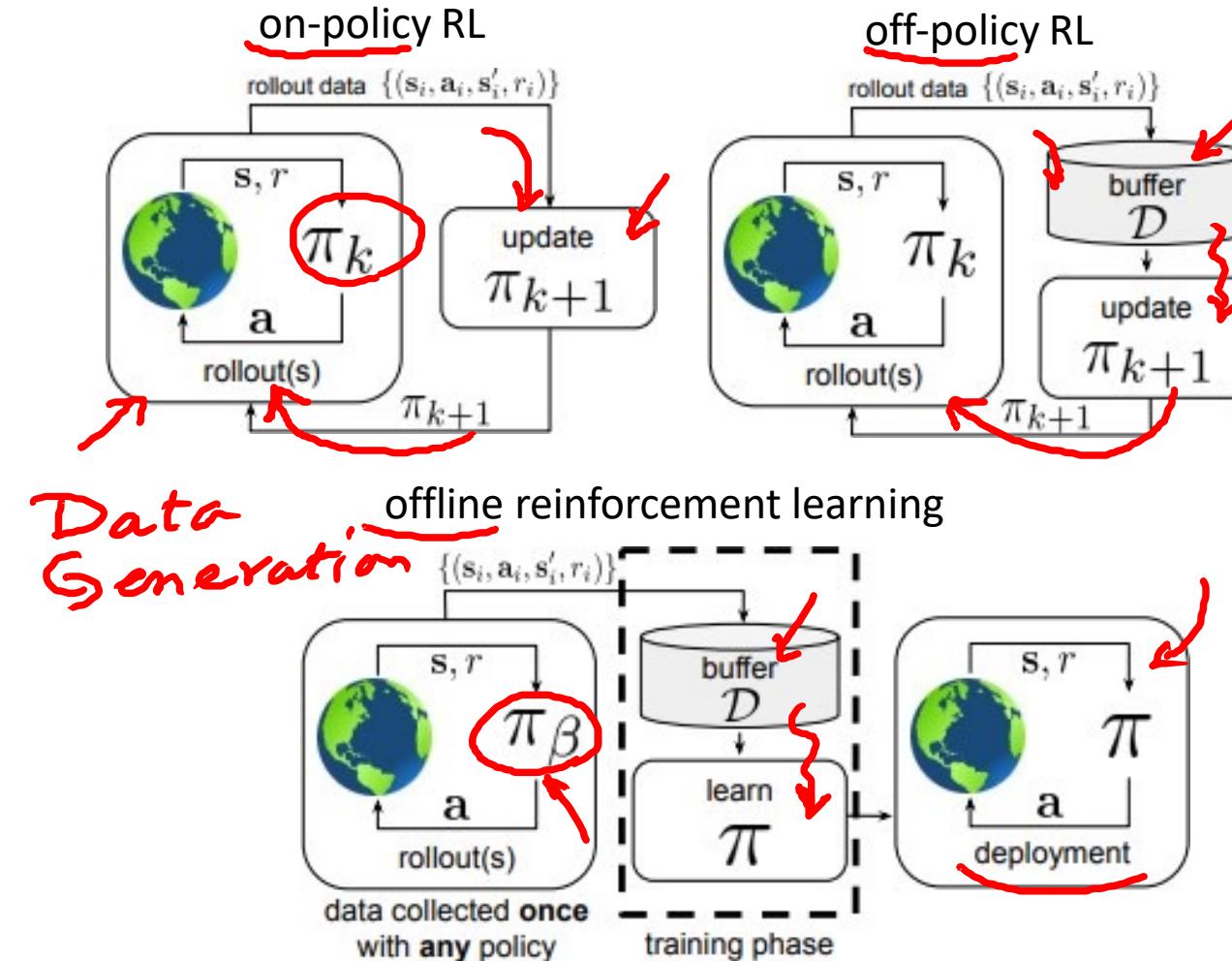
# What makes modern machine learning work?



# Can we develop data-driven RL methods?



# What does offline RL mean?



Formally:

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

$$s \sim d^{\pi_\beta}(s)$$

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

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

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

generally **not** known

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

# Types of offline RL problems

off-policy evaluation (OPE):

given  $\mathcal{D}$ , estimate  $J(\pi) = E_{\pi} \left[ \sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) \right]$

$$\pi_{\beta} \neq \pi$$

1

$$\mathcal{D} = \{(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i)\}$$

$$\mathbf{s} \sim d^{\pi_{\beta}}(\mathbf{s})$$

$$\mathbf{a} \sim \pi_{\beta}(\mathbf{a}|\mathbf{s})$$

$$\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})$$

$$r \leftarrow r(\mathbf{s}, \mathbf{a})$$

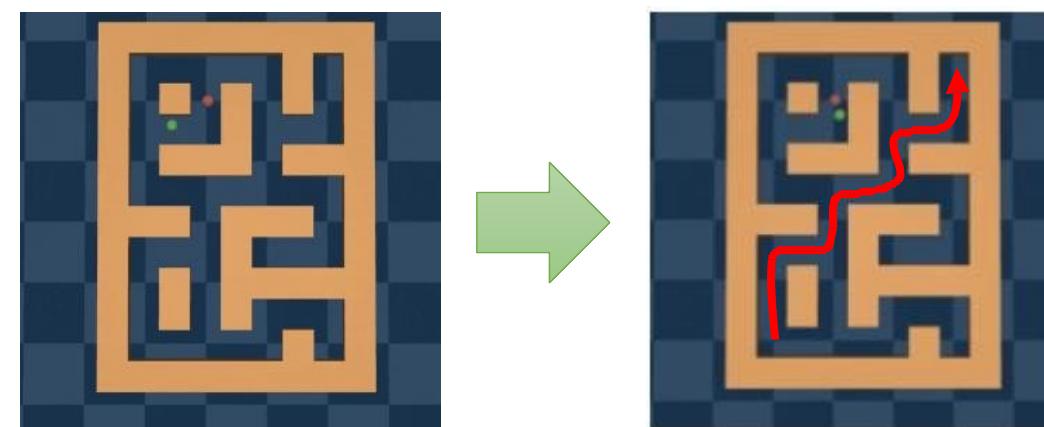
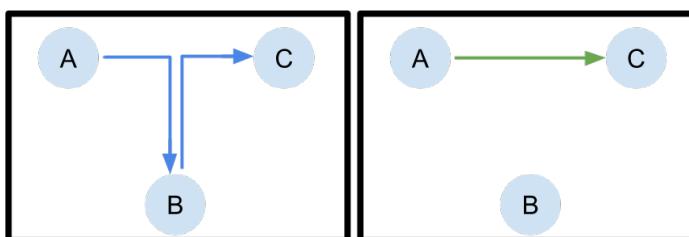
offline reinforcement learning: (a.k.a. batch RL, sometimes fully off-policy RL)

given  $\mathcal{D}$ , learn the best possible policy  $\pi_{\theta}$

not necessarily obvious what this means

# How is this even possible?

1. Find the “good stuff” in a dataset full of good and bad behaviors
2. Generalization: good behavior in one place may suggest good behavior in another place
3. “Stitching”: parts of good behaviors can be recombined



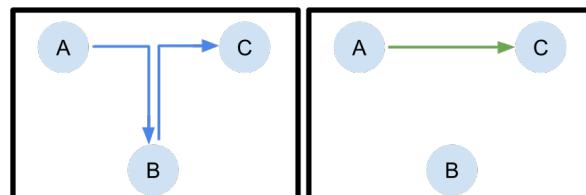
# What do we expect offline RL methods to do?

**Bad intuition:** it's like imitation learning

Though it can be shown to be **provably** better than imitation learning even with optimal data, under some structural assumptions!

See: Kumar, Hong, Singh, Levine. Should I Run Offline Reinforcement Learning or Behavioral Cloning?

**Better intuition:** get order from chaos

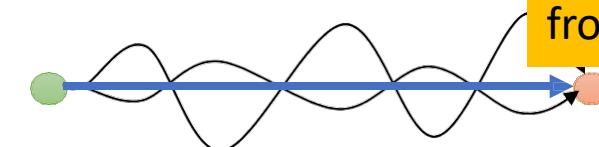
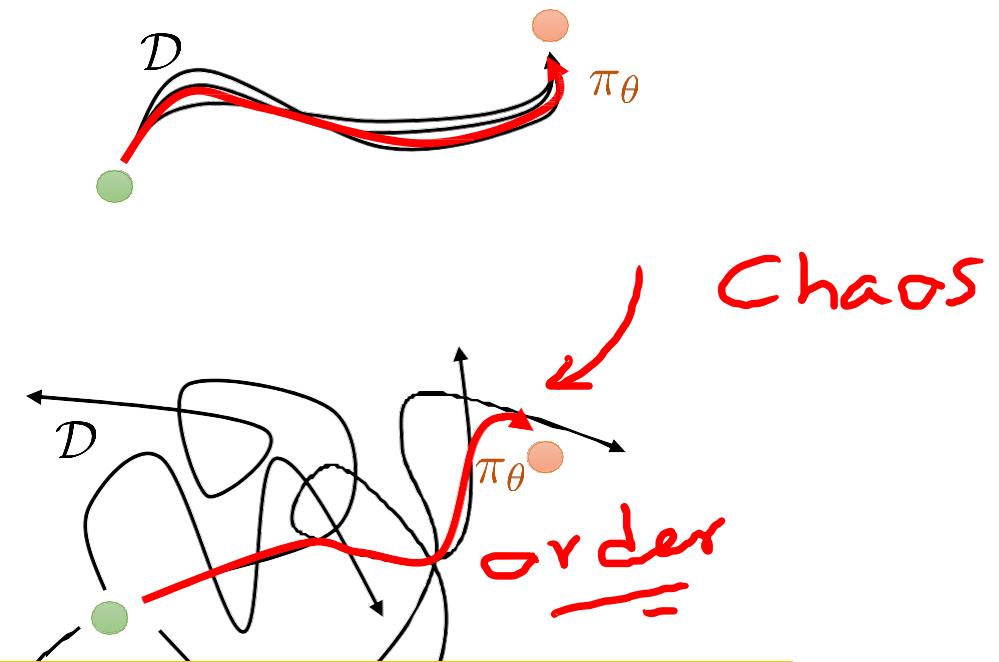


**But this is just the clearest example!**

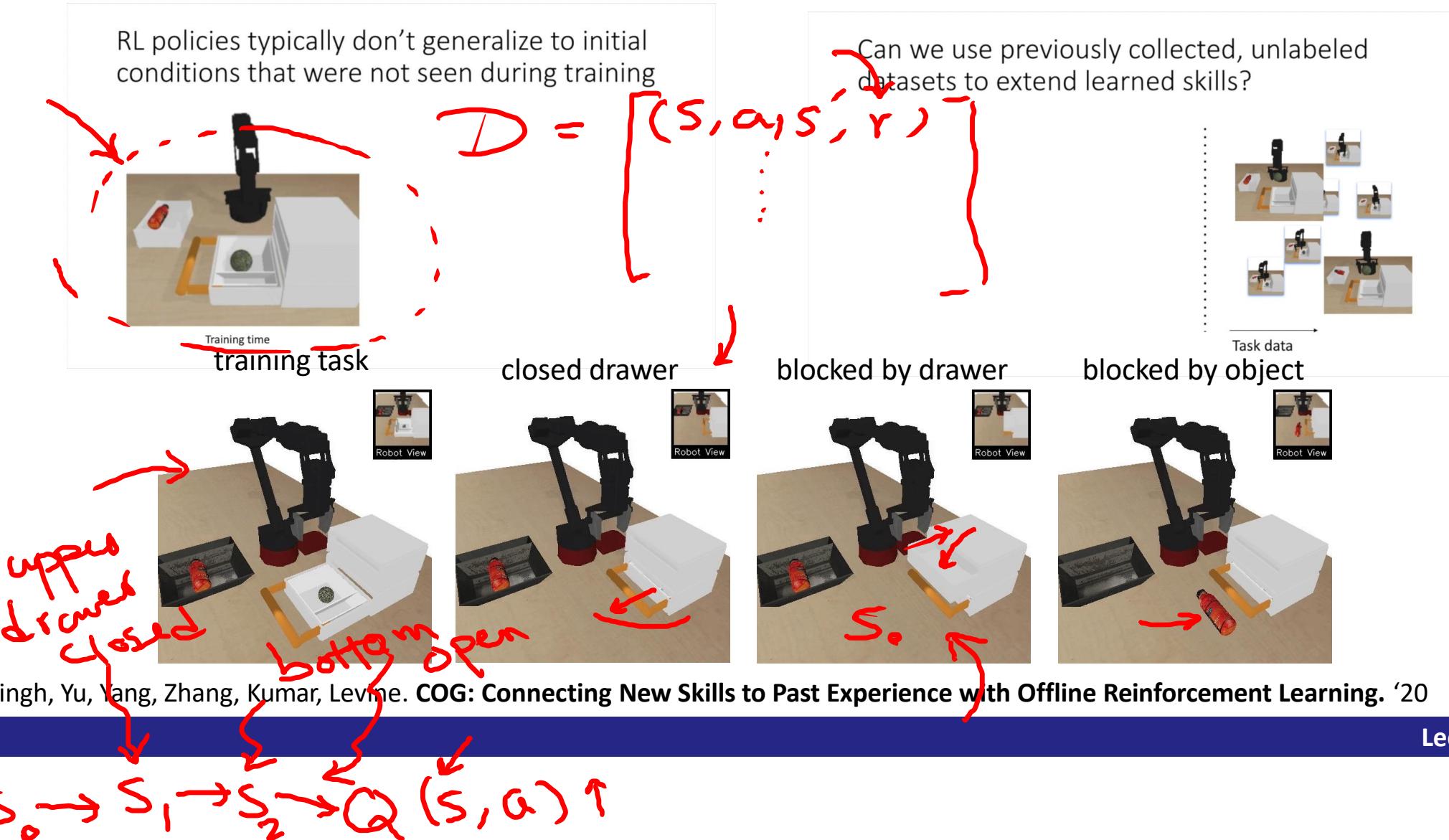
“Micro-scale” stitching:

“Macro-scale” stitching

If we have algorithms that properly perform dynamic programming, we can take this idea much further and get near-optimal policies from highly suboptimal data

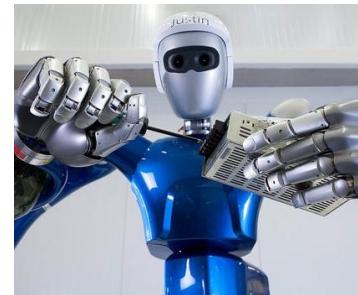
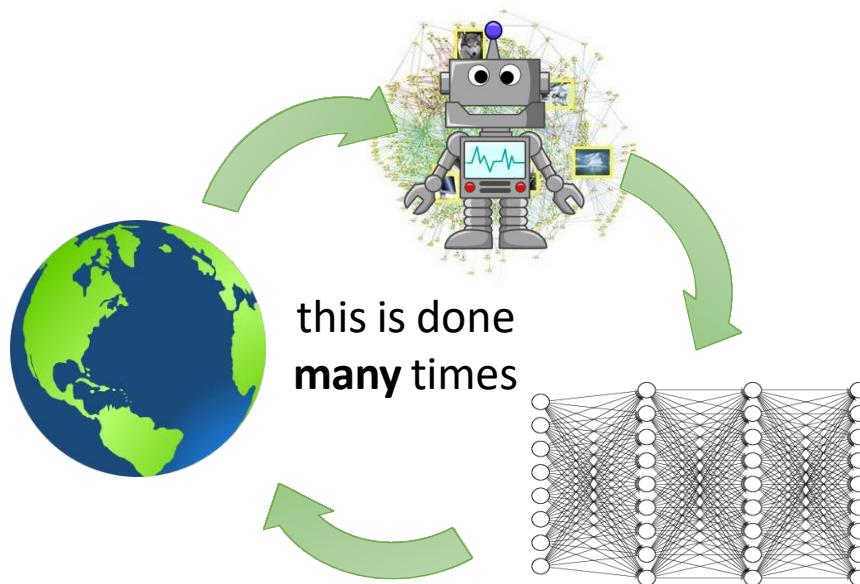


# A vivid example

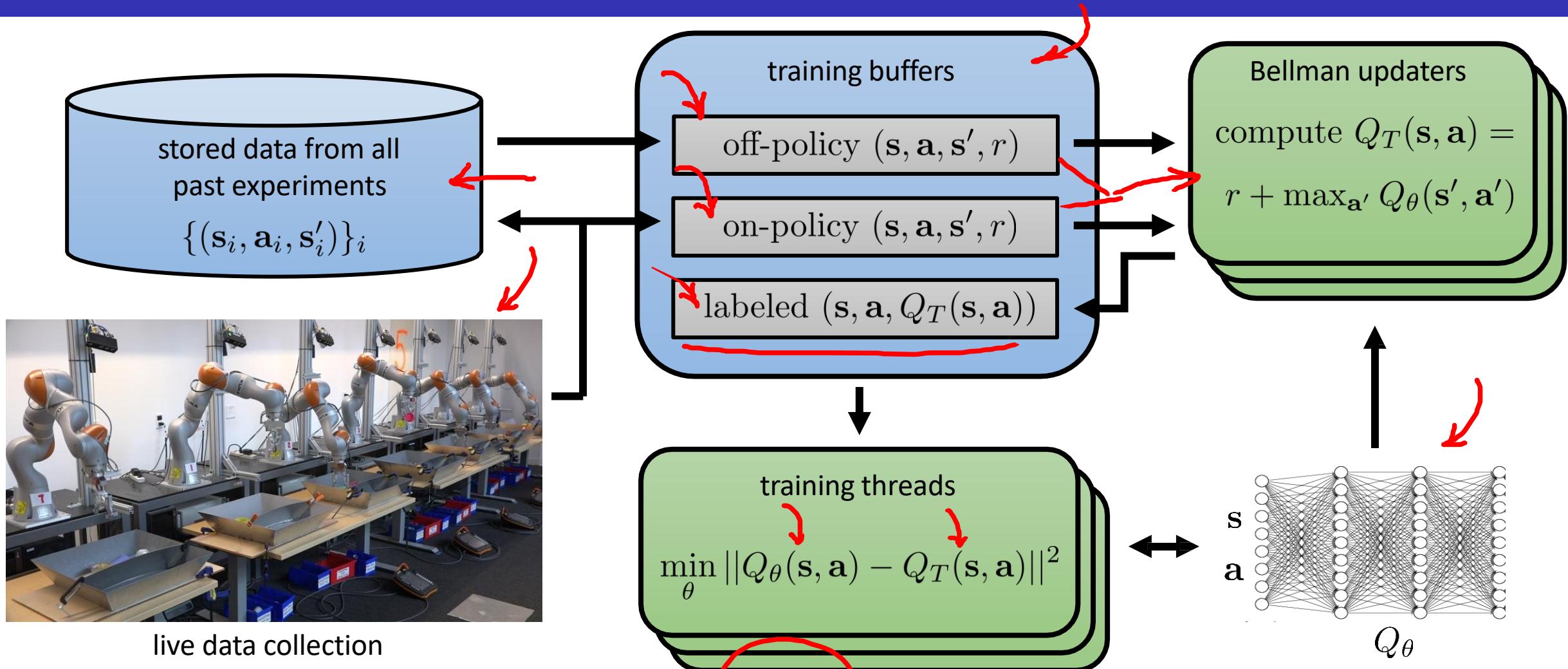


Singh, Yu, Yang, Zhang, Kumar, Levine. COG: Connecting New Skills to Past Experience with Offline Reinforcement Learning. '20

# Why should we care?



# Does it work?



# Does it work?



2x

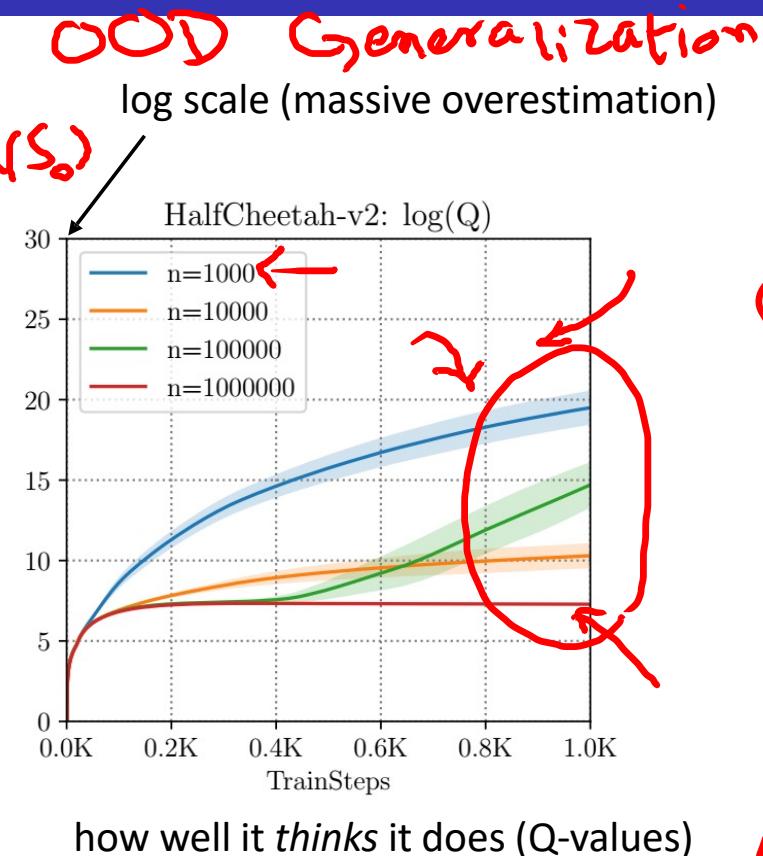
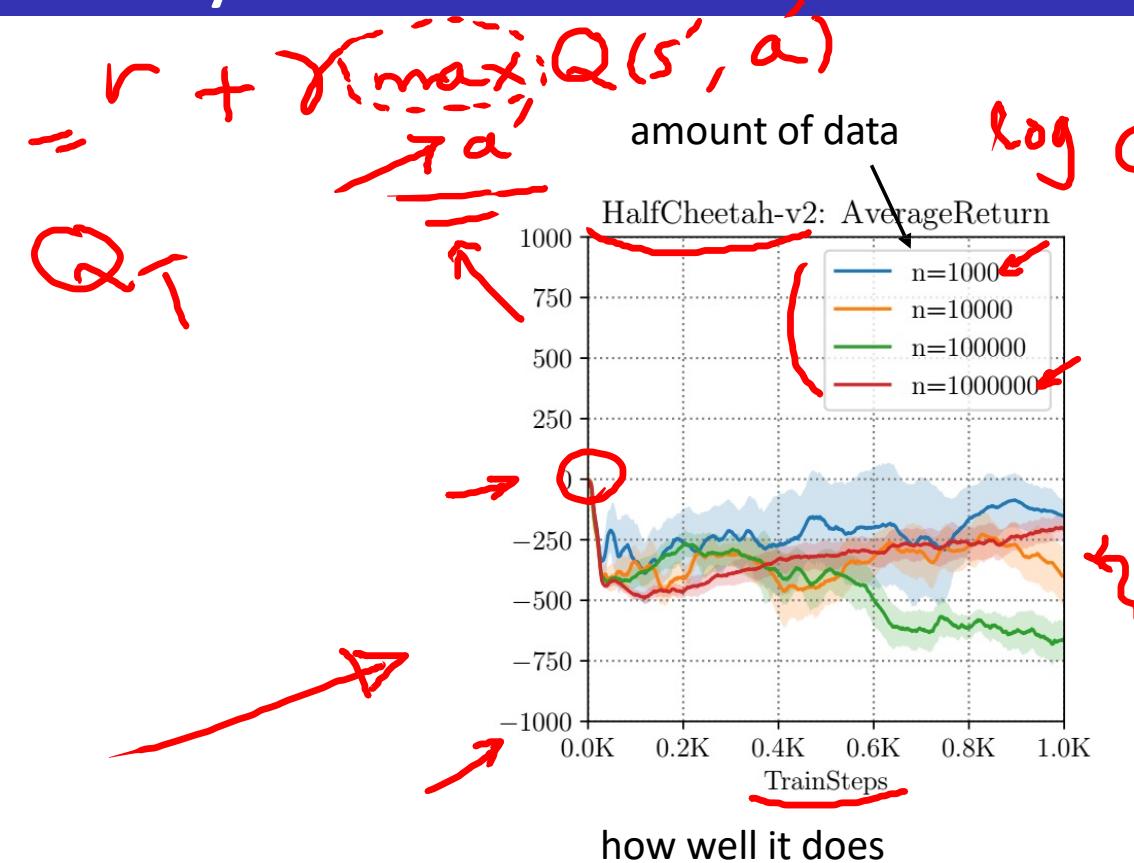


4x speed



Method	Dataset	Success	Failure
Offline QT-Opt	580k offline	87%	13%
Finetuned QT-Opt	580k offline + 28k online	96%	4%

# Why is offline RL hard?

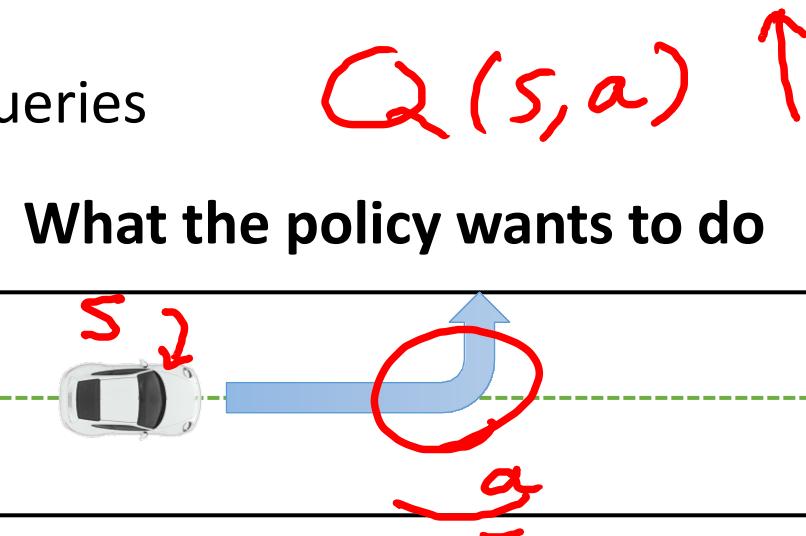


- ① act. freq. in D taken into ctcc.
- ② Don't Take

③ Take multiple steps.

# Why is offline RL hard?

Fundamental problem: counterfactual queries



Is this good? Bad?  
How do we know if  
we didn't see it in  
the data?

**Online RL** algorithms don't have to handle this, because they can simply **try** this action and see what happens

**Offline RL** methods must somehow account for these unseen ("out-of-distribution") actions, ideally in a safe way

...while still making use of generalization to come up with behaviors that are better than the best thing seen in the data!

# Distribution shift in a nutshell

Example empirical risk minimization (ERM) problem:

$$\theta \leftarrow \arg \min_{\theta} E_{\mathbf{x} \sim p(\mathbf{x}), y \sim p(y|\mathbf{x})} [(f_{\theta}(\mathbf{x}) - y)^2]$$

given some  $\mathbf{x}^*$ , is  $f_{\theta}(\mathbf{x}^*)$  correct?

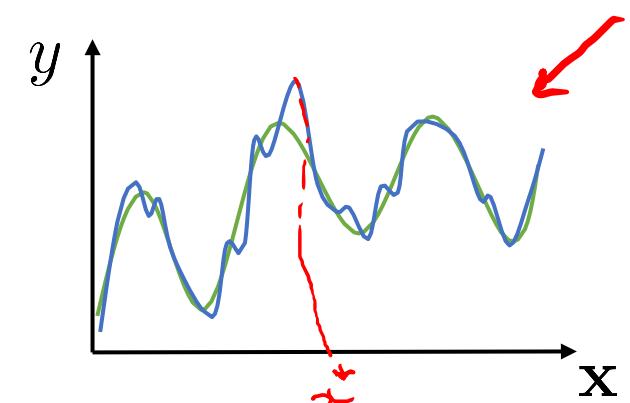
$$E_{\mathbf{x} \sim p(\mathbf{x}), y \sim p(y|\mathbf{x})} [(f_{\theta}(\mathbf{x}) - y)^2] \text{ is low}$$

$$E_{\mathbf{x} \sim \bar{p}(\mathbf{x}), y \sim p(y|\mathbf{x})} [(f_{\theta}(\mathbf{x}) - y)^2] \text{ is not, for general } \bar{p}(\mathbf{x}) \neq p(\mathbf{x})$$

what if  $\mathbf{x}^* \sim p(\mathbf{x})$ ? not necessarily...

usually we are not worried – neural nets generalize well!

what if we pick  $\mathbf{x}^* \leftarrow \arg \max_{\mathbf{x}} f_{\theta}(\mathbf{x})$ ?



# Where do we suffer from distribution shift?

$$Q(\mathbf{s}, \mathbf{a}) \leftarrow r(\mathbf{s}, \mathbf{a}) + \max_{\mathbf{a}'} Q(\mathbf{s}', \mathbf{a}') \\ Q(\mathbf{s}, \mathbf{a}) \leftarrow r(\mathbf{s}, \mathbf{a}) + E_{\mathbf{a}' \sim \pi_{\text{new}}} [Q(\mathbf{s}', \mathbf{a}')] \\ y(\mathbf{s}, \mathbf{a})$$

what is the objective?

$$\min_Q E_{(\mathbf{s}, \mathbf{a}) \sim \pi_\beta(\mathbf{s}, \mathbf{a})} [(Q(\mathbf{s}, \mathbf{a}) - y(\mathbf{s}, \mathbf{a}))^2]$$

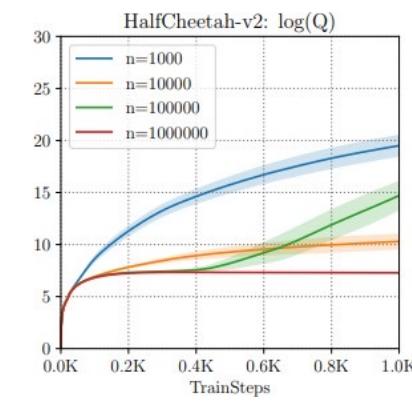
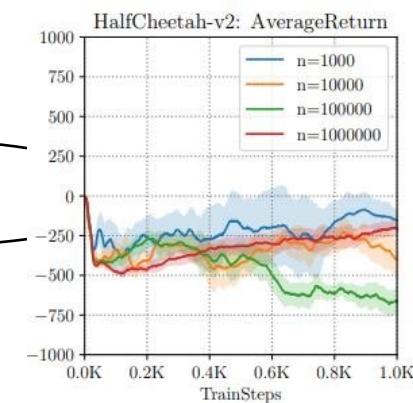


expect good accuracy when  $\pi_\beta(\mathbf{a}|\mathbf{s}) = \pi_{\text{new}}(\mathbf{a}|\mathbf{s})$

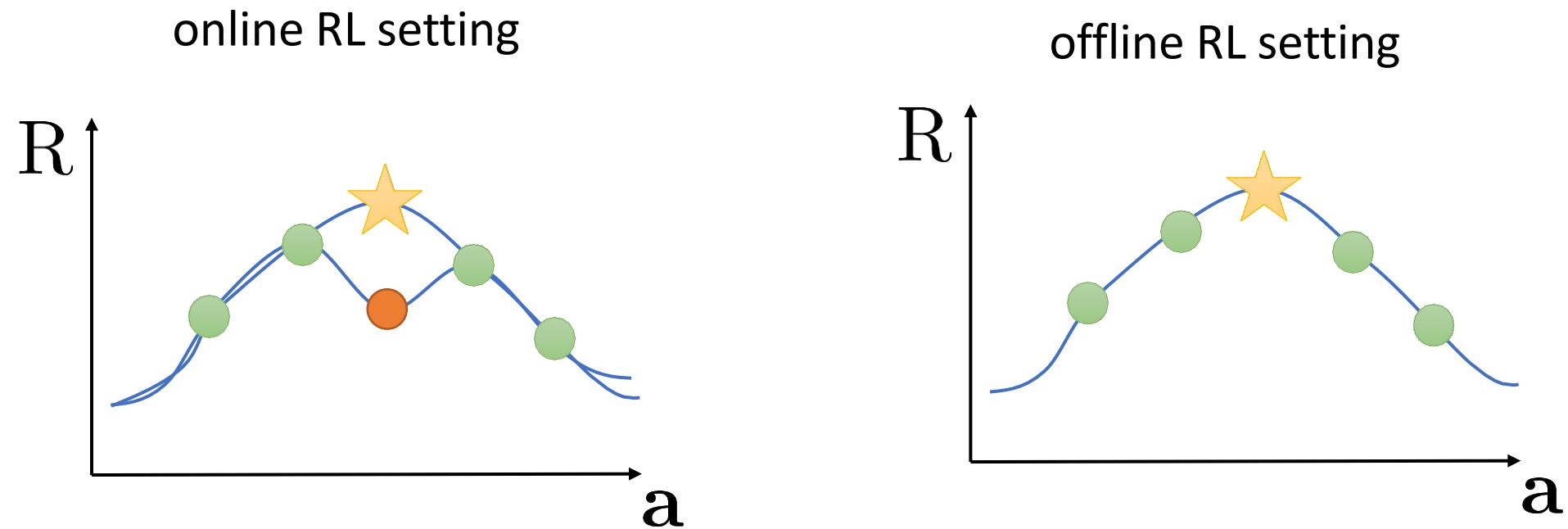
how often does *that* happen?

even worse:  $\pi_{\text{new}} = \arg \max_{\pi} E_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})}[Q(\mathbf{s}, \mathbf{a})]$

(what if we pick  $\mathbf{x}^* \leftarrow \arg \max_{\mathbf{x}} f_{\theta}(\mathbf{x})$ ?)



# Issues with generalization are not corrected



Existing challenges with sampling error and function approximation error in standard RL become **much more severe** in offline RL

## Policy Constraint Methods



# How do prior methods address this?

$$Q(s, a) \leftarrow r(s, a) + E_{a' \sim \pi_{\text{new}}} [Q(s', a')]$$
$$\pi_{\text{new}}(a|s) = \arg \max_{\pi} E_{a \sim \pi(a|s)} [Q(s, a)] \text{ s.t. } D_{\text{KL}}(\pi \| \pi_{\beta}) \leq \epsilon$$

This solves distribution shift, right?

No more erroneous values?

**Issue 1:** we usually don't know the behavior policy  $\pi_{\beta}(a|s)$

- human-provided data
- data from hand-designed controller
- data from many past RL runs
- all of the above

**Issue 2:** this is both *too pessimistic* and *not pessimistic enough*

“policy constraint” method

very old idea (but it had no single name?)

Todorov et al. [passive dynamics in linearly-solvable MDPs]

Kappen et al. [KL-divergence control, etc.]

trust regions, covariant policy gradients, natural policy gradients, etc.

used in some form in recent papers:

Fox et al. '15 (“Taming the Noise...”)

Fujimoto et al. '18 (“Off Policy...”)

Jaques et al. '19 (“Way Off Policy...”)

Kumar et al. '19 (“Stabilizing...”)

Wu et al. '19 (“Behavior Regularized...”)

# Explicit policy constraint methods

What kinds of constraints can we use?

KL-divergence:  $D_{\text{KL}}(\pi \parallel \pi_\beta)$

+ easy to implement (more on this later)

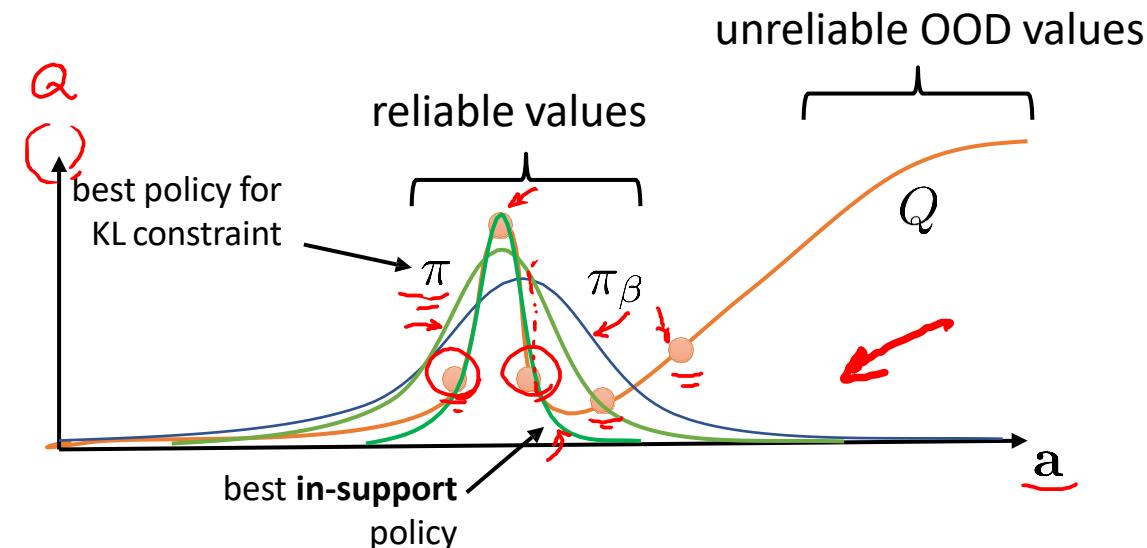
- not necessarily what we want

support constraint:  $\pi(a|s) \geq 0$  only if  $\pi_\beta(a|s) \geq \epsilon$

can approximate with MMD

- significantly more complex to implement

+ much closer to what we really want



For more information, see:

Levine, Kumar, Tucker, Fu. **Offline Reinforcement Learning: Tutorial, Review, and Perspectives on Open Problems.** '20

Kumar, Fu, Tucker, Levine. **Stabilizing Off-Policy Q-Learning via Bootstrapping Error Reduction.** '19

Wu, Tucker, Nachum. **Behavior Regularized Offline Reinforcement Learning.** '19

# Explicit policy constraint methods

How do we implement constraints?

1. Modify the actor objective

$$\theta \leftarrow \arg \max_{\theta} E_{s \sim D} [E_{a \sim \pi_{\theta}(a|s)} [Q(s, a)]]$$

$$\theta \leftarrow \arg \max_{\theta} E_{s \sim D} [E_{a \sim \pi_{\theta}(a|s)} [Q(s, a) + \lambda \log \pi_{\beta}(a|s)] + \lambda \mathcal{H}(\pi(a|s))]$$

$$D_{KL}(\pi \| \pi_{\beta}) = E_{\pi} [\log \pi(a|s) - \log \pi_{\beta}(a|s)] = -E_{\pi} [\log \pi_{\beta}(a|s)] - \mathcal{H}(\pi)$$

2. Modify the reward function

$$IE_{\pi} [\log \frac{\pi}{\pi_{\beta}}]$$

simple modification to directly penalize divergence  
also accounts for **future** divergence

$$\bar{r}(s, a) = r(s, a) - D(\pi, \pi_{\beta})$$

Lagrange multiplier

easy to compute and differentiate  
for Gaussian or categorical policies



See: Wu, Tucker, Nachum. **Behavior Regularized Offline Reinforcement Learning.** '19

generally, the best modern offline RL methods do not do either of these things

# Implicit policy constraint methods

$$\pi_{\text{new}}(\mathbf{a}|\mathbf{s}) = \arg \max_{\pi} E_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})}[Q(\mathbf{s}, \mathbf{a})] \text{ s.t. } D_{\text{KL}}(\pi \| \pi_{\beta}) \leq \epsilon$$

*Q(s, a)*

$\pi^*(\mathbf{a}|\mathbf{s}) = \frac{1}{Z(\mathbf{s})} \pi_{\beta}(\mathbf{a}|\mathbf{s}) \exp\left(\frac{1}{\lambda} A^{\pi}(\mathbf{s}, \mathbf{a})\right)$

approximate via **weighted** max likelihood!

$\pi_{\text{new}}(\mathbf{a}|\mathbf{s}) = \arg \max_{\pi} E_{(\mathbf{s}, \mathbf{a}) \sim \pi_{\beta}} \left[ \log \pi(\mathbf{a}|\mathbf{s}) \frac{1}{Z(\mathbf{s})} \exp\left(\frac{1}{\lambda} A^{\pi_{\text{old}}}(\mathbf{s}, \mathbf{a})\right) \right]$

samples from dataset  $\mathbf{a} \sim \pi_{\beta}(\mathbf{a}|\mathbf{s})$

$w(\mathbf{s}, \mathbf{a}) = \min_{\pi} D(\pi_{\pi^*}, \pi^*)$

straightforward to show via duality

See also: Peters et al. (REPS)

$\pi^*(\mathbf{s}/\mathbf{s}) \log \pi(\mathbf{a}/\mathbf{s})$

$\pi_{\beta}(\mathbf{s}/\mathbf{s}) \log \pi_{\beta}(\mathbf{a}/\mathbf{s})$

$\exp(-Q(\mathbf{s}, \mathbf{a}))$

$\log \pi(\mathbf{a}/\mathbf{s})$

Peng\*, Kumar\*, Levine. **Advantage-Weighted Regression.** '19

# Implicit policy constraint methods

$$\mathcal{L}_C(\phi) = E_{(s,a,s') \sim D} \left[ (Q_\phi(s, a) - (r(s, a) + \gamma E_{a' \sim \pi_\theta(a'|s')} [Q_\phi(s', a')]) )^2 \right]$$

Bellman Err.  
minimization

$$\mathcal{L}_A(\theta) = -E_{(s,a) \sim \pi_\beta} \left[ \log \pi_\theta(a|s) \frac{1}{Z(s)} \exp \left( \frac{1}{\lambda} A^{\pi_{\text{old}}}(s, a) \right) \right]$$

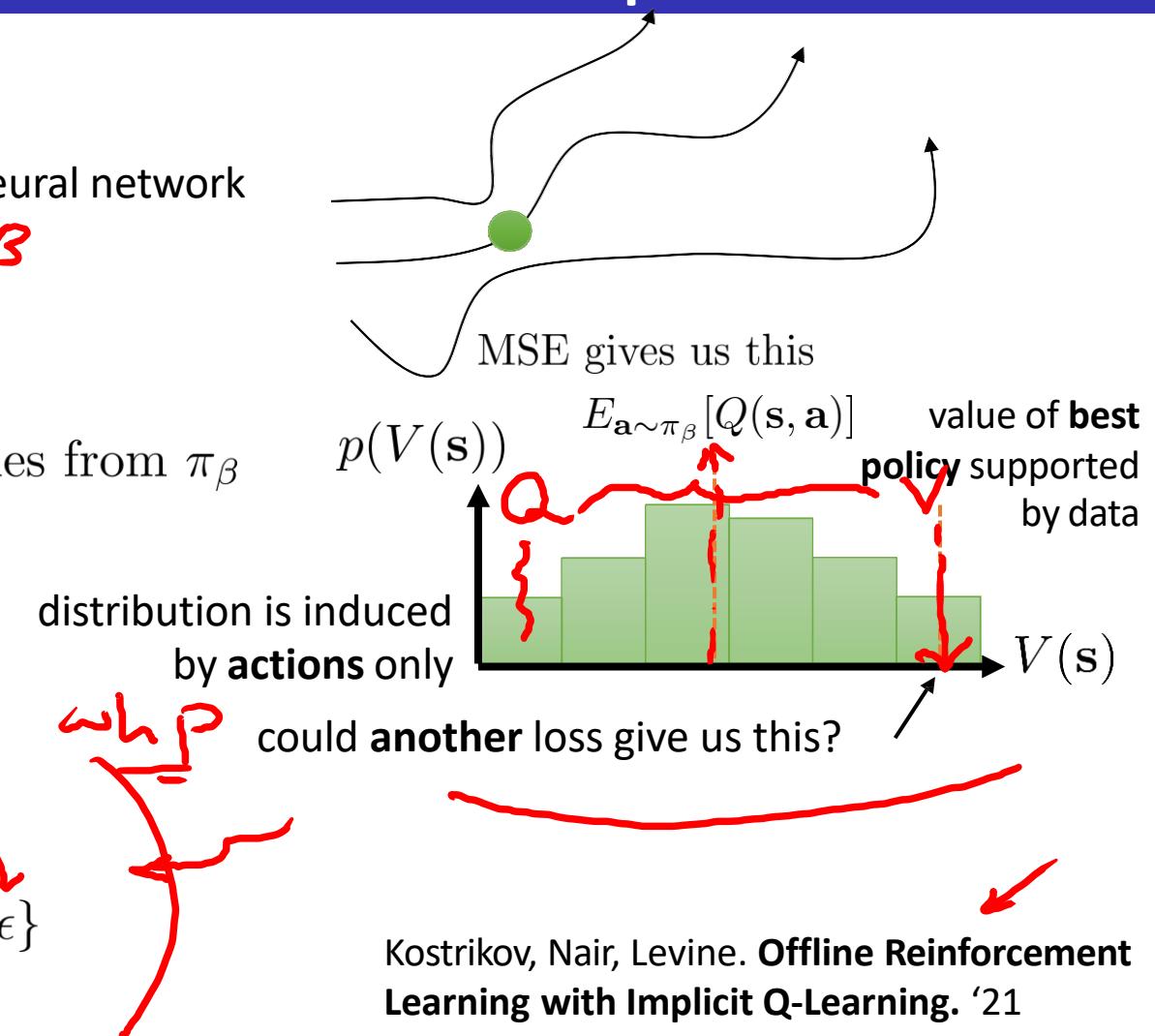
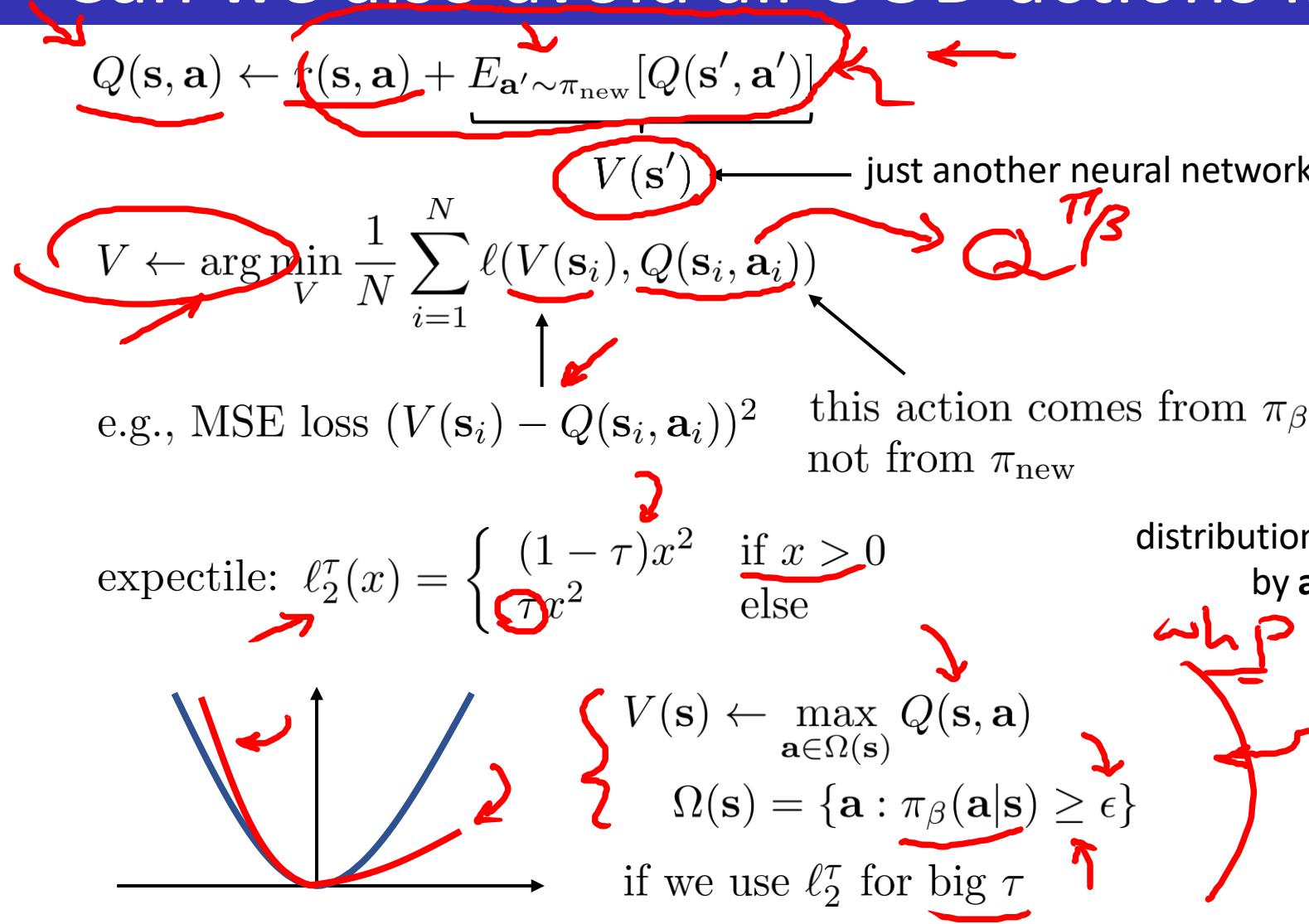
C 1.  $\phi \leftarrow \phi - \alpha \nabla_\phi \mathcal{L}_C(\phi)$

C 2.  $\theta \leftarrow \theta - \alpha \nabla_\theta \mathcal{L}_A(\theta)$

C  $Q(s, a) \leftarrow r(s, a) + E_{a' \sim \pi_{\text{new}}} [Q(s', a')]$

C  $\pi_{\text{new}}(a|s) = \arg \max_\pi E_{a \sim \pi(a|s)} [Q(s, a)] \text{ s.t. } D_{\text{KL}}(\pi \| \pi_\beta) \leq \epsilon$

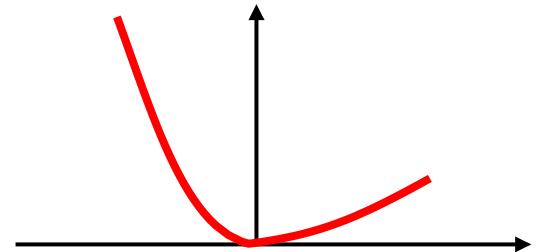
# Can we also avoid all OOD actions in the Q update?



# Implicit Q-learning (IQL)

$$Q(\mathbf{s}, \mathbf{a}) \leftarrow r(\mathbf{s}, \mathbf{a}) + V(\mathbf{s}')$$

$$V \leftarrow \arg \min_V \frac{1}{N} \sum_{i=1}^N \ell_2^\tau(V(\mathbf{s}_i), Q(\mathbf{s}_i, \mathbf{a}_i))$$



$$V(\mathbf{s}) \leftarrow \max_{\mathbf{a} \in \Omega(\mathbf{s})} Q(\mathbf{s}, \mathbf{a})$$

$$\Omega(\mathbf{s}) = \{\mathbf{a} : \pi_\beta(\mathbf{a}|\mathbf{s}) \geq \epsilon\}$$



if we use  $\ell_2^\tau$  for big  $\tau$

$$Q(\mathbf{s}, \mathbf{a}) \leftarrow r(\mathbf{s}, \mathbf{a}) + \max_{\mathbf{a}' \in \Omega(\mathbf{s}')} Q(\mathbf{s}', \mathbf{a}')$$

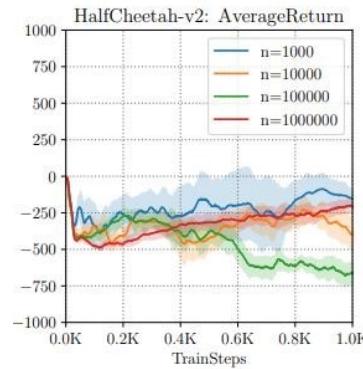
“implicit” policy

$$\pi_{\text{new}}(\mathbf{a}|\mathbf{s}) = \delta(\mathbf{a} = \arg \max_{\mathbf{a} \in \Omega(\mathbf{s})} Q(\mathbf{s}, \mathbf{a}))$$

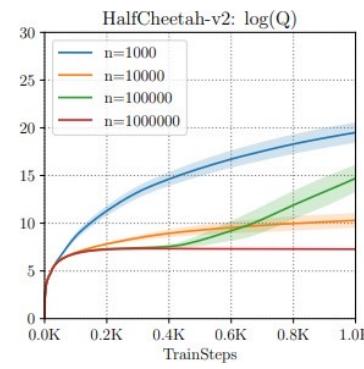
Now we can do value function updates without ever risking out-of-distribution actions!

## Conservative Q-Learning

# Conservative Q-learning (CQL)

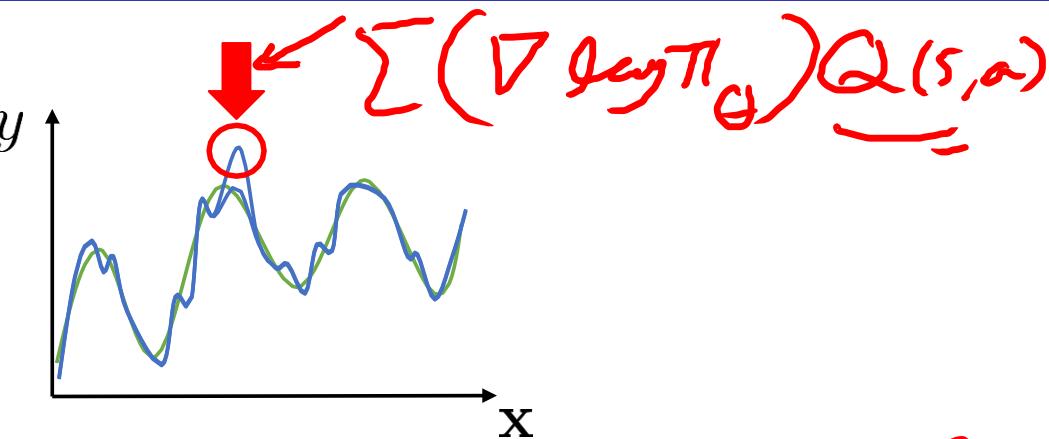


how well it does



how well it *thinks* it does (Q-values)

$$\hat{Q}^\pi = \arg \min_Q \max_\mu E_{s \sim D, a \sim \mu(a|s)} [Q(s, a)] \quad \text{regular objective}$$



$$\pi(s) = \arg \max_a Q(s, a)$$

term to push down big Q-values

$$+ E_{(s, a, s') \sim D} [(Q(s, a) - (r(s, a) + E_\pi [Q(s', a')]))^2]$$

can show that  $\hat{Q}^\pi \leq Q^\pi$  for large enough  $\alpha$

w.h.p.

true Q-function

$\alpha$

$$\hat{Q}(s, a) = Q(s, a) + \underbrace{\dots}_{\leq 0} + \underbrace{\dots}_{\geq 0}$$

# Conservative Q-learning (CQL)

A *better* bound:

$$\hat{Q}^\pi = \arg \min_Q \max_\mu \alpha E_{s \sim D, a \sim \mu(a|s)} [Q(s, a)] - \alpha E_{(s, a) \sim D} [Q(s, a)] + E_{(s, a, s') \sim D} [(Q(s, a) - (r(s, a) + E_\pi [Q(s', a')]))^2]$$

$\left. \right\} \mathcal{L}_{\text{CQL}}(\hat{Q}^\pi)$

Annotations in red:

- always pushes Q-values down (points to the first term  $\alpha E_{s \sim D, a \sim \mu(a|s)} [Q(s, a)]$ )
- push up on  $(s, a)$  samples in data (points to the second term  $\alpha E_{(s, a) \sim D} [Q(s, a)]$ )
- $(Q(s, a) - (r(s, a) + E_\pi [Q(s', a')]))^2$  (points to the third term  $E_{(s, a, s') \sim D} [(Q(s, a) - (r(s, a) + E_\pi [Q(s', a')]))^2]$ )

no longer guaranteed that  $\hat{Q}^\pi(s, a) \leq Q^\pi(s, a)$  for all  $(s, a)$

but guaranteed that  $E_{\pi(a|s)}[\hat{Q}^\pi(s, a)] \leq E_{\pi(a|s)}[Q^\pi(s, a)]$  for all  $s \in D$

# Conservative Q-learning (CQL)

- 
1. Update  $\hat{Q}^\pi$  w.r.t.  $\mathcal{L}_{\text{CQL}}(\hat{Q}^\pi)$  using  $\mathcal{D}$
  2. Update policy  $\pi$

if actions are discrete:

$$\pi(\mathbf{a}|\mathbf{s}) = \begin{cases} 1 & \text{if } \mathbf{a} = \arg \max_{\mathbf{a}} \hat{Q}(\mathbf{s}, \mathbf{a}) \\ 0 & \text{otherwise} \end{cases}$$

if actions are continuous:

$$\theta \leftarrow \theta + \alpha \nabla_\theta \sum_i E_{\mathbf{a} \sim \pi_\theta(\mathbf{a}|\mathbf{s}_i)} [\hat{Q}(\mathbf{s}_i, \mathbf{a})]$$

# Conservative Q-learning (CQL)

$$\hat{Q}^\pi = \arg \min_Q \max_\mu \alpha E_{\mathbf{s} \sim D, \mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} [Q(\mathbf{s}, \mathbf{a})] - \alpha E_{(\mathbf{s}, \mathbf{a}) \sim D} [Q(\mathbf{s}, \mathbf{a})] + \beta \mathcal{R}(\mu) + E_{(\mathbf{s}, \mathbf{a}, \mathbf{s}') \sim D} \left[ (Q(\mathbf{s}, \mathbf{a}) - (r(\mathbf{s}, \mathbf{a}) + E_\pi [Q(\mathbf{s}', \mathbf{a}')]))^2 \right]$$

$\left. \quad \quad \quad \right\} \mathcal{L}_{\text{CQL}}(\hat{Q}^\pi)$

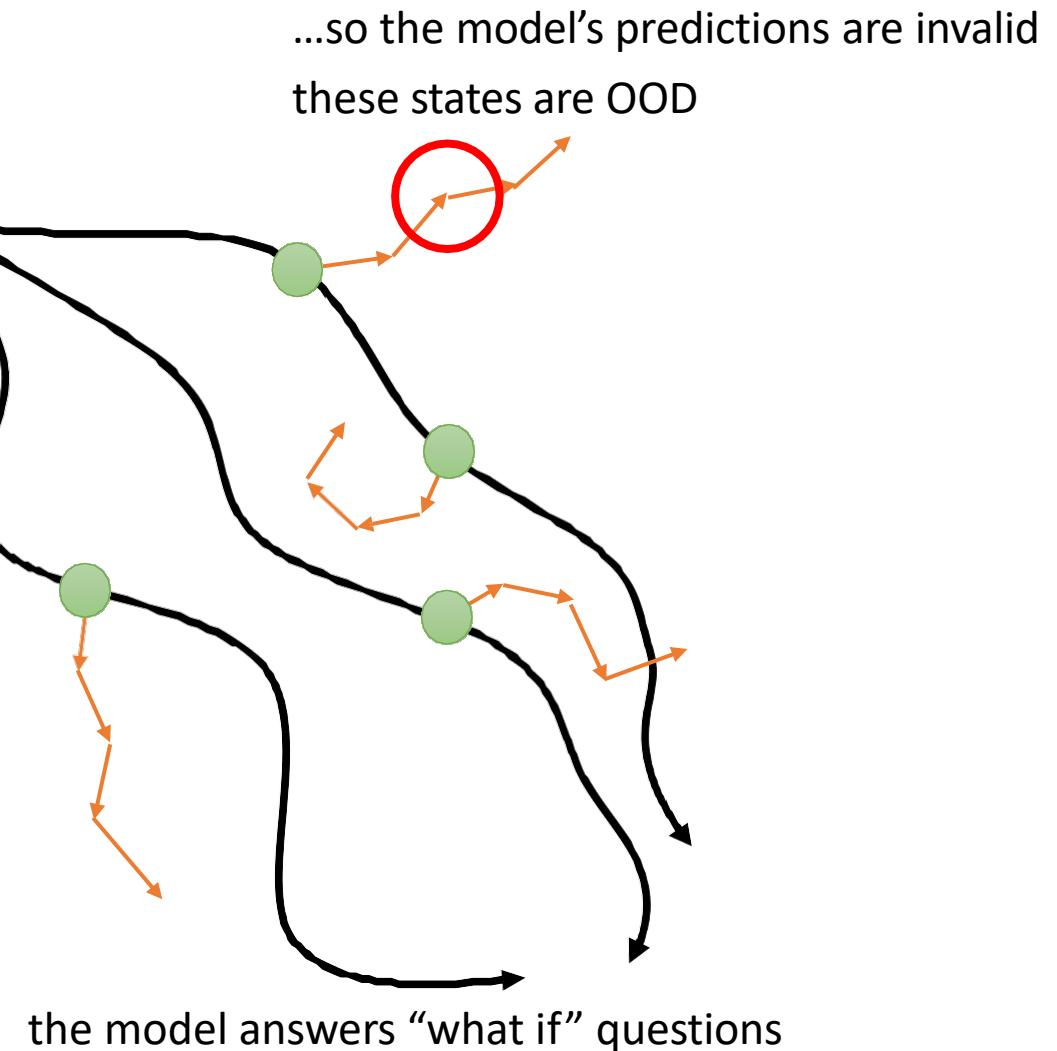
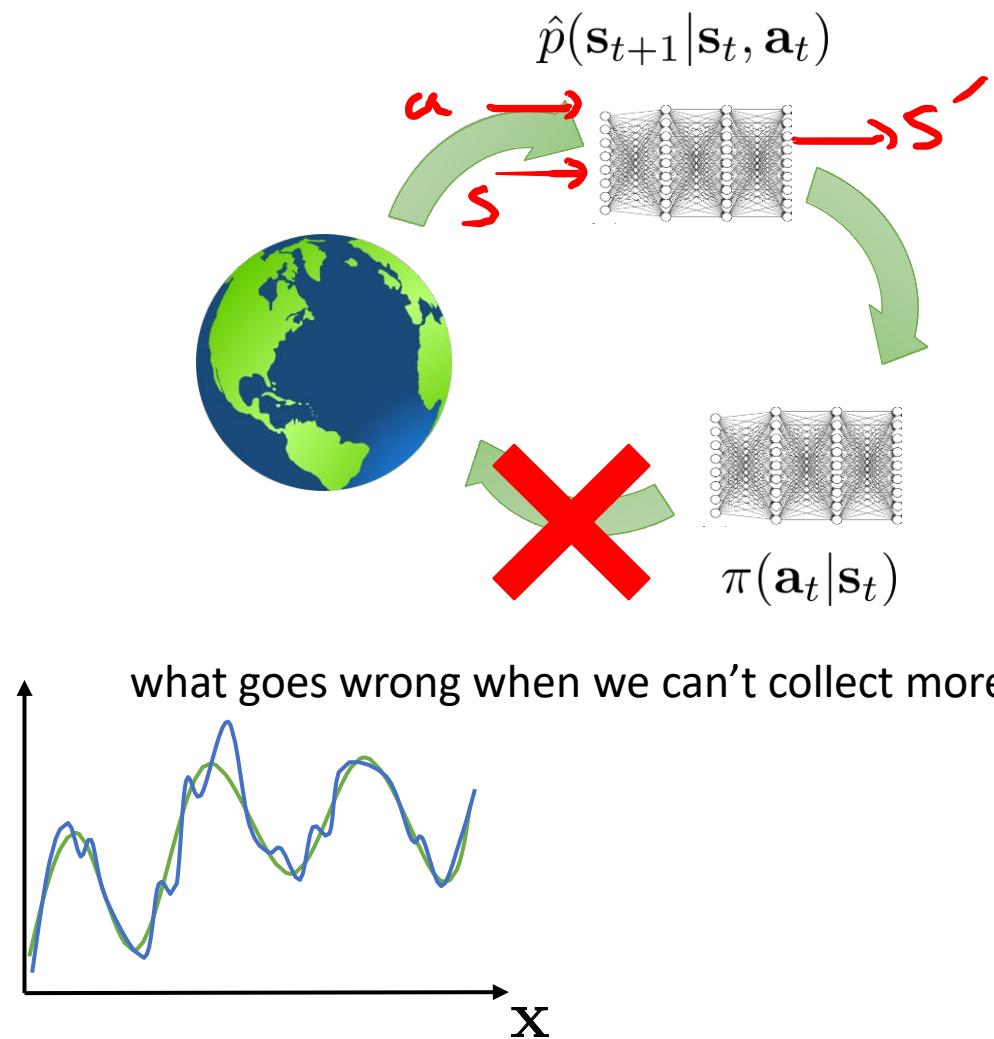
$\mu \propto \exp(\beta Q(s, a))$

regularization

common choice:  $\mathcal{R} = E_{\mathbf{s} \sim D} [\mathcal{H}(\mu(\cdot|\mathbf{s}))]$  maximum entropy regularization

## Model-Based Offline RL

# How does model-based RL work?



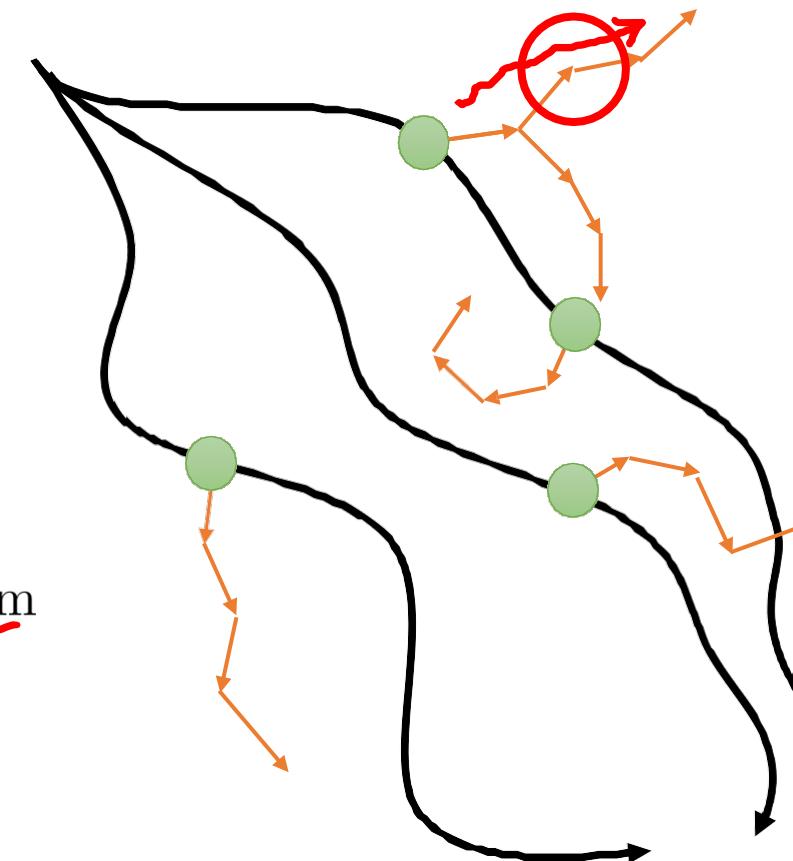
# Model-Based Offline RL

solution: “punish” the policy for exploiting

$$\tilde{r}(s, a) = r(s, a) - \lambda u(s, a)$$

uncertainty penalty

...and then use any existing model-based RL algorithm



Yu\*, Thomas\*, Yu, Ermon, Zou, Levine, Finn, Ma. **MOPO: Model-Based Offline Policy Optimization.** '20

See also: Kidambi et al., **MOReL : Model-Based Offline Reinforcement Learning.** '20 (concurrent)

# Conservative Model-Based RL

**Basic idea:** just like CQL minimizes Q-value of policy actions, we can minimize Q-value of model state-action tuples

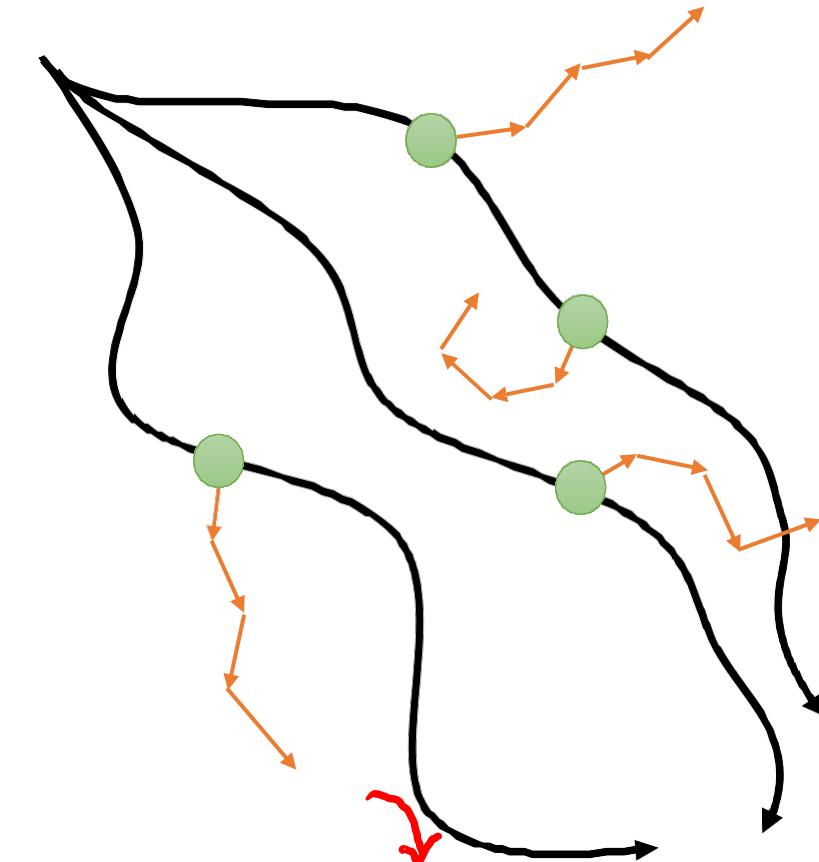
state-action tuples from the model

$$\begin{aligned}\hat{Q}^{k+1} \leftarrow \arg \min_Q & \beta (\mathbb{E}_{s,a \sim \rho(s,a)} [Q(s,a)] - \mathbb{E}_{s,a \sim \mathcal{D}} [Q(s,a)]) \\ & + \frac{1}{2} \mathbb{E}_{s,a,s' \sim d_f} \left[ (Q(s,a) - \hat{\mathcal{B}}^\pi \hat{Q}^k(s,a))^2 \right].\end{aligned}\quad (4)$$

*Dyna Style*

**Intuition:** if the model produces something that looks clearly different from real data, it's easy for the Q-function to make it look bad

$s, a \mid s'$



# Summary, Applications, Open Questions

# Which offline RL algorithm do I use?

If you want to only train offline...

Conservative Q-learning	+ just one hyperparameter	+ well understood and widely tested
Implicit Q-learning	+ more flexible (offline + online)	- more hyperparameters

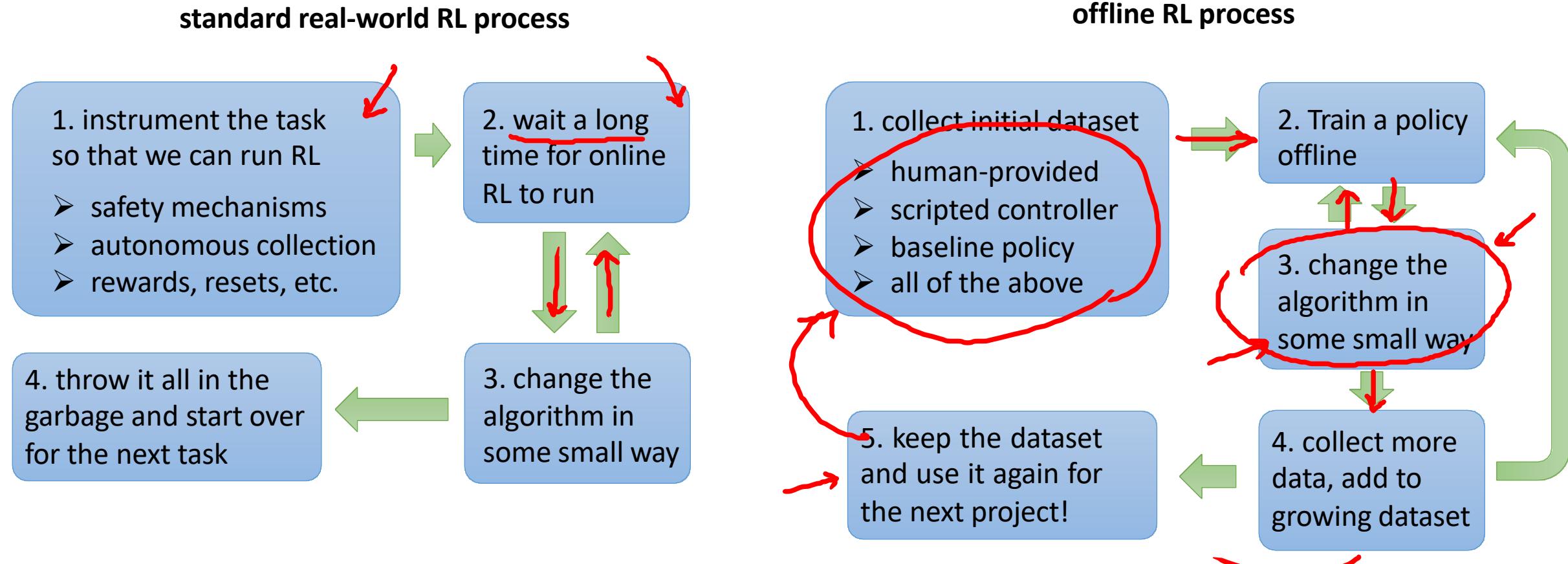
If you want to only train offline and finetune online

Advantage-weighted actor-critic (AWAC)	+ widely used and well tested
Implicit Q-learning	+ seems to perform much better!

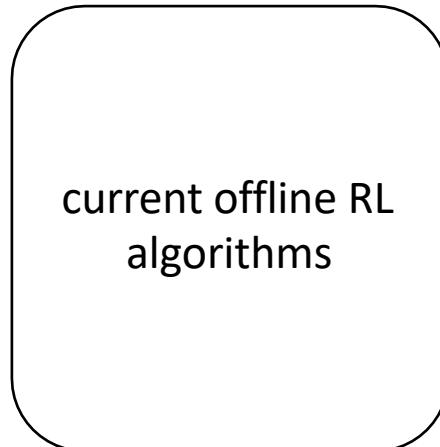
If you have a good way to train models in your domain

COMBO	+ similar properties as CQL, but benefits from models
	- not always easy to train a good model in your domain!

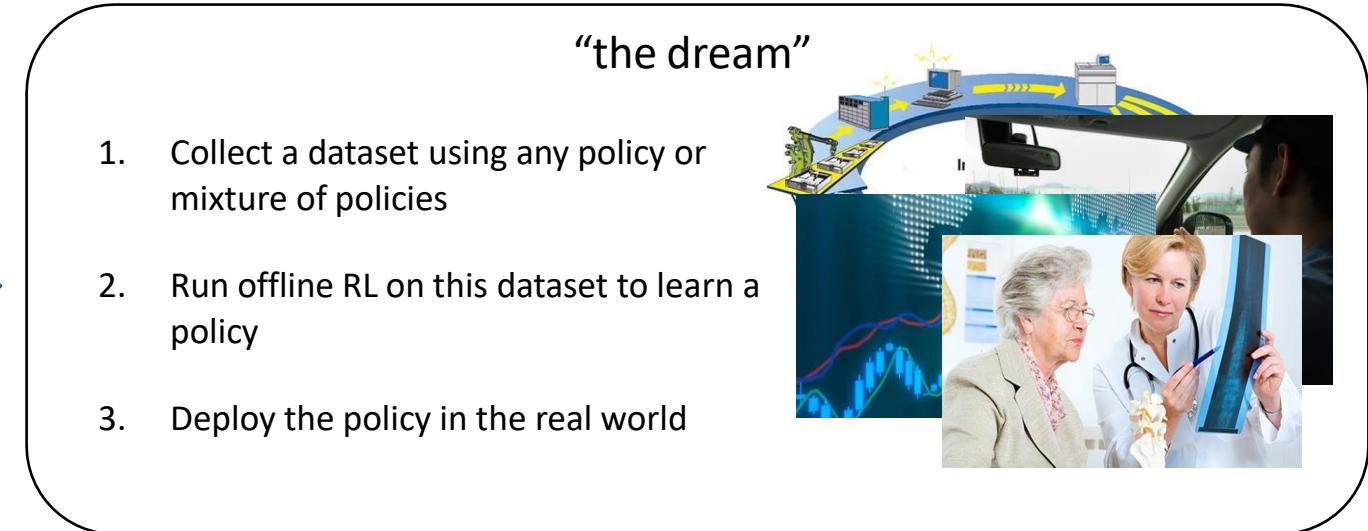
# The power of offline RL



# Takeaways, conclusions, future directions



“the gap”



- An offline RL **workflow**
  - Supervised learning workflow: train/test split
  - Offline RL workflow: ???
- Statistical **guarantees**
  - Biggest challenge: distributional shift/counterfactuals
  - Can we make any guarantees?
- Scalable methods, large-scale applications
  - Dialogue systems
  - Data-driven navigation and driving