

# Robot Learning

Offline Reinforcement Learning

Prof. Songhwai Oh

ECE, SNU

## Handling extrapolation error

- Off-Policy Deep Reinforcement Learning without Exploration
- Stabilizing Off-Policy Q-Learning via Bootstrapping Error Reduction
- Conservative Q-learning for offline reinforcement learning
- Implicit Constraint Approach for Offline Multi-Agent Reinforcement Learning

## Data augmentation

- S4RL: Surprisingly Simple Self-Supervision for Offline Reinforcement Learning

## Model-based

- MOREL: Model-Based Offline Reinforcement Learning

## Supervised learning

- Decision Transformer: Reinforcement Learning via Sequence Modeling
- Trajectory Transformer
- Learning to generalize across long-horizon tasks from human demonstrations

Fujimoto, Scott, David Meger, and Doina Precup. "Off-policy deep reinforcement learning without exploration." ICML, 2019.

# **OFF-POLICY DEEP REINFORCEMENT LEARNING WITHOUT EXPLORATION**

- Problem
  - In batch reinforcement learning (=offline RL), off-policy algorithms (DQN, DDPG) do not learn well.
  - This is due to an extrapolation error.
  - The cause of this problem is that the distribution of data induced by the policy is different from the distribution of batch data.
  
- Proposed method
  - Propose batch-constrained deep Q-learning (BCQ).
  - Use a state-conditioned generative model to produce actions similar to those in the batch.
  - Can learn successfully in a setting without interaction.

- Extrapolation Error
  - Error caused by the mismatch between the dataset and state-action visitations of the current policy
  - The value estimate  $Q(s, a)$  is affected by the extrapolation error during the value update
- Cause of extrapolation error
  - 1. Absent Data**
    - The estimate of  $Q(s', a')$  may be arbitrarily bad because there is not enough data near  $s', a'$ .
  - 2. Model Bias**
    - Occurs due to the difference between the transitions of the true MDP and the transition distribution within the data.
  - 3. Training Mismatch**
    - Occurs when distribution under the current policy and data distribution are different
- In case of on-policy setting, it is possible to reduce error by exploring overestimated part due to extrapolation error, but it is impossible in offline setting.

- To avoid extrapolation error, a policy should induce a similar state-action to the batch (batch-constrained).
- Batch-constrained policies are trained to select actions with respect to three objectives:
  1. Minimize the distance of selected actions to the data in the batch. (most important)
  2. Lead to states where familiar data can be observed.
  3. Maximize the value function.
- BCQ approaches the notion of batch-constrained through a generative model (produce likely actions under the batch).

- Generative model for a batch (in this paper, use VAE) :  $G_\omega(s)$
- Perturbation model :  $\xi_\phi(s, a, \Phi)$ 
  - Outputs an adjustment to an action  $a$  in the range  $[-\Phi, \Phi]$
- Policy  $\pi$  :  $\pi(s) = \underset{a_i + \xi_\phi(s, a_i, \Phi)}{\operatorname{argmax}} Q_\theta(s, a_i + \xi_\phi(s, a_i, \Phi)),$   
 $\{a_i \sim G_\omega(s)\}_{i=1}^n.$
- Perturbation model training: actions are sampled from a generative model.

$$\phi \leftarrow \underset{\phi}{\operatorname{argmax}} \sum_{(s,a) \in \mathcal{B}} Q_\theta(s, a + \xi_\phi(s, a, \Phi))$$

- To penalize high variance in uncertainty region, use clipped double Q-Learning.

$$r + \gamma \max_{a_i} \left[ \lambda \min_{j=1,2} Q_{\theta'_j}(s', a_i) + (1 - \lambda) \max_{j=1,2} Q_{\theta'_j}(s', a_i) \right]$$

---

## Algorithm 1 BCQ

---

**Input:** Batch  $\mathcal{B}$ , horizon  $T$ , target network update rate  $\tau$ , mini-batch size  $N$ , max perturbation  $\Phi$ , number of sampled actions  $n$ , minimum weighting  $\lambda$ .

Initialize Q-networks  $Q_{\theta_1}, Q_{\theta_2}$ , perturbation network  $\xi_\phi$ , and VAE  $G_\omega = \{E_{\omega_1}, D_{\omega_2}\}$ , with random parameters  $\theta_1, \theta_2, \phi, \omega$ , and target networks  $Q_{\theta'_1}, Q_{\theta'_2}, \xi_{\phi'}$  with  $\theta'_1 \leftarrow \theta_1, \theta'_2 \leftarrow \theta_2, \phi' \leftarrow \phi$ .

**for**  $t = 1$  **to**  $T$  **do**

    Sample mini-batch of  $N$  transitions  $(s, a, r, s')$  from  $\mathcal{B}$

$\mu, \sigma = E_{\omega_1}(s, a), \quad \tilde{a} = D_{\omega_2}(s, z), \quad z \sim \mathcal{N}(\mu, \sigma)$

$\omega \leftarrow \operatorname{argmin}_\omega \sum (a - \tilde{a})^2 + D_{\text{KL}}(\mathcal{N}(\mu, \sigma) \parallel \mathcal{N}(0, 1))$

    Sample  $n$  actions:  $\{a_i \sim G_\omega(s')\}_{i=1}^n$

    Perturb each action:  $\{a_i = a_i + \xi_\phi(s', a_i, \Phi)\}_{i=1}^n$

    Set value target  $y$  (Eqn. 13)

$\theta \leftarrow \operatorname{argmin}_\theta \sum (y - Q_\theta(s, a))^2$

$\phi \leftarrow \operatorname{argmax}_\phi \sum Q_{\theta_1}(s, a + \xi_\phi(s, a, \Phi)), a \sim G_\omega(s)$

    Update target networks:  $\theta'_i \leftarrow \tau\theta + (1 - \tau)\theta'_i$

$\phi' \leftarrow \tau\phi + (1 - \tau)\phi'$

**end for**

---

Kumar, Aviral, Justin Fu, Matthew Soh, George Tucker, and Sergey Levine.  
"Stabilizing off-policy Q-learning via bootstrapping error reduction." NeurIPS, 2019.

# STABILIZING OFF-POLICY Q-LEARNING VIA BOOTSTRAPPING ERROR REDUCTION

- The proposed method, **Bootstrapping Error Accumulation Reduction (BEAR)**, restricts the policy to ensure that it outputs actions that lie in the support of the training distribution.
- The existing methods (e.g., BCQ) implicitly constrain the distribution of the learned policy to be close to the behavior policy.  
→ It is overly restrictive.  
For example, if the behavior policy is close to uniform, the learned policy will behave randomly, even if the data is sufficient.
- BEAR enforces the learned policy  $\pi(a|s)$  to have a positive density **only where the density of the behavior policy  $\beta(a|s)$  is more than a threshold.**  
(i.e.,  $\forall a, \beta(a|s) \leq \varepsilon \implies \pi(a|s) = 0$ )

- **Support (policy) set**  $\Pi$ : the set of restricted policies
- Upper bound on the performance of distribution-constrained Q-iteration:

$$\lim_{k \rightarrow \infty} \mathbb{E}_{\rho_0} [|V^{\pi^k}(s) - V^*(s)|] \leq \frac{\gamma}{(1-\gamma)^2} \left[ \underbrace{C(\Pi) \mathbb{E}_{\mu} [\max_{\pi \in \Pi} \mathbb{E}_{\pi} [\delta(s, a)]]}_{\text{current Bellman error}} + \frac{1-\gamma}{\gamma} \underbrace{\alpha(\Pi)}_{\text{suboptimality constant}} \right]$$

concentrability coefficient
suboptimality constant  
: measures how far  $\Pi$  is from the training data distribution
: measures how far  $\pi^*$  is from  $\Pi$

This bound formalizes the tradeoff between keeping policies chosen during backups close to the data and keeping the set  $\Pi$  large enough to capture well-performing policies.

→ BEAR uses the following support set:  $\Pi_{\epsilon} = \{\pi \mid \pi(a|s) = 0 \text{ whenever } \beta(a|s) < \epsilon\}$   
 $\beta$ : the behavior policy

→ Therefore, **using the support set  $\Pi_{\epsilon}$**  gives us a single lever,  $\epsilon$ , which **simultaneously trades off the value of  $C(\Pi_{\epsilon})$  and  $\alpha(\Pi_{\epsilon})$** .

- BEAR uses a set of Q-functions:  $\hat{Q}_1, \dots, \hat{Q}_K$
- Then, the policy is updated to maximize the conservative estimate of Q-values within  $\Pi_\epsilon$ .

$$\pi_\phi(s) := \max_{\pi \in \Pi_\epsilon} \mathbb{E}_{a \sim \pi(\cdot|s)} \left[ \min_{j=1, \dots, K} \hat{Q}_j(s, a) \right]$$

- Since the behavior policy  $\beta$  is unknown, they use the maximum mean discrepancy (MMD) approximation to measure distance between two distributions using samples.

$$\text{MMD}^2(\{x_1, \dots, x_n\}, \{y_1, \dots, y_m\}) = \frac{1}{n^2} \sum_{i, i'} k(x_i, x_{i'}) - \frac{2}{nm} \sum_{i, j} k(x_i, y_j) + \frac{1}{m^2} \sum_{j, j'} k(y_j, y_{j'})$$

k: universal kernel

In BEAR, k is Laplacian or Gaussian Kernel

- Finally, the optimization problem in the policy improvement step becomes

$$\pi_\phi := \max_{\pi \in \Delta_{|S|}} \mathbb{E}_{s \sim \mathcal{D}} \mathbb{E}_{a \sim \pi(\cdot|s)} \left[ \min_{j=1, \dots, K} \hat{Q}_j(s, a) \right] \quad \text{s.t.} \quad \mathbb{E}_{s \sim \mathcal{D}} [\text{MMD}(\mathcal{D}(s), \pi(\cdot|s))] \leq \epsilon$$

→ Equation 1

- Equation 1:

$$\pi_\phi := \max_{\pi \in \Delta_{|S|}} \mathbb{E}_{s \sim \mathcal{D}} \mathbb{E}_{a \sim \pi(\cdot|s)} \left[ \min_{j=1, \dots, K} \hat{Q}_j(s, a) \right] \quad \text{s.t.} \quad \mathbb{E}_{s \sim \mathcal{D}} [\text{MMD}(\mathcal{D}(s), \pi(\cdot|s))] \leq \varepsilon$$

---

## Algorithm 1 BEAR Q-Learning (BEAR-QL)

---

- input** : Dataset  $\mathcal{D}$ , target network update rate  $\tau$ , mini-batch size  $N$ , sampled actions for MMD  $n$ , minimum  $\lambda$
- 1: Initialize Q-ensemble  $\{Q_{\theta_i}\}_{i=1}^K$ , actor  $\pi_\phi$ , Lagrange multiplier  $\alpha$ , target networks  $\{Q_{\theta'_i}\}_{i=1}^K$ , and a target actor  $\pi_{\phi'}$ , with  $\phi' \leftarrow \phi, \theta'_i \leftarrow \theta_i$
  - 2: **for**  $t$  in  $\{1, \dots, N\}$  **do**
  - 3:   Sample mini-batch of transitions  $(s, a, r, s') \sim \mathcal{D}$
  - Q-update:**
  - 4:   Sample  $p$  action samples,  $\{a_i \sim \pi_{\phi'}(\cdot|s')\}_{i=1}^p$
  - 5:   Define  $y(s, a) := \max_{a_i} [\lambda \min_{j=1, \dots, K} Q_{\theta'_j}(s', a_i) + (1 - \lambda) \max_{j=1, \dots, K} Q_{\theta'_j}(s', a_i)]$
  - 6:    $\forall i, \theta_i \leftarrow \arg \min_{\theta_i} (Q_{\theta_i}(s, a) - (r + \gamma y(s, a)))^2$
  - Policy-update:**
  - 7:   Sample actions  $\{\hat{a}_i \sim \pi_\phi(\cdot|s)\}_{i=1}^m$  and  $\{a_j \sim \mathcal{D}(s)\}_{j=1}^n$ ,  $n$  preferably an intermediate integer(1-10)
  - 8:   Update  $\phi, \alpha$  by minimizing Equation 1 by using dual gradient descent with Lagrange multiplier  $\alpha$
  - 9:   **Update Target Networks:**  $\theta'_i \leftarrow \tau \theta_i + (1 - \tau) \theta'_i; \phi' \leftarrow \tau \phi + (1 - \tau) \phi'$
  - 10: **end for**
-

Kumar, Aviral, Aurick Zhou, George Tucker, and Sergey Levine. "**Conservative Q-learning for offline reinforcement learning.**" NeurIPS 2020.

# **CONSERVATIVE Q-LEARNING FOR OFFLINE REINFORCEMENT LEARNING**

- Offline RL often fails due to *extrapolation errors*
  - Overestimate Q-value on actions not appearing in dataset
- In order to prevent overestimation, conservative q-learning (CQL) penalizes the value of out-of-distribution state-action pairs.
- By doing so, the agent's policy tends to follow the actions explored by the behavior policy.

- Learning the lower-bound of Q-function by minimizing the expected Q-value under a particular distribution,

$$\mu(s, a) = d^{\pi_\beta}(s)\mu(a|s), \quad \pi_\beta: \text{behavior policy}$$

$$\hat{Q}^{k+1} \leftarrow \arg \min_Q \underbrace{\alpha \mathbb{E}_{s \sim \mathcal{D}, a \sim \mu(a|s)} [Q(s, a)]}_{\text{Penalizing term}} + \frac{1}{2} \mathbb{E}_{s, a \sim \mathcal{D}} \left[ \left( Q(s, a) - \hat{B}^\pi \hat{Q}^k(s, a) \right)^2 \right]$$

Penalizing term

- When  $\mu(a|s) = \pi_\beta(a|s)$ , the resulting Q-function  $\lim_{k \rightarrow \infty} \hat{Q}_k(s, a)$  lower-bounds  $Q^\pi$  at all  $s \in \mathcal{D}, a \in \mathcal{A}$ .

**Theorem 3.1.** For any  $\mu(a|s)$  with  $\text{supp } \mu \subset \text{supp } \hat{\pi}_\beta$ , with probability  $\geq 1 - \delta$ ,  $\hat{Q}^\pi$  (the Q-function obtained by iterating Equation [1]) satisfies:

$$\forall s \in \mathcal{D}, a, \quad \hat{Q}^\pi(s, a) \leq Q^\pi(s, a) - \alpha \left[ (I - \gamma P^\pi)^{-1} \frac{\mu}{\hat{\pi}_\beta} \right](s, a) + \left[ (I - \gamma P^\pi)^{-1} \frac{C_{r, T, \delta} R_{\max}}{(1 - \gamma)\sqrt{|\mathcal{D}|}} \right](s, a).$$

Thus, if  $\alpha$  is sufficiently large, then  $\hat{Q}^\pi(s, a) \leq Q^\pi(s, a), \forall s \in \mathcal{D}, a$ . When  $\hat{B}^\pi = B^\pi$ , any  $\alpha > 0$  guarantees  $\hat{Q}^\pi(s, a) \leq Q^\pi(s, a), \forall s \in \mathcal{D}, a \in \mathcal{A}$ .

- To make  $\hat{Q}^\pi$  be a tighter lower-bound, we introduce an additional Q-value maximization term, under the data distribution.

$$\hat{Q}^{k+1} \leftarrow \arg \min_Q \alpha \cdot (\mathbb{E}_{\mathbf{s} \sim \mathcal{D}, \mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} [Q(\mathbf{s}, \mathbf{a})] - \mathbb{E}_{\mathbf{s} \sim \mathcal{D}, \mathbf{a} \sim \hat{\pi}_\beta(\mathbf{a}|\mathbf{s})} [Q(\mathbf{s}, \mathbf{a})]) + \frac{1}{2} \mathbb{E}_{\mathbf{s}, \mathbf{a}, \mathbf{s}' \sim \mathcal{D}} \left[ \left( Q(\mathbf{s}, \mathbf{a}) - \hat{B}^\pi \hat{Q}^k(\mathbf{s}, \mathbf{a}) \right)^2 \right]. \quad (2)$$

- Then,  $\lim_{k \rightarrow \infty} \hat{Q}_k(s, a)$  is not pointwise lower-bound of  $Q^\pi$ , however,  $\hat{V}^\pi(s)$  becomes a lower-bound of  $V^\pi(s)$

**Theorem 3.2** (Equation 2 results in a tighter lower bound). *The value of the policy under the Q-function from Equation 2,  $\hat{V}^\pi(s) = \mathbb{E}_{\pi(\mathbf{a}|\mathbf{s})}[\hat{Q}^\pi(s, \mathbf{a})]$ , lower-bounds the true value of the policy obtained via exact policy evaluation,  $V^\pi(s) = \mathbb{E}_{\pi(\mathbf{a}|\mathbf{s})}[Q^\pi(s, \mathbf{a})]$ , when  $\mu = \pi$ , according to:*

$$\forall \mathbf{s} \in \mathcal{D}, \hat{V}^\pi(\mathbf{s}) \leq V^\pi(\mathbf{s}) - \alpha \left[ (I - \gamma P^\pi)^{-1} \mathbb{E}_\pi \left[ \frac{\pi}{\hat{\pi}_\beta} - 1 \right] \right] (\mathbf{s}) + \left[ (I - \gamma P^\pi)^{-1} \frac{C_{r,T,\delta} R_{\max}}{(1-\gamma)\sqrt{|\mathcal{D}|}} \right] (\mathbf{s}).$$

Thus, if  $\alpha > \frac{C_{r,T} R_{\max}}{1-\gamma} \cdot \max_{\mathbf{s} \in \mathcal{D}} \frac{1}{\sqrt{|\mathcal{D}(\mathbf{s})|}} \cdot \left[ \sum_{\mathbf{a}} \pi(\mathbf{a}|\mathbf{s}) \left( \frac{\pi(\mathbf{a}|\mathbf{s})}{\hat{\pi}_\beta(\mathbf{a}|\mathbf{s})} - 1 \right) \right]^{-1}$ ,  $\forall \mathbf{s} \in \mathcal{D}, \hat{V}^\pi(\mathbf{s}) \leq V^\pi(\mathbf{s})$ , with probability  $\geq 1 - \delta$ . When  $\hat{B}^\pi = B^\pi$ , then any  $\alpha > 0$  guarantees  $\hat{V}^\pi(\mathbf{s}) \leq V^\pi(\mathbf{s}), \forall \mathbf{s} \in \mathcal{D}$ .

- Now we can define CQL problem as,

$$\min_Q \max_{\mu} \alpha \left( \mathbb{E}_{\mathbf{s} \sim \mathcal{D}, \mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} [Q(\mathbf{s}, \mathbf{a})] - \mathbb{E}_{\mathbf{s} \sim \mathcal{D}, \mathbf{a} \sim \hat{\pi}_{\beta}(\mathbf{a}|\mathbf{s})} [Q(\mathbf{s}, \mathbf{a})] \right) + \frac{1}{2} \mathbb{E}_{\mathbf{s}, \mathbf{a}, \mathbf{s}' \sim \mathcal{D}} \left[ \left( Q(\mathbf{s}, \mathbf{a}) - \hat{B}^{\pi_k} \hat{Q}^k(\mathbf{s}, \mathbf{a}) \right)^2 \right] + \mathcal{R}(\mu) \quad (\text{CQL}(\mathcal{R})).$$

where  $R(\mu)$  is a policy regularizer.

- For example, for maximum entropy regularization,  $R(\mu) = -D_{KL}(\mu | \text{Unif}(a))$ , CQL becomes

$$\min_Q \alpha \mathbb{E}_{\mathbf{s} \sim \mathcal{D}} \left[ \log \sum_{\mathbf{a}} \exp(Q(\mathbf{s}, \mathbf{a})) - \mathbb{E}_{\mathbf{a} \sim \hat{\pi}_{\beta}(\mathbf{a}|\mathbf{s})} [Q(\mathbf{s}, \mathbf{a})] \right] + \frac{1}{2} \mathbb{E}_{\mathbf{s}, \mathbf{a}, \mathbf{s}' \sim \mathcal{D}} \left[ \left( Q - \hat{B}^{\pi_k} \hat{Q}^k \right)^2 \right].$$

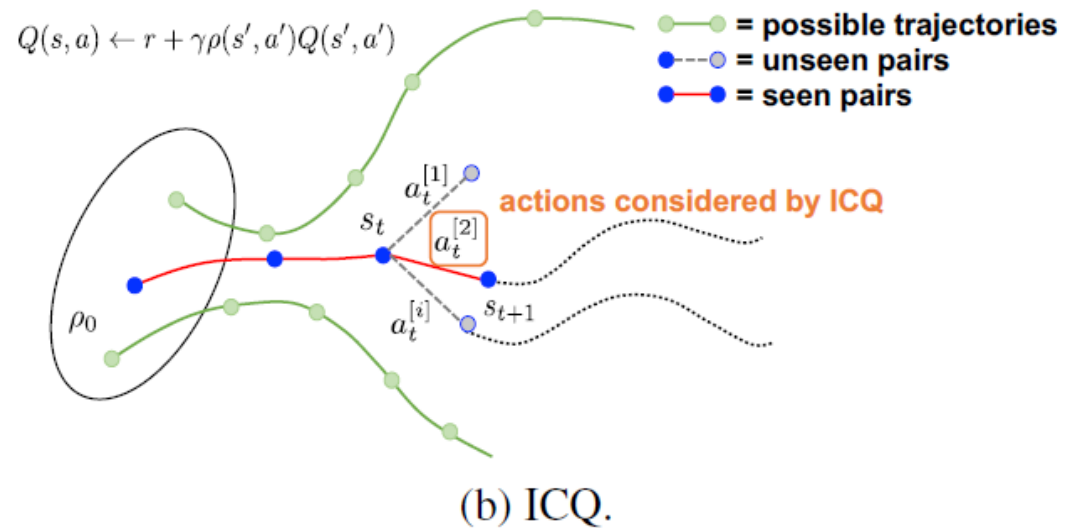
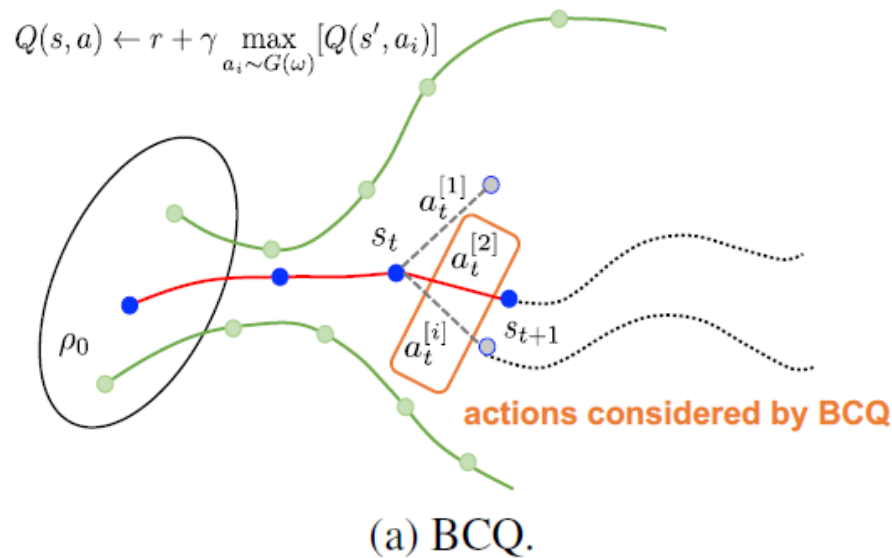
Yang, Yiqin, Xiaoteng Ma, Li Chenghao, Zewu Zheng, Qiyuan Zhang, Gao Huang, Jun Yang, and Qianchuan Zhao. "**Believe what you see: Implicit constraint approach for offline multi-agent reinforcement learning.**" NeurIPS 2021.

## **BELIEVE WHAT YOU SEE: IMPLICIT CONSTRAINT APPROACH FOR OFFLINE MULTI-AGENT REINFORCEMENT LEARNING**

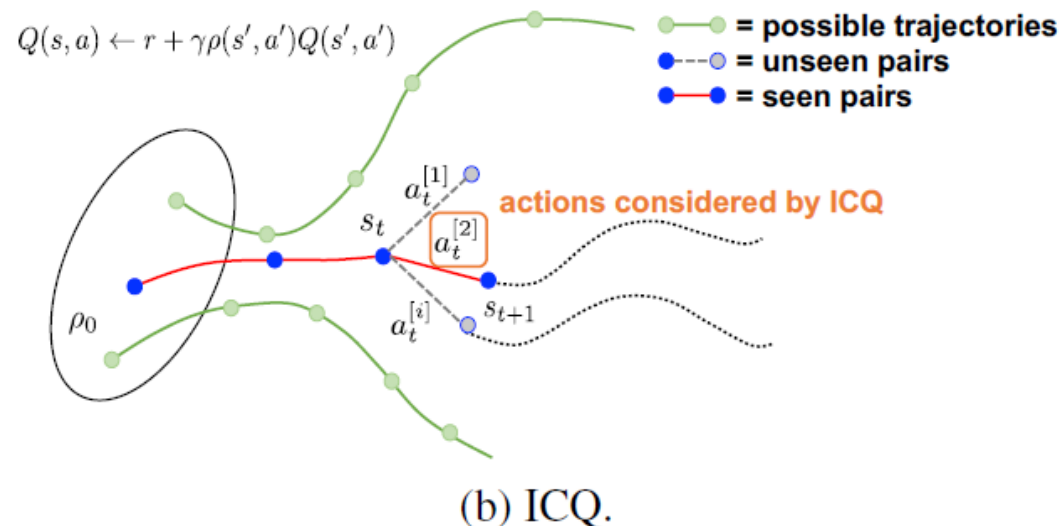
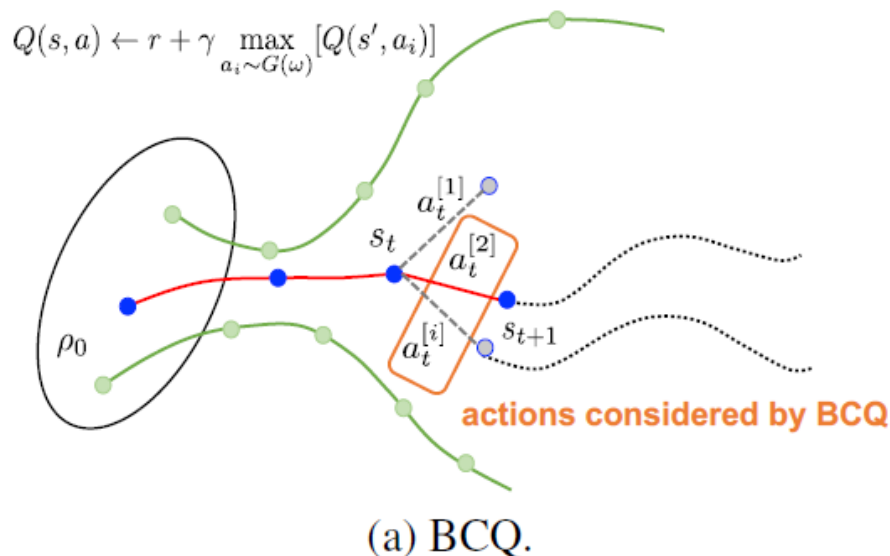
- *Extrapolation error* is the greatest obstacle of offline RL.

$$Q(s, a) \leftarrow r + \gamma \mathbb{E}_{a' \sim \pi} [Q(s', a')]$$

Can be inaccurate for unseen  $a'$



- Implicit constraint Q-learning (ICQ) aims to estimate Q-value using only the *seen* state-action pairs.



- Implicit constraint Q-learning (ICQ) utilizes only the seen state-action pairs using *importance sampling*.

$$(\mathcal{T}^\pi Q)(\tau, a) = Q(\tau, a) + \mathbb{E}_{\tau'} [r + \gamma \mathbb{E}_{a' \sim \mu} [\rho(\tau', a') Q(\tau', a')]] - Q(\tau, a),$$

Where,  $\rho(\tau', a') = \frac{\pi(a'|\tau')}{\mu(a'|\tau')}$  ( $\mu$ : behavior policy used for collecting data)

- How to estimate  $\rho(\tau, a)$
- Policy optimization with the behavior regularized constraint.

$$\pi_{k+1} = \arg \max \mathbb{E}_{a \sim \pi(\cdot | \tau)} [Q^{\pi_k}(\tau, a)], \quad \text{s.t.} \quad D_{\text{KL}}(\pi \parallel \mu)[\tau] \leq \epsilon.$$

- Then the optimal policy  $\pi^*$  becomes,

$$\pi_{k+1}^*(a | \tau) = \frac{1}{Z(\tau)} \mu(a | \tau) \exp \left( \frac{Q^{\pi_k}(\tau, a)}{\alpha} \right)$$

- And we can estimate  $\rho(\tau, a)$  as,

$$\rho(\tau, a) = \frac{\pi_{k+1}^*(a | \tau)}{\mu(a | \tau)} = \frac{1}{Z(\tau)} \exp \left( \frac{Q^{\pi_k}(\tau, a)}{\alpha} \right)$$

- Implicit Constraint Q-learning operator:

$$\mathcal{T}_{\text{ICQ}} Q(\tau, a) = r + \gamma \mathbb{E}_{a' \sim \mu} \left[ \frac{1}{Z(\tau')} \exp \left( \frac{Q(\tau', a')}{\alpha} \right) Q(\tau', a') \right]$$

- CTDE paradigm : Centralized Training and Decentralized Execution
- Value decomposition assumption:
  - Global q-value is expressed as a linear combination of the q-values of each agent.

$$Q^\pi(\tau, \mathbf{a}) = \sum_i w^i(\tau) Q^i(\tau^i, a^i) + b(\tau),$$

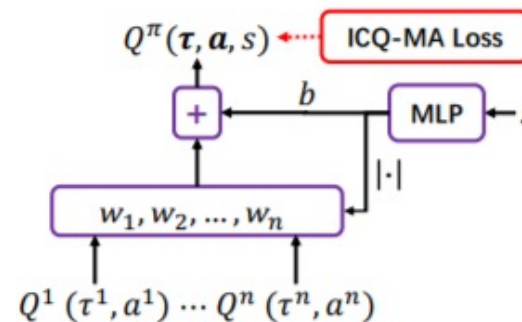


Figure 3: Mixer Network.

- Then, the joint policy can be expressed as a combination of the optimal policies of the all agent (for  $i=1,2,\dots,n$ )

**Theorem 3.** Assuming the joint action-value function is linearly decomposed, we can decompose the multi-agent joint-policy under implicit constraint as follows

$$\pi = \arg \max_{\pi^1, \dots, \pi^n} \sum_i \mathbb{E}_{\tau^i, a^i \sim \mathcal{B}} \left[ \frac{1}{Z^i(\tau^i)} \log(\pi^i(a^i | \tau^i)) \exp \left( \frac{w^i(\tau) Q^i(\tau^i, a^i)}{\alpha} \right) \right], \quad (51)$$

where  $Z^i(\tau^i) = \sum_{\tilde{a}^i} \mu^i(\tilde{a}^i | \tau^i) \exp \left( \frac{1}{\alpha} w^i(\tau) Q^i(\tau^i, \tilde{a}^i) \right)$  is the normalizing partition function.

- Q-function objective:

$$\mathcal{J}_Q(\phi, \psi) = \mathbb{E}_{\mathcal{B}} \left[ \sum_{t \geq 0} (\gamma \lambda)^t \left( r_t + \gamma \frac{1}{Z(\boldsymbol{\tau}_{t+1})} \exp \left( \frac{Q(\boldsymbol{\tau}_{t+1}, \mathbf{a}_{t+1})}{\alpha} \right) Q(\boldsymbol{\tau}_{t+1}, \mathbf{a}_{t+1}) - Q(\boldsymbol{\tau}_t, \mathbf{a}_t) \right) \right]^2 \quad (19)$$

where  $Q(\boldsymbol{\tau}_{t+1}, \mathbf{a}_{t+1}) = \sum_i w^i(\boldsymbol{\tau}_{t+1}; \psi') Q^i(\tau_{t+1}^i, a_{t+1}^i; \phi'_i) - b(\boldsymbol{\tau}_{t+1}; \psi')$ .

- Policy gradient:

$$\mathcal{J}_{\pi}(\theta) = \sum_i \mathbb{E}_{\tau^i, a^i \sim \mathcal{B}} \left[ -\frac{1}{Z^i(\tau^i)} \log(\pi^i(a^i | \tau^i; \theta_i)) \exp \left( \frac{w^i(\boldsymbol{\tau}) Q^i(\tau^i, a^i)}{\alpha} \right) \right]$$

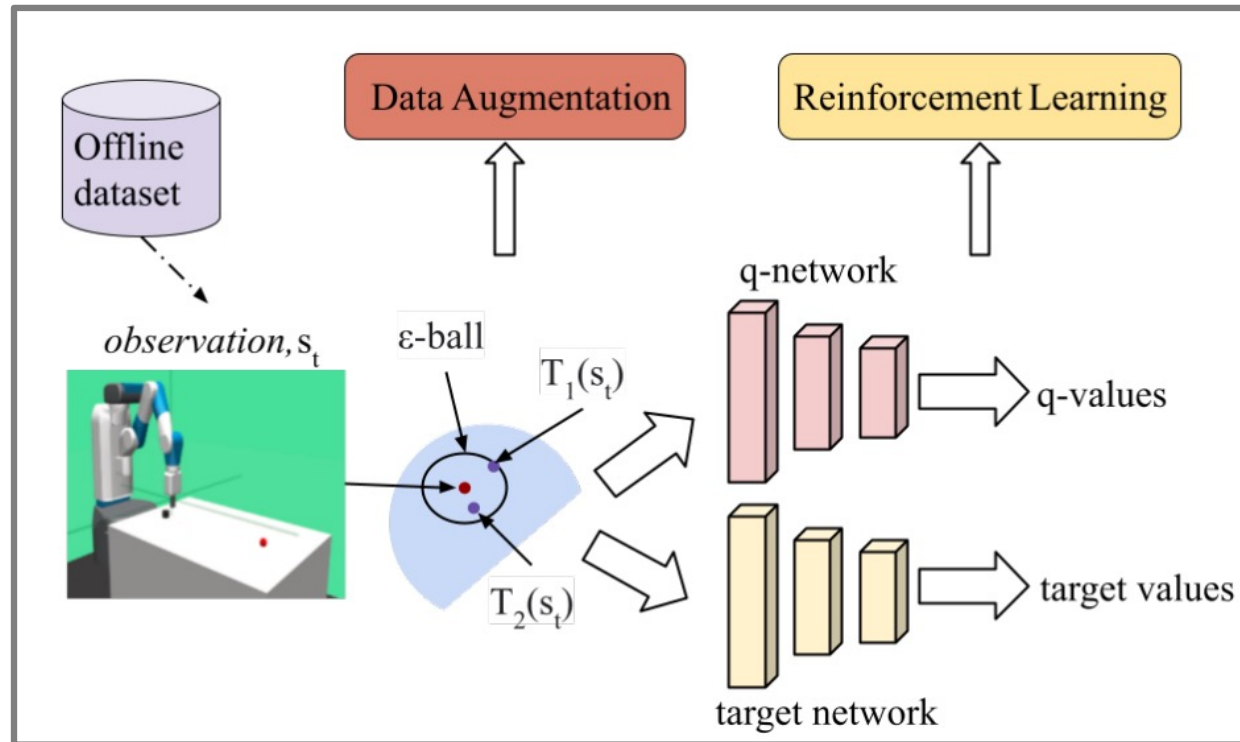
- Multi-agent value estimation with  $\lambda$ -return

$$(\mathcal{T}_{\text{ICQ}}^{\lambda} Q)(\boldsymbol{\tau}, \mathbf{a}) \triangleq Q(\boldsymbol{\tau}, \mathbf{a}) + \mathbb{E}_{\mu} \left[ \sum_{t \geq 0} (\gamma \lambda)^t (r_t + \gamma \rho(\boldsymbol{\tau}_{t+1}, \mathbf{a}_{t+1}) Q(\boldsymbol{\tau}_{t+1}, \mathbf{a}_{t+1}) - Q(\boldsymbol{\tau}_t, \mathbf{a}_t)) \right]$$

Sinha, Samarth, Ajay Mandlekar, and Animesh Garg. "**S4RL: Surprisingly simple self-supervision for offline reinforcement learning in robotics.**" Conference on Robot Learning (CoRL), 2021.

## **S4RL: SURPRISINGLY SIMPLE SELF-SUPERVISION FOR OFFLINE REINFORCEMENT LEARNING IN ROBOTICS**

- Existing offline RL methods mainly focus on reducing **Q-value estimation errors due to out-of-distribution** from offline data.
- However, due to the use of parameterized neural networks, **approximation errors** are also introduced.
- In this paper, a **data augmentation method** for offline RL is proposed to reduce the approximation errors
- “Existing offline RL methods + proposed data augmentation” outperforms baselines in the D4RL experiments.



## Overview of the proposed method

- A data augmentation transformation is denoted by:  $\mathcal{T}(\tilde{s}_t | s_t)$  where  $s_t \sim \mathcal{D}$ .
- The authors propose seven different data augmentation transformations and compare the results of each transformation method.

- Proposed Q-learning objective:

$$\min_Q \mathbb{E}_{s_t, a_t \sim \mathcal{D}} \left[ r_t + \gamma \frac{1}{i} \sum_i Q(\mathcal{T}_i(\tilde{s}_{t+1}|s_{t+1}), a_{t+1}) - \frac{1}{i} \sum_i Q(\mathcal{T}_i(\tilde{s}_t|s_t), a_t) \right]$$

- Q values and target values are used as the average of each augmented state value.
  - In the experiments, however, use only one augmentation transformation.
- Proposed candidates for augmentation transformation:
    - zero-mean Gaussian noise**  $\tilde{s}_t \leftarrow s_t + \epsilon$ , where  $\epsilon \in \mathcal{N}(0, \sigma I)$
    - zero-mean Uniform noise**  $\tilde{s}_t \leftarrow s_t + \epsilon$ , where  $\epsilon \in \mathcal{U}(-\alpha, \alpha)$
    - random amplitude scaling**  $\tilde{s}_t \leftarrow s_t * \epsilon$ , where  $\epsilon \in \mathcal{U}(\alpha, \beta)$
    - dimension dropout**  $\tilde{s}_t \leftarrow s_t \cdot \mathbf{1}$ ,  
where  $\mathbf{1}$  is a vector of 1s with one 0 randomly sampled from a Bernoulli
    - state-switch** where we flip the value of 2 randomly selected dimensions
    - state mix-up**  $\tilde{s}_t \leftarrow \lambda s_t + (1 - \lambda)s_{t+1}$ , where  $\lambda \sim \text{Beta}(\alpha, \alpha)$
    - adversarial state training**  $\tilde{s}_t \leftarrow s_t + \epsilon \nabla_{s_t} \mathbb{J}_Q(\pi(s_t))$

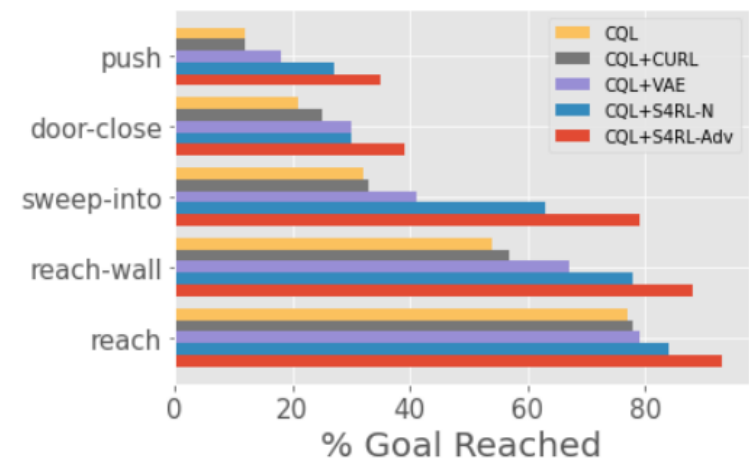
- Existing offline RL methods:
  - Conservative Q-learning (CQL)
  - Behavior regularized actor-critic (BRAC)
- Proposed methods: “[offline RL] + S4RL (transformation method)”  
(“CQL+S4RL (Adv)”: update CQL with the adversarial state training transformation)
- Baseline augmentation (self-supervision) methods:
  - Contrastive unsupervised reinforcement learning (CURL): use another data augmentation transformation.
  - Variational auto-encoder (VAE): feed the augmented states to the auto-encoder to learn generalizable representations.
- Offline Dataset:
  - Use D4RL dataset.  
(“cheetah-medium”: environment – cheetah, data – collect using a medium-level policy)

# Experiment: results

Task Name	Normal	CQL						
		+S4RL (N)	+S4RL (L)	+S4RL (Amp-Scale)	+S4RL (Dim-Drop)	+S4RL (State-Switch)	+S4RL (MixUp)	+S4RL (Adv)
cheetah-random	35.4	<b>52.3</b>	<b>50.3</b>	46.4	45.5	40.3	45.2	<b>53.9</b>
cheetah-medium	44.4	<b>48.8</b>	<b>47.0</b>	42.5	45.6	41.2	46.2	<b>48.6</b>
cheetah-medium-replay	42.0	<b>51.4</b>	<b>50.8</b>	43.1	49.6	41.0	46.2	<b>51.7</b>
cheetah-medium-expert	62.4	<b>79.0</b>	<b>78.5</b>	71.2	66.1	68.3	73.1	<b>78.1</b>
hopper-random	<b>10.8</b>	<b>10.8</b>	10.7	<b>10.8</b>	9.3	9.5	<b>11.0</b>	10.7
hopper-medium	58.0	78.9	<b>80.6</b>	60.3	45.3	54.6	<b>79.2</b>	<b>81.3</b>
hopper-medium-replay	29.5	<b>35.4</b>	35.2	28.6	20.4	20.9	<b>35.6</b>	<b>36.8</b>
hopper-medium-expert	111.0	<b>115.2</b>	112.3	102.6	98.3	108.6	<b>113.5</b>	<b>117.9</b>
walker-random	7.0	<b>24.9</b>	<b>20.5</b>	6.9	3.2	9.9	10.2	<b>25.1</b>
walker-medium	79.2	<b>93.6</b>	<b>93.9</b>	81.6	65.3	65.3	88.6	<b>93.1</b>
walker-medium-replay	21.1	<b>30.3</b>	<b>31.9</b>	25.1	8.5	9.3	27.6	<b>35.0</b>
walker-medium-expert	98.7	<b>112.2</b>	104.1	100.2	86.8	103.8	<b>104.3</b>	<b>107.1</b>
<b>average-score</b>	49.96	<b>60.57</b>	<b>59.65</b>	51.60	45.31	47.73	56.73	<b>61.6</b>
<b>average-ranking</b>	5.83	<b>1.92</b>	<b>3.00</b>	5.00	6.42	6.25	3.17	<b>1.83</b>

→ show that adversarial state is the best transformation method.

Domain	Task Name	Normal	CQL					BRAC-v			
			+CURL (N)	+VAE (N)	+S4RL (N)	+S4RL (Mix-Up)	+S4RL (Adv)	Normal	+S4RL (N)	+S4RL (Mix-Up)	+S4RL (Adv)
AntMaze	antmaze-umaze	74.0	70.3	86.1	<b>91.3</b>	86.7	<b>94.1</b>	70.0	76.5	66.4	<b>80.3</b>
	antmaze-umaze-diverse	84.0	81.4	85.1	<b>87.8</b>	<b>86.9</b>	<b>88.0</b>	70.0	<b>81.3</b>	65.4	<b>80.9</b>
	antmaze-medium-play	<b>61.2</b>	<b>60.9</b>	<b>62.1</b>	<b>61.9</b>	<b>61.3</b>	<b>61.6</b>	0.0	0.0	0.0	0.0
	antmaze-medium-diverse	53.7	45.4	59.9	78.1	62.3	<b>82.3</b>	0.0	0.0	0.0	0.0
	antmaze-large-play	15.8	12.3	15.2	<b>24.4</b>	16.7	<b>25.1</b>	0.0	0.0	0.0	0.0
	antmaze-large-diverse	14.9	8.3	17.3	<b>27.0</b>	15.9	<b>26.2</b>	0.0	0.0	0.0	0.0
Gym	cheetah-random	35.4	43.1	35.7	<b>52.3</b>	45.2	<b>53.9</b>	31.2	<b>35.6</b>	34.3	<b>36.1</b>
	cheetah-medium	44.4	44.8	45.6	<b>48.8</b>	46.2	48.6	46.3	<b>49.1</b>	46.1	46.0
	cheetah-medium-replay	42.0	36.5	41.9	51.4	46.2	<b>51.7</b>	<b>47.7</b>	<b>47.9</b>	<b>47.0</b>	<b>47.9</b>
	cheetah-medium-expert	62.4	65.6	70.7	<b>79.0</b>	73.1	78.1	41.9	<b>52.1</b>	46.7	<b>53.4</b>
	hopper-random	<b>10.8</b>	<b>10.8</b>	<b>10.8</b>	<b>10.8</b>	<b>11.0</b>	10.7	<b>12.2</b>	<b>12.9</b>	10.9	<b>12.1</b>
	hopper-medium	58.0	61.9	66.4	<b>78.9</b>	<b>79.2</b>	<b>81.3</b>	31.1	45.1	41.9	<b>48.0</b>
	hopper-medium-replay	29.5	30.1	<b>34.8</b>	<b>35.4</b>	<b>35.6</b>	<b>36.8</b>	<b>0.6</b>	<b>1.1</b>	<b>0.9</b>	<b>0.4</b>
	hopper-medium-expert	<b>111.0</b>	<b>109.7</b>	<b>114.3</b>	<b>113.5</b>	<b>112.3</b>	117.9	0.8	<b>59.0</b>	49.5	55.3
	walker-random	7.0	6.6	6.5	<b>24.9</b>	10.2	<b>25.1</b>	1.9	9.5	2.4	<b>14.1</b>
	walker-medium	79.2	81.3	80.8	<b>93.6</b>	88.6	<b>93.1</b>	81.1	<b>85.4</b>	80.9	<b>86.7</b>
	walker-medium-replay	21.1	24.5	26.4	30.3	27.6	<b>35.0</b>	<b>0.9</b>	<b>1.1</b>	<b>0.9</b>	<b>1.1</b>
walker-medium-expert	98.7	103.0	99.6	<b>112.2</b>	104.3	107.1	81.6	90.4	89.7	<b>94.5</b>	
Adroit	pen-human	37.5	33.4	39.7	<b>44.4</b>	41.6	<b>51.2</b>	0.6	<b>11.3</b>	0.8	<b>14.1</b>
	pen-cloned	39.2	40.1	41.3	<b>57.1</b>	44.5	<b>58.2</b>	-2.5	<b>10.6</b>	0.5	<b>9.9</b>
	hammer-human	4.4	<b>6.1</b>	<b>6.0</b>	<b>5.9</b>	<b>6.3</b>	<b>6.3</b>	0.2	<b>3.2</b>	<b>3.3</b>	<b>5.4</b>
	hammer-cloned	2.1	1.9	2.1	<b>2.7</b>	<b>2.4</b>	<b>2.9</b>	0.3	0.1	0.3	1.2
	door-human	9.9	10.3	14.3	27.0	16.7	<b>35.3</b>	-0.3	<b>2.8</b>	0.2	4.2
	door-cloned	0.4	0.4	0.7	<b>2.1</b>	0.3	0.8	-0.1	<b>0.3</b>	-0.1	<b>0.5</b>
	relocate-human	<b>0.2</b>	<b>0.2</b>	<b>0.2</b>	<b>0.2</b>	<b>0.2</b>	<b>0.2</b>	-0.3	<b>0.1</b>	<b>0.1</b>	<b>0.2</b>
	relocate-cloned	<b>-0.1</b>	<b>-0.1</b>	-0.3	<b>-0.1</b>	<b>-0.1</b>	<b>-0.1</b>	-0.3	<b>-0.1</b>	<b>-0.1</b>	<b>-0.1</b>
Franka	kitchen-complete	43.8	54.4	42.1	77.1	60.4	<b>88.1</b>	0.0	0.0	0.0	<b>12.3</b>
	kitchen-partial	49.8	59.6	49.9	74.8	59.0	<b>83.6</b>	0.0	0.0	0.0	<b>13.9</b>

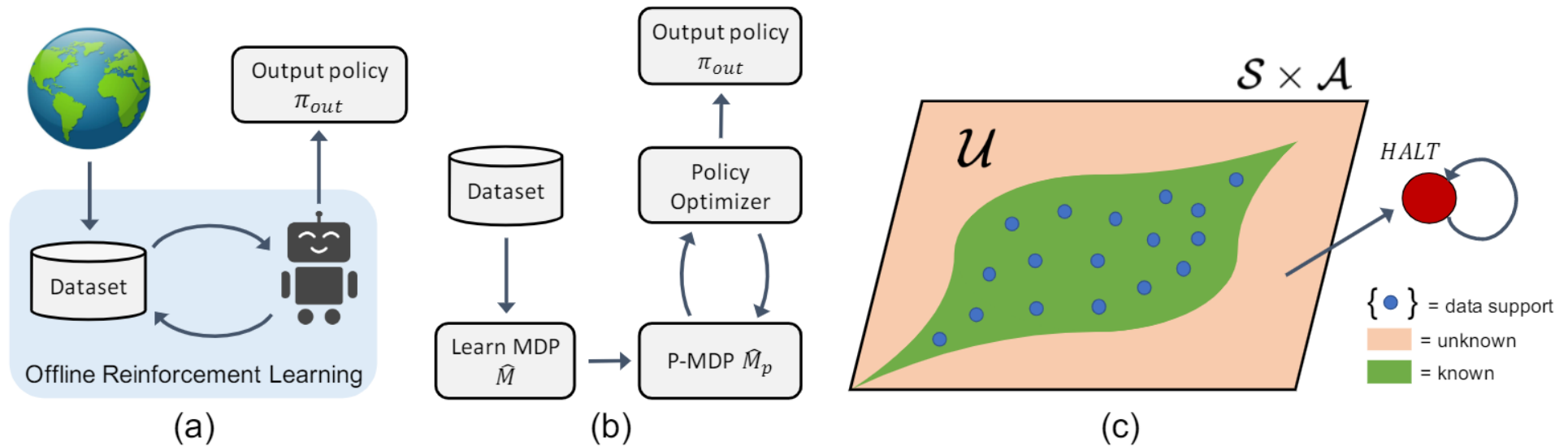


→ CQL+S4RL-Adv outperforms all other baselines.

Rahul Kidambi, Aravind Rajeswaran, Praneeth Netrapalli, and Thorsten Joachims.  
"MOReL: Model-based offline reinforcement learning." NeurIPS 2020.

# **MOReL: MODEL-BASED OFFLINE REINFORCEMENT LEARNING**

- Prior works in offline RL have been confined to model-free RL approaches.
- Model-based RL (MBRL) algorithms are highly sample efficient, but the learned model is unlikely to be globally accurate due to lack of offline data.
- In this paper, a pessimistic MDP (P-MDP) is proposed, where a policy is penalized when entering “unknown” regions.
- Theoretically, a policy trained in P-MDP has the upper bound for the sub-optimality.
- The proposed method, MOREL, obtains SOTA results in 12 out of 20 environments.



Overview of the proposed method ((a), (b): offline RL, (c): P-MDP)

## Definition 1

- Learned MDP:  $\hat{\mathcal{M}} = \{S, A, r, \hat{P}, \hat{\rho}_0, \gamma\}$
- Unknown state-action detector (USAD):

$$U^\alpha(s, a) = \begin{cases} FALSE \text{ (i.e. Known)} & \text{if } D_{TV}(\hat{P}(\cdot|s, a), P(\cdot|s, a)) \leq \alpha \text{ can be guaranteed} \\ TRUE \text{ (i.e. Unknown)} & \text{otherwise} \end{cases}$$

## Definition 2

- Pessimistic-MDP (P-MDP):  $\hat{\mathcal{M}}_p := \{S \cup \text{HALT}, A, r_p, \hat{P}_p, \hat{\rho}_0, \gamma\}$

$$\hat{P}_p(s'|s, a) = \begin{cases} \delta(s' = \text{HALT}) & \text{if } U^\alpha(s, a) = \text{TRUE} \\ & \text{or } s = \text{HALT} \\ \hat{P}(s'|s, a) & \text{otherwise} \end{cases} \quad r_p(s, a) = \begin{cases} -\kappa & \text{if } s = \text{HALT} \\ r(s, a) & \text{otherwise} \end{cases}$$

## Implementation

- Dynamic model learning:

$\hat{P}(\cdot|s, a) \equiv \mathcal{N}(f_\phi(s, a), \Sigma)$ , with  $f_\phi(s, a) = s + \sigma_\Delta \text{MLP}_\phi((s - \mu_s)/\sigma_s, (a - \mu_a)/\sigma_a)$ ,  
 $\mu_s, \sigma_s, \mu_a, \sigma_a$  are the mean and standard deviations of states/actions in  $\mathcal{D}$ ,  
 $\sigma_\Delta$  is the standard deviation of state differences, i.e.  $\Delta = s' - s, (s, s') \in \mathcal{D}$ .

- Unknown state-action detector (USAD):

Train an ensemble dynamic model and define the disc function as:

$$\text{disc}(s, a) = \max_{i,j} \|f_{\phi_i}(s, a) - f_{\phi_j}(s, a)\|_2$$

Then  $U_{\text{practical}}(s, a) = \begin{cases} \text{FALSE (i.e. Known)} & \text{if } \text{disc}(s, a) \leq \text{threshold} \\ \text{TRUE (i.e. Unknown)} & \text{if } \text{disc}(s, a) > \text{threshold} \end{cases}$

- Policy learning: use natural policy gradient (NPG)

## Definition 3

- Hitting time  $T_{(s,a)}^\pi$ : the first time action  $a$  is taken at state  $s$  by  $\pi$ .  

$$\rightarrow T_S^\pi \stackrel{\text{def}}{=} \min_{(s,a) \in \mathcal{S}} T_{(s,a)}^\pi$$

**Theorem 1.** (Policy value with pessimism) The value of any policy  $\pi$  on the original MDP  $\mathcal{M}$  and its  $(\alpha, R_{\max})$ -pessimistic MDP  $\hat{\mathcal{M}}_p$  satisfies:

$$J_{\hat{\rho}_0}(\pi, \hat{\mathcal{M}}_p) \geq J_{\rho_0}(\pi, \mathcal{M}) - \frac{2R_{\max}}{1-\gamma} \cdot D_{TV}(\rho_0, \hat{\rho}_0) - \frac{2\gamma R_{\max}}{(1-\gamma)^2} \cdot \alpha - \frac{2R_{\max}}{1-\gamma} \cdot \mathbb{E} \left[ \gamma^{T_{\mathcal{U}}^\pi} \right], \text{ and}$$

$$J_{\hat{\rho}_0}(\pi, \hat{\mathcal{M}}_p) \leq J_{\rho_0}(\pi, \mathcal{M}) + \frac{2R_{\max}}{1-\gamma} \cdot D_{TV}(\rho_0, \hat{\rho}_0) + \frac{2\gamma R_{\max}}{(1-\gamma)^2} \cdot \alpha,$$

where  $T_{\mathcal{U}}^\pi$  denotes the hitting time of unknown states  $\mathcal{U} \stackrel{\text{def}}{=} \{(s, a) : U^\alpha(s, a) = \text{TRUE}\}$  by  $\pi$  on  $\mathcal{M}$ .

**Corollary 3.** (Upper bound) Suppose the dataset  $\mathcal{D}$  is large enough so that  $\alpha = D_{TV}(\rho_0, \hat{\rho}_0) = 0$ . Then, the output  $\pi_{out}$  of Algorithm 1 satisfies:

$$J_{\rho_0}(\pi^*, \mathcal{M}) - J_{\rho_0}(\pi_{out}, \mathcal{M}) \leq \epsilon_\pi + \frac{2R_{\max}}{1-\gamma} \cdot \mathbb{E} \left[ \gamma^{T_{\mathcal{U}}^{\pi^*}} \right] \leq \epsilon_\pi + \frac{2R_{\max}}{(1-\gamma)^2} \cdot d^{\pi^*, \mathcal{M}}(\mathcal{U})$$

→ give an **upper bound of the sub-optimality gap** between the optimal policy in the true MDP and P-MDP.

- Environments:
  - Hopper, Half-Cheetah, Walker-2d, and Ant in MuJoCo.
- Dataset:
  - First train a policy  $\pi_p$  to obtain reward sums 1000, 4000, 1000, and 1000 for the four environments.
  - Pure: collecting data using  $\pi_p$ .
  - `eps-1`, `eps-3`, `gauss-1`, `gauss-3` : collecting data using  $\pi_p$  with noise.
- Baselines:
  - Batch-constrained Q learning (BCQ)
  - Bootstrapping error accumulation reduction (BEAR)
  - Behaviour regularized actor-critic (BRAC)

# Experiment: results

Environment: Ant-v2					
Algorithm	BCQ [15]	BEAR [16]	BRAC [18]	Best Baseline	MOReL (Ours)
Pure	1921	2100	<u>2839</u>	2839	<b>3663±247</b>
Eps-1	1864	1897	<u>2672</u>	2672	<b>3305±413</b>
Eps-3	1504	2008	<u>2602</u>	2602	<b>3008±231</b>
Gauss-1	1731	2054	<u>2667</u>	2667	<b>3329±270</b>
Gauss-3	1887	2018	2640	2661	<b>3693±33</b>

Environment: Hopper-v2					
Algorithm	BCQ [15]	BEAR [16]	BRAC [18]	Best Baseline	MOReL (Ours)
Pure	1543	0	2291	2774	<b>3642±54</b>
Eps-1	1652	1620	2282	2360	<b>3724±46</b>
Eps-3	1632	2213	1892	2892	<b>3535±91</b>
Gauss-1	1599	1825	<u>2255</u>	2255	<b>3653±52</b>
Gauss-3	1590	1720	1458	2097	<b>3648±148</b>

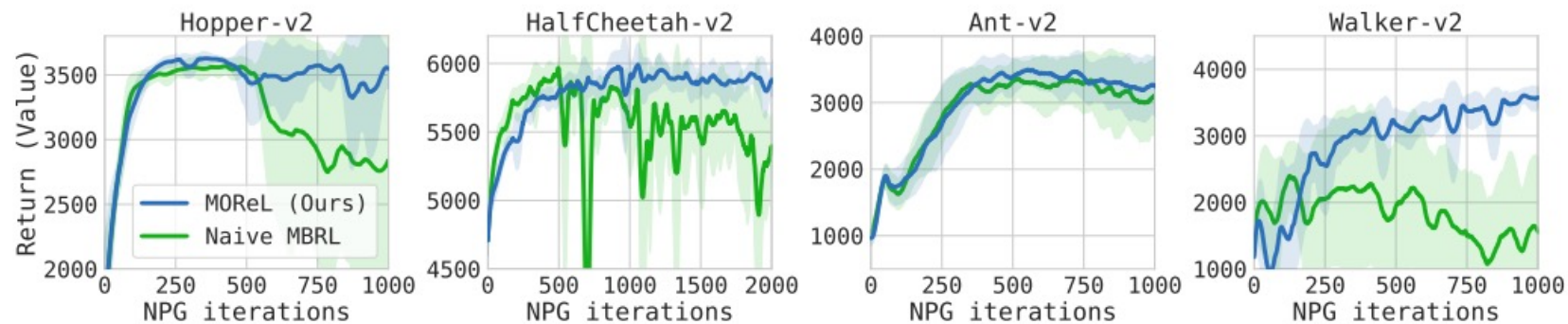
  

Environment: HalfCheetah-v2					
Algorithm	BCQ [15]	BEAR [16]	BRAC [18]	Best Baseline	MOReL (Ours)
Pure	5064	5325	6207	<b>6209</b>	<b>6028±192</b>
Eps-1	5693	5435	<b>6307</b>	<b>6307</b>	5861±192
Eps-3	5588	5149	<u>6263</u>	<b>6359</b>	5869±139
Gauss-1	5614	5394	<b>6323</b>	<b>6323</b>	6026±74
Gauss-3	5837	5329	<b>6400</b>	<b>6400</b>	5892±128

Environment: Walker-v2					
Algorithm	BCQ [15]	BEAR [16]	BRAC [18]	Best Baseline	MOReL (Ours)
Pure	2095	2646	2694	2907	<b>3709±159</b>
Eps-1	1921	2695	3241	<b>3490</b>	2899±588
Eps-3	1953	2608	<b>3255</b>	<b>3255</b>	<b>3186±92</b>
Gauss-1	2094	2539	2893	3193	<b>4027±314</b>
Gauss-3	1734	2194	<b>3368</b>	<b>3368</b>	<b>2828±589</b>

→ MOReL shows the best performance in 12 of 20 environments and exceeds the performance of the data-collecting policy.



Results on MDP (naive) vs P-MDP (ours)

→ P-MDP is essential to prevent performance degradation.

Chen, Lili, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Misha Laskin, Pieter Abbeel, Aravind Srinivas, and Igor Mordatch. "**Decision transformer: Reinforcement learning via sequence modeling.**" NeurIPS 2021.

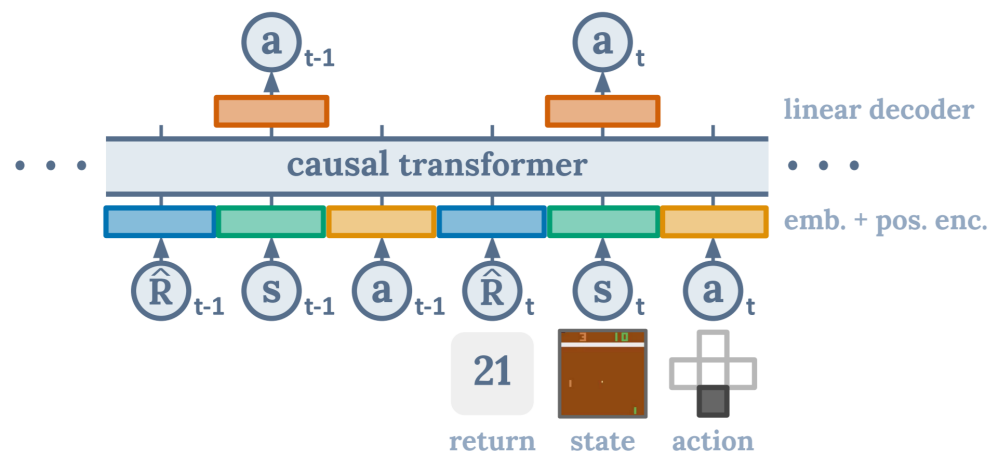
# **DECISION TRANSFORMER: REINFORCEMENT LEARNING VIA SEQUENCE MODELING**

- Instead of TD-learning, train transformer models on collected experience using a sequence modeling objective
- Autoregressive transformer model (GPT) is used.
- Advantages
  - Can bypass the need for bootstrapping for long term credit assignment.
  - Avoids the need for discounting future rewards.
  - Perform credit assignment directly via self-attention
  - No other methods are needed to solve the problems of offline RL (error propagation and value overestimation).

- Trajectory representation
  - Use returns-to-go to allow the model to generate actions based on future desired returns.

$$\text{returns-to-go } \hat{R}_t = \sum_{t'=t}^T r_{t'} \quad \tau = \left( \hat{R}_1, s_1, a_1, \hat{R}_2, s_2, a_2, \dots, \hat{R}_T, s_T, a_T \right)$$

- Architecture
  - Put tokens (return-to-go, state, action) of K time steps as input.
  - Use GPT model which predicts future action tokens via autoregressive modeling
  - The action of each time step is the output of the model.
- Training
  - Sample minibatches of sequence length K from the trajectory dataset
  - Use cross entropy (discrete) or mean-squared error (continuous) as loss.



---

**Algorithm 1** Decision Transformer Pseudocode (for continuous actions)

---

```
# R, s, a, t: returns-to-go, states, actions, or timesteps
# transformer: transformer with causal masking (GPT)
# embed_s, embed_a, embed_R: linear embedding layers
# embed_t: learned episode positional embedding
# pred_a: linear action prediction layer

# main model
def DecisionTransformer(R, s, a, t):
    # compute embeddings for tokens
    pos_embedding = embed_t(t) # per-timestep (note: not per-token)
    s_embedding = embed_s(s) + pos_embedding
    a_embedding = embed_a(a) + pos_embedding
    R_embedding = embed_R(R) + pos_embedding

    # interleave tokens as (R_1, s_1, a_1, ..., R_K, s_K)
    input_embeds = stack(R_embedding, s_embedding, a_embedding)

    # use transformer to get hidden states
    hidden_states = transformer(input_embeds=input_embeds)

    # select hidden states for action prediction tokens
    a_hidden = unstack(hidden_states).actions

    # predict action
    return pred_a(a_hidden)

# training loop
for (R, s, a, t) in dataloader: # dims: (batch_size, K, dim)
    a_preds = DecisionTransformer(R, s, a, t)
    loss = mean((a_preds - a)**2) # L2 loss for continuous actions
    optimizer.zero_grad(); loss.backward(); optimizer.step()

# evaluation loop
target_return = 1 # for instance, expert-level return
R, s, a, t, done = [target_return], [env.reset()], [], [1], False
while not done: # autoregressive generation/sampling
    # sample next action
    action = DecisionTransformer(R, s, a, t)[-1] # for cts actions
    new_s, r, done, _ = env.step(action)

    # append new tokens to sequence
    R = R + [R[-1] - r] # decrement returns-to-go with reward
    s, a, t = s + [new_s], a + [action], t + [len(R)]
    R, s, a, t = R[-K:], ... # only keep context length of K
```

---

- How well does Decision Transformer model the distribution of returns?

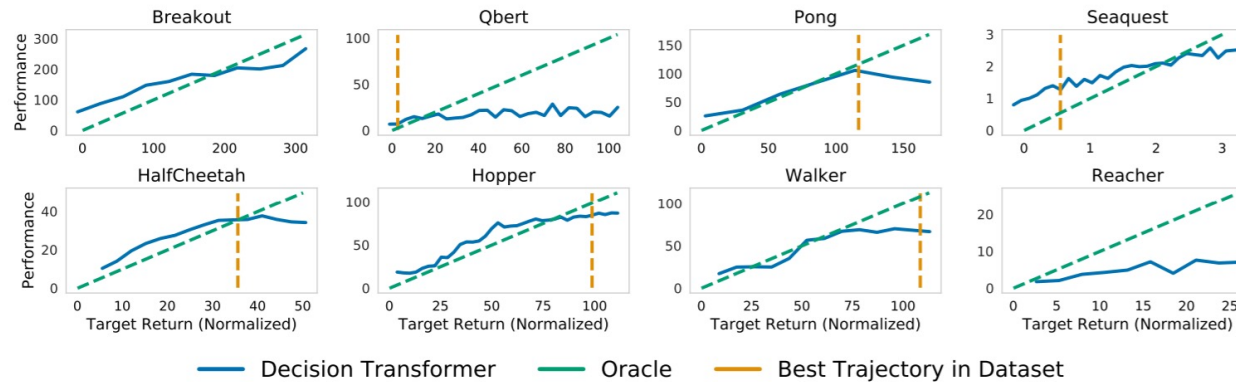


Figure 4: Sampled (evaluation) returns accumulated by Decision Transformer when conditioned on the specified target (desired) returns. **Top:** Atari. **Bottom:** D4RL medium-replay datasets.

It shows that it is possible to generate a trajectory with a desired reward sum.

- Does Decision Transformer perform well in sparse reward settings?

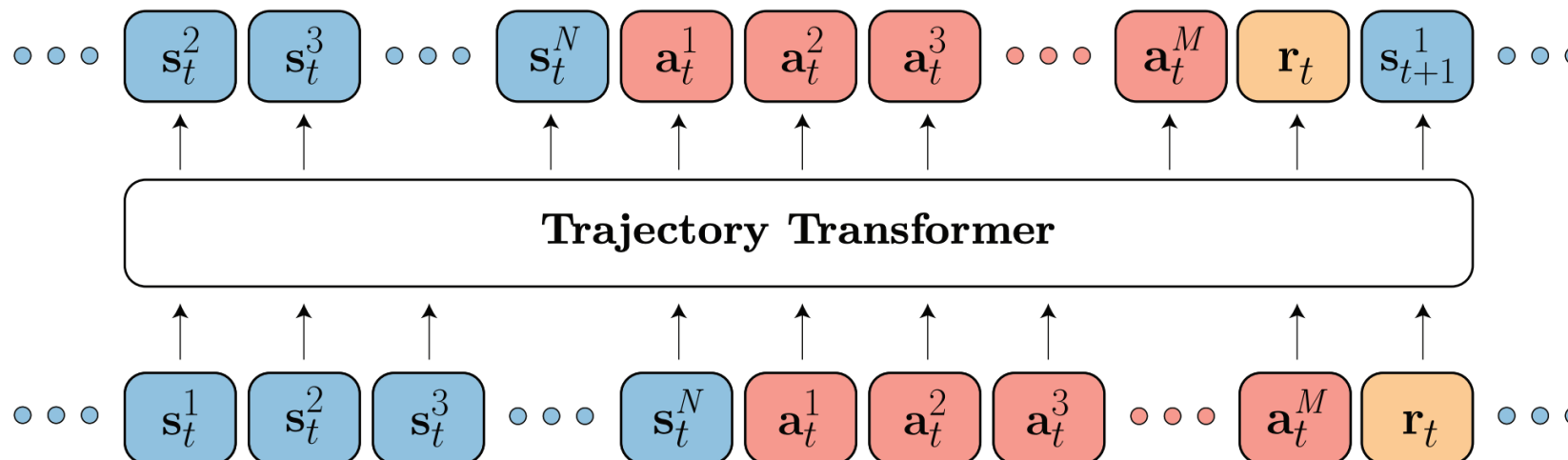
Dataset	Environment	Delayed (Sparse)		Agnostic		Original (Dense)	
		DT (Ours)	CQL	BC	%BC	DT (Ours)	CQL
Medium-Expert	Hopper	<b>107.3 ± 3.5</b>	9.0	59.9	102.6	107.6	111.0
Medium	Hopper	60.7 ± 4.5	5.2	63.9	<b>65.9</b>	67.6	58.0
Medium-Replay	Hopper	<b>78.5 ± 3.7</b>	2.0	27.6	70.6	82.7	48.6

Table 7: Results for D4RL datasets with delayed (sparse) reward. Decision Transformer (DT) and imitation learning are minimally affected by the removal of dense rewards, while CQL fails.

Compared to CQL (another offline RL algorithm), it works well in sparse reward setting.

Michael Janner, Qiyang Li, and Sergey Levine, "**Offline Reinforcement Learning as One Big Sequence Modeling Problem**", Advances in Neural Information Processing Systems (NeurIPS), 2021.

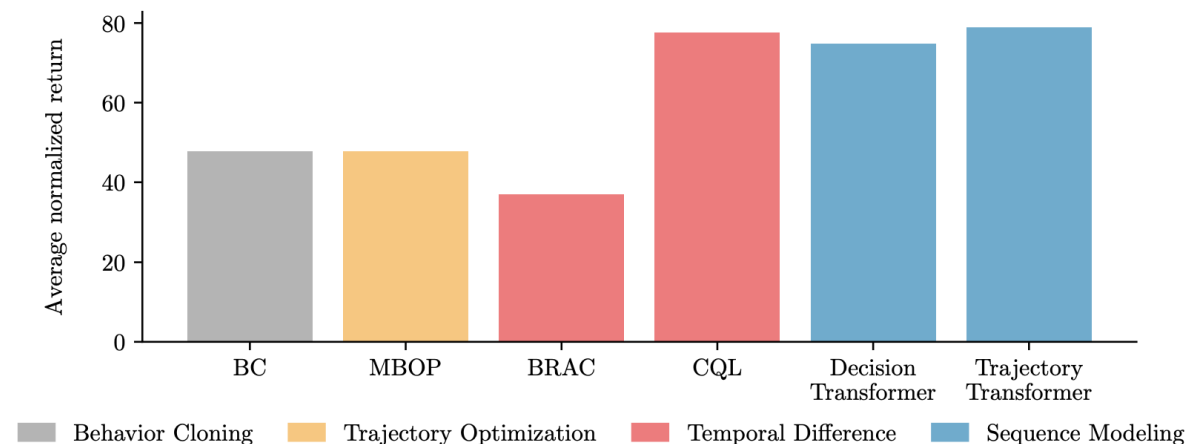
# TRAJECTORY TRANSFORMER



**Figure 1 (Architecture)** The Trajectory Transformer trains on sequences of (autoregressively discretized) states, actions, and rewards. Planning with the Trajectory Transformer mirrors the sampling procedure used to generate sequences from a language model.

# Experimental Results

Dataset	Environment	BC	MBOP	BRAC	CQL	DT	TT (uniform)	TT (quantile)
Med-Expert	HalfCheetah	59.9	105.9	41.9	91.6	86.8	40.8 $\pm$ 2.3	95.0 $\pm$ 0.2
Med-Expert	Hopper	79.6	55.1	0.9	105.4	107.6	106.0 $\pm$ 0.28	110.0 $\pm$ 2.7
Med-Expert	Walker2d	36.6	70.2	81.6	108.8	108.1	91.0 $\pm$ 2.8	101.9 $\pm$ 6.8
Medium	HalfCheetah	43.1	44.6	46.3	44.0	42.6	44.0 $\pm$ 0.31	46.9 $\pm$ 0.4
Medium	Hopper	63.9	48.8	31.3	58.5	67.6	67.4 $\pm$ 2.9	61.1 $\pm$ 3.6
Medium	Walker2d	77.3	41.0	81.1	72.5	74.0	81.3 $\pm$ 2.1	79.0 $\pm$ 2.8
Med-Replay	HalfCheetah	4.3	42.3	47.7	45.5	36.6	44.1 $\pm$ 0.9	41.9 $\pm$ 2.5
Med-Replay	Hopper	27.6	12.4	0.6	95.0	82.7	99.4 $\pm$ 3.2	91.5 $\pm$ 3.6
Med-Replay	Walker2d	36.9	9.7	0.9	77.2	66.6	79.4 $\pm$ 3.3	82.6 $\pm$ 6.9
<b>Average</b>		<b>47.7</b>	<b>47.8</b>	<b>36.9</b>	<b>77.6</b>	<b>74.7</b>	<b>72.6</b>	<b>78.9</b>

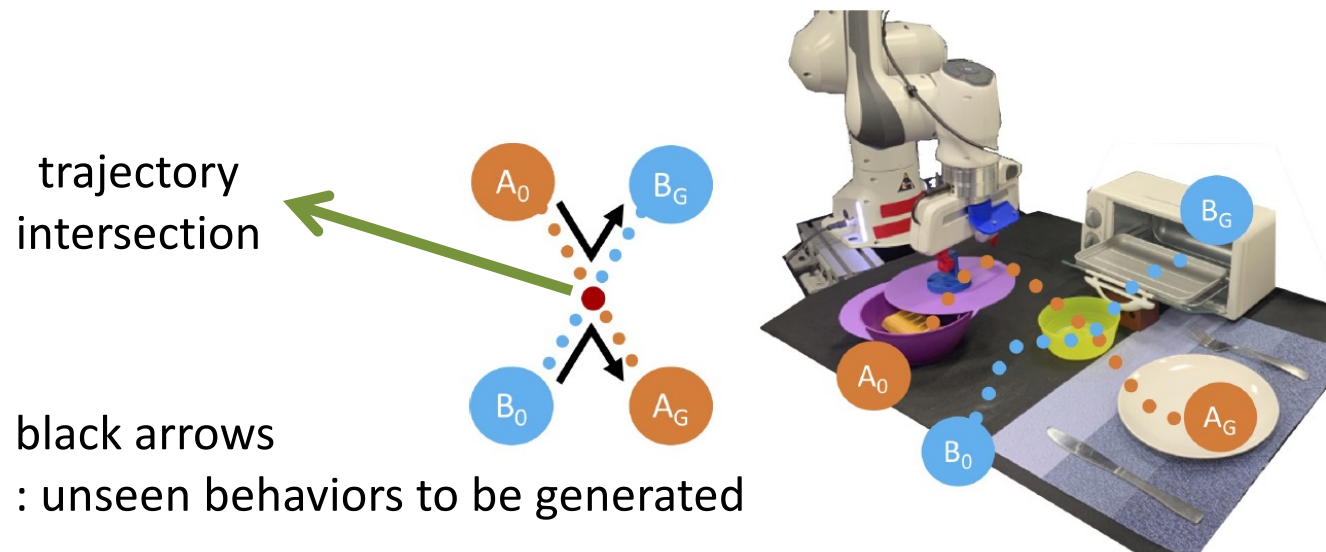


Mandlekar, Ajay, Danfei Xu, Roberto Martín-Martín, Silvio Savarese, and Li Fei-Fei. "GTI: Learning to generalize across long-horizon tasks from human demonstrations." Robotics: Science and Systems (RSS), 2020.

# GTI: LEARNING TO GENERALIZE ACROSS LONG-HORIZON TASKS FROM HUMAN DEMONSTRATIONS

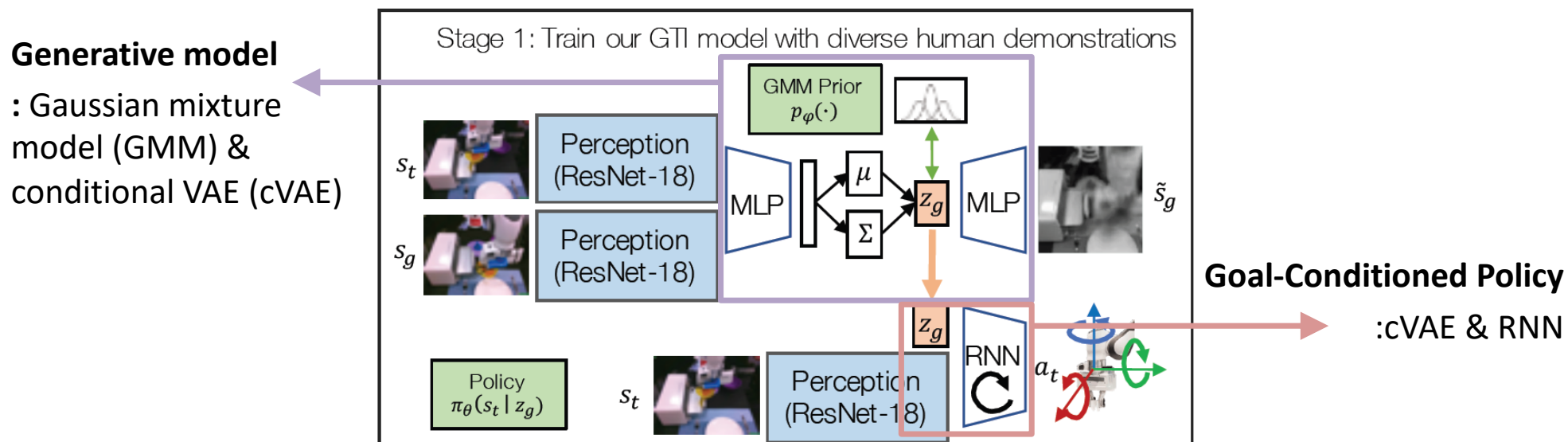
- Learning from offline data is an effective and safe technique but learning policies that **generalize beyond the demonstrated behaviors** is still an open challenge.
- This paper proposes Generalization Through Imitation (GTI), a two-stage offline imitation learning algorithm.
- GTI enable robots to
  - 1) learn complex real world manipulation tasks efficiently from **a small number of human demonstrations**,
  - 2) **synthesize new behaviors not contained** in the offline data.

- In the **first stage** of GTI, they train a **stochastic policy** that leverages trajectory intersections to have the capacity to compose behaviors from different demonstrations together.

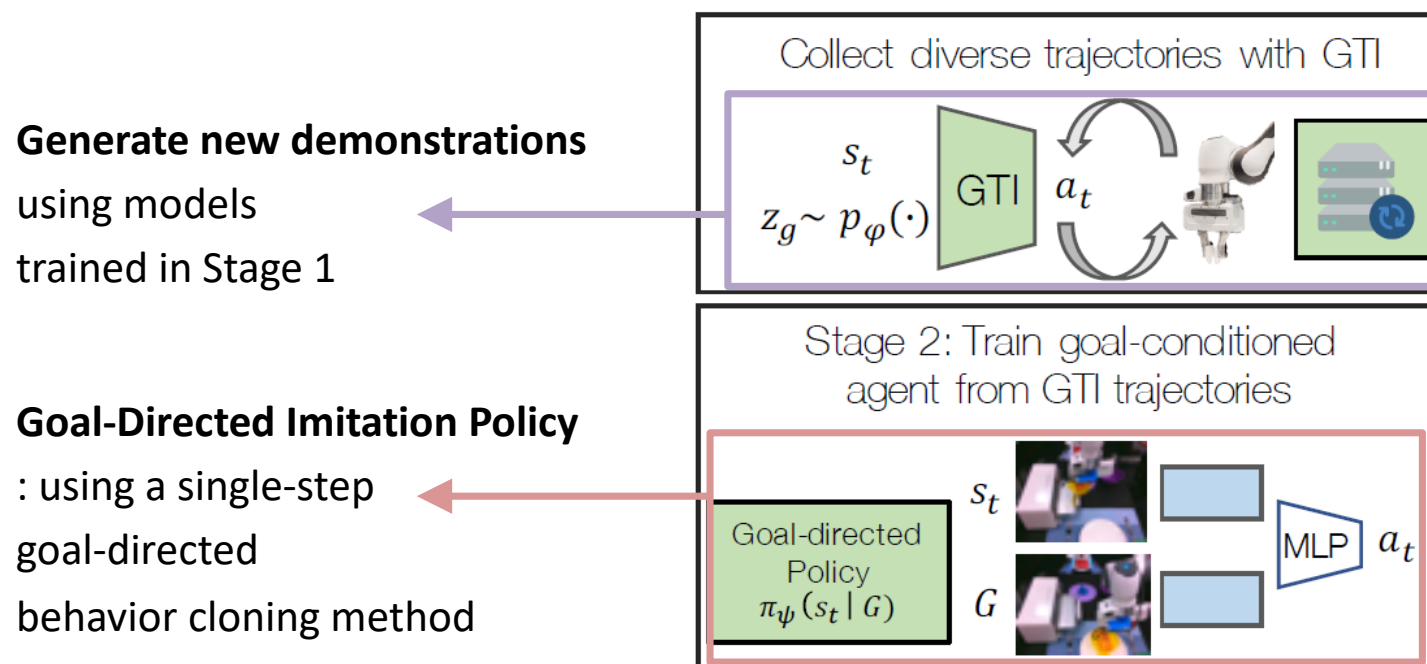


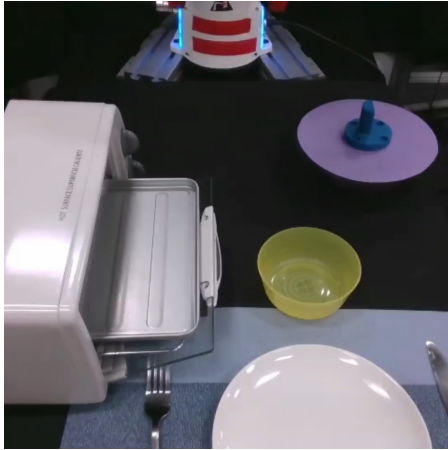
- In the **second stage** of GTI, they **collect a small set of rollouts** from the unconditioned stochastic policy of the first stage and **train a goal-directed agent** to generalize to novel start and goal state.

- In Stage 1, GTI decomposes imitation learning into two subproblems.
  - learn a **generative model** from demonstrations to predict possible future states conditioned on a current state.
  - train a **goal-conditioned imitation policy** from the demonstrations using predicted states as goals.



- In Stage 2, GTI uses trained models in Stage 1 to generate new demonstrations through compositional generalization.
- Then, GTI trains a new policy using newly collected behaviors in a goal directed manner.





**Demonstration 1**

- Start: Bread in Container
- End: Bread on Plate



**Demonstration 2**

- Start: Bread on Table
- End: Bread in Oven

**Demonstrations**  
(offline data)

**Train GTI**



**Unseen Trajectory**

- Start: Bread in Container
- End: Bread in Oven