# Versatile Offline Imitation from Observations and Examples via Regularized State Occupancy Matching

Jason Yecheng Ma, Andrew Shen, Dinesh Jayaraman, Osbert Bastani

University of Pennsylvania

# Offline Imitation Learning from Observations



Reinforcement Learning with Online Interactions

Offline Reinforcement Learning

- (Few) expert observations: $\mathcal{D}^E = \{(s_0, ..., s_T)\}$

- Offline non-expert dataset: $\mathcal{D}^O = \{(s, a, s')\}$

- Objective: $\mathrm{D_{KL}}(d^\pi(s) \| d^E(s))$

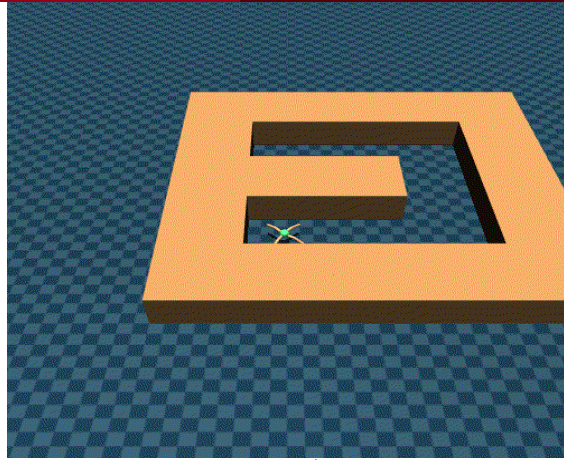- Offline IL: $\pi = \mathcal{A}(\mathcal{D}^E, \mathcal{D}^O)$

How can we leverage small number of expert observations and large amount of unlabeled offline data to achieve offline imitation learning?
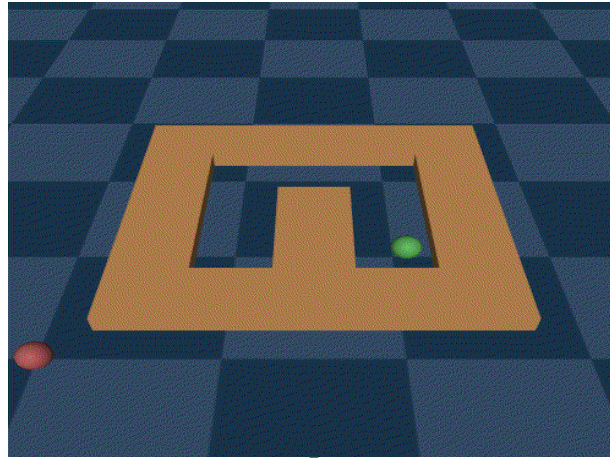
# Objective: State Occupancy Matching

$$\min_{\pi} \mathrm{D}_{\mathrm{KL}}(d^{\pi}(s) \| d^{E}(s))$$

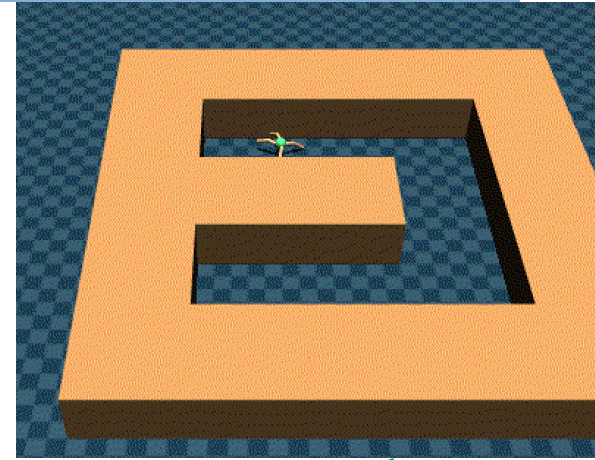"distribution of task-relevant states the policy visits"

Penn Engineering

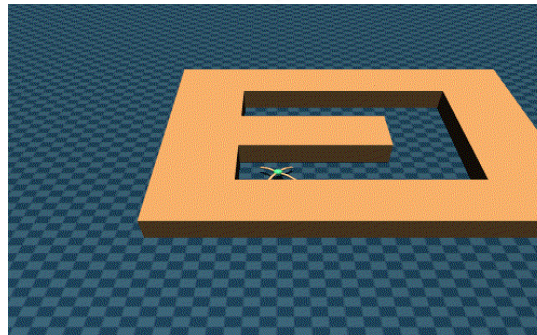# Versatility of State-Occupancy Matching



Observations

Mismatched Experts

Examples of Success

# Imitation Learning from Examples

$$\min_{\pi} \mathrm{D}_{\mathrm{KL}}(d^{\pi}(s) \| d^{E}(s))$$

"state distribution of a teleporting expert"

Dynamics-abiding imitator



"Teleporting" expert

# Offline Imitation Learning from via State Matching

- Formulated as a <u>state-matching</u> problem:

$$\mathrm{D_{KL}}(d^{\pi}(s)\|d^{E}(s)) = \mathbb{E}_{s\sim \boxed{d^{\pi}}}\left[\log\frac{d^{\pi}(s)}{d^{E}(s)}\right]$$

- Key challenge is that this objective requires samples from the policy we are optimizing

- Difficult to do offline without access to the environment!

Penn Engineering

# Regularized State-Occupancy Matching

Under some mild assumptions, for any $f$-divergence such that $D_f \geq D_{KL}$,

$$D_{KL}(d^\pi(s) \| d^E(s)) \leq \mathbb{E}_{s \sim d^\pi} \left[ \log \left( \frac{d^O(s)}{d^E(s)} \right) \right] + D_f(d^\pi(s, a) \| d^O(s, a))$$

"reward function": encourages visiting expert states; learned using a discriminator!
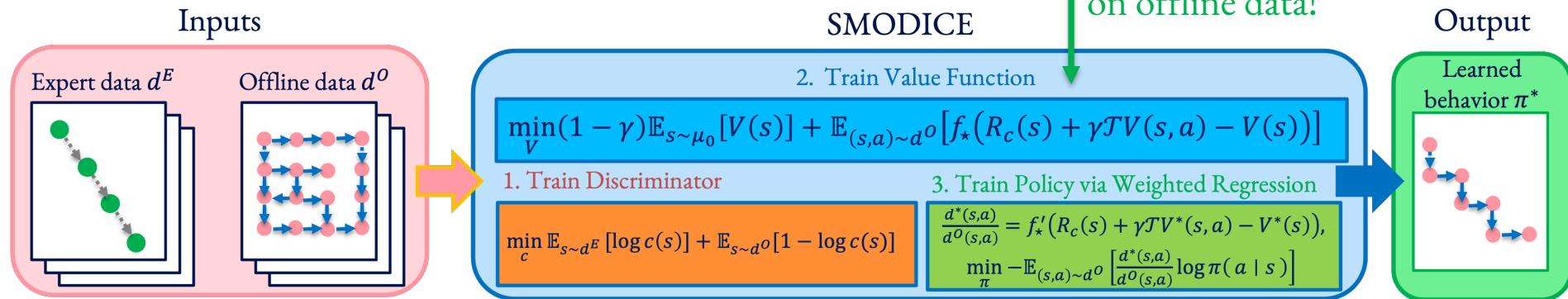
"constraint": encourages staying close to the offline dataset

Still requires on-policy samples!

# State Matching Offline DIstribution Correction Estimation (SMODICE)

$$\mathrm{D}_{\mathrm{KL}}(d^\pi(s)\|d^E(s)) \le \mathbb{E}_{s \sim d^\pi}\left[\log\left(\frac{d^O(s)}{d^E(s)}\right)\right] + \mathrm{D}_f(d^\pi(s,a)\|d^O(s,a))$$

Dual problem depends only on offline data!

Inputs

Expert data $d^E$

Offline data $d^O$

SMODICE

**2. Train Value Function**

$$\min_V (1-\gamma)\mathbb{E}_{s \sim \mu_0}[V(s)] + \mathbb{E}_{(s,a) \sim d^O}\left[f_\star\left(R_c(s) + \gamma\mathcal{T}V(s,a) - V(s)\right)\right]$$

**1. Train Discriminator**

$$\min_c \mathbb{E}_{s \sim d^E}[\log c(s)] + \mathbb{E}_{s \sim d^O}[1 - \log c(s)]$$

**3. Train Policy via Weighted Regression**

$$\frac{d^*(s,a)}{d^O(s,a)} = f'_\star\left(R_c(s) + \gamma\mathcal{T}V^*(s,a) - V^*(s)\right),$$

$$\min_\pi -\mathbb{E}_{(s,a)\sim d^O}\left[\frac{d^*(s,a)}{d^O(s,a)}\log\pi(a \mid s)\right]$$

Output

Learned behavior $\pi^*$

☺ 3 Disjoint Optimization Steps
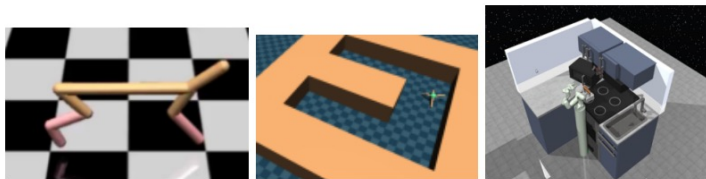
# Experiments

Penn Engineering

# Questions

1. Can SMODICE effectively learn from <u>observations</u>?

2. How robust is SMODICE to <u>mismatched experts ?</u>

3. Can SMODICE learn from <u>examples of success outcomes</u>?

Penn Engineering

# 1. Offline Imitation Learning from Observations

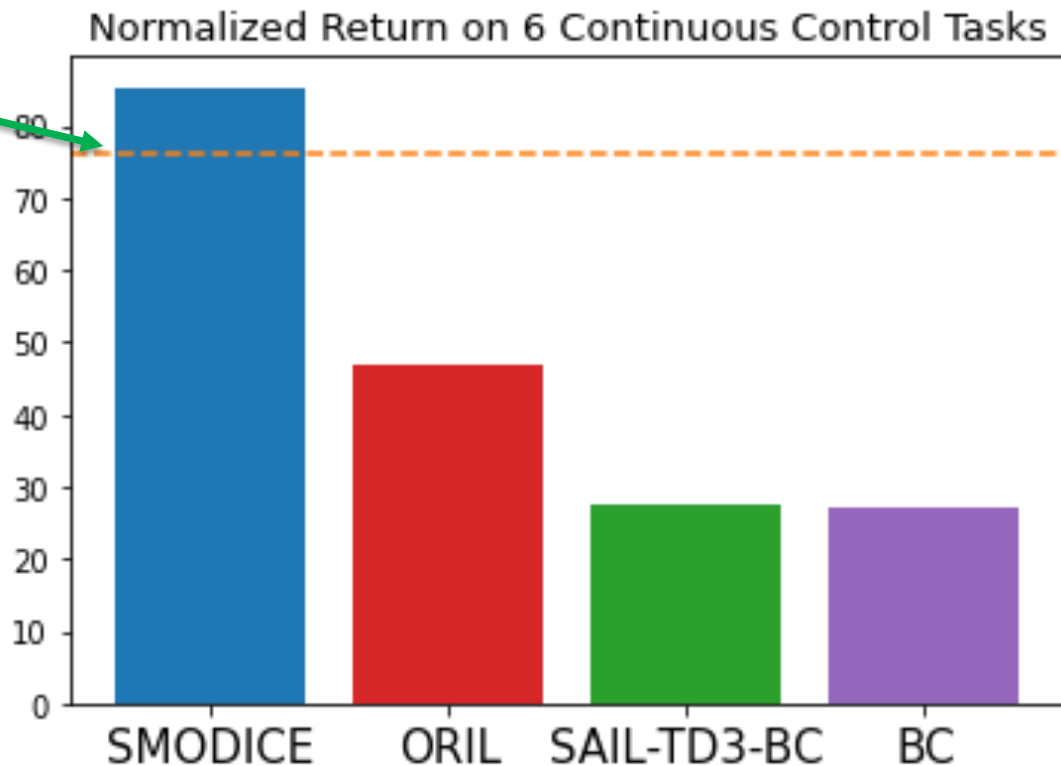Outperforms state-of-art with privileged action information!



(a) Mujoco   (b) AntMaze   (c) Franka Kitchen
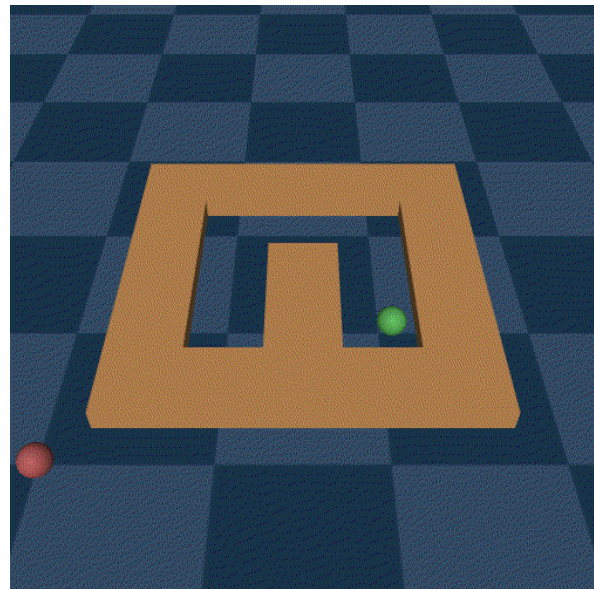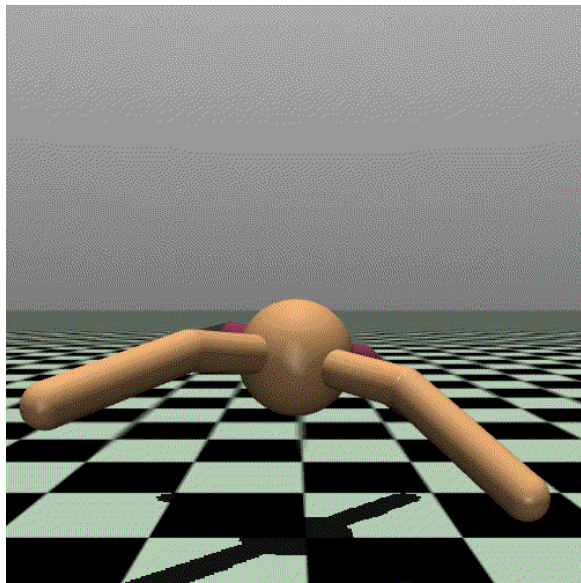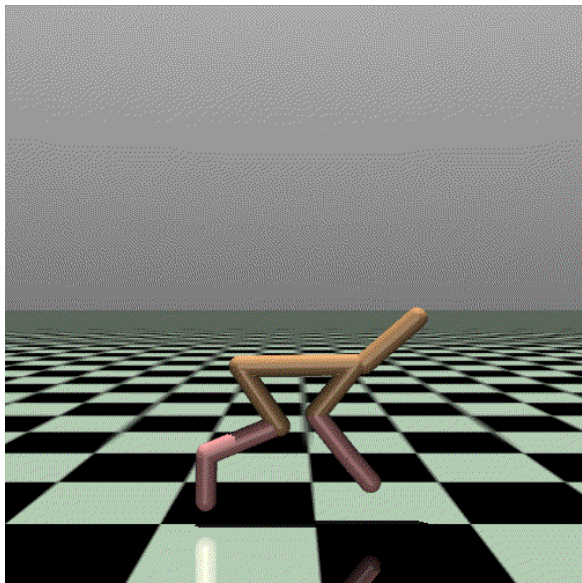
Figure 2. Illustrations of the evaluation environments.



Normalized Return on 6 Continuous Control Tasks

Penn Engineering

# 2. Offline IL from Mismatched Experts

Penn Engineering

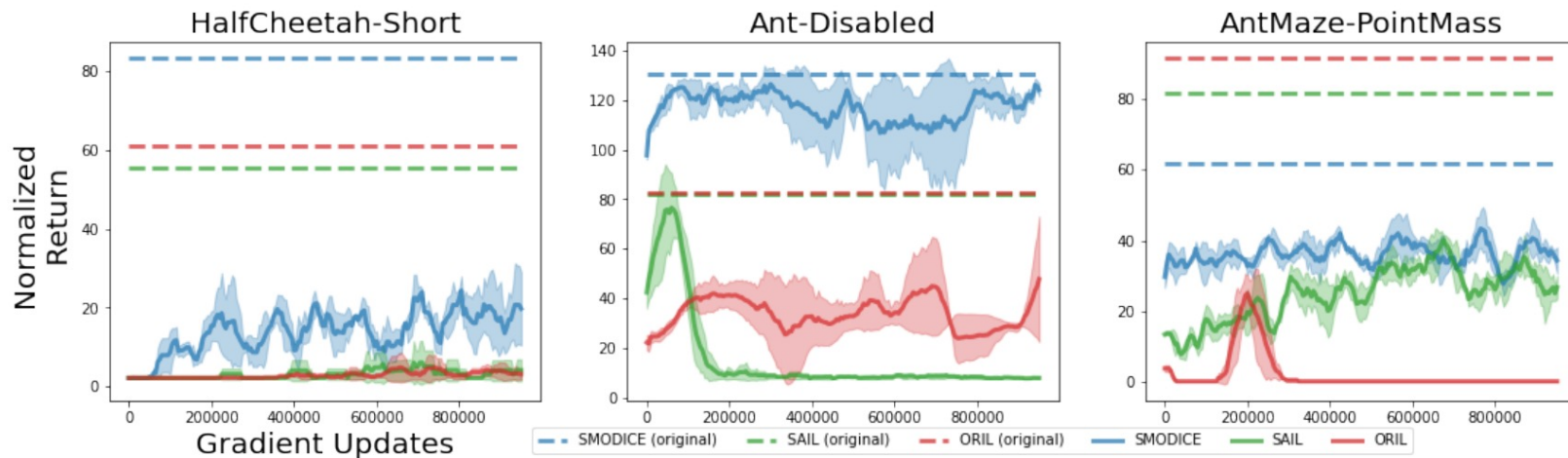# 2. Offline IL from Mismatched Experts



Figure 5. **Offline imitation learning from heterogeneous experts results.**

SMODICE is most robust to mismatched experts!

# 3. Offline IL from Examples



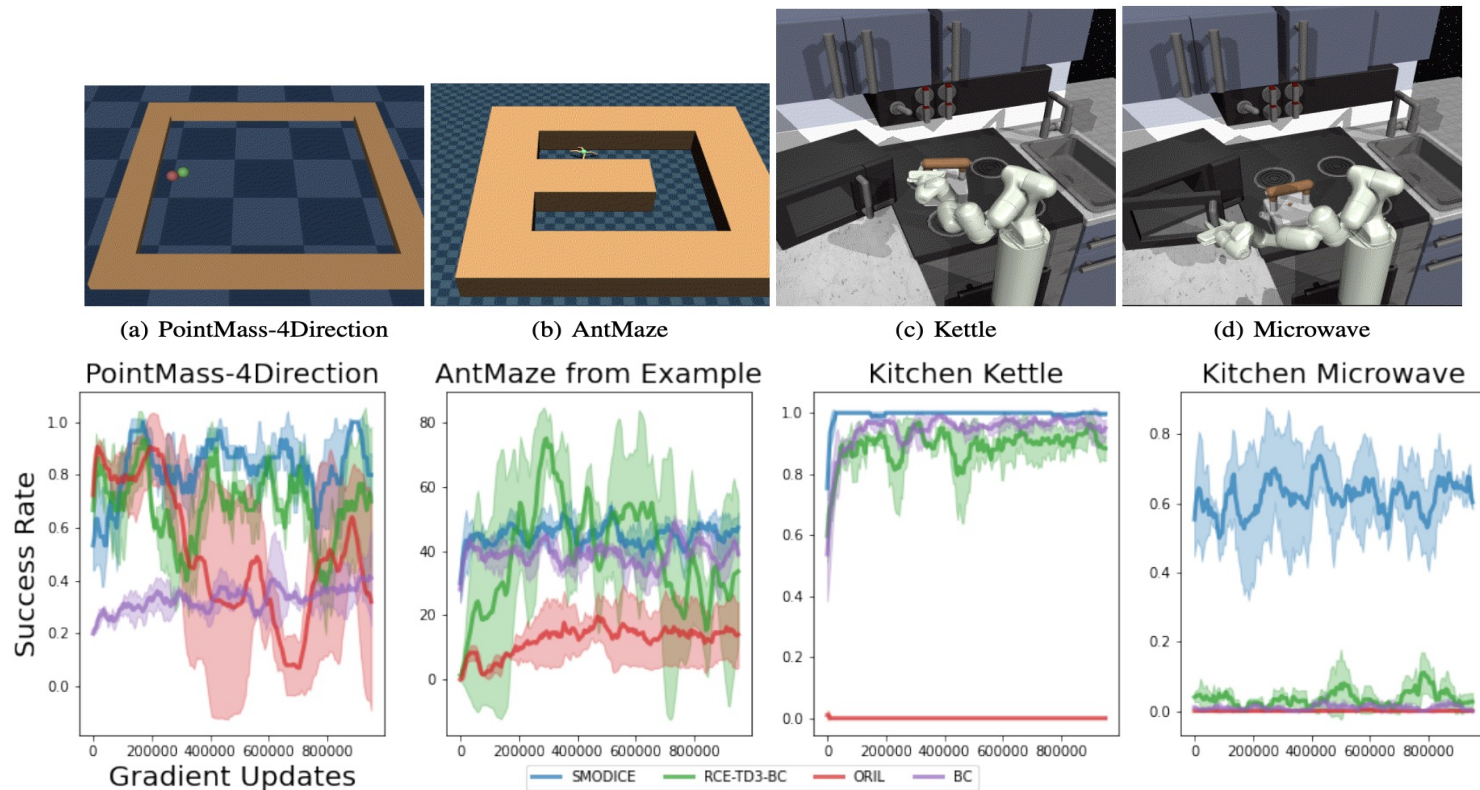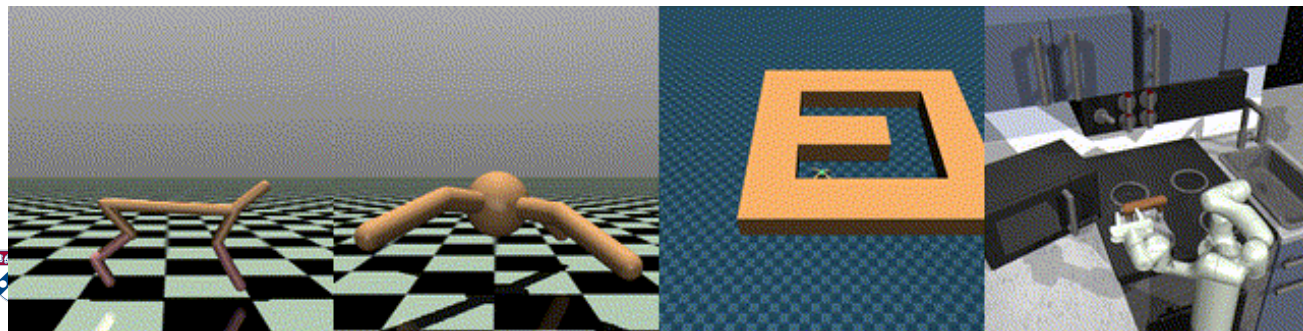(a) PointMass-4Direction  (b) AntMaze  (c) Kettle  (d) Microwave

Figure 6. **Offline imitation learning from examples results.**

# Conclusion

- State-matching as a framework for versatile offline IL

- SMODICE: A Regression-Based Offline IL Algorithm

- State-of-art results in all three settings without any

  hyperparameter tuning!

Project Website:
https://sites.google.com/view/smodice/home

Jason Yecheng Ma, Andrew Shen, Dinesh Jayaraman, Osbert Bastani
University of Pennsylvania