

Ishan Durugkar

Estimation and Control of Visitation Distributions for Reinforcement Learning

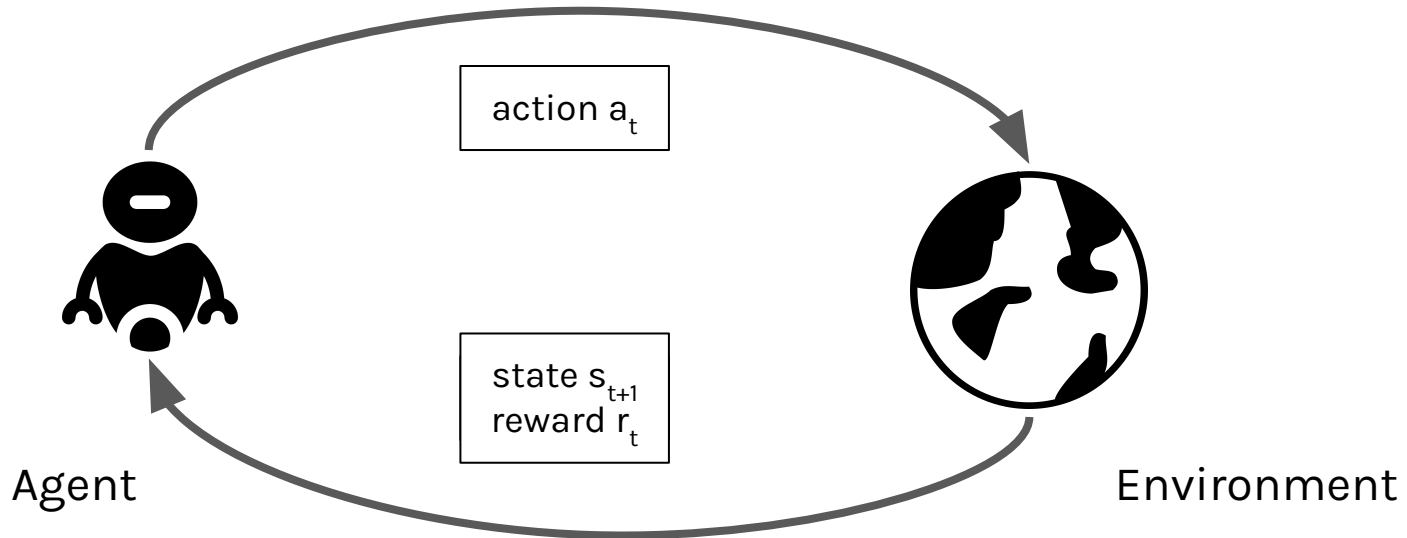
Committee:

Peter Stone, Qiang Liu, Philipp Krähenbühl

Scott Niekum, Marc Bellemare



Reinforcement Learning (RL)^[1]



$$\text{Return: } R(s) = \sum_{t=0}^{\infty} [\gamma^t r(s_t, a_t, s_{t+1}) | s_0 = s]$$

$$\text{Policy: } \pi : \mathcal{S} \mapsto \Delta(\mathcal{A})$$

[1] Reinforcement Learning: An Introduction, Sutton and Barto, 2018

Successes



[1] TD-Gammon, a self-teaching backgammon program, achieves master-level play; Tesauro, G.; Neural Computation 1994

[2] Mastering the game of go without human knowledge; Silver et al.; Nature 2017

[3] Human-level control through deep reinforcement learning; Mnih et al.; Nature 2015

[4] Outracing champion Gran Turismo drivers with deep reinforcement learning; Wurman et al.; Nature 2022

Ishan Durugkar, UT Austin

Challenges



Google's Robot Farm^[1]

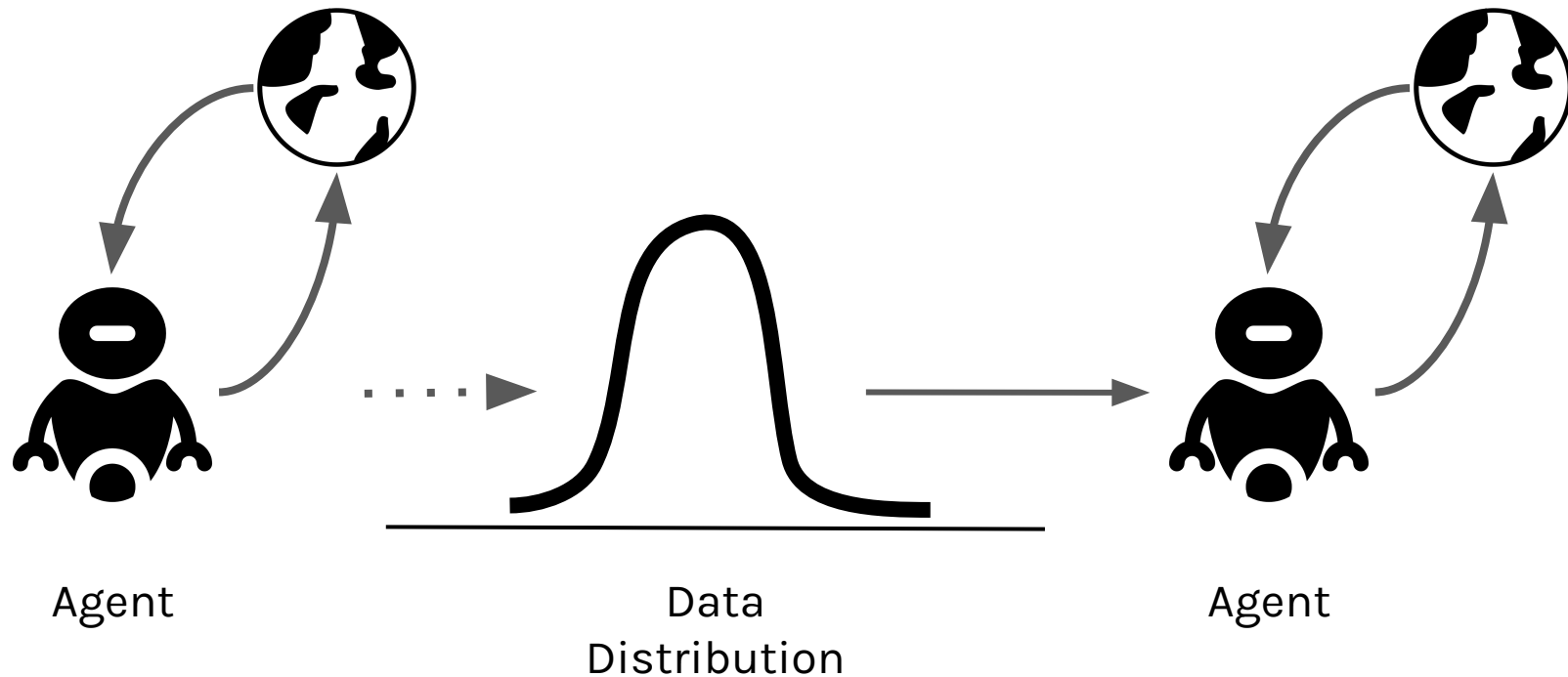


Autonomous Driving^[2]

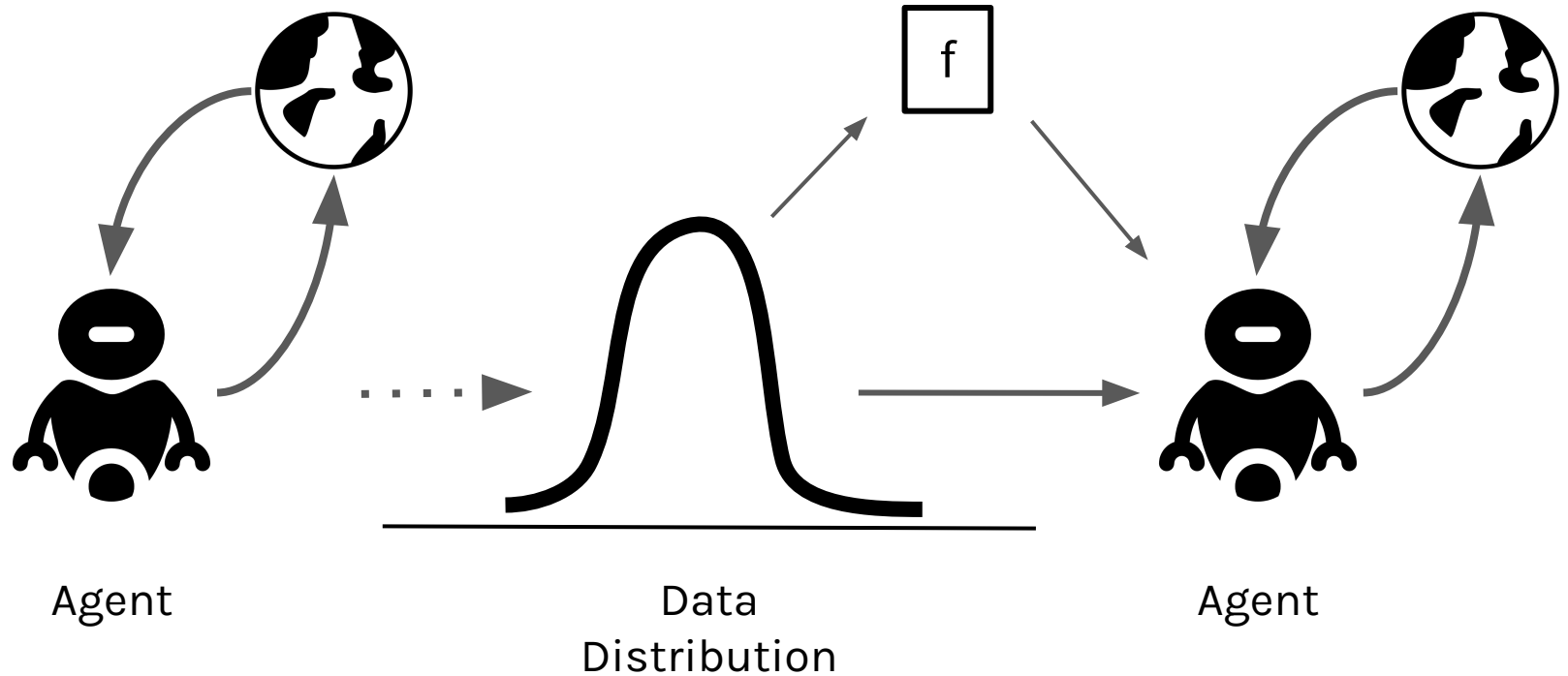
[1] Learning Hand-Eye Coordination for Robotic Grasping with Deep Learning and Large-Scale Data Collection; Levine et al.; ArXiv 2016

[2] Reward (Mis)design for Autonomous Driving; Knox et al.; ArXiv 2021

Why is RL Difficult?



Estimate and Control the Data Distribution



The Thesis Question

How can methods for estimating and controlling an agent's visitation distribution be beneficially incorporated into RL algorithms?

Distribution Matching for RL

How can methods for estimating and controlling an agent's visitation distribution be beneficially incorporated into RL algorithms?

Policy

Transitions

States

Definitions

Policy

$$\pi : \mathcal{S} \mapsto \Delta(\mathcal{A})$$

State visitation distribution

$$\rho_{\pi}(s) := \mathbb{E}_{s_0 \sim \rho_0} \left[(1 - \gamma) \sum_{t=0}^{\infty} \gamma^t P(s_t = s | \pi) \right]$$

Transition visitation distribution

$$\rho_{\pi}(s, a, s') := \mathbb{E}_{s_0 \sim \rho_0} \left[(1 - \gamma) \pi(a|s) P(s'|s, a) \sum_{t=0}^{\infty} \gamma^t P(s_t = s | \pi) \right]$$

Overview

1. Overcoming policy sampling error	2. Simulator grounding as imitation from observations (IfO)	3. An algorithm for sim-to-real transfer
4. Time-step metric for estimating Wasserstein distance in MDPs	5. Learning a goal conditioned policy by minimizing W_1 distance	6. Multi-agent coordination

Policy

Transitions

States

Overview

1. Overcoming policy sampling error Chapter 4	2. Simulator grounding as imitation from observations (IfO) Chapter 3	3. An algorithm for sim-to-real transfer
4. Time-step metric for estimating Wasserstein distance in MDPs Chapter 5	5. Learning a goal conditioned policy by minimizing W_1 distance	6. Multi-agent coordination Chapter 6

Policy

Transitions

States

Overview - Completed Before Proposal

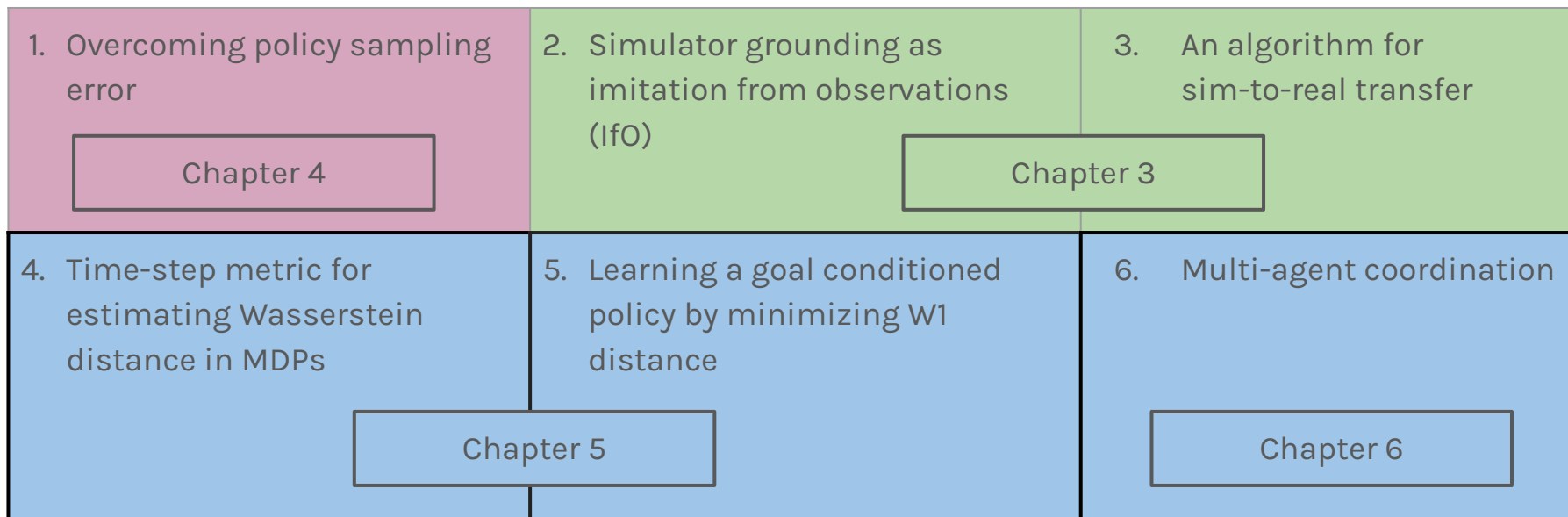
<p>1. Overcoming policy sampling error</p> <p>Chapter 4</p>	<p>2. Simulator grounding as imitation from observations (IfO)</p> <p>Chapter 3</p>	<p>3. An algorithm for sim-to-real transfer</p>
<p>4. Time-step metric for estimating Wasserstein distance in MDPs</p> <p>Chapter 5</p>	<p>5. Learning a goal conditioned policy by minimizing W_1 distance</p>	<p>6. Multi-agent coordination</p> <p>Chapter 6</p>

Policy

Transitions

States

Overview - After Proposal



Policy

Transitions

States

Distribution Estimation and Control for RL

Actions

- Reducing Sampling Error in Batch Temporal Difference Learning; Pavse, B., **Durugkar, I.**, Hanna, J., Stone, P.; ICML 2021

Transitions

- An Imitation from Observation Approach to Transfer Learning with Dynamics Mismatch; *Desai, S., ***Durugkar, I.**, *Karnan, H., Warnell, G., Hanna, J. and Stone, P; NeurIPS 2020

States

- Adversarial Intrinsic Motivation for Reinforcement Learning; **Durugkar, I.**, Tec, M., Niekum S., Stone, P.; NeurIPS 2021
- DM²: Distributed Multi-Agent Reinforcement Learning by Distribution Matching; *Wang, C., ***Durugkar, I.**, *Liebman, E., Stone, P.; AAI 2023

* - joint first authors

Overview

1. Overcoming policy sampling error	2. Simulator grounding as imitation from observations (IfO)	3. An algorithm for sim-to-real transfer
4. Time-step metric for estimating Wasserstein distance in MDPs	5. Learning a goal conditioned policy by minimizing W_1 distance	6. Multi-agent coordination

Policy

Transitions

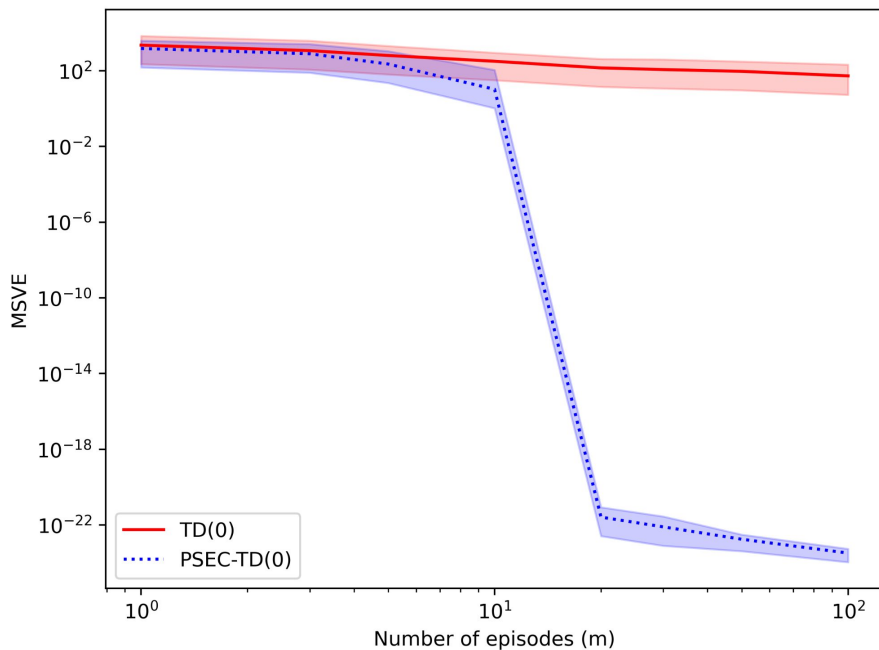
States

Policy Sampling Error Corrected (PSEC) - TD Learning

- Batch temporal difference (TD) learning will have some sampling error
- Contribution^[1]: Estimating the maximum likelihood policy implied by the dataset allows me to eliminate sampling error
- Analysis shows that PSEC-TD(0) converges to a fixed point with no policy sampling error

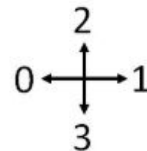
[1] Reducing sampling error in batch temporal difference learning; Pavse, B., **Durugkar, I.**, Hanna, J. and Stone, P.; ICML 2020

Results - Grid World



- 4 x 4 grid world, deterministic transitions
- Tabular representation
- equiprobable policy being evaluated
- PSEC-TD uses correction on TD error

0	1	2	3
4	5	6	7
8	9	10	11
12	13	14	15



Potential Future Work - Mixed Batches

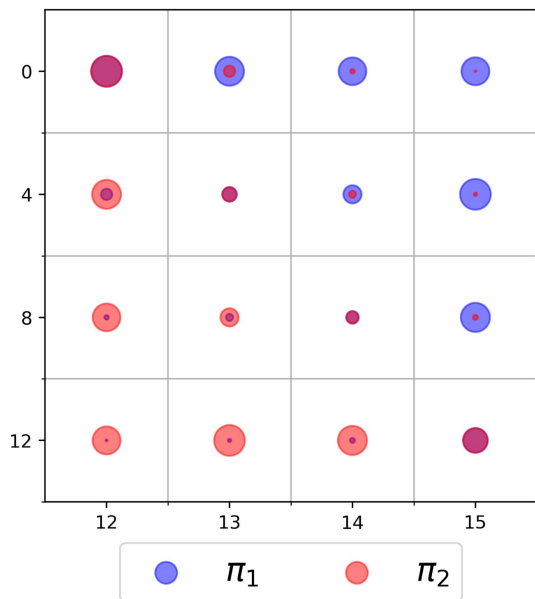
- How to deal with batches made up of data from different policies?
- Behavior policy estimation stays the same
- What evaluation policy to use?

$$\pi_{mix}(a|s) = \frac{\sum_{i=1}^K \lambda_i \rho_{\pi_i}(s, a)}{\sum_{i=1}^K \lambda_i \rho_{\pi_i}(s)}$$

where the batch was obtained by executing K policies π_1, \dots, π_K

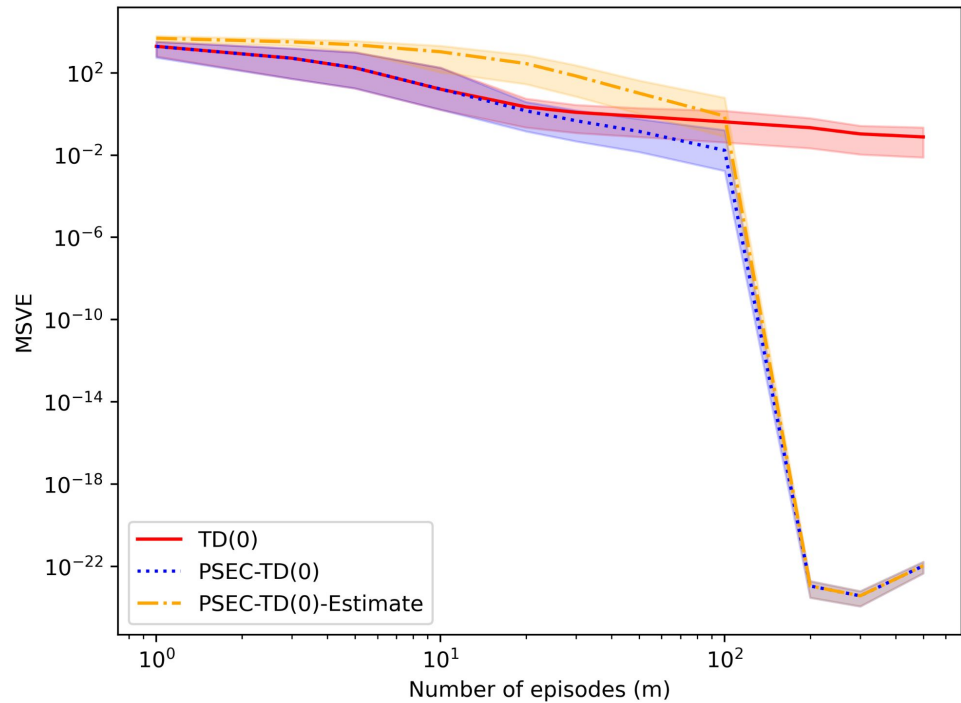
and policy i is executed with likelihood λ_i

Experiment



- 4x4 grid world
- Two policies, data collected from them equally
- Size of bubbles shows relative likelihood of visitation under corresponding policy

Results - Mixed Batches



- Two policies with data 50% from each policy
- Evaluation policy calculated with DP
- Takes more episodes to eliminate policy sampling error

Policy Sampling Error – Summary

- Estimating the empirical distribution of the policy can help eliminate sampling error
- Analysis shows that PSEC-TD(0) converges to a more desirable fixed point compared to TD(0)
- Experiments show that PSEC-TD(0) eliminates sampling error
- Introduce a potential avenue for future work

Overview

1. Overcoming policy sampling error	2. Simulator grounding as imitation from observations (IfO)	3. An algorithm for sim-to-real transfer
4. Time-step metric for estimating Wasserstein distance in MDPs	5. Learning a goal conditioned policy by minimizing W_1 distance	6. Multi-agent coordination

Policy

Transitions

States

Simulator Grounding and GARAT ^[1]

- Transfer with dynamics mismatch seen through the transition distributions induced
- Contribution 1: Show that simulator grounding via grounded action transformation (GAT)^[2] is equivalent to imitation from observations (IfO) where the expert is the target environment (real world)
- Contribution 2: Derive an adversarial distribution matching algorithm, generative adversarial reinforced action transformation (GARAT), to train the action transformation function

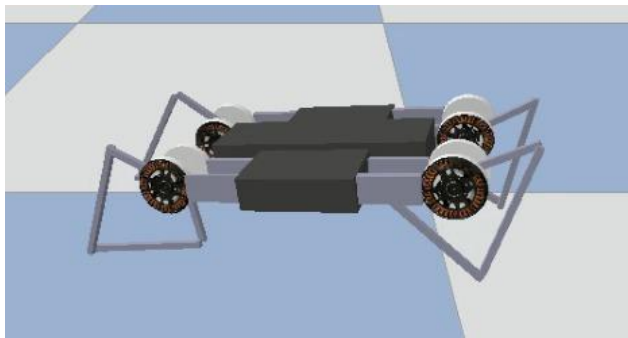
[1] An Imitation from Observation Approach to Transfer Learning with Dynamics Mismatch; *Desai, S., ***Durugkar, I.**, *Karnan, H., Warnell, G., Hanna, J. and Stone, P; NeurIPS 2020

* - joint first authors

[2] Grounded action transformation; Hanna et al.; AAAI 2017

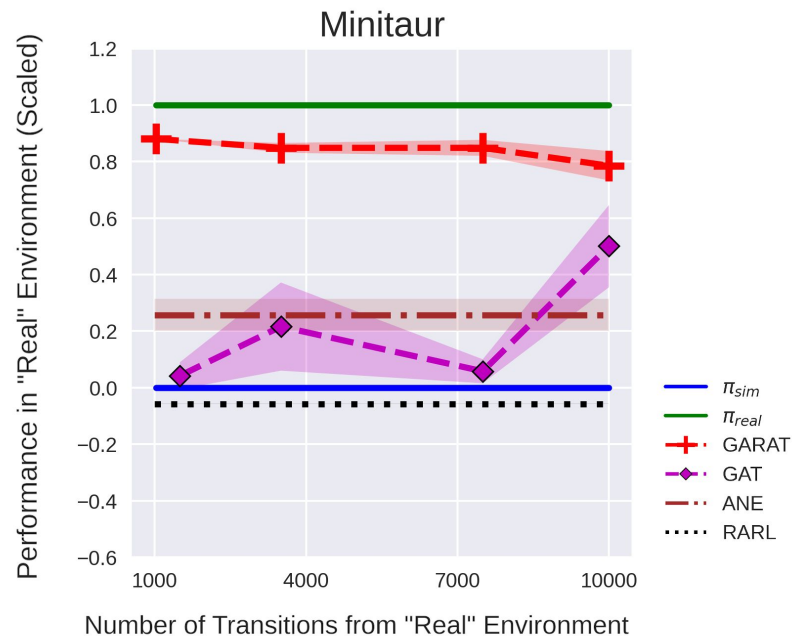
Results - Evaluating Transfer

- Transfer between two simulators for Minitaur^[1]
- Baselines trained for 1 million time-steps
- Results scaled to set performance of π_{sim} to 0 and π_{real} to 1



Minitaur domain

[1] Sim-to-Real: Learning Agile Locomotion For Quadruped Robots, Tan et al., RSS 2018



ANE - Noise and the reality gap: The use of simulation in evolutionary robotics, Morán et al., Advances in Artificial Life, 1995

RARL - Robust adversarial reinforcement learning, Pinto et al., ICML 2017

Overview

1. Overcoming policy sampling error	2. Simulator grounding as imitation from observations (IfO)	3. An algorithm for sim-to-real transfer
4. Time-step metric for estimating Wasserstein distance in MDPs	5. Learning a goal conditioned policy by minimizing W_1 distance	6. Multi-agent coordination

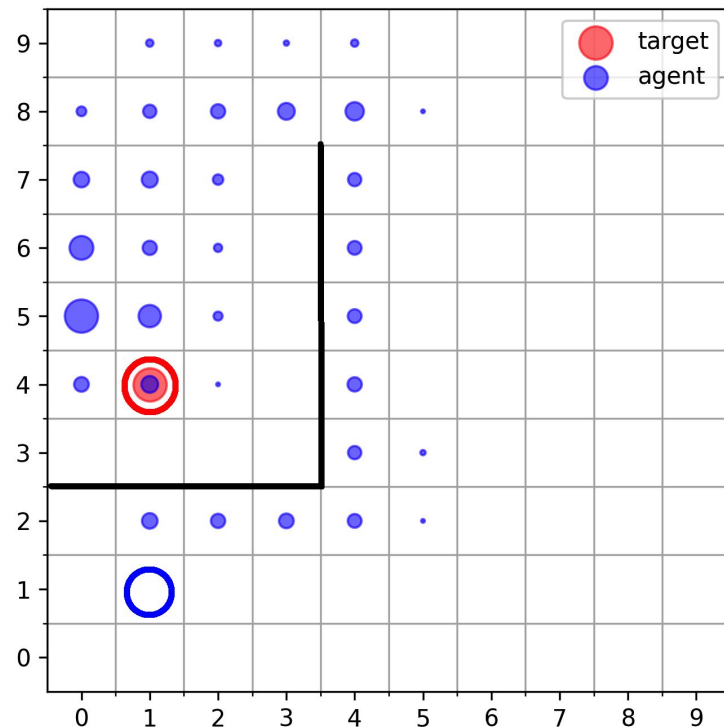
Policy

Transitions

States

Minimize distribution mismatch for Goal-conditioned RL

- **Goal-conditioned RL**^[1]: Agent needs to reach a goal given to it at the beginning of its episode.
- Blue circle - start state
- Red circle - goal state
- Target distribution can be specified as **Dirac distribution at the goal**
- Agent needs to **minimize mismatch of its state visitation distribution** to this target distribution.



[1] Learning to achieve goals; Kaelbling; IJCAI 1993

Minimize distribution mismatch in Goal-conditioned RL^[1]

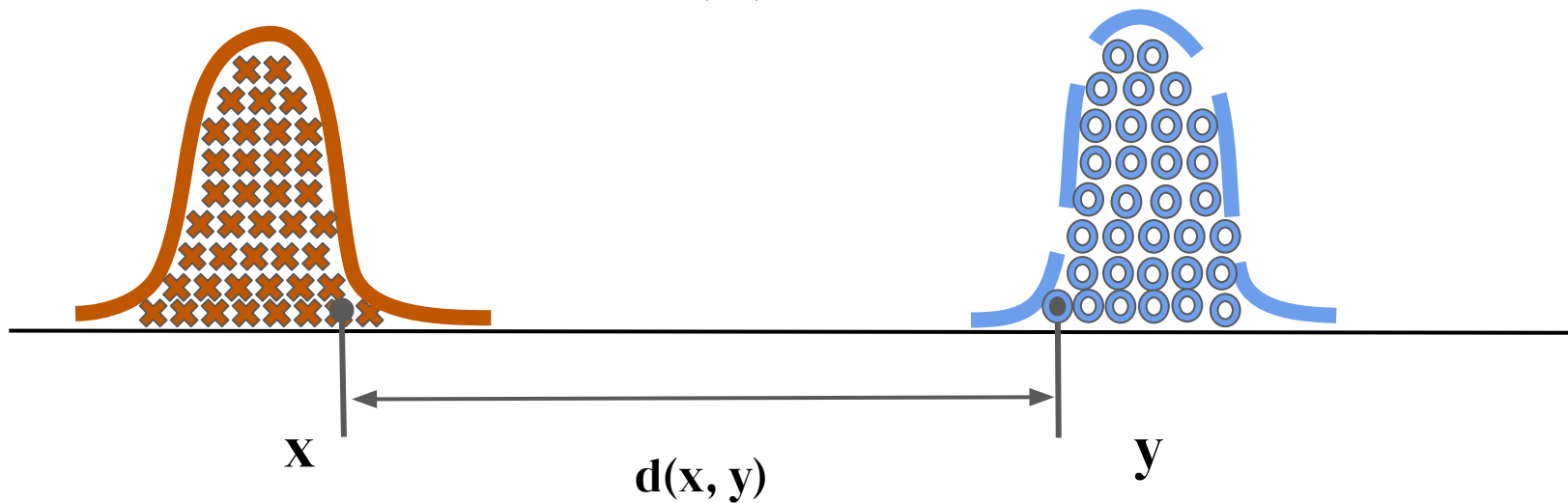
- Contribution 1: Study use of Wasserstein distance to minimize state visitation distribution mismatch. Propose use of time-step metric as ground metric for Wasserstein distance
- Contribution 2: Propose an adversarial procedure, adversarial intrinsic motivation (AIM) to learn a reward function which results in a policy that minimizes Wasserstein distance to a goal.

[1] Adversarial intrinsic motivation for reinforcement learning; **Durugkar, I.**, Tec, M., Niekum S., and Stone, P.; NeurIPS 2021

Wasserstein Distance

- Distance between distributions (say μ and ν)

$$W_d^p(\mu, \nu) := \inf_{\zeta \in Z(\mu, \nu)} \mathbb{E} [d(x, y)^p]^{\frac{1}{p}}$$



Wasserstein Distance using Kantorovich Duality

- If estimating Wasserstein-1 distance, the dual form can be used

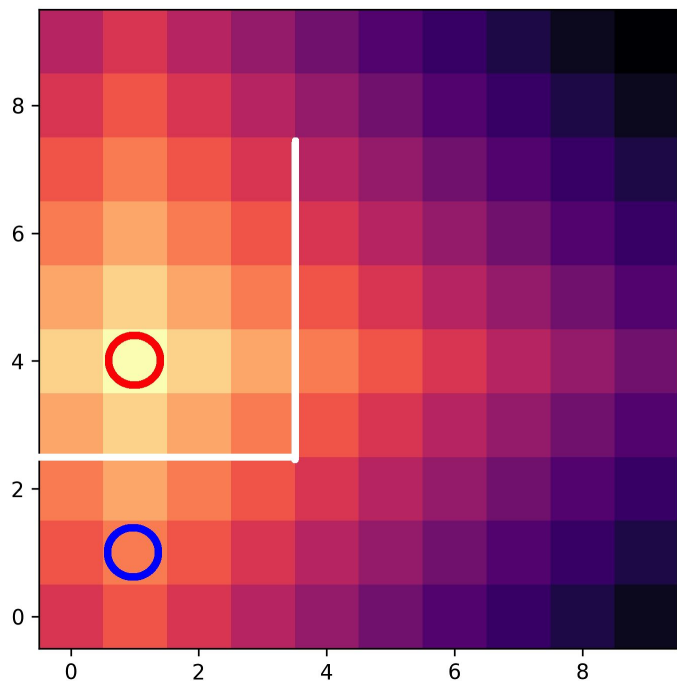
$$W_d^1(\mu, \nu) = \sup_{\text{Lip}(f) \leq 1} \mathbb{E}_{y \sim \nu} [f(y)] - \mathbb{E}_{x \sim \mu} [f(x)]$$

- The potential function f needs to be 1-Lipschitz w.r.t. metric d

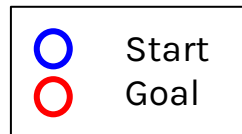
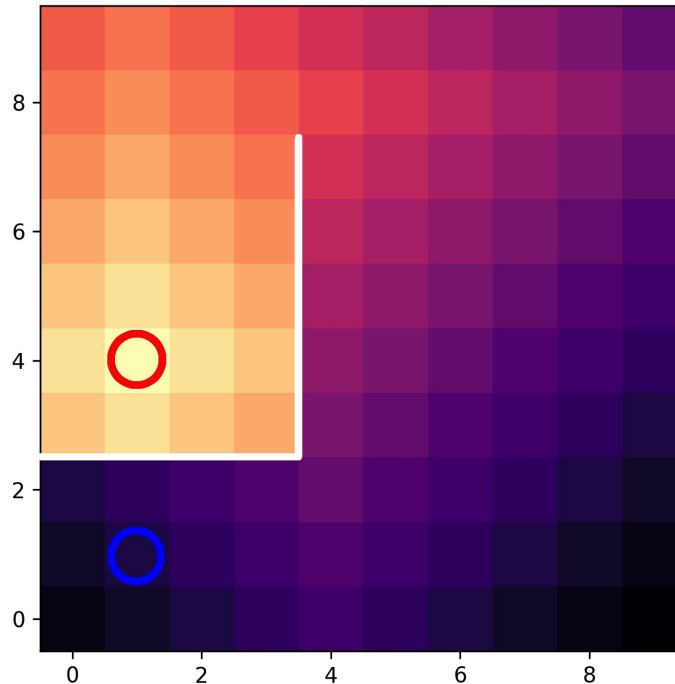


Why the Ground Metric Matters

Manhattan distance



Dynamics-based distance



Wasserstein Distance using Kantorovich Duality

- In most previous work, d is assumed to be L2 distance between features.
- I propose the use of the time-step metric for d in MDPs

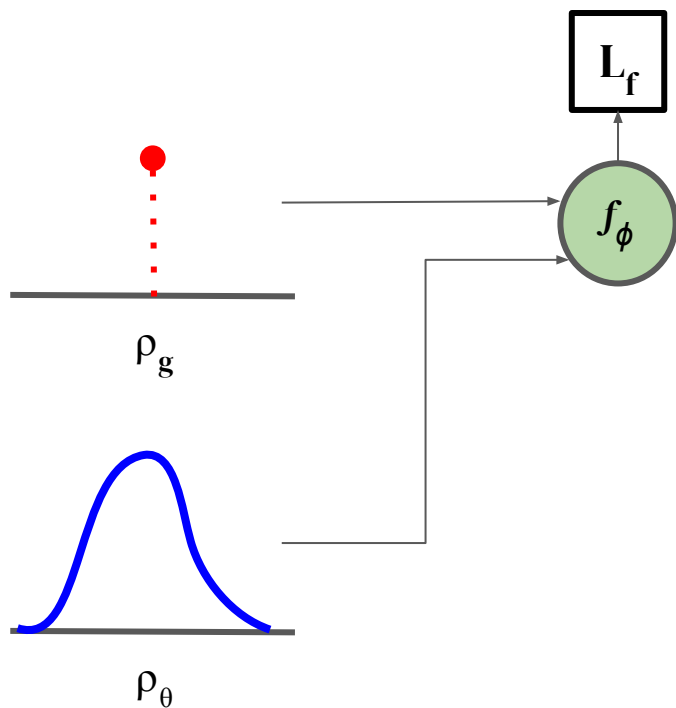
$$d_T^\pi(x, y) = \mathbb{E} [T(y|\pi, x)]$$

- Lipschitz continuity can be enforced as follows:

$$\mathbb{E}_{s' \sim \pi, P} [\|f(s) - f(s')\|] \leq 1 \quad \forall s \in \mathcal{S}$$

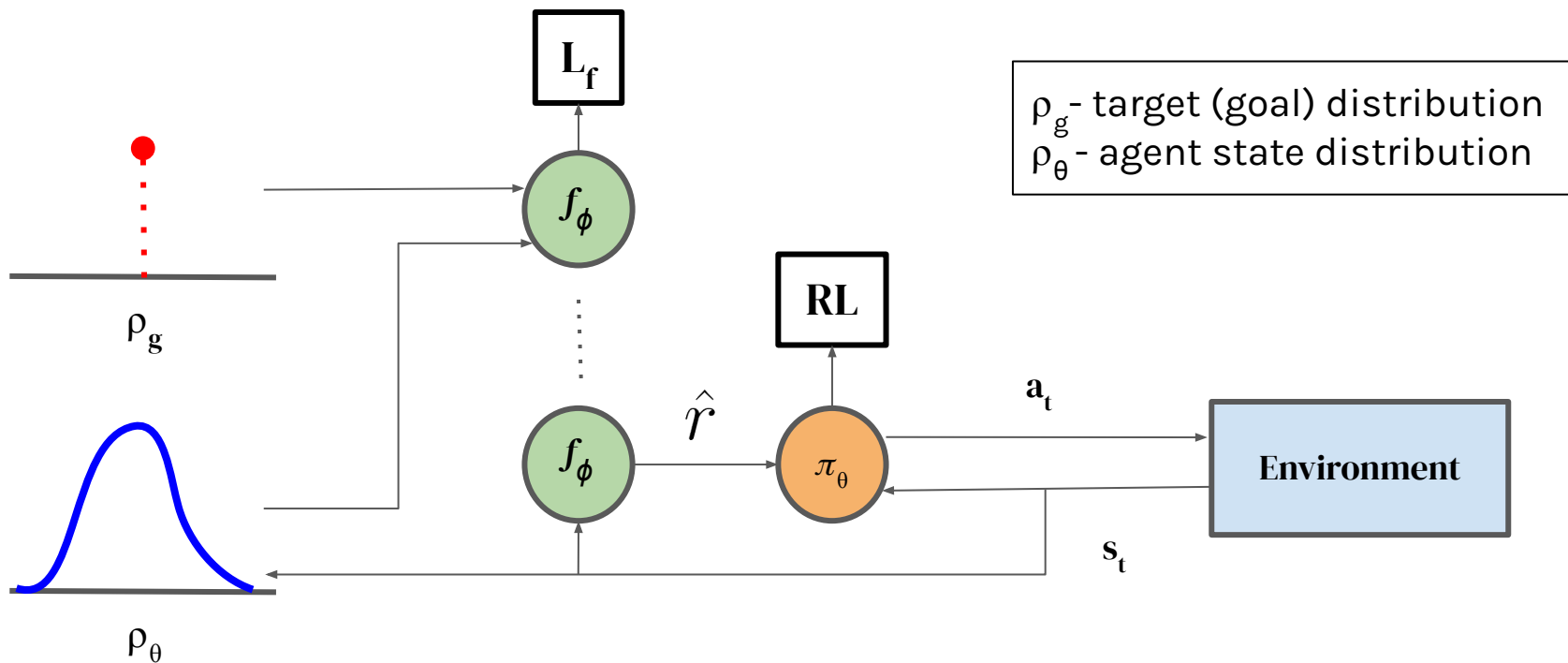
- In practice, we enforce it using samples from the environment

Adversarial Intrinsic Motivation (AIM)



ρ_g - target (goal) distribution
 ρ_θ - agent state distribution

Adversarial Intrinsic Motivation (AIM)



Analysis

Analysis - Discounted Setting

Comparing the optimal policy (π^*) and the policy that minimizes Wasserstein distance to goal (π ♦)

Proposition 4: A lower bound on the value of any state under a policy π can be expressed in terms of the time-step distance from that state to the goal:

$$v^\pi(s|s_g) = \mathbb{E} \left[\gamma^{T(s_g|\pi, s)} \right] \geq \gamma^{d_T^\pi(s, s_g)} \quad \forall s \in \mathcal{S}$$

Analysis - Discounted Setting

Comparing the optimal policy (π^*) and the policy that minimizes Wasserstein distance to goal (π^\blacklozenge)

Proposition 4: A lower bound on the value of any state under a policy π can be expressed in terms of the time-step distance from that state to the goal:

$$v^\pi(s|s_g) = \mathbb{E} \left[\gamma^{T(s_g|\pi,s)} \right] \geq \gamma^{d_T^\pi(s,s_g)} \quad \forall s \in \mathcal{S}$$

Theorem 5: If the transition dynamics are deterministic, the policy that minimizes the Wasserstein distance over the time-step metric in a goal-conditioned MDP is the optimal policy.

Analysis - Undiscounted Setting

- Undiscounted setting ($\gamma = 1$), with reward function

$$r(s_t, a_t, s_{t+1} | s_g) := \begin{cases} 0 & \text{if } s_{t+1} = \bar{s} \\ -1 & \text{otherwise} \end{cases}$$

- Assume agent reaches goal state from any start state within T steps

Analysis - Undiscounted Setting

- Undiscounted setting ($\gamma = 1$), with reward function

$$r(s_t, a_t, s_{t+1}|s_g) := \begin{cases} 0 & \text{if } s_{t+1} = \bar{s} \\ -1 & \text{otherwise} \end{cases}$$

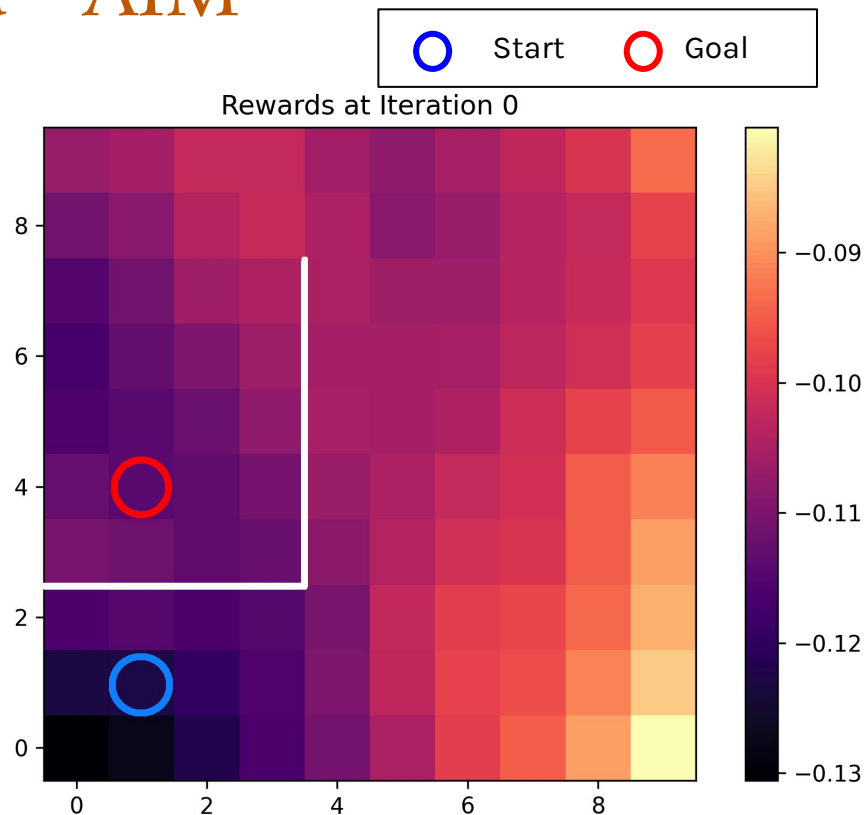
- Assume agent reaches goal state from any start state within T steps

Proposition 8: Assuming non-zero measure for all states s under the agent's state visitation distribution ρ_π , and considering s_g as the given goal state, the difference in potentials $f(s) - f(s_g) = v_\pi(s|s_g)$

Experiments: Grid World - AIM

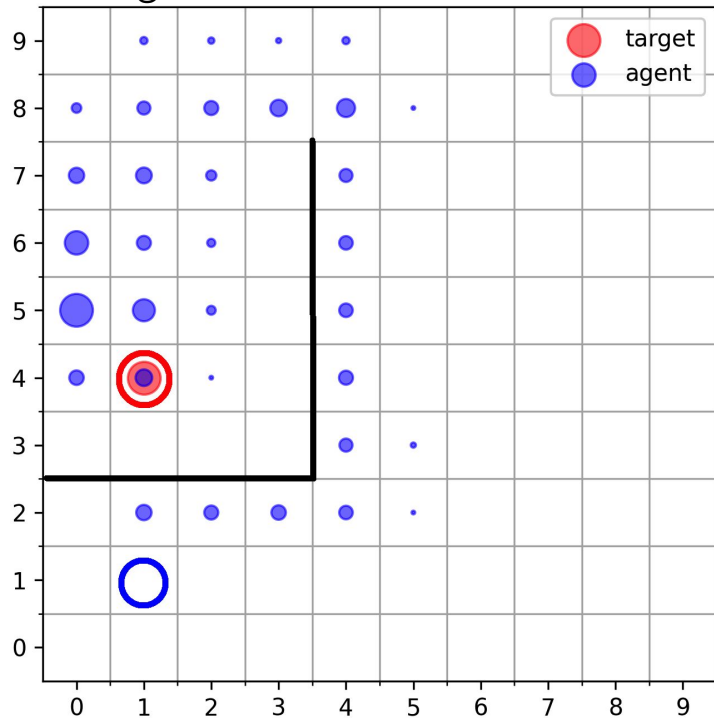
- 10 x 10 grid world, 4 actions, deterministic transitions
- Bold white lines are walls that agent cannot cross
- Features - (x, y) coordinates of agent state
- Agent algorithm - soft Q-learning^[1]
- Every iteration involves data collection, 5 potential function update steps, and 10 Q-function update steps

[1] Reinforcement learning with deep energy-based policies; Haarnoja et al.; ICML 2017

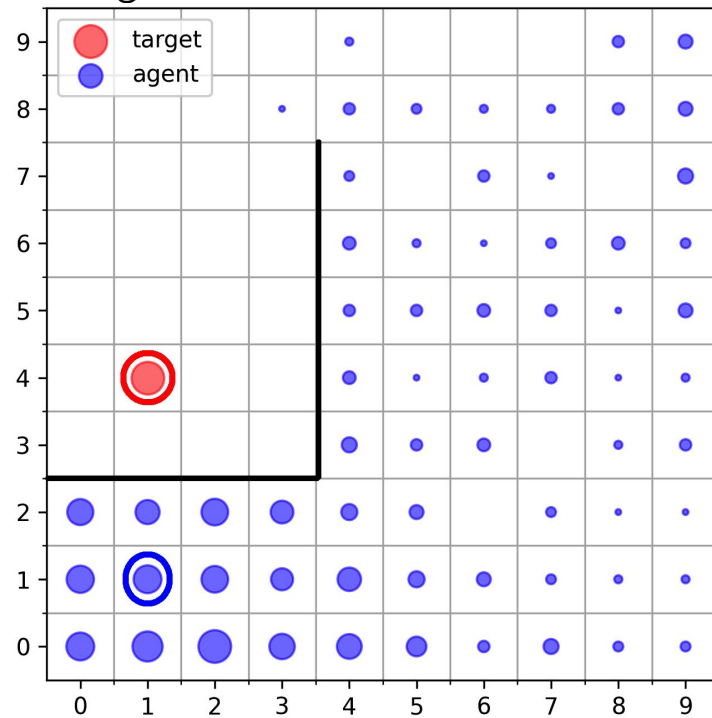


Experiments: Grid World

Agent with AIM - 500 iterations

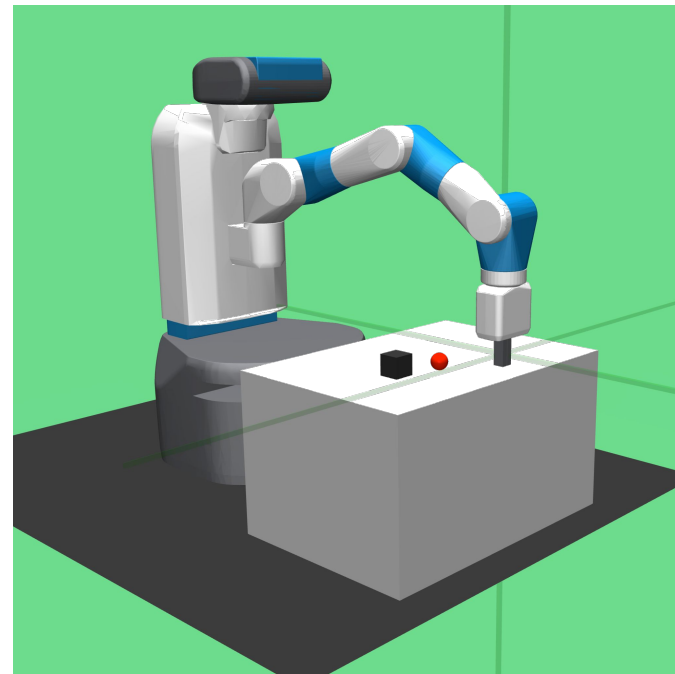


Agent without AIM - 500 iterations



Experiments: Fetch Domain

- MuJoCo gym environment
- Continuous state and action space
- Various tasks: Reach, Push, Slide, and Pick and Place
- AIM combined with HER^[1] (AIM + HER)
- Policy trained with TD3^[2]



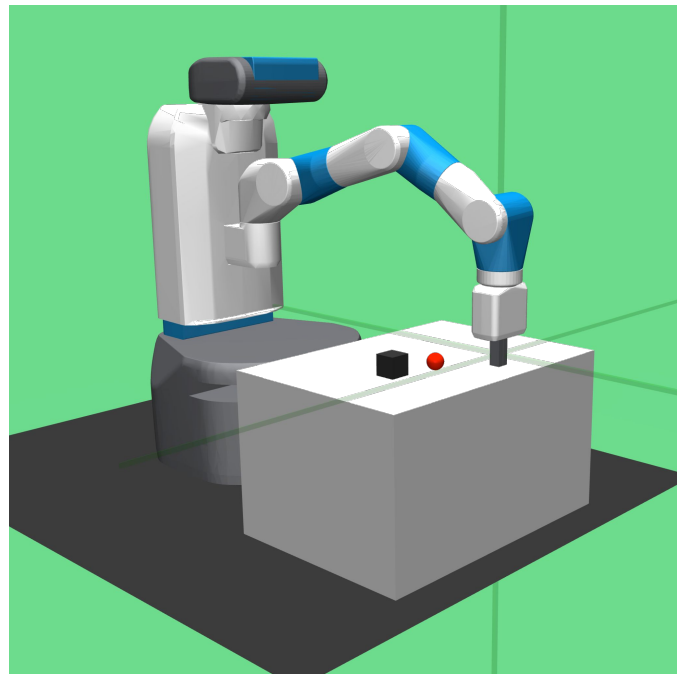
[1] Hindsight experience replay; Andrychowicz et al.; NeurIPS 2017

[2] Addressing function approximation error in actor-critic methods; Fujimoto et al.; ICML 2018

Experiments: Fetch Domain

Baselines:

- Only sparse reward (R + HER)
- Exact distance to goal (-L2 + HER)
 - Oracle reward
- Distance learned via regression from MC rollouts^[1] (MC + R + HER)
- General exploration bonus^[2] (RND + R + HER)
- GAIL^[3] reward from hindsight trajectories (GAIL + R + HER)

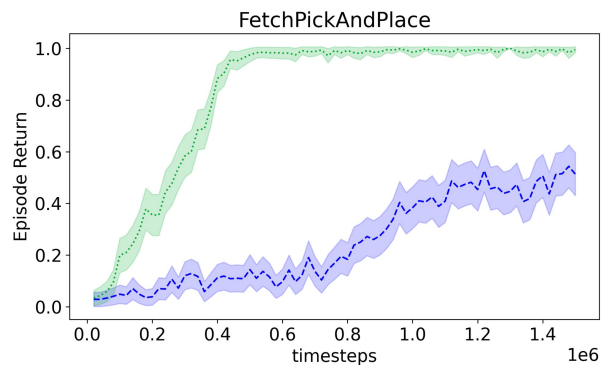
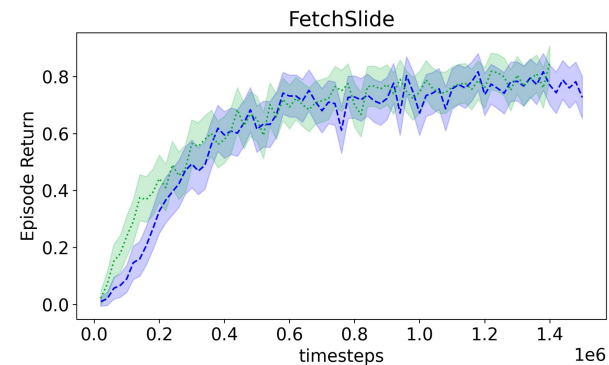
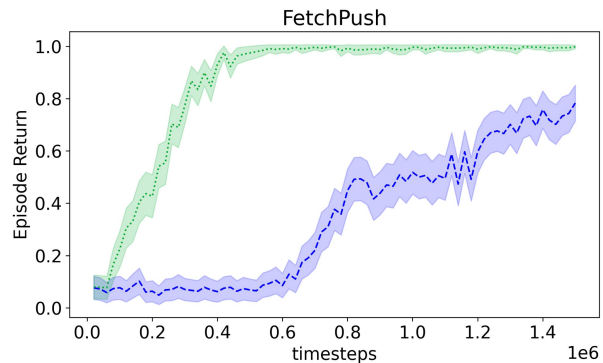
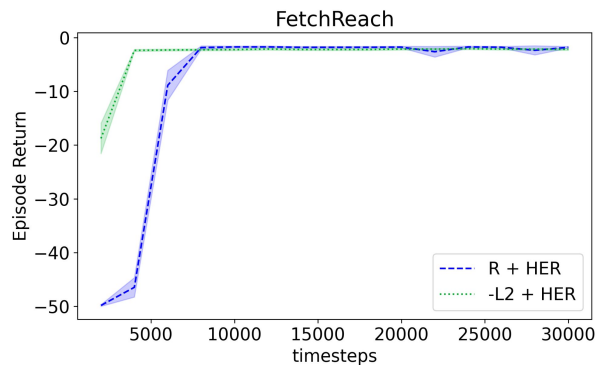


[1] Dynamical distance learning for semi-supervised and unsupervised skill discovery; Hartikainen et al.; ICLR 2020

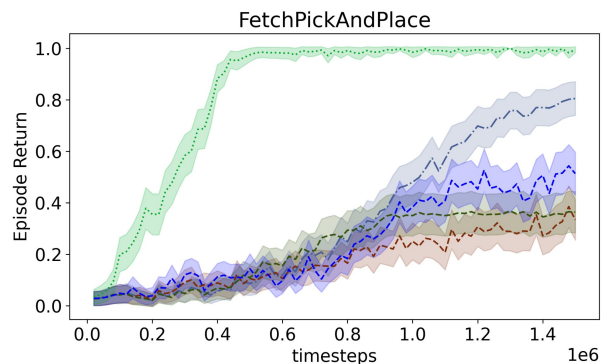
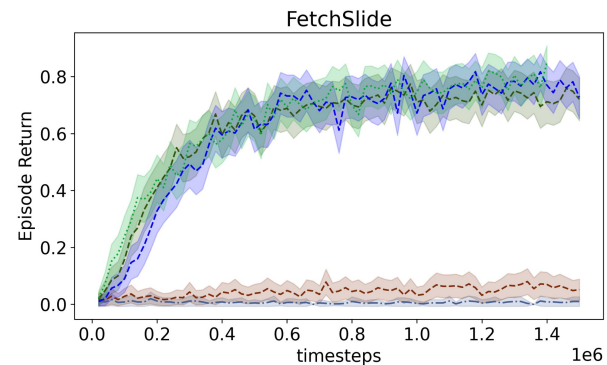
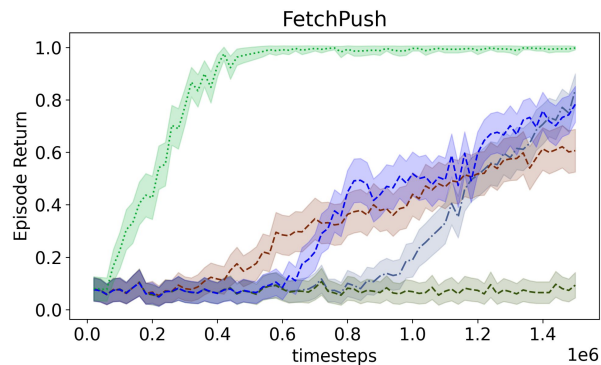
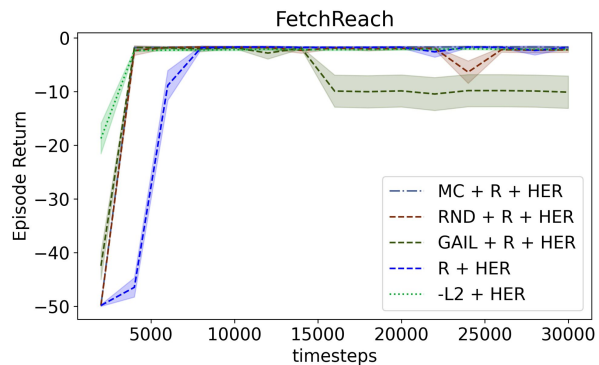
[2] Exploration by random network distillation; Burda et al.; ICLR 2019

[3] Generative adversarial imitation learning; Ho and Ermon; NeurIPS 2016

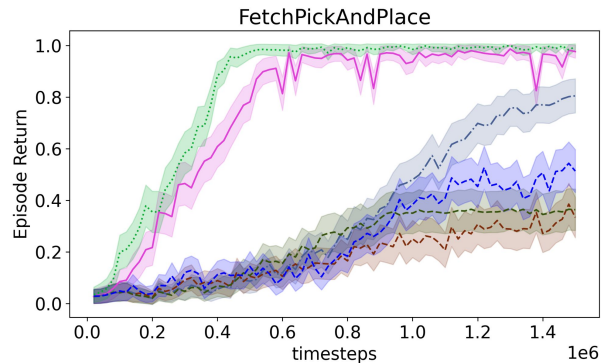
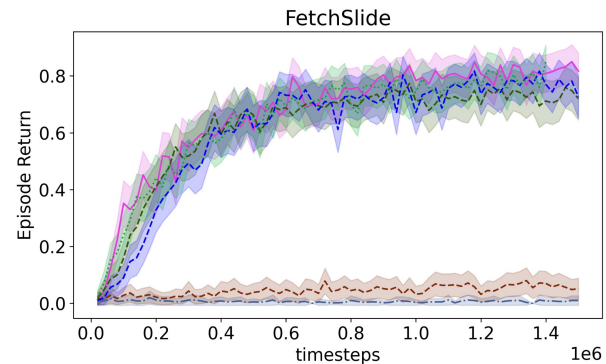
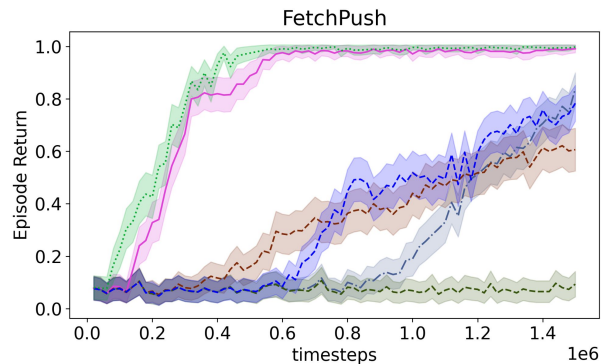
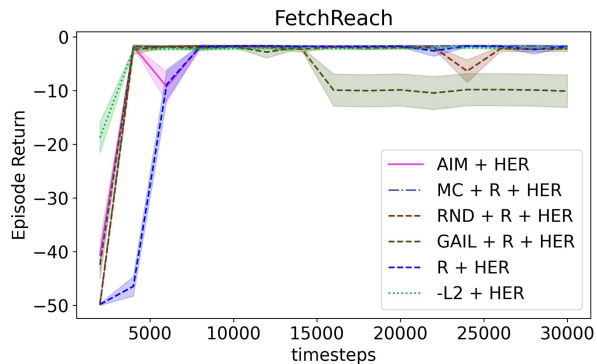
Experiments: Fetch Domain



Experiments: Fetch Domain



Experiments: Fetch Domain



Adversarial Intrinsic Motivation – Summary

- Considering the goal-conditioned RL problem through a perspective of distribution mismatch minimization.
- Requires use of the Wasserstein distance.
- Introduce a novel regularization objective for estimating Wasserstein distance in MDPs
- Compare learning under AIM with learning with sparse reward in goal-conditioned RL
- Experimental validation

Overview

1. Overcoming policy sampling error	2. Simulator grounding as imitation from observations (IfO)	3. An algorithm for sim-to-real transfer
4. Time-step metric for estimating Wasserstein distance in MDPs	5. Learning a goal conditioned policy by minimizing W_1 distance	6. Multi-agent coordination

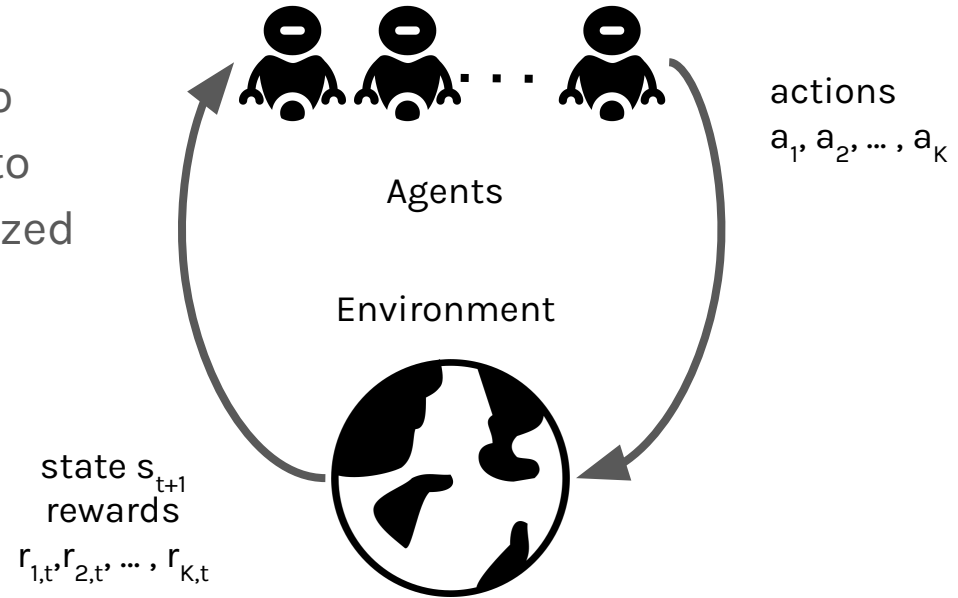
Policy

Transitions

States

Multi-agent Coordination via Distribution Matching^[1]

- Contribution: Distribution matching as a novel method to present a coordination signal to agents learning in a decentralized manner



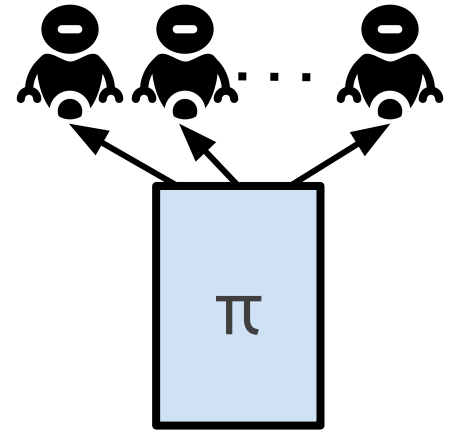
[1] DM²: Decentralized Multi-Agent Reinforcement Learning via Distribution Matching; *Wang, C., *Durugkar, I., *Liebman, E., Stone P.; AAAI 2023

* - joint first authors

Ishan Durugkar, UT Austin

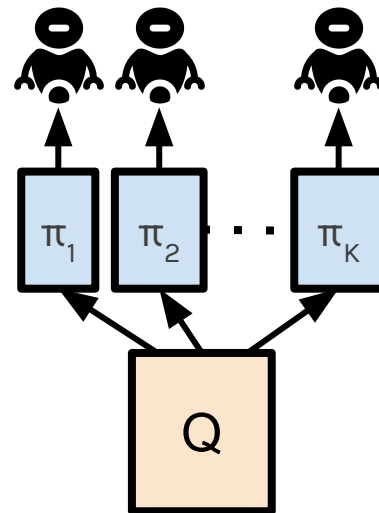
Motivation

- Multi-agent reinforcement learning (MARL) is challenging – agents learning simultaneously makes the environment nonstationary
- Strategies:
 - Fully centralized learning



Motivation

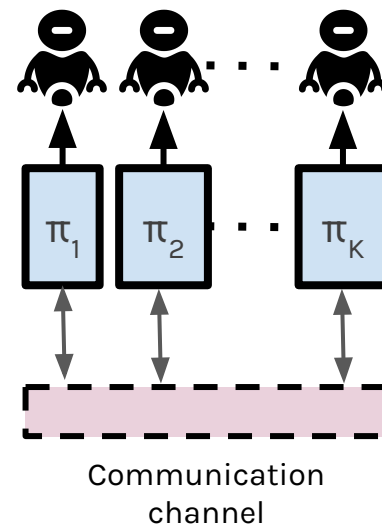
- Multi-agent reinforcement learning (MARL) is challenging – agents learning simultaneously makes the environment nonstationary
- Strategies:
 - Fully centralized learning
 - Centralized training, decentralized execution (CTDE) ^[1]



[1] Sunehag et al., Value Decomposition Networks for Cooperative Multiagent learning, AAMAS 2018.

Motivation

- Multi-agent reinforcement learning (MARL) is challenging – agents learning simultaneously makes the environment nonstationary
- Strategies:
 - Fully centralized learning
 - Centralized training, decentralized execution (CTDE) ^[1]
 - Decentralized learning + communication ^[2]

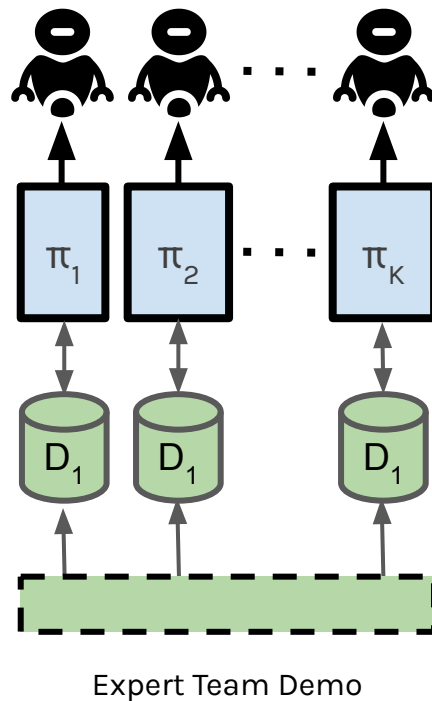


[1] Sunehag et al., Value Decomposition Networks for Cooperative Multiagent learning, AAMAS 2018.

[2] Jaques et al., Social Influence as Intrinsic Motivation for Multi-Agent Deep Reinforcement Learning, ICML 2019.

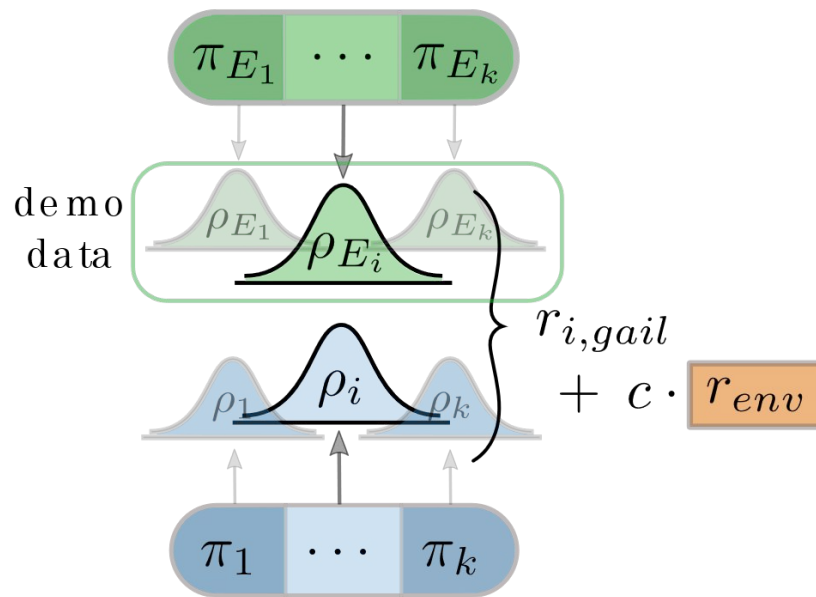
DM²: Decentralized MARL with Distribution Matching

- Control agent visitation distributions to induce coordination between agents learning independently
- The target distribution acts as coordination signal



DM²: Decentralized MARL with Distribution Matching

- Individual agents distribution matching to target distributions induced by demonstrations from coordinated expert demonstrations
- Distribution matching reward combined with task reward



Analysis

Theorem 7: Each agent maximizing its individual return over the individual distribution matching rewards r_{ϕ_i} will converge to the joint expert policy π_E

Analysis

Theorem 7: Each agent maximizing its individual return over the individual distribution matching rewards r_{ϕ_i} will converge to the joint expert policy π_E

If the expert policies are optimal with respect to the shared task, then π_E is a Nash equilibrium for rewards that are a linear combination of the task and distribution matching reward.

Experimental Setting

- StarCraft II Multi-Agent Challenge^[1] tasks
 - 5m vs 6m (5v6)
 - 3s vs 4z (3sv4z)
- Baselines w/environment reward alone
 - IPPO (decentralized)
 - QMIX^[2] (CTDE)
 - R-MAPPO^[3] (CTDE)
- Distribution Matching Baseline: DM² w/SIL^[4]

[1] Samvelyan et al., The StarCraft Multi-Agent Challenge, AAMAS 2019.

[2] Rashid et al., Qmix: Monotonic Value Function Factorisation for Deep Multi-agent Reinforcement Learning, ICML 2018.

[3] Yu et al., The Surprising Effectiveness of PPO in Cooperative, Multi-Agent Games, ArXiv 2021.

[4] Oh et al., Self-Imitation Learning, ICML 2018.

Experimental Setting

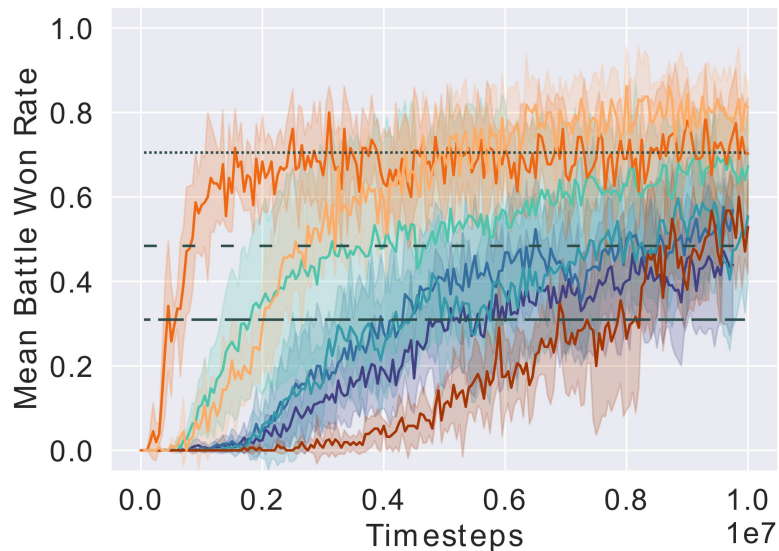
- MARL algorithm: Independent PPO (IPPO)^[1]
- Demonstrations from K experts
 - State-only demonstrations sampled from saved IPPO **and** QMIX checkpoints
- Per-agent reward function:

$$r_{i,mix} = r_{env} + r_{i,GAIL} * c$$

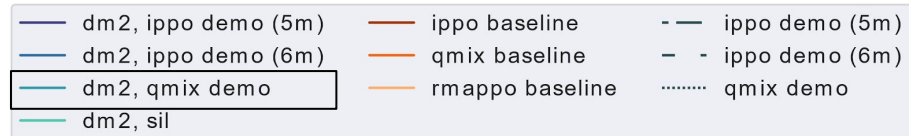
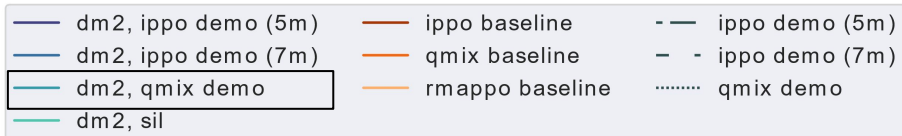
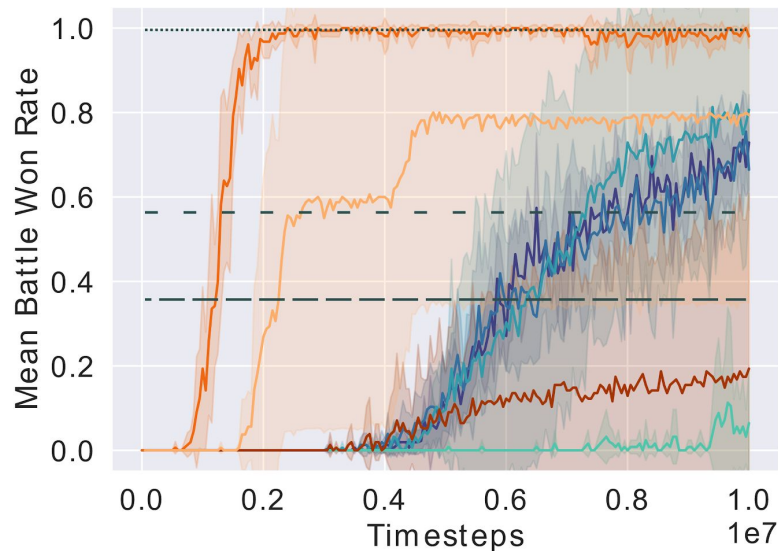
[1] Yu et al., The Surprising Effectiveness of PPO in Cooperative, Multi-Agent Games, ArXiv 2021.

Results – DM²

5v6



3sv4z



Multi-agent Coordination – Summary

- Controlling the state visitation distributions of individual agents can be a strategy for multi-agent coordination
- Can speed up learning for tasks, and improve upon performance of target distributions

Overview

1. Overcoming policy sampling error	2. Simulator grounding as imitation from observations (IfO)	3. An algorithm for sim-to-real transfer
4. Time-step metric for estimating Wasserstein distance in MDPs	5. Learning a goal conditioned policy by minimizing W_1 distance	6. Multi-agent coordination

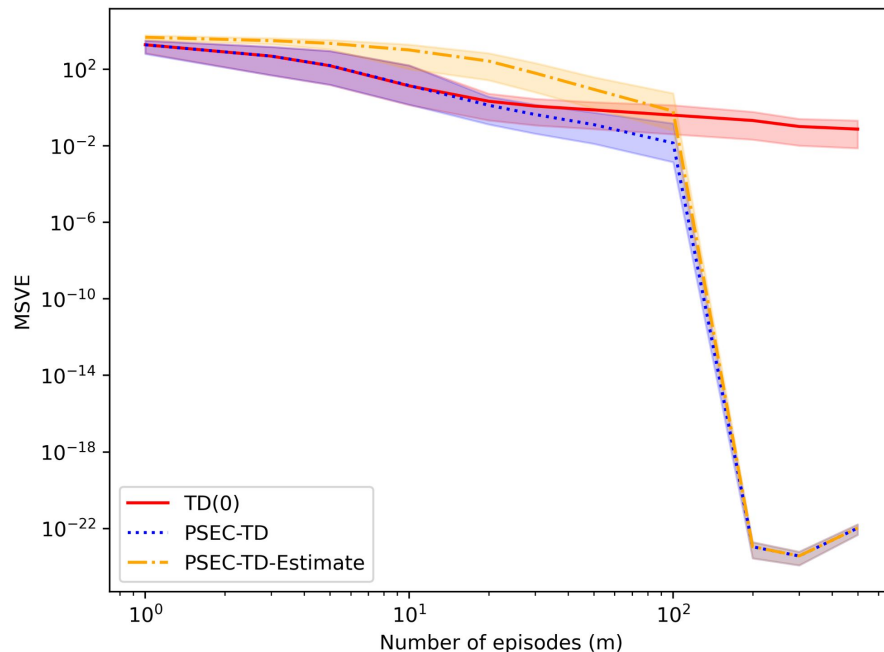
Policy

Transitions

States

Future Work - Estimation

- PSEC-TD(0) showed that it can eliminate sampling error
- When evaluating a batch of data from a mixture of policies, initial experiments indicate that with enough data, PSEC might work as expected.
- More investigation is needed



Future Work - Estimation

- Estimation - long term research:
 - Combination of distributional RL and successor features
 - Impact of other distribution estimation techniques (diffusion, density estimation, and others)

Future Work – Minimizing Distribution Mismatch

- GARAT showed that learning an action transformation function can be seen as a distribution mismatch problem
- What other problems can benefit similarly?
- Short term avenue:
 - Other objectives for behavioral cloning^[1]
- Long term avenue:
 - Learning a dynamics model of the environment

[1] ABC: Adversarial Behavioral Cloning for Offline Mode-seeking Imitation Learning; Hudson, E., **Durugkar, I.**, Warnell, G., Stone, P.; ArXiv 2022

Future Work – Extending AIM

- Use of the Wasserstein distance to measure distance between distributions
- Considering RL problems as controlling visitation distributions
- Short term:
 - Exploration
 - Beyond goal-conditioned RL
 - Skill learning
- Long term:
 - Distribution control for general reward functions

Future Work – Distribution Control in MARL

- DM^2 has opened doors for the kind of impact distribution control can have in MARL
- Potential avenues:
 - Better coordination techniques
 - Beyond cooperative tasks
 - Bootstrapping K -expert demonstrations to N agents ($N > K$)

Related Work

- Generative Adversarial Nets; Goodfellow et al.; NeurIPS 2014

Imitation Learning

- Generative adversarial imitation learning; Ho and Ermon; NeurIPS 2016

Off-Policy Evaluation

- Breaking the curse of horizon: infinite-horizon off-policy estimation; Liu et al.; NeurIPS 2018
- DualDICE: Behavior-agnostic estimation of discounted stationary distribution corrections; Nachum et al.; NeurIPS 2019

Distributional RL

- A distributional perspective on reinforcement learning; Bellemare et al.; ICML 2017

Exploration

- Provably efficient maximum entropy exploration; Hazan et al.; ICML 2019
- Efficient exploration via state marginal matching; Lee et al.; ArXiv 2019

Summary

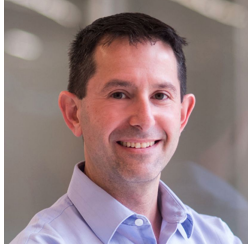
How can methods for estimating and controlling an agent's visitation distribution be beneficially incorporated into RL algorithms?

- Variety of problems benefit from estimating or controlling visitation distributions
- Various actionable insights and broader implications for future work

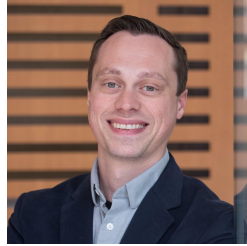
Thank you!



Ishan
Durugkar



Peter Stone



Scott
Niekum



Garrett
Warnell



Josiah
Hanna



Elad
Liebman



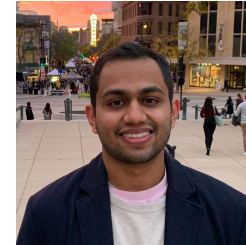
Mauricio
Tec



Siddharth
Desai



Haresh
Karnan



Brahma
Pavse



Caroline
Wang

Questions?

How can methods for estimating and controlling an agent's visitation distribution be beneficially incorporated into RL algorithms?

1. Overcoming policy sampling error	2. Simulator grounding as imitation from observations (IfO)	3. An algorithm for sim-to-real transfer
4. Time-step metric for estimating Wasserstein distance in MDPs	5. Learning a goal conditioned policy by minimizing W_1 distance	6. Multi-agent coordination