

Imitation Learning from Observation

PhD Defense

Faraz Torabi

Supervisor: Peter Stone

University of Texas at Austin

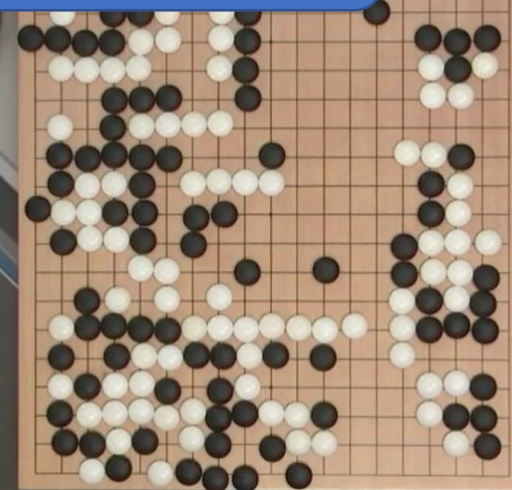
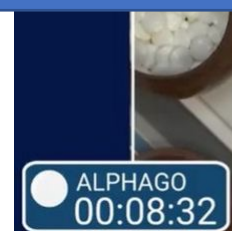
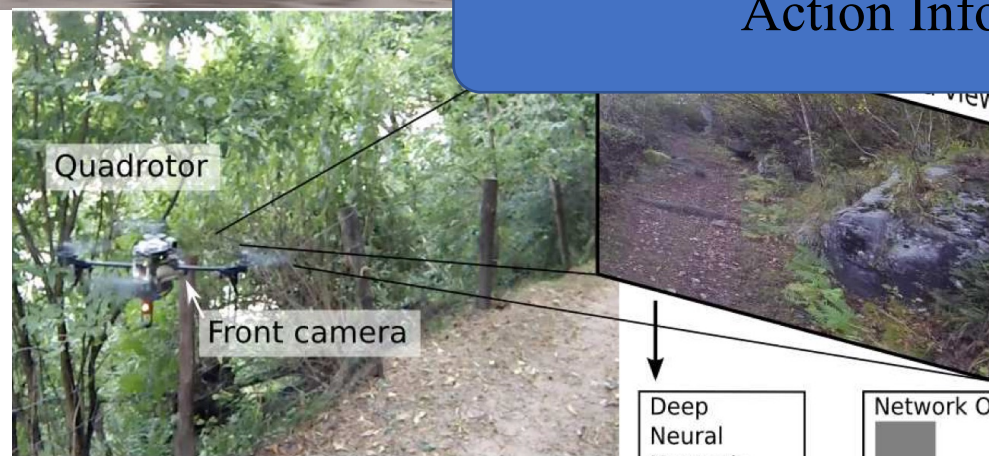
faraztrb@cs.utexas.edu



Success of Imitation Learning



Action Information is needed.



Humans Mostly Learn by Observation



Example: Watching YouTube

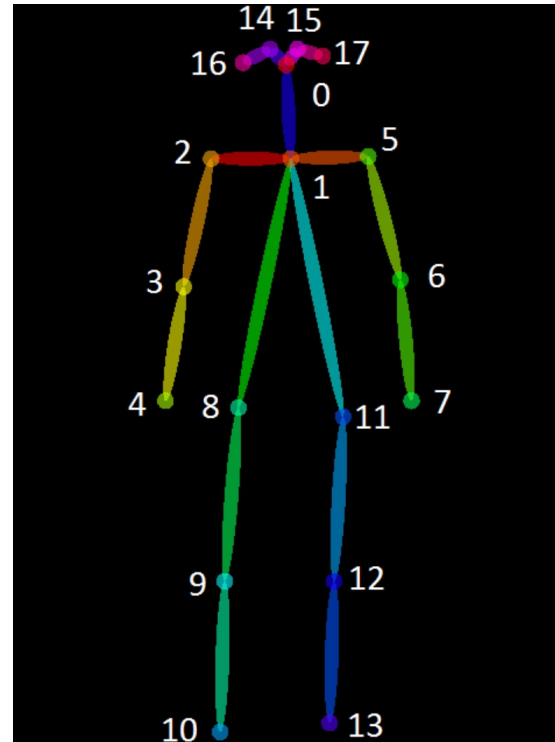


Imitation Learning from Observation



Visual observations

Perception →



States

Control →

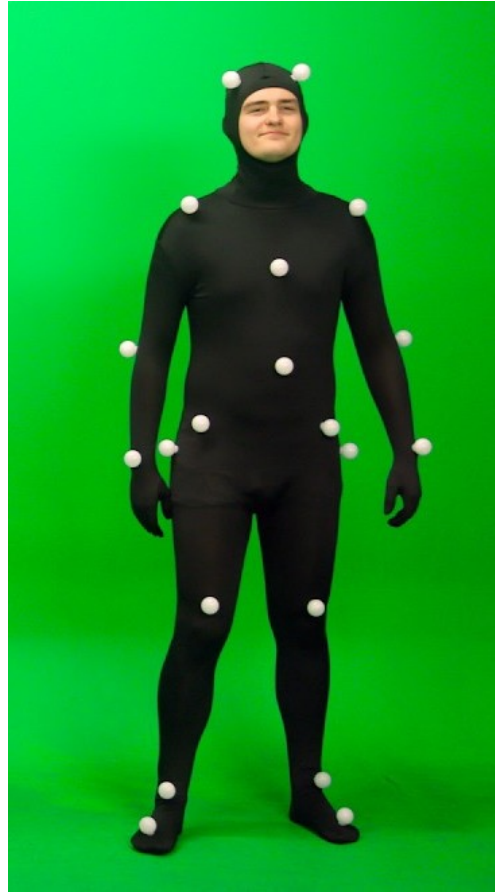


Policy

Perception Module



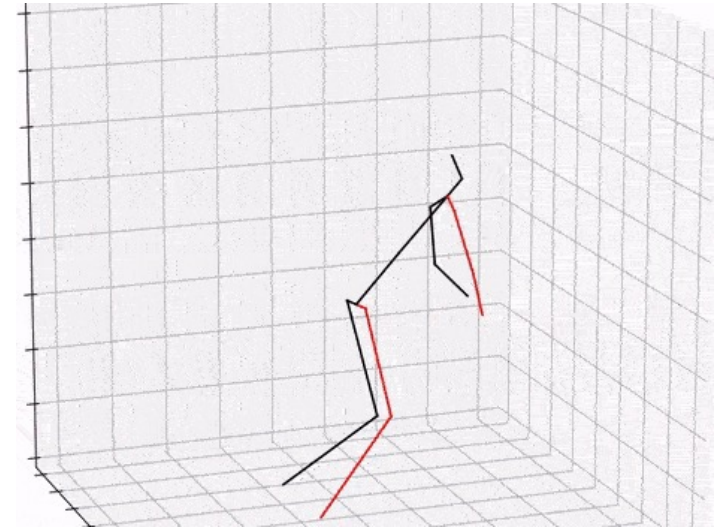
Sensors



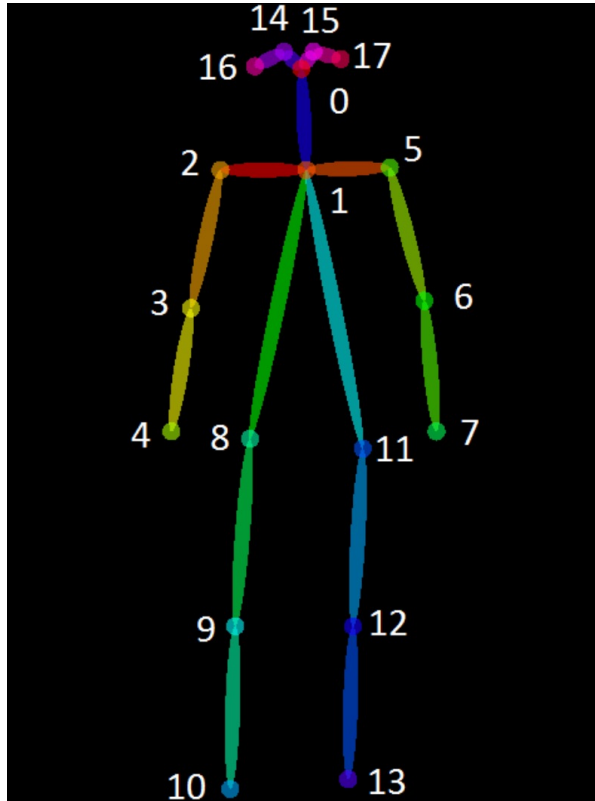
Motion Capture



Pose Estimation



Focus of My Research is on the ...



Control

Research Question

Behavioral Cloning from
Observation (BCO)
[IJCAI 2018]

Generative Adversarial
Imitation from Observation
(GAIfO)
[AAMAS 2019, ICML
Workshop]

Data-Efficient Adversarial
Learning for Imitation from
Observation (DEALIO)
[Under Review]

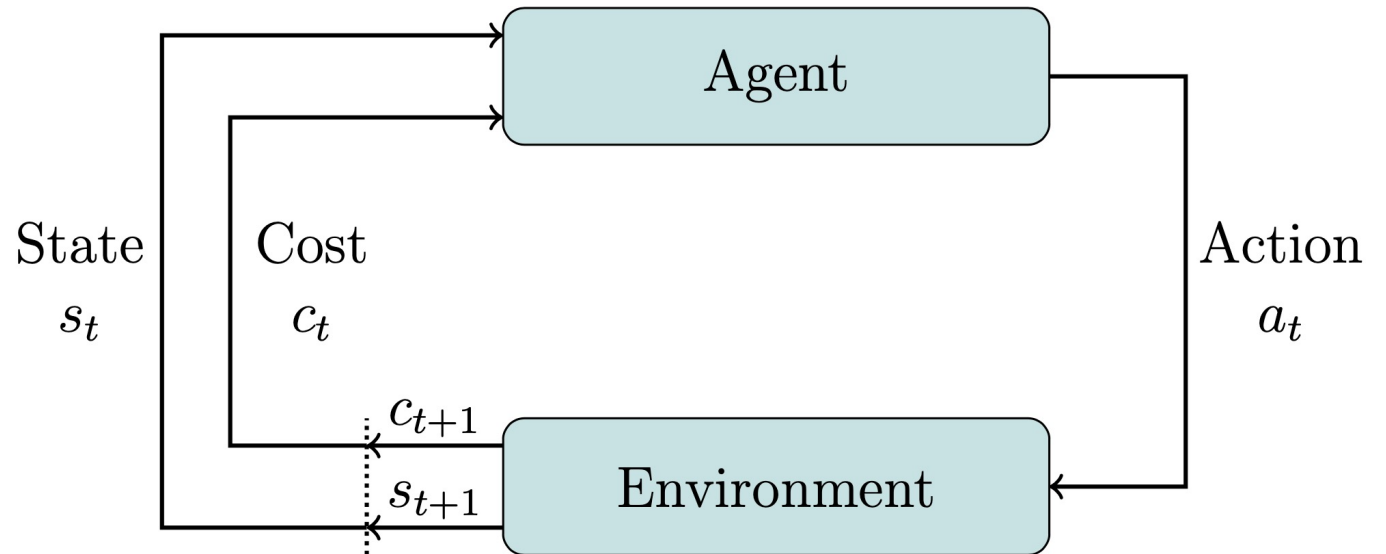
Reinforced Inverse Dynamics
Modeling (RIDM)
[RAL, IROS 2020]

Visual Extension of GAIfO
with Self-observation
[ICML Workshop]

Visual Extension of GAIfO
with Proprioceptive Information
[IJCAI 2020]

Reinforcement Learning

- Goal:
 - Learn how to make decisions by minimizing the cumulative cost feedback.
- $M = \langle S, A, P, c \rangle$
 - S : Set of states
 - A : Set of actions
 - P : Transition function
 - c : Cost function
- Learn a policy $\pi: S \rightarrow A$



Reinforcement Learning

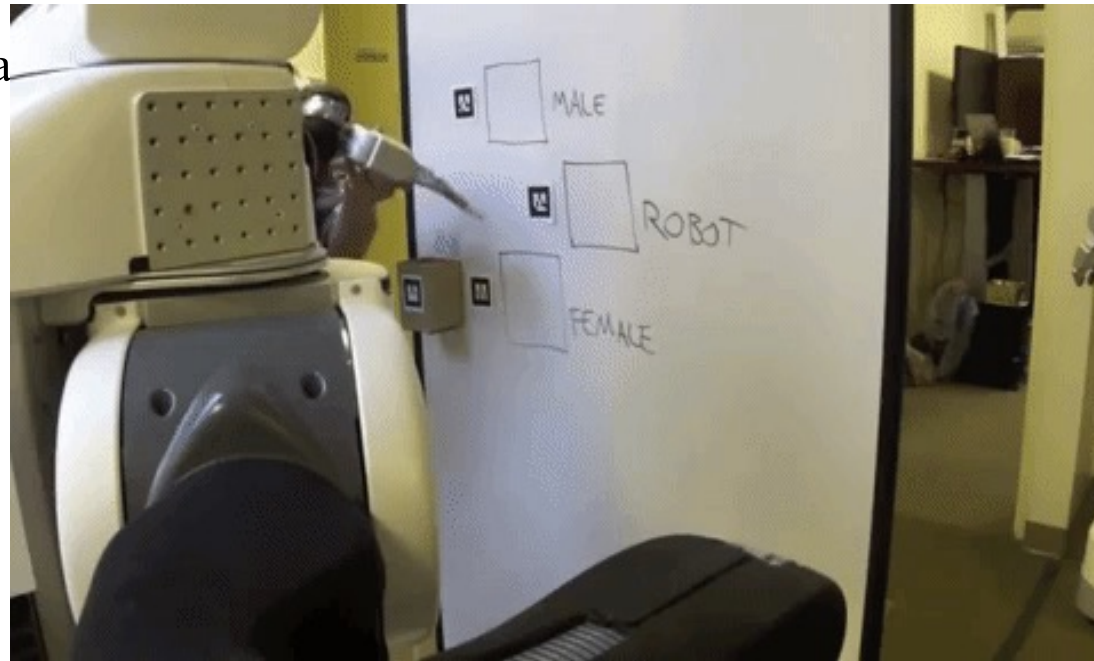
- Algorithms:
 - Model-based:
 - Known model (planning): LQR, MCTS, etc.
 - Unknown model: PILQR, PLATO, etc.
 - Model-free:
 - Policy-based: Reinforce, TRPO, PPO, etc.
 - Value-based: SARSA, Q-learning, etc.

Imitation Learning

- Goal:
 - Learn how to make decisions by trying to imitate another agent.
- $M \setminus c$:
 - Provided: $\tau_i^e = \{(s_0^e, a_0^e), (s_1^e, a_1^e), \dots, (s_N^e, a_N^e)\}_i$
 - Learn: $\pi: S_t \rightarrow a_t$
- Algorithms:
 - Behavioral Cloning (BC)
 - Inverse Reinforcement Learning (IRL)
 - Adversarial Imitation Learning (AIL)

Imitation Learning

- Observations of other agent (demonstrations) consist of state-action pairs.
- Limitation:
 - Precludes using a large number of demonstrations that are not given.



Scott Niekum et al. "Learning and generalization of complex tasks from unstructured demonstrations". In: Intelligent Robots and Systems (IROS), 2012

Imitation Learning from Observation

- Goal: from state-only demonstrations of
 - Learn how to perform a task ~~by visually observing~~ an expert.
- $M \setminus c$:
 - Provided: ~~$\tau_i^e = \{o_0^e, o_1^e, \dots, o_N^e\}_i$~~ $\tau_i^e = \{s_0^e, s_1^e, \dots, s_N^e\}_i$
 - Learn: $\pi : s_t \rightarrow a_t$



In what ways can autonomous agents learn to imitate experts using
state-only observations?

Behavioral Cloning from
Observation (BCO)
[IJCAI 2018]

Generative Adversarial
Imitation from Observation
(GAIfO)
[AAMAS 2019, ICML
Workshop]

Data-Efficient Adversarial
Learning for Imitation from
Observation (DEALIO)
[Under Review]

Reinforced Inverse Dynamics
Modeling (RIDM)
[RAL, IROS 2020]

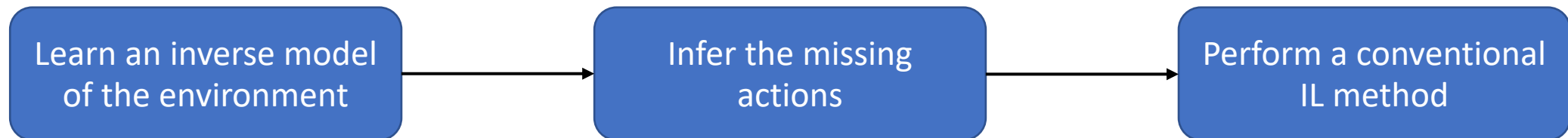
Visual Extension of GAIfO
with Self-observation
[ICML Workshop]

Visual Extension of GAIfO
with Proprioceptive Information
[IJCAI 2020]

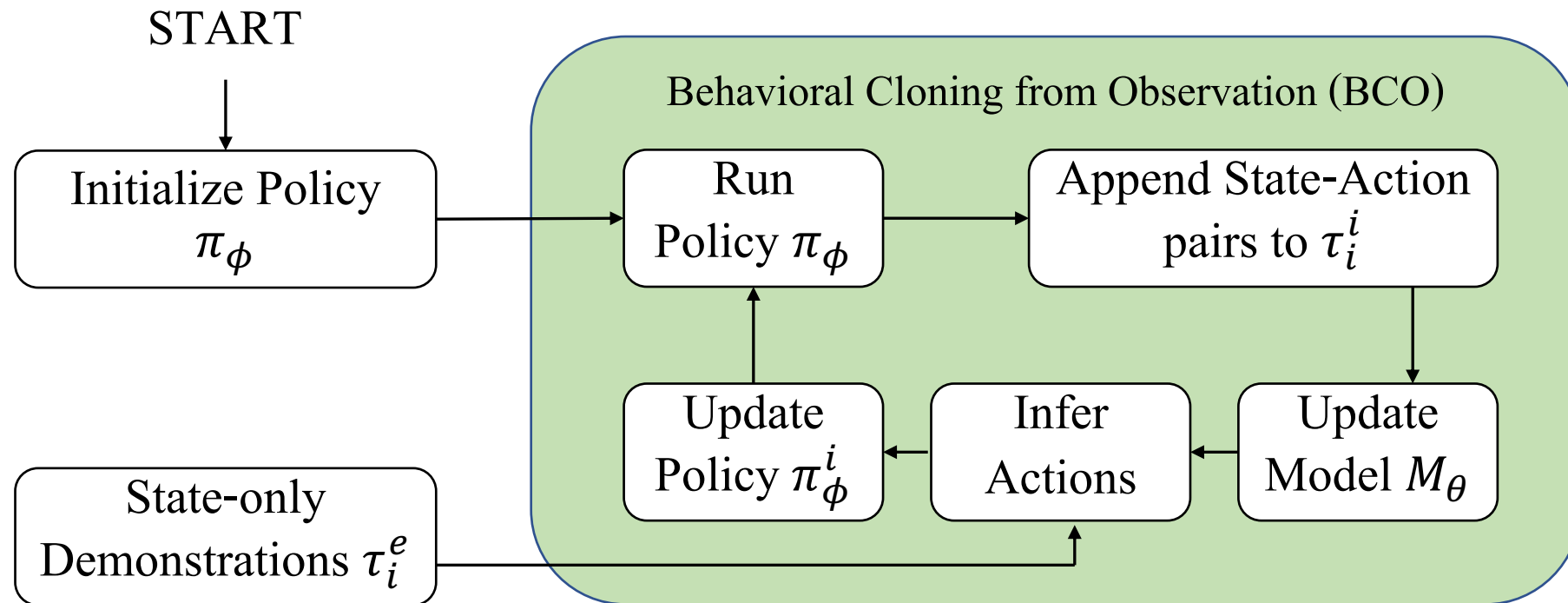
Behavioral Cloning from Observation

- Goal:
 - Propose a *Model-based Algorithm* for Imitation from Observation.
- Imitation Learning (IL): $\tau_i^e = \{(s_0^e, a_0^e), (s_0^e, a_0^e), \dots, (s_N^e, a_N^e)\}_i$
- Imitation from Observation (IfO): $\tau_i^e = \{(s_0^e, ?), (s_1^e, ?), \dots, (s_N^e, ?)\}_i$

A Model-based Approach

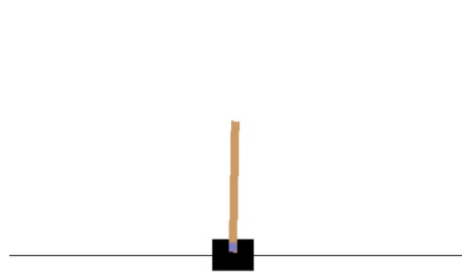


Behavioral Cloning from Observation

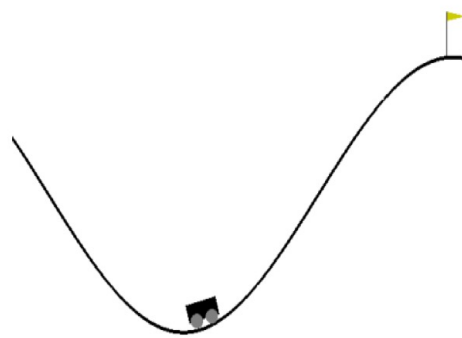


Experiments

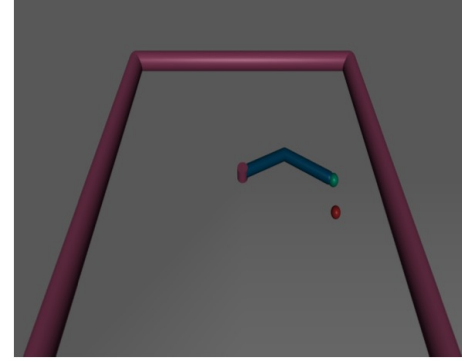
- Tasks:



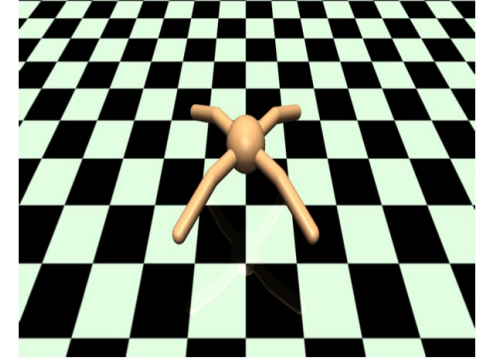
(a) CartPole



(b) MountainCar



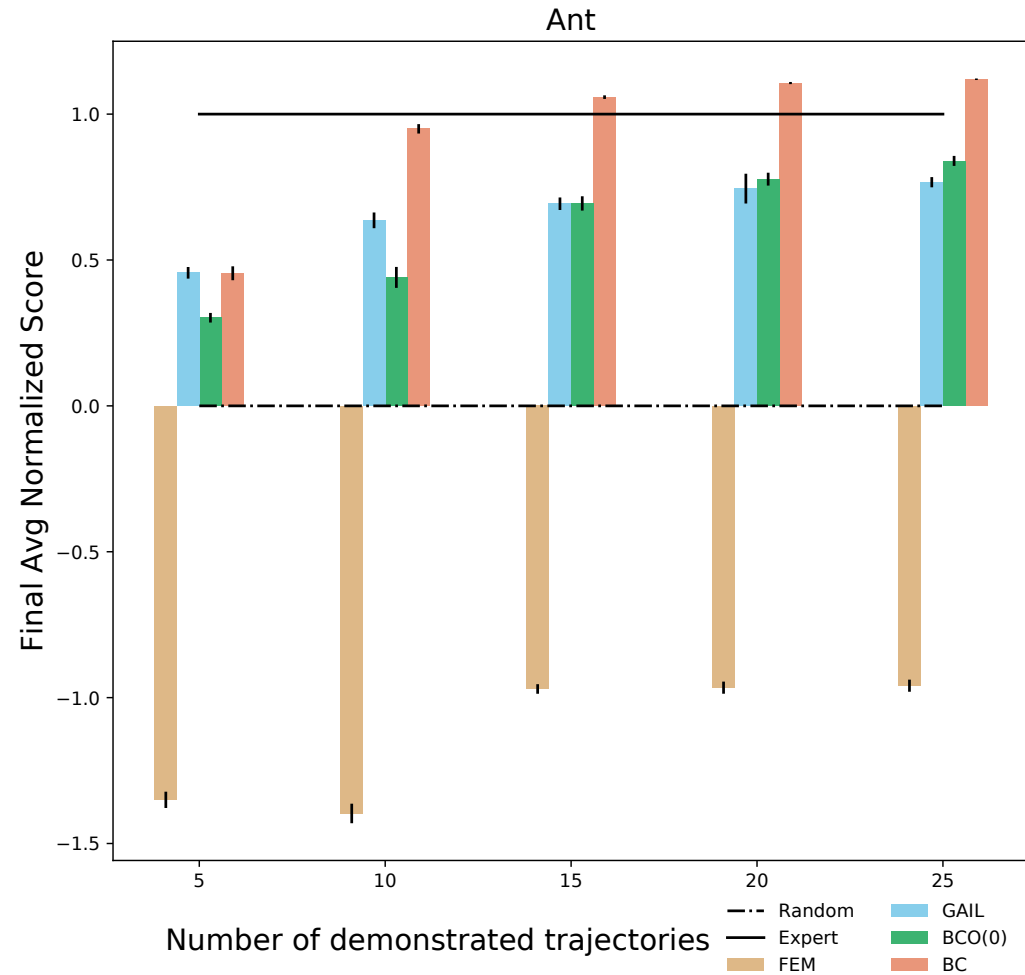
(c) Reacher



(d) Ant

Experiments

- Ant



In what ways can autonomous agents learn to imitate experts using
state-only observations?

Behavioral Cloning from
Observation (BCO)
[IJCAI 2018]

Generative Adversarial
Imitation from Observation
(GAIfO)
[AAMAS 2019, ICML
Workshop]

Data-Efficient Adversarial
Learning for Imitation from
Observation (DEALIO)
[Under Review]

Reinforced Inverse Dynamics
Modeling (RIDM)
[RAL, IROS 2020]

Visual Extension of GAIfO
with Self-observation
[ICML Workshop]

Visual Extension of GAIfO
with Proprioceptive Information
[IJCAI 2020]

Generative Adversarial Imitation from Observation

- Goal:
 - Propose a *Model-free Algorithm* for Imitation from Observation.
- IfO problem:

$$RL \circ IRLfO_{\psi}(\pi^e) = \operatorname{argmin}_{\pi \in \Pi} \operatorname{argmax}_{c \in R^{S \times S}} -\psi(c) + (\min_{\pi \in \Pi} E_{\pi}[c(s, s')]) - E_{\pi^e}[c(s, s')]$$

- Which is a composition of:

$$IRLFO_{\psi}(\pi^e) = \operatorname{argmax}_{c \in R^{S \times S}} -\psi(c) + (\min_{\pi \in \Pi} E_{\pi}[c(s, s')]) - E_{\pi^e}[c(s, s')]$$

$$RL(\tilde{c}) = \operatorname{argmin}_{\pi \in \Pi} E_{\pi}[\tilde{c}(s, s')]$$

- It is equivalent to solving:

$$\operatorname{argmin}_{\pi \in \Pi} \operatorname{argmax}_{D \in (0,1)^{S \times S}} [\log(\pi(s, s')) + E_{\pi}[\log(1 - D(s, s'))]]$$

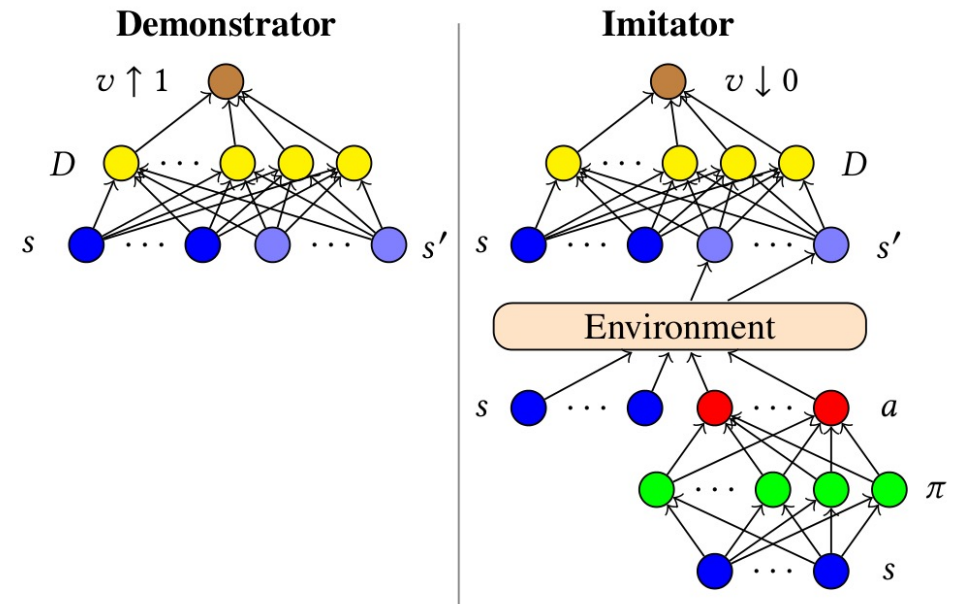
Difficult to Solve

- $c(s, s')$: Cost as a function of state transition
- π^e : Expert policy
- Π : Set of all possible policies
- $\psi(c)$: Regularizer

Generative Adversarial Imitation from Observation

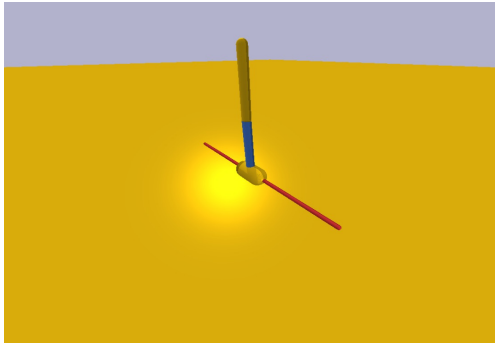
- Algorithm:
 - Initialize π_ϕ and D_θ
 - While π_ϕ improves do:
 - Execute π_ϕ and store state transitions $\tau_i^i = \{(s^i)\}_i$
 - Update D_θ using loss:
$$-(E_\pi[\log(D(s, s'))] + E_{\pi^e}[\log(1 - D(s, s'))])$$
 - Update π_ϕ by performing TRPO updates with cost function:

$$E_\pi[\log(D(s, s'))]$$

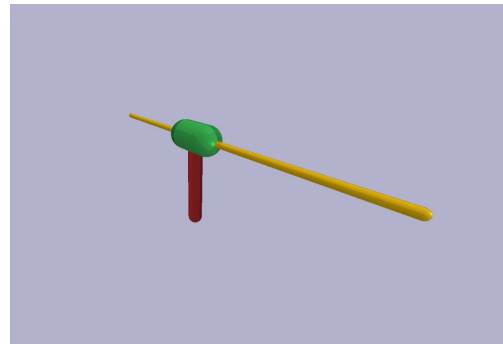


Experiments

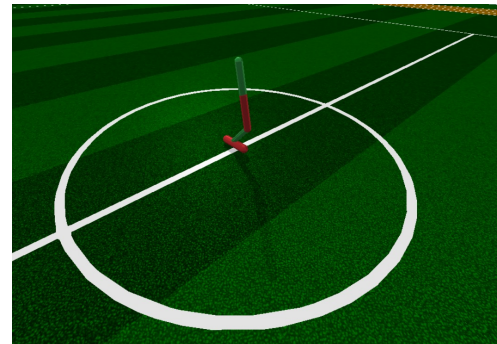
- Tasks:



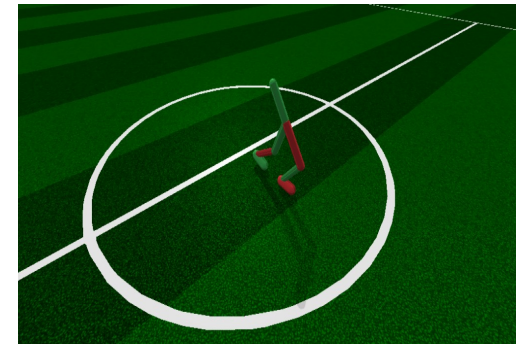
InvertedDoublePendulum



InvertedPendulumSwingup



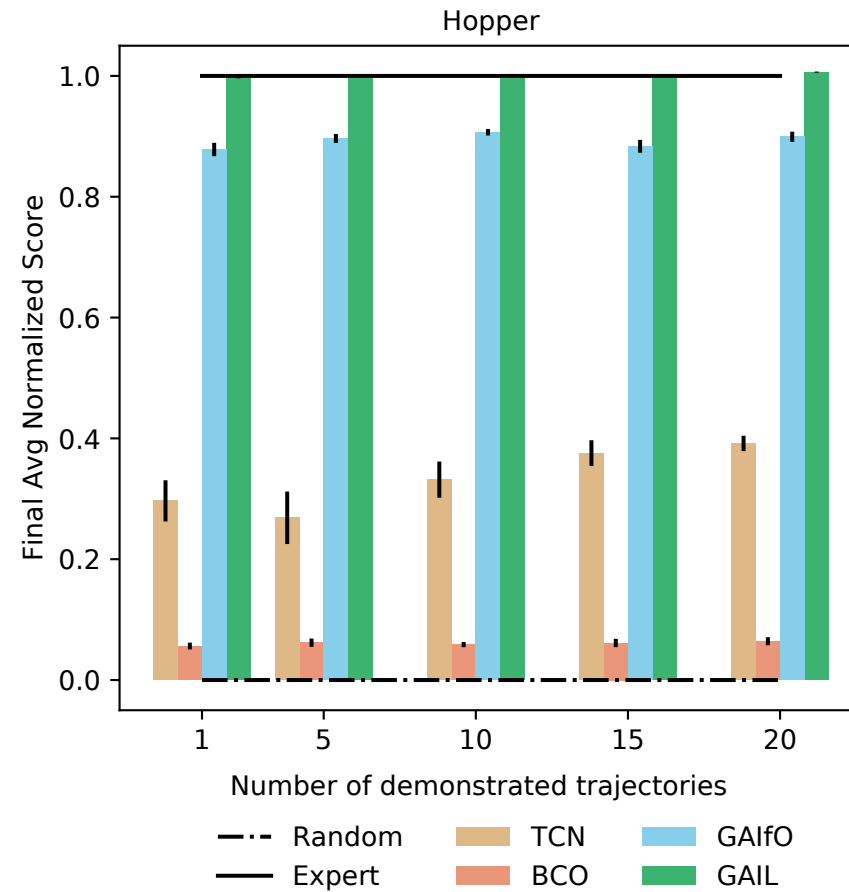
Hopper



Walker2D

Experiments

- Hopper:



Theoretical Contribution

- IfO problem:

$$RL \circ IRLfO_{\psi}(\pi^e) = \operatorname{argmin}_{\pi \in \Pi} \operatorname{argmax}_{c \in R^{S \times S}} -\psi(c) + (\min_{\pi \in \Pi} E_{\pi}[c(s, s')]) - E_{\pi^e}[c(s, s')]$$

- Equivalent to:

Difficult to Solve

$$\operatorname{argmin}_{\pi \in \Pi} \operatorname{argmax}_{D \in (0,1)^{S \times S}} E_{\pi}[\log(D(s, s'))] + E_{\pi^e}[\log(1 - D(s, s'))]$$

How are they equivalent?

Proposition 5.4.1. $RL \circ IRLfO_{\psi}(\pi^e)$ and $\operatorname{argmin}_{\pi \in \Pi} \psi^*(\rho_{\pi}^s - \rho_{\pi^e}^s)$ induce policies that have the same state transition occupancy measure, $\rho_{\tilde{\pi}}^s$.

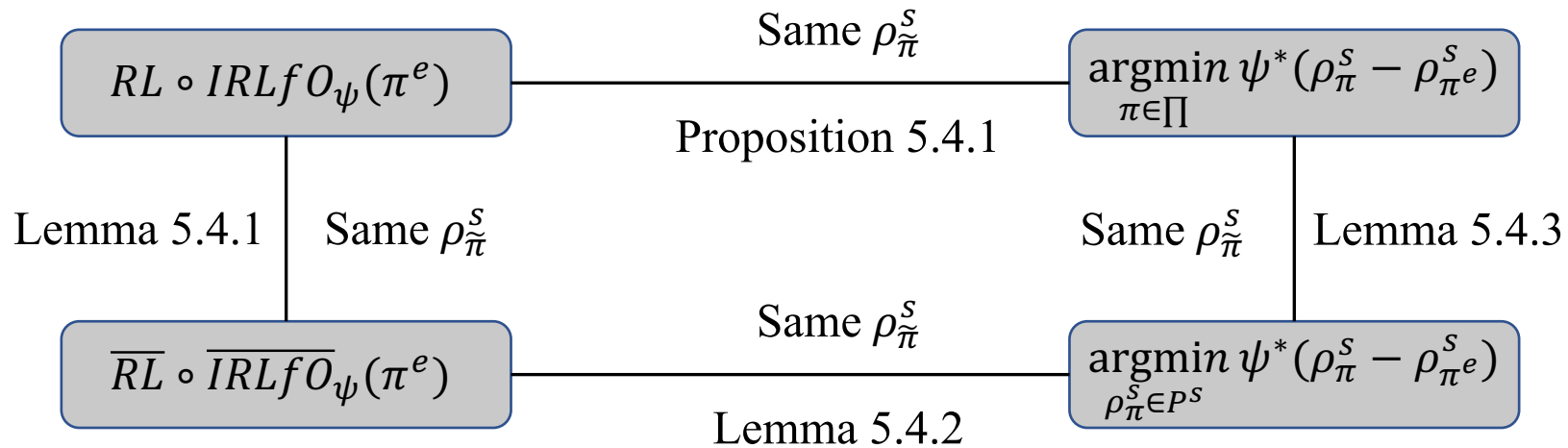
Where $\rho_{\pi}^s(s_i, s_j) = \sum_a P(s_j | s_i, a) \pi(a | s_i) \sum_{t=0}^{\infty} \gamma^t P(s_t = s_i | \pi)$

Proposition 5.4.1. $RL \circ IRLfO_\psi(\pi^e)$ and $\operatorname{argmin}_{\pi \in \Pi} \psi^*(\rho_\pi^s - \rho_{\pi^e}^s)$ induce policies that have the same state transition occupancy measure, $\rho_{\tilde{\pi}}^s$.

$$RL \circ IRLfO_\psi(\pi^e) = \operatorname{argmin}_{\pi \in \Pi} \operatorname{argmax}_{c \in R^{S \times S}} -\psi(c) + (\min_{\pi \in \Pi} E_\pi[c(s, s')]) - E_{\pi^e}[c(s, s')]$$

$$\overline{RL} \circ \overline{IRLfO}_\psi(\pi^e) = \operatorname{argmin}_{\rho_\pi^s \in P^S} \operatorname{argmax}_{c \in R^{S \times S}} -\psi(c) + (\min_{\rho_\pi^s \in P^S} \sum_{s, s'} \rho_\pi^s(s, s') c(s, s')) - \sum_{s, s'} \rho_{\pi^e}^s(s, s') c(s, s')$$

Proof:



Theoretical Contribution

- New IfO problem:

$$\operatorname{argmin}_{\pi \in \Pi} \psi^*(\rho_{\pi}^s - \rho_{\pi^e}^s)$$

- Specifying ψ :

$$\psi(c) = \begin{cases} E_{\pi^e}[g(c(s, s')))] & \text{if } c < 0 \\ +\infty & \text{otherwise} \end{cases} \quad \text{where} \quad g(x) = \begin{cases} -x - \log(1 - e^x) & \text{if } x < 0 \\ +\infty & \text{otherwise} \end{cases}$$

- The optimization problem becomes [proposition A.1.1]:

$$\operatorname{argmin}_{\pi \in \Pi} \operatorname{argmax}_{D \in (0,1)^{S \times S}} E_{\pi}[\log(D(s, s'))] + E_{\pi^e}[\log(1 - D(s, s'))]$$

Similar to the Generative Adversarial Loss

In what ways can autonomous agents learn to imitate experts using
state-only observations?

Behavioral Cloning from
Observation (BCO)
[IJCAI 2018]

Generative Adversarial
Imitation from Observation
(GAIfO)
[AAMAS 2019, ICML
Workshop]

Data-Efficient Adversarial
Learning for Imitation from
Observation (DEALIO)
[Under Review]

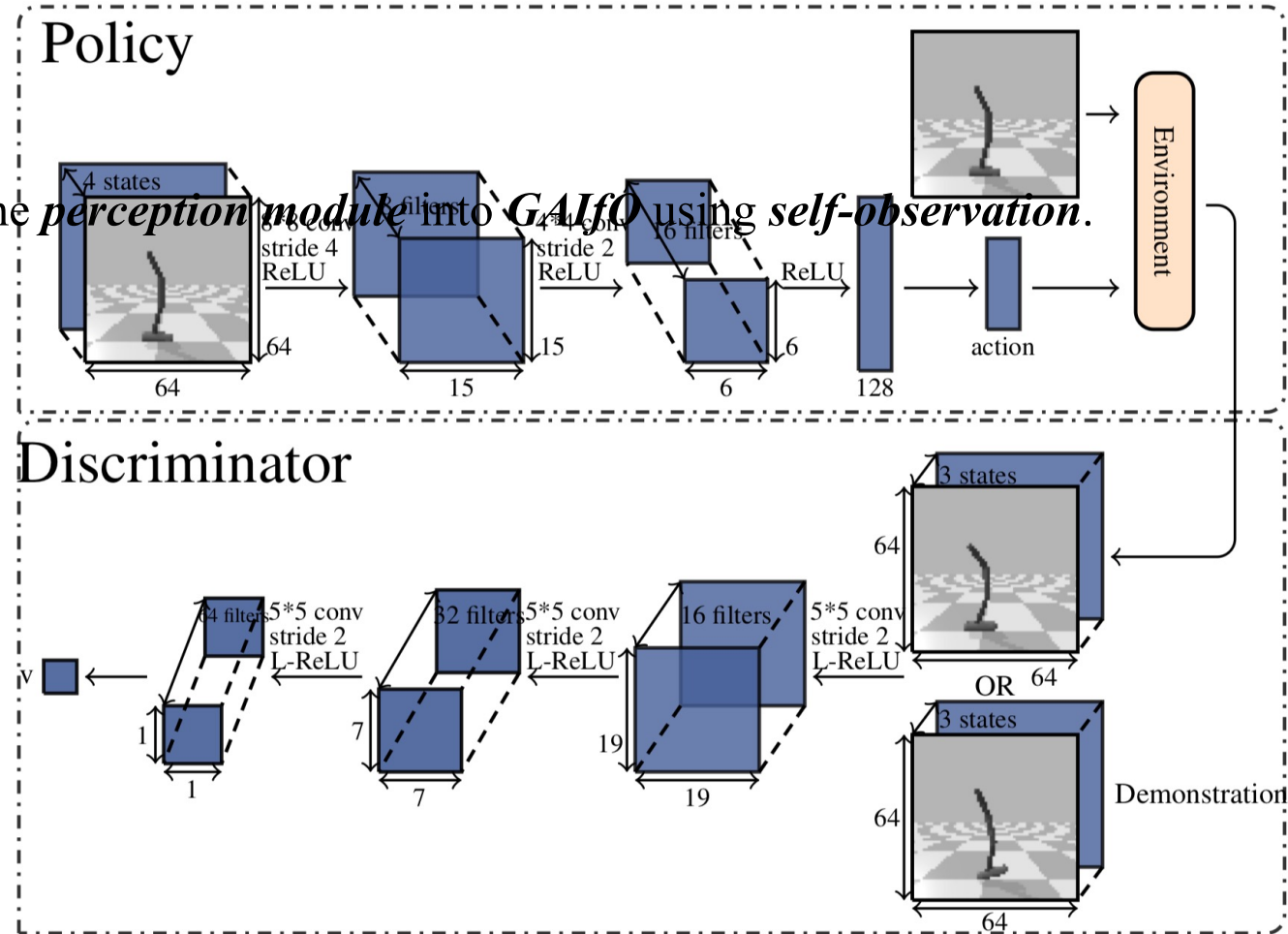
Reinforced Inverse Dynamics
Modeling (RIDM)
[RAL, IROS 2020]

Visual Extension of GAIfO
with Self-observation
[ICML Workshop]

Visual Extension of GAIfO
with Proprioceptive Information
[IJCAI 2020]

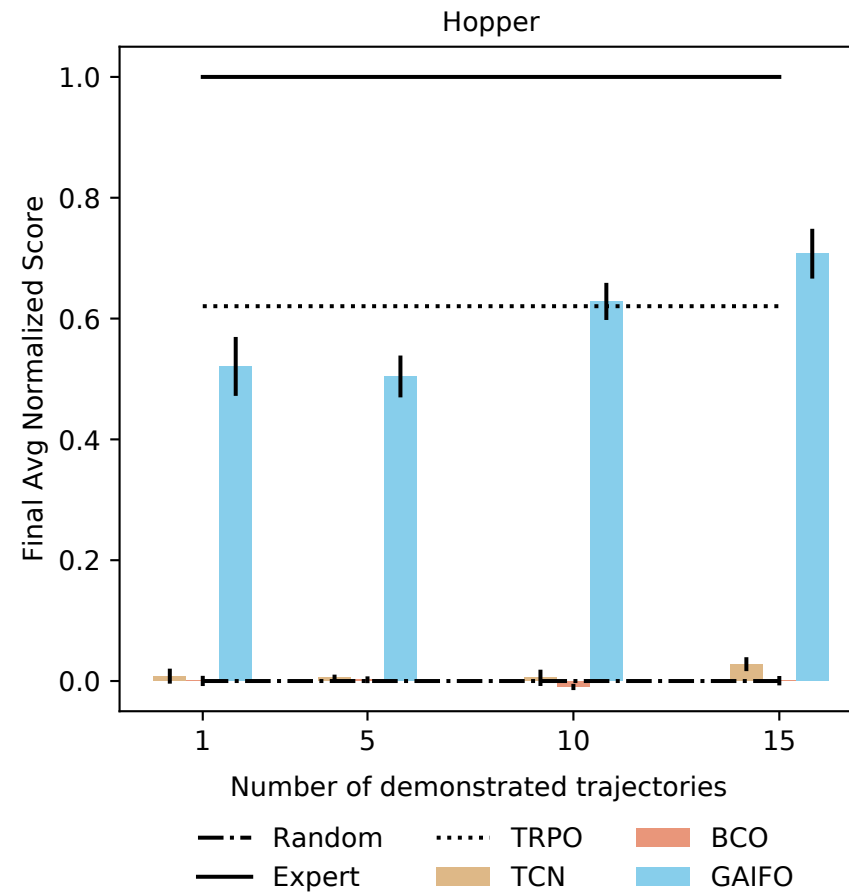
GAIfO with Self-observation

- Goal:
 - Incorporate the *perception module* into GAIfO using *self-observation*.



Experiments

- Hopper:



In what ways can autonomous agents learn to imitate experts using
state-only observations?

Behavioral Cloning from
Observation (BCO)
[IJCAI 2018]

Generative Adversarial
Imitation from Observation
(GAIfO)
[AAMAS 2019, ICML
Workshop]

Data-Efficient Adversarial
Learning for Imitation from
Observation (DEALIO)
[Under Review]

Reinforced Inverse Dynamics
Modeling (RIDM)
[RAL, IROS 2020]

Visual Extension of GAIfO
with Self-observation
[ICML Workshop]

Visual Extension of GAIfO
with Proprioceptive Information
[IJCAI 2020]

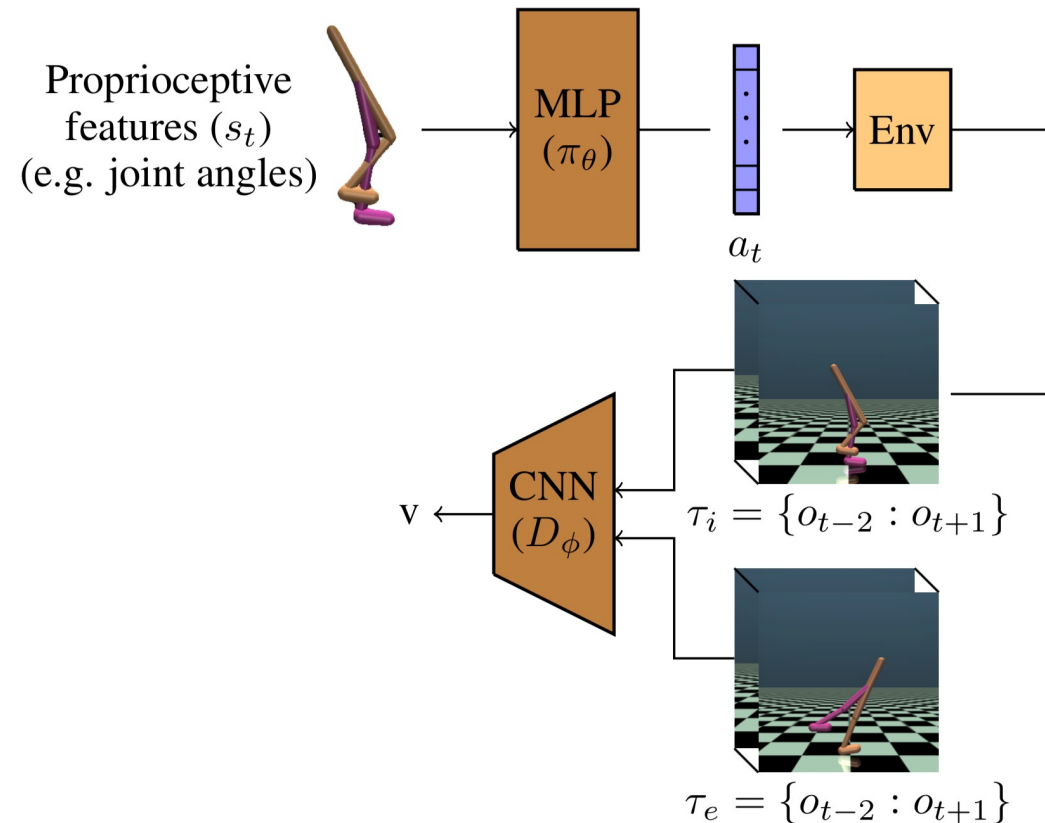
GAlfO with Proprioceptive Information

- Goal:
 - To improve the *performance* and *sample-complexity* of *GAlfO with self-observation*.
- Hypothesis:

Leveraging *proprioceptive information* will help with both issues

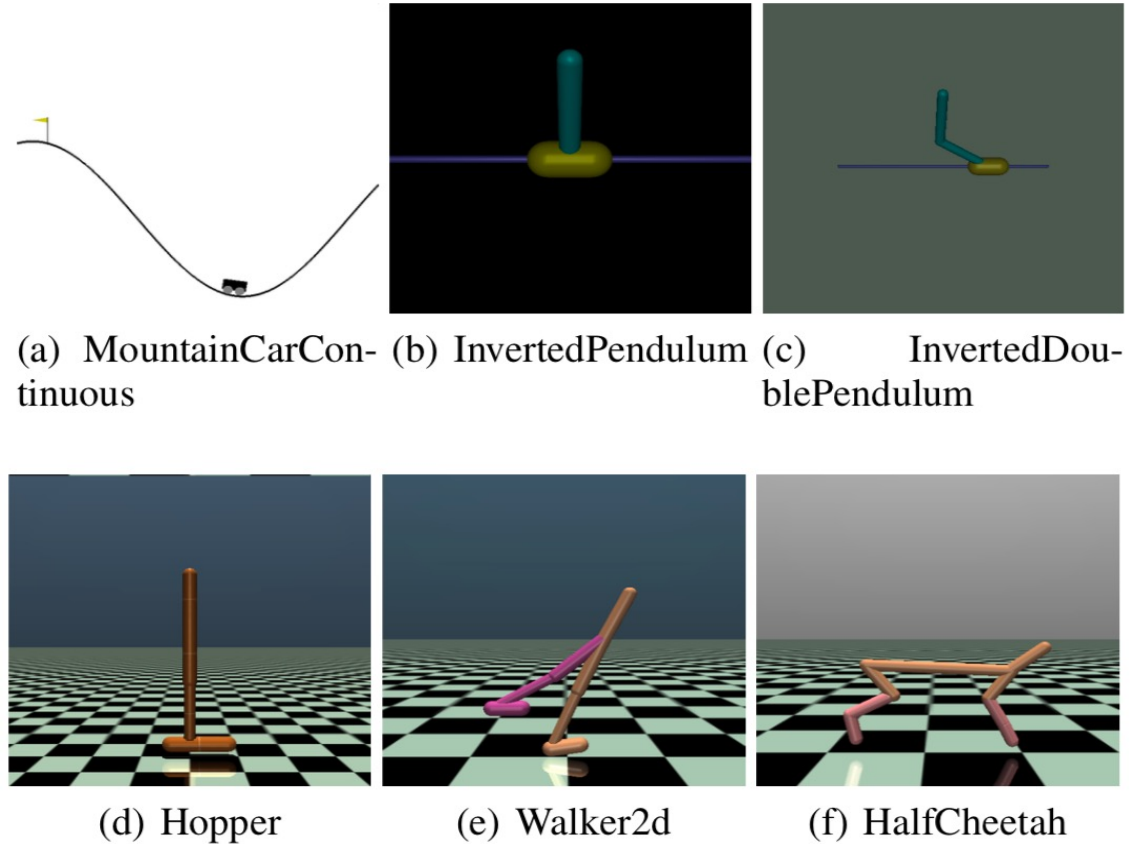
GAfO with Proprioceptive Information

- Propose an algorithm that uses both proprioceptive and visual information in order to:
 - Improve performance
 - Improve sample complexity



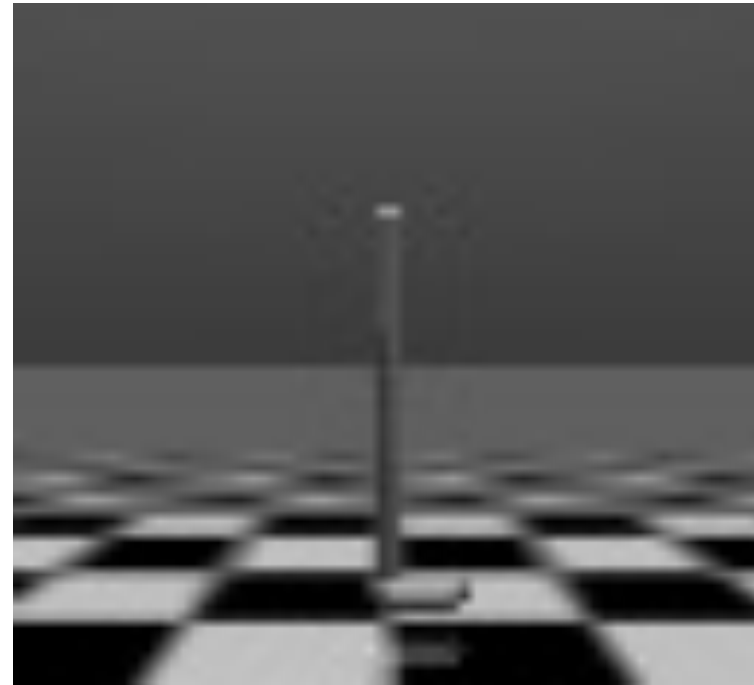
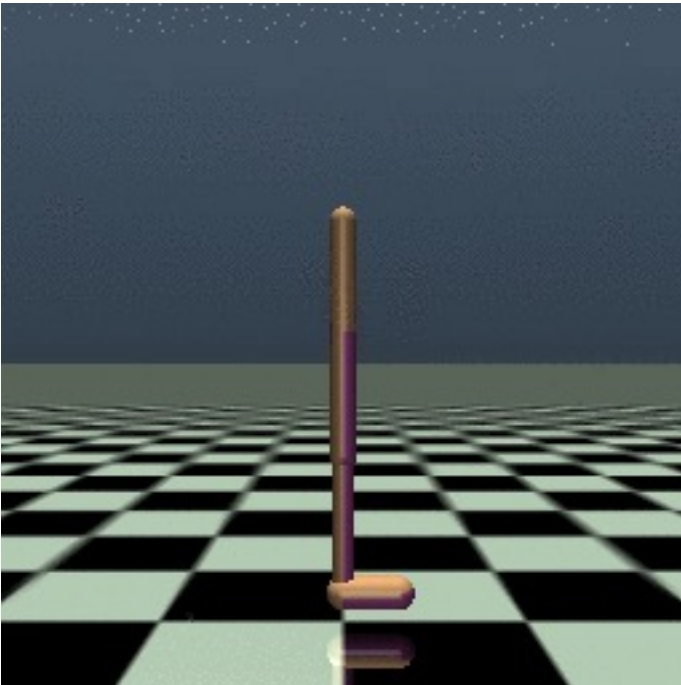
Experiments

- Tasks:
 - OpenAI Gym Environments
- Visual Demonstrations:
 - 64*64 grayscale frames



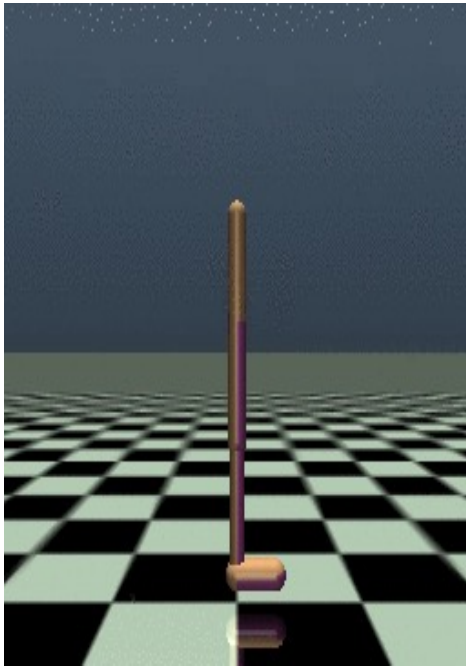
Experiments

- Walker2D Expert:

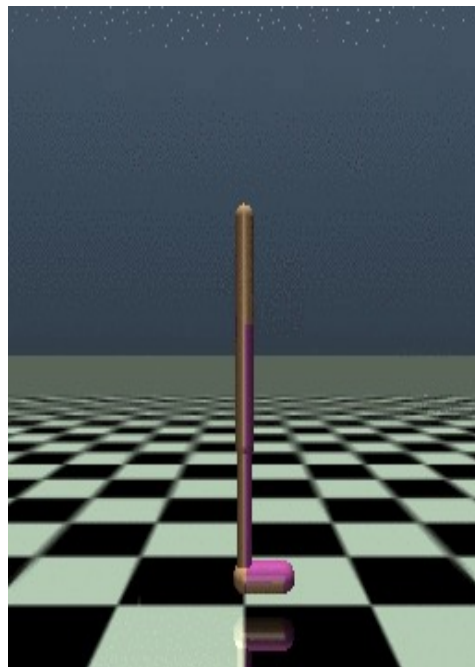


Experiments

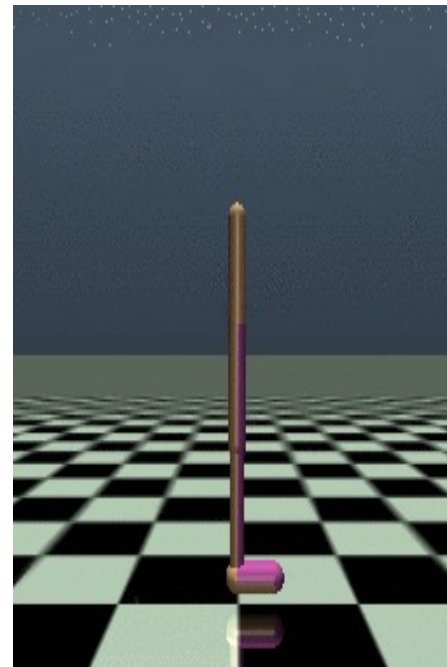
- Walker2D



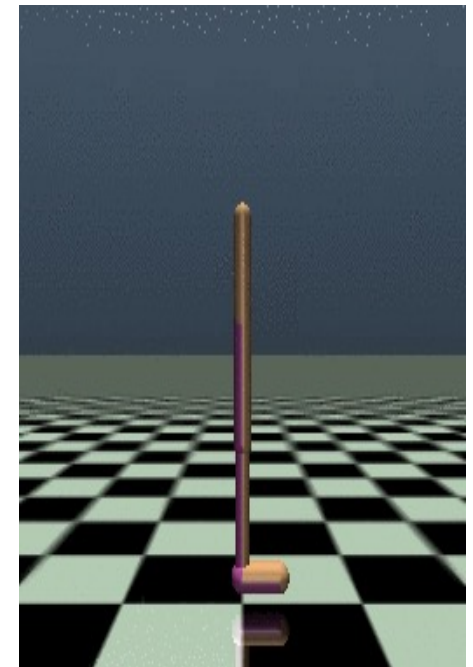
Expert



Iteration 0



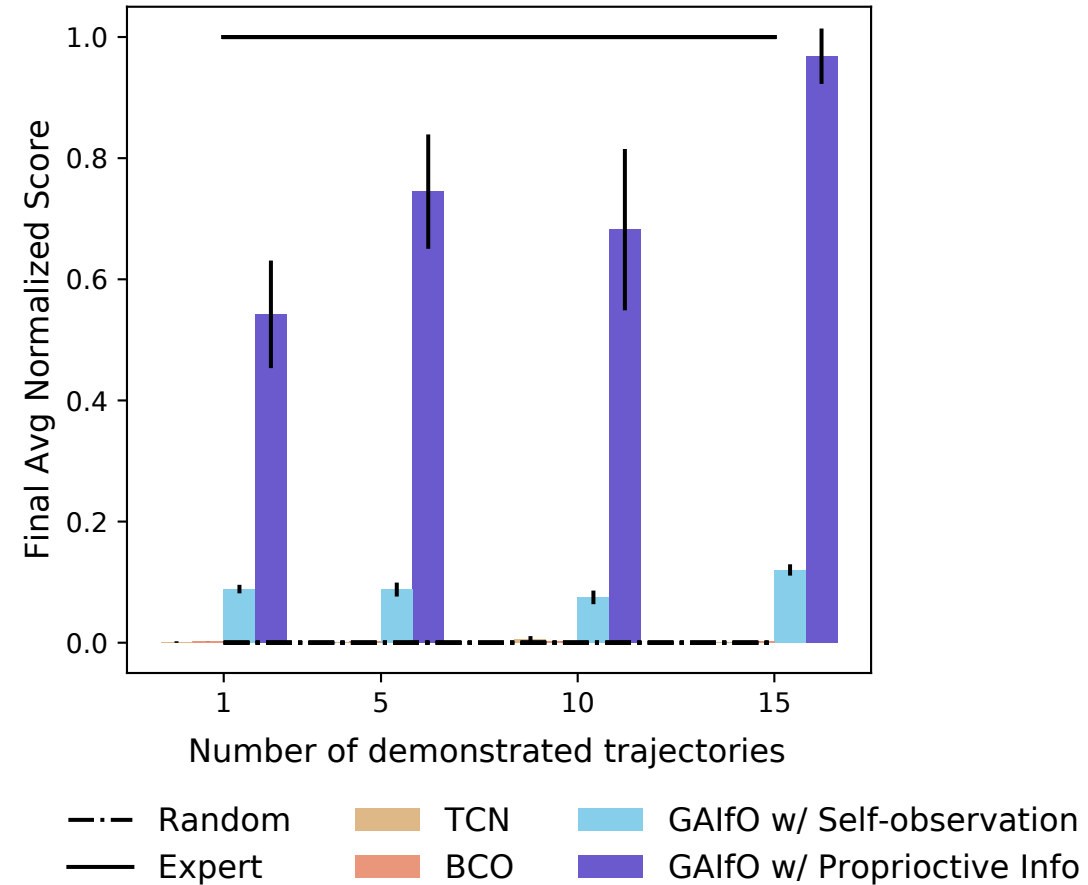
Iteration 100



Iteration 1942

Experiments

- Walker2D



In what ways can autonomous agents learn to imitate experts using
state-only observations?

Behavioral Cloning from
Observation (BCO)
[IJCAI 2018]

Generative Adversarial
Imitation from Observation
(GAIfO)
[AAMAS 2019, ICML
Workshop]

Data-Efficient Adversarial
Learning for Imitation from
Observation (DEALIO)
[Under Review]

Reinforced Inverse Dynamics
Modeling (RIDM)
[RAL, IROS 2020]

Visual Extension of GAIfO
with Self-observation
[ICML Workshop]

Visual Extension of GAIfO
with Proprioceptive Information
[IJCAI 2020]

Motivation

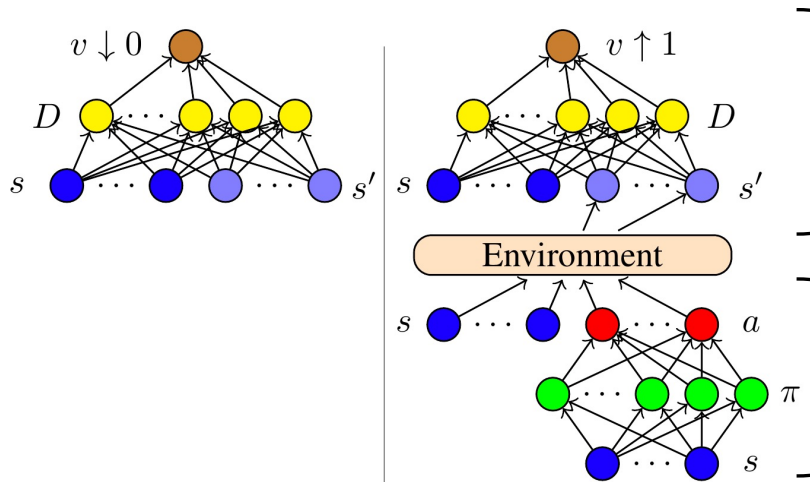
- Goal:
 - Improve *sample complexity* of *GAIfo* to enable *application on physical robots*.
- Integrating
 - Sample-efficient RL updates from PILQR [1] with
 - High-performing GAIfo algorithm for IfO.

[1] Chebotar, Yevgen, et al. "Combining model-based and model-free updates for trajectory-centric reinforcement learning." *International conference on machine learning*. PMLR, 2017.

PILQR

- Combines:
 - iterative Linear Quadratic Regulator (iLQR)
 - And, Path Integral Policy Improvement (PI²)
- iLQR constraints:
 - Linear dynamics $s_{t+1} = F_t \begin{bmatrix} s_t \\ a_t \end{bmatrix} + f_t$
 - Quadratic cost function $c(s_t, a_t) = \frac{1}{2} \begin{bmatrix} s_t \\ a_t \end{bmatrix}^T C_t \begin{bmatrix} s_t \\ a_t \end{bmatrix} + \begin{bmatrix} s_t \\ a_t \end{bmatrix}^T c_t$
- PILQR:
 - Twice differentiable-cost function
 - iLQR on quadratic approximation of the cost
 - PI² policy update on the residual cost
 - Returns a Gaussian controller $p(a|s)$

~~GAIfO~~ DEALIO



Demonstrator

Imitator

Some constraints on the cost

~~Used as cost function~~

Twice-differentiable

A function of both states and actions

PILQR

~~PPO~~ used to update

• We consider:

$$c(s_t, a_t) = \frac{1}{2} \begin{bmatrix} s_t \\ a_t \end{bmatrix}^T C_t \begin{bmatrix} s_t \\ a_t \end{bmatrix} + \begin{bmatrix} s_t \\ a_t \end{bmatrix}^T c_t + c c_t$$

Quadratic approximation

DEALIO

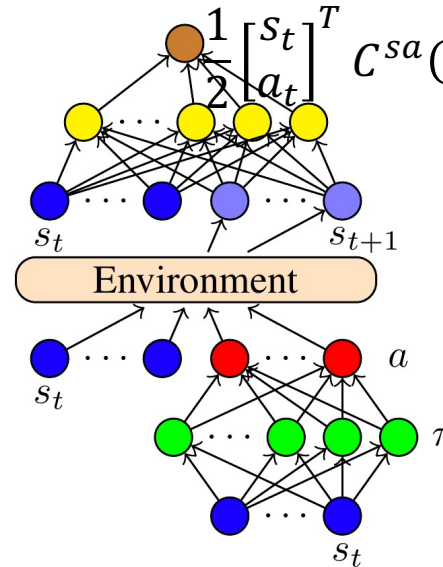
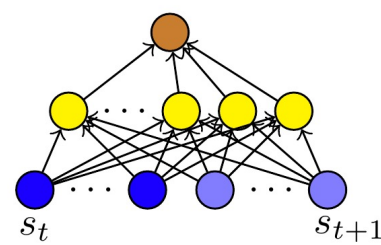
$$c(s_t, a_t) = \frac{1}{2} \begin{bmatrix} s_t \\ a_t \end{bmatrix}^T C_t \begin{bmatrix} s_t \\ a_t \end{bmatrix} + \begin{bmatrix} s_t \\ a_t \end{bmatrix}^T c_t + cc_t$$

$$D_\theta(s_t, s_{t+1}) = \frac{1}{2} \begin{bmatrix} s_t \\ s_{t+1} \end{bmatrix}^T C^{ss}(s_t, s_{t+1}) \begin{bmatrix} s_t \\ s_{t+1} \end{bmatrix} + \begin{bmatrix} s_t \\ s_{t+1} \end{bmatrix}^T c^{ss}(s_t, s_{t+1}) + cc^{ss}(s_t, s_{t+1}) = F_t \begin{bmatrix} s_t \\ a_t \end{bmatrix} + f_t$$

\uparrow $v \downarrow 0$ \uparrow $v \uparrow 1$

$C^{ss}(s_t, s_{t+1}), c^{ss}(s_t, s_{t+1}), cc^{ss}(s_t, s_{t+1})$

• Substituting s_{t+1} in $D_\theta(s_t, s_{t+1})$ to find



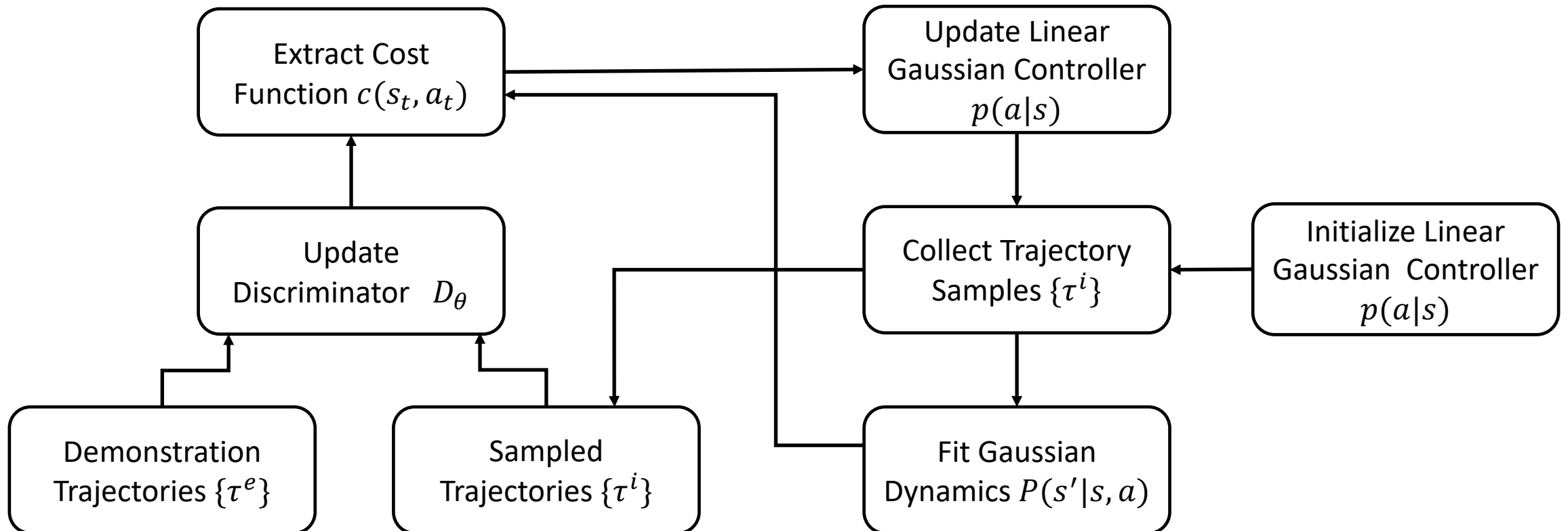
$$C_t = C^{sa}(\bar{s}_t, \bar{s}_{t+1})$$

$$c_t = c^{sa}(\bar{s}_t, \bar{s}_{t+1})$$

$$cc_t = cc^{ss}(\bar{s}_t, \bar{s}_{t+1})$$

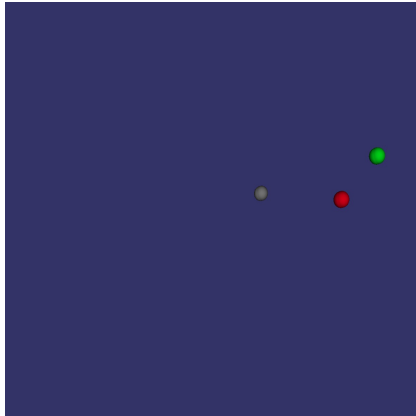
DEALIO

$$D_{\theta}(s_t, s_{t+1}) = \frac{1}{2} \begin{bmatrix} s_t \\ s_{t+1} \end{bmatrix}^T \begin{bmatrix} C_t & c_t \\ c_t & cc_t \end{bmatrix} \begin{bmatrix} s_t \\ s_{t+1} \end{bmatrix} + \frac{1}{2} \begin{bmatrix} s_t \\ s_{t+1} \end{bmatrix}^T \begin{bmatrix} C_t & c_t \\ c_t & cc_t \end{bmatrix} \begin{bmatrix} s_t \\ s_{t+1} \end{bmatrix} + cc_t$$

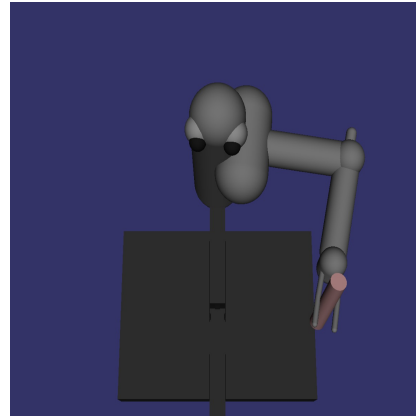


Experiments

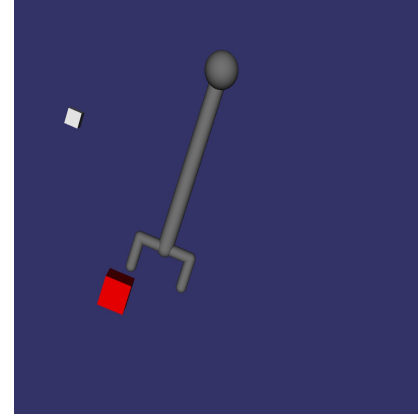
- MuJoCo Simulation Domains:



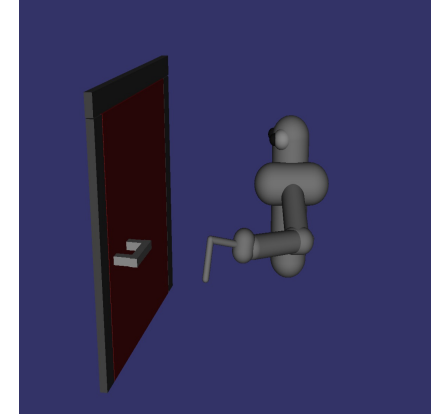
Disc



PegInsertion



GripperPusher



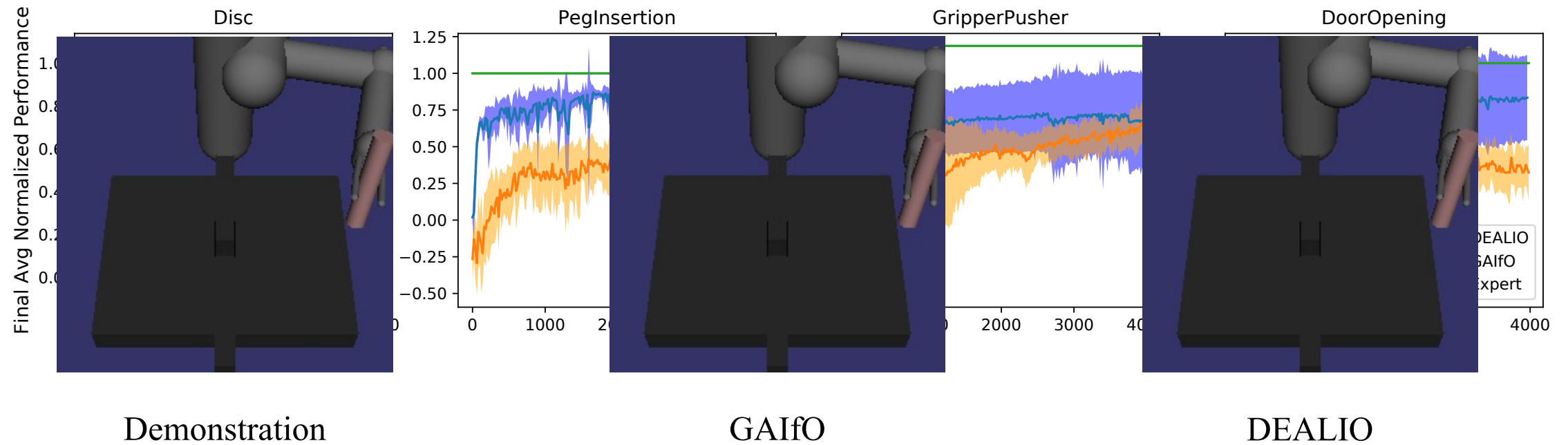
DoorOpening

- Hypothesis:

- DEALIO is able to learn tasks efficiently compared to GAIfo
- DEALIO is able to perform better compared to GAIfo

Experiments

- Our experiments on MuJoCo show DEALIO is faster in learning and has higher performance compared to GAIfo.
- PegInsertion:



In what ways can autonomous agents learn to imitate experts using
state-only observations?

Behavioral Cloning from
Observation (BCO)
[IJCAI 2018]

Generative Adversarial
Imitation from Observation
(GAIfO)
[AAMAS 2019, ICML
Workshop]

Data-Efficient Adversarial
Learning for Imitation from
Observation (DEALIO)
[Under Review]

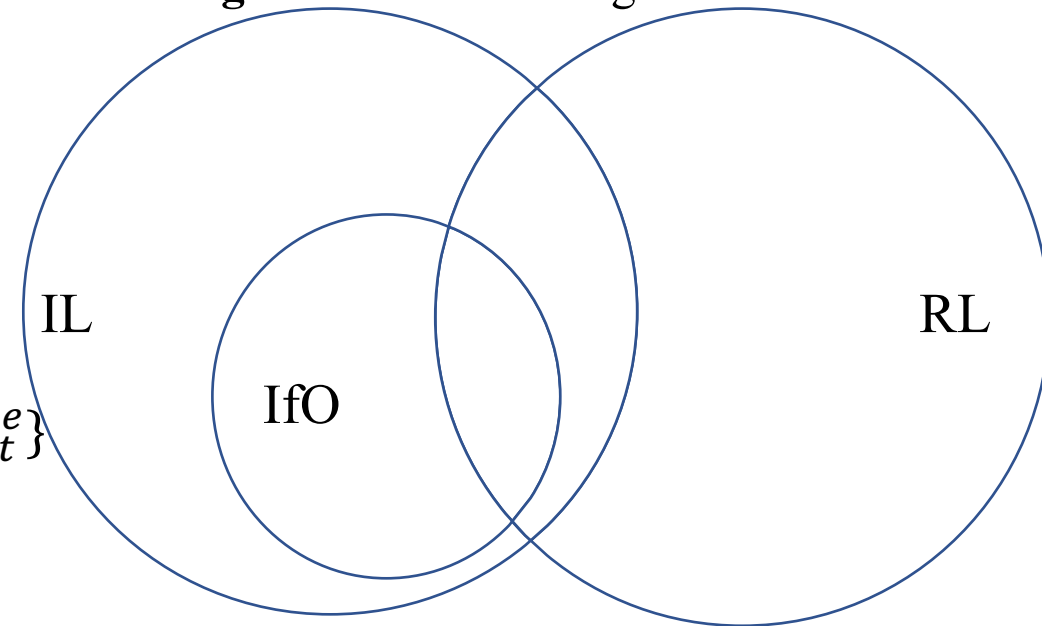
Reinforced Inverse Dynamics
Modeling (RIDM)
[RAL, IROS 2020]

Visual Extension of GAIfO
with Self-observation
[ICML Workshop]

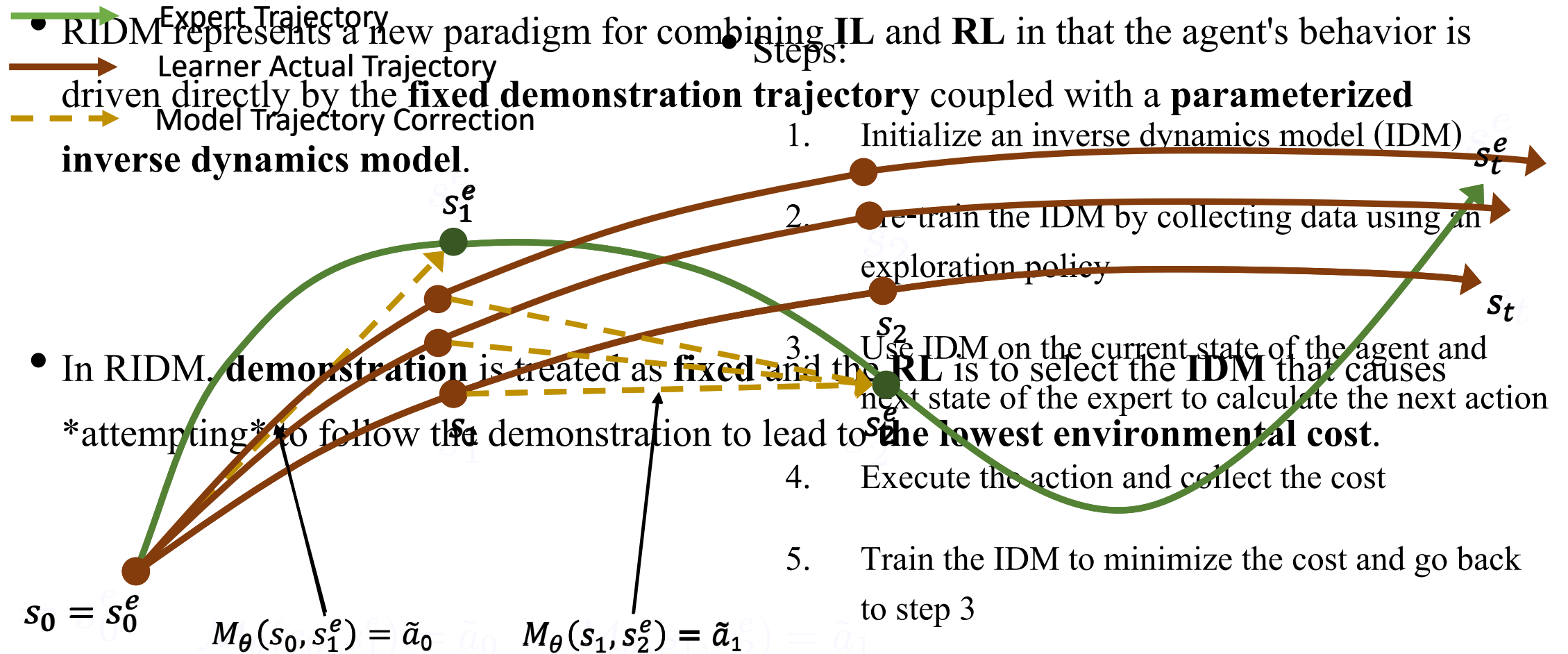
Visual Extension of GAIfO
with Proprioceptive Information
[IJCAI 2020]

Motivation

- Goal:
 - Combine “**imitation from observation**” and “**reinforcement learning**” to enable learning when:
 - The demonstrator is sub-optimal
 - Not many demonstration trajectories are available
- **Given:**
 - A single state-only (sub-optimal) demonstration: $D^e = \{s_t^e\}$
 - A cost function: c_{env}
- **Learn:**
 - A policy to perform the task



RIDM: Reinforced Inverse Dynamics Modeling



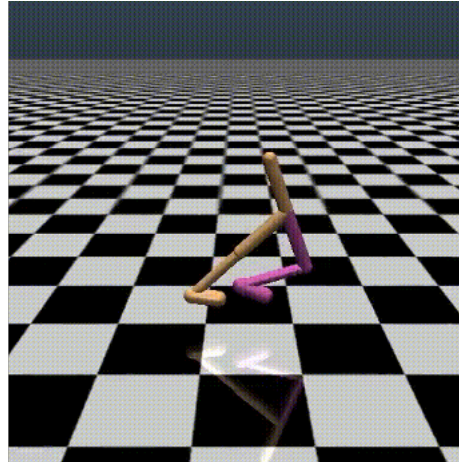
Experiments

- Robot control domains:

- MuJoCo Simulator
- SimSpark Simulator
- UR5 Arm Robot

- Hypothesis:

- RIDM is able to learn tasks efficiently with comparable performance compared to the demonstrator
- If the demonstrator is sub-optimal, RIDM is potentially able to outperform the demonstrator



Experiments

- Computationally challenging:
 - Hopper: 4.5 days
 - Nao's Fast Walk: 2.5 days

Parallelized Over 50 Machines



2-3 hours each experiment



Experiments: SimSpark Robot Soccer

- Used in 3D Simulation RoboCup
- Developing skills such as walk and kick is challenging
- Tasks:
 - Fast Walk
 - Long Kick
- Demonstrators:
 - FUT-K
 - FC Portugal
- Demonstrators are sub-optimal with respect to the designed cost.



Experiments: SimSpark Robot Soccer

- Our experiments on SimSpark 3D simulator show learned behavior outperforms the suboptimal experts.
- ~~EastgWalk~~ ((FUT-K)):

Learned
Behavior:



Demo:

- ~~EastgWalk~~ ((FC Portugal)):

Learned
Behavior:



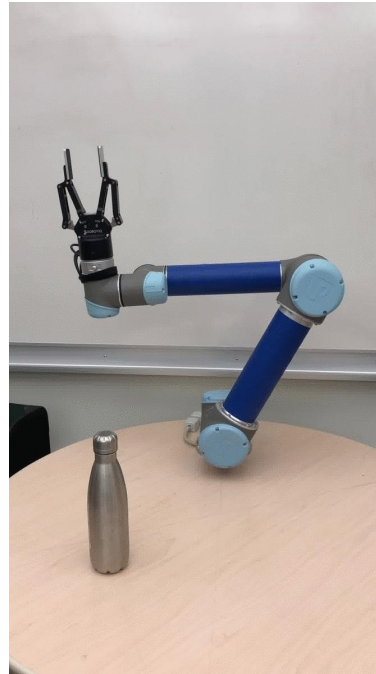
Demo:

Experiments: UR5 Arm Robot

- Our experiments on UR5 robots show learned behavior outperforms robot's default PID performance.
- Pushing Task (10x):



Demonstration



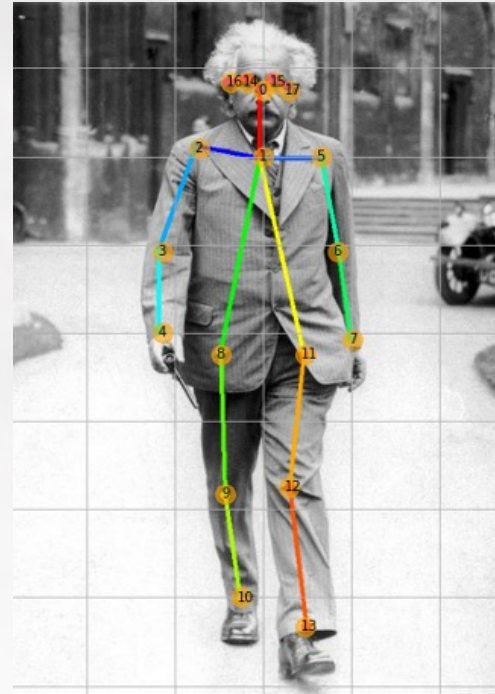
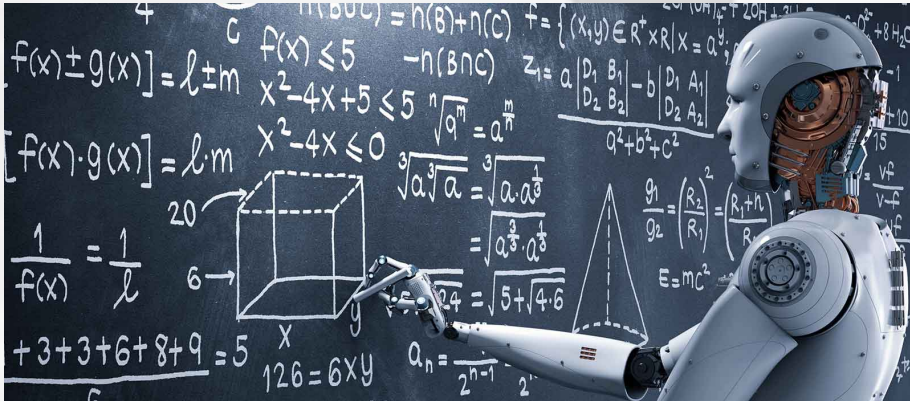
Default UR5 PID
Imitation Learning from Observation



Learned Behavior

Future Work

- Perception
- Application to Physical Robots
- Fully-intelligent Agents



Future Work

- Perception Challenges

Integration of Perception and Control

Embodiment Mismatch

CycleGAN [Zhu et al. 2017]

Pix2pix [Isola et al. 2017]

Dual GAN [Yi et al. 2017]

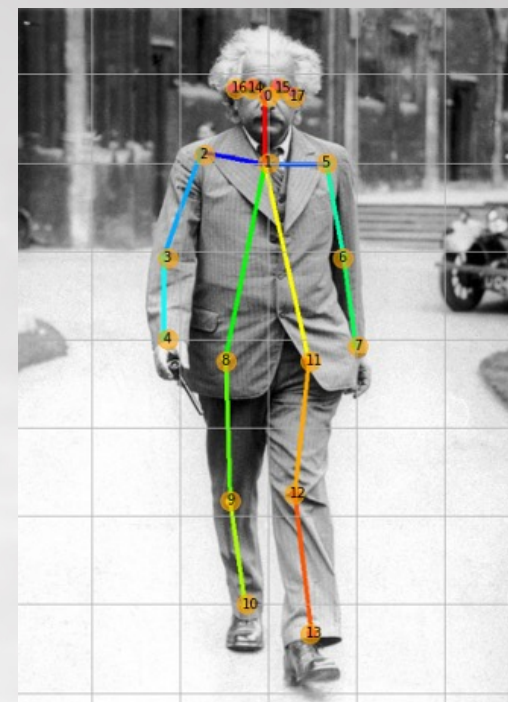
Disco GAN [Kim et al. 2017]

Viewpoint Mismatch

Pose-estimation [Cao et al. 2017,

Wang et al. 2019]

Keypoint Detection



Future Work

- Application to Physical Robots



Sample-efficiency

Safe

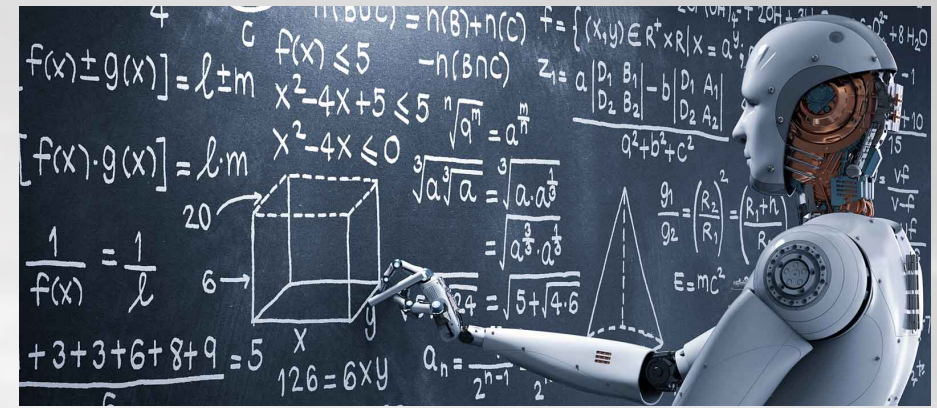
Future Work

- Fully-intelligent Agents

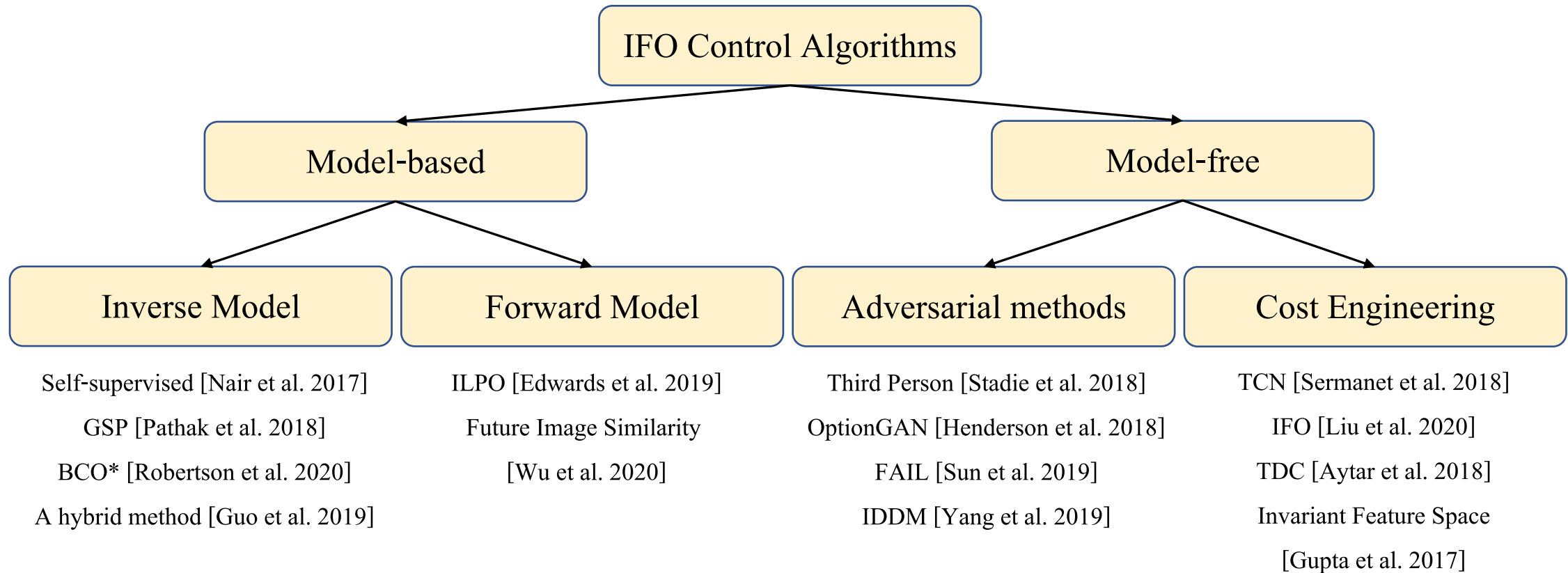
Reinforcement learning

Imitation Learning

Imitation from Observation



Related Work



Summary

- Area A:
 - Equivalency of solving the model-free IfO problem and solving the GANs like optimization problem.
- Area B:
 - Implementation of the introduced algorithms.
 - Training the models on hundreds of machines.
 - Extensive hyperparameter search for each algorithm.
- Area C:
 - Modeling the human ability of imitation from observation.
 - Application of the developed algorithms to simulated and physical robots.

Acknowledgements



Peter Stone



Garrett Warnell



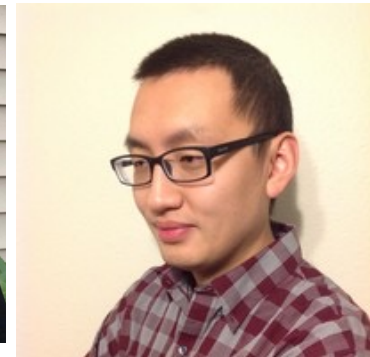
Patrick MacAlpine



Josiah Hanna



Brahma Pavse

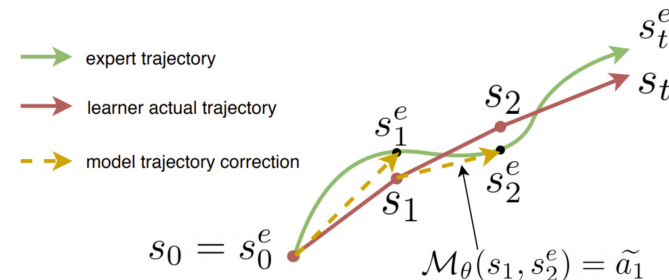
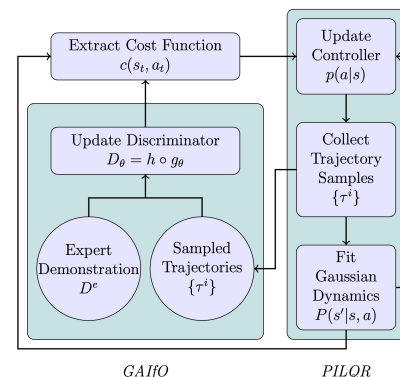
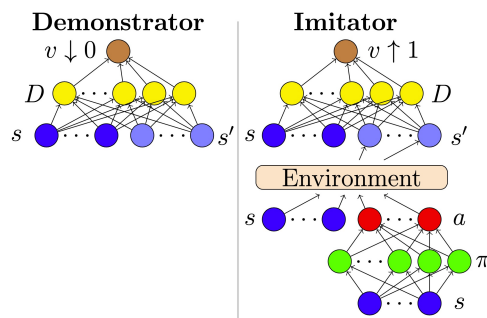
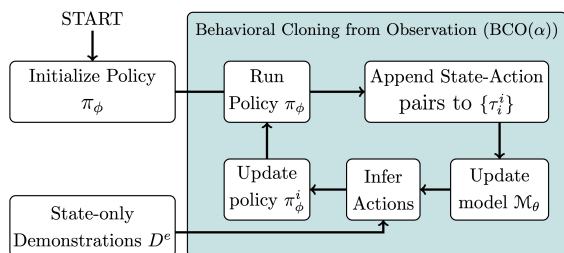


Ruohan Zhang



Sean Geiger

In what ways can autonomous agents learn to imitate experts using state-only observations?

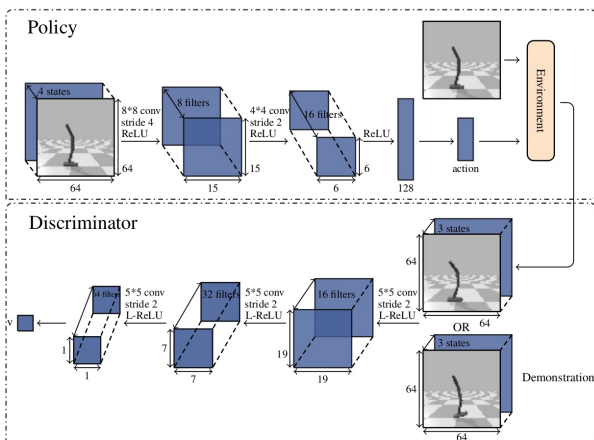


Behavioral Cloning from Observation (BCO)
[IJCAI 2018]

Generative Adversarial Imitation from Observation (GAIfO)
[AAMAS 2019, ICML Workshop]

Data-Efficient Adversarial Learning for Imitation from Observation (DEALIO)
[Under Review]

Reinforced Inverse Dynamics Modeling (RIDM)
[RAL, IROS 2020]



Visual Extension of GAIfO with Self-observation
[ICML Workshop]

Visual Extension of GAIfO with Proprioceptive Information
[IJCAI 2020]

