

Imitation Learning

Goal:

- Learn how to make decisions by trying to imitate another agent.

Imitation Learning

Goal:

- Learn how to make decisions by trying to imitate another agent.

Conventional Imitation Learning:

- Observations of other agent (demonstrations) consist of state-action pairs.¹

¹Niekum et al., "Learning and generalization of complex tasks from unstructured demonstrations".

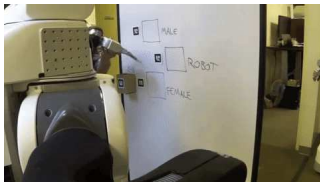
Imitation Learning

Goal:

- Learn how to make decisions by trying to imitate another agent.

Conventional Imitation Learning:

- Observations of other agent (demonstrations) consist of state-action pairs.¹



Imitation Learning

Goal:

- Learn how to make decisions by trying to imitate another agent.

Conventional Imitation Learning:

- Observations of other agent (demonstrations) consist of state-action pairs.¹

Challenge:

- Precludes using a large amount of demonstration data where action sequences are not given (e.g. YouTube videos).

¹Niekum et al., "Learning and generalization of complex tasks from unstructured demonstrations".

Imitation Learning

Algorithms:

Imitation Learning

Algorithms:

- Behavioral Cloning:

Imitation Learning

Algorithms:

- Behavioral Cloning:
 - ▶ End to End Learning for Self-Driving Cars.²

²Zhang and Cho, "Query-Efficient Imitation Learning for End-to-End Simulated Driving."

Imitation Learning

Algorithms:

- Behavioral Cloning:
 - ▶ End to End Learning for Self-Driving Cars.²
- Inverse Reinforcement Learning:

²Zhang and Cho, "Query-Efficient Imitation Learning for End-to-End Simulated Driving."

Imitation Learning

Algorithms:

- Behavioral Cloning:
 - ▶ End to End Learning for Self-Driving Cars.²
- Inverse Reinforcement Learning:
 - ▶ Generative Adversarial Imitation Learning.³
 - ▶ Guided Cost Learning.⁴

²Zhang and Cho, "Query-Efficient Imitation Learning for End-to-End Simulated Driving."

³Ho and Ermon, "Generative adversarial imitation learning".

⁴Finn, Levine, and Abbeel, "Guided cost learning: Deep inverse optimal control via policy optimization".

Imitation from Observation

Goal:

- Learn how to perform a task given state-only demonstrations.



Imitation from Observation

Goal:

- Learn how to perform a task given state-only demonstrations.

Imitation from Observation

Goal:

- Learn how to perform a task given state-only demonstrations.

Formulation:

- Given:
 - ▶ $D_{demo} = (s_0, s_1, \dots)$
- Learn:
 - ▶ $\pi : \mathcal{S} \rightarrow \mathcal{A}$

Imitation from Observation

Previous work:

Imitation from Observation

Previous work:

- Time Contrastive Networks (TCN).⁵
- Imitation from observation: Learning to imitate behaviors from raw video via context translation.⁶
- Learning invariant feature spaces to transfer skills with reinforcement learning.⁷

⁵Sermanet et al., "Time-contrastive networks: Self-supervised learning from multi-view observation".

⁶Liu et al., "Imitation from observation: Learning to imitate behaviors from raw video via context translation".

⁷Gupta et al., "Learning invariant feature spaces to transfer skills with reinforcement learning".

Imitation from Observation

Previous work:

- Time Contrastive Networks (TCN).⁵
- Imitation from observation: Learning to imitate behaviors from raw video via context translation.⁶
- Learning invariant feature spaces to transfer skills with reinforcement learning.⁷

Concentrate on perception; require time-aligned demonstrations.

⁵Sermanet et al., "Time-contrastive networks: Self-supervised learning from multi-view observation".

⁶Liu et al., "Imitation from observation: Learning to imitate behaviors from raw video via context translation".

⁷Gupta et al., "Learning invariant feature spaces to transfer skills with reinforcement learning".

Efficient Robot Skill Learning

- Motivation: RoboCup
- Sim2Real: Grounded Simulation Learning
- Imitation Learning from Observation:
 - ▶ **Model-based approach:** BCO
 - ▶ Model-free approach: GAlfO

Model-based Approach

- Imitation Learning: $D_{demo} = \{(s_0, a_0), (s_1, a_1), \dots\}$

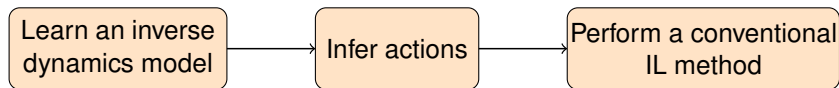
Model-based Approach

- Imitation Learning: $D_{demo} = \{(s_0, a_0), (s_1, a_1), \dots\}$
- Imitation from Observation: $D_{demo} = \{(s_0, ?), (s_1, ?), \dots\}$

Model-based Approach

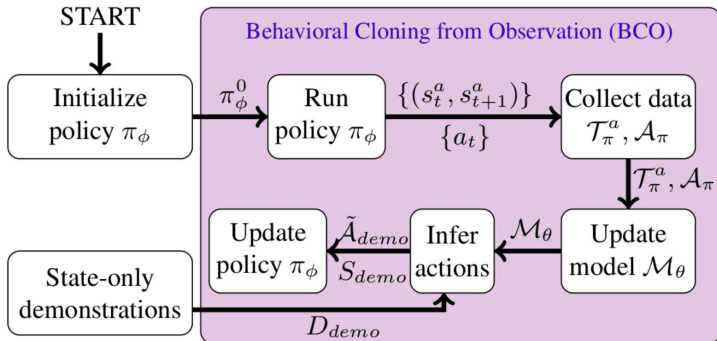
- Imitation Learning: $D_{demo} = \{(s_0, a_0), (s_1, a_1), \dots\}$
- Imitation from Observation: $D_{demo} = \{(s_0, ?), (s_1, ?), \dots\}$

Model-based Approach:



Behavioral Cloning from Observation (BCO)

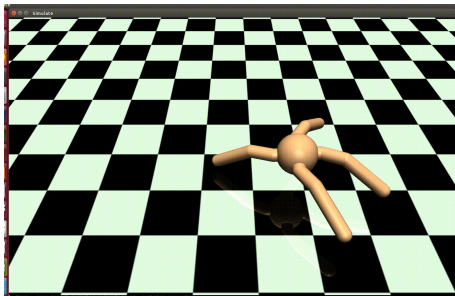
Algorithm:



Behavioral Cloning from Observation (BCO)

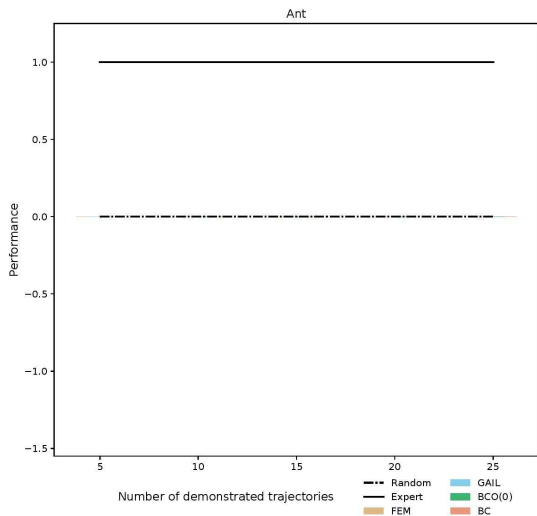
Experimental Results:

- Domain:
 - ▶ Mujoco domain "Ant" with 111 dimensional state space and 8 dimensional action space.



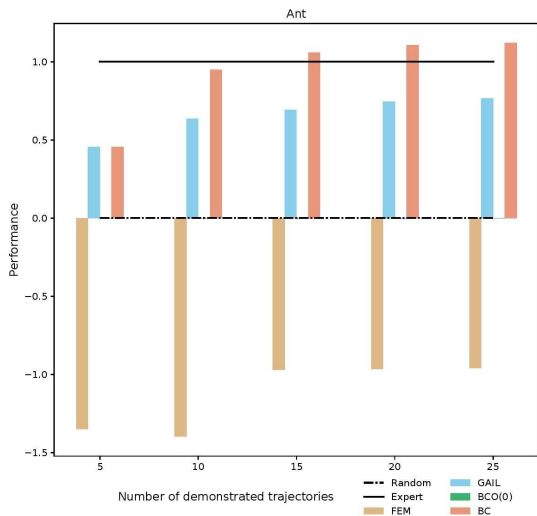
Behavioral Cloning from Observation (BCO)

Experimental Results:



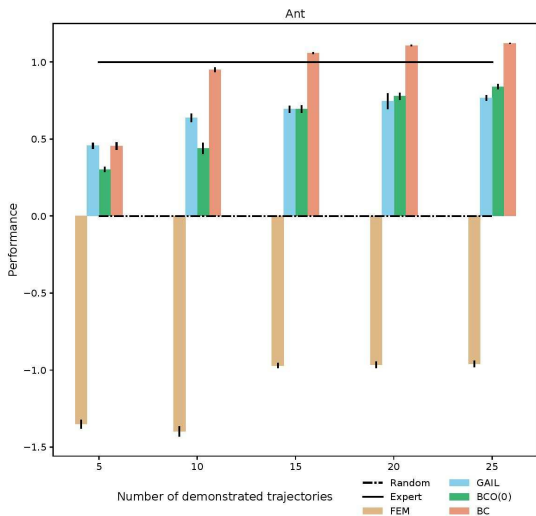
Behavioral Cloning from Observation (BCO)

Experimental Results:



Behavioral Cloning from Observation (BCO)

Experimental Results:



Behavioral Cloning from Observation (BCO(α))

Issue:

- Inverse dynamics model is learned using a random policy.

Behavioral Cloning from Observation (BCO(α))

Issue:

- Inverse dynamics model is learned using a random policy.

Solution: BCO(α)

Behavioral Cloning from Observation (BCO(α))

Issue:

- Inverse dynamics model is learned using a random policy.

Solution: BCO(α)

- Update the model with the learned policy.

Behavioral Cloning from Observation (BCO(α))

Issue:

- Inverse dynamics model is learned using a random policy.

Solution: BCO(α)

- Update the model with the learned policy.
- Parameter α controls tradeoff between performance and environment interactions

Behavioral Cloning from Observation (BCO(α))

Issue:

- Inverse dynamics model is learned using a random policy.

Solution: BCO(α)

- Update the model with the learned policy.
- Parameter α controls tradeoff between performance and environment interactions
 - ▶ $\alpha = 0$: no post-demonstration interaction.

Behavioral Cloning from Observation (BCO(α))

Issue:

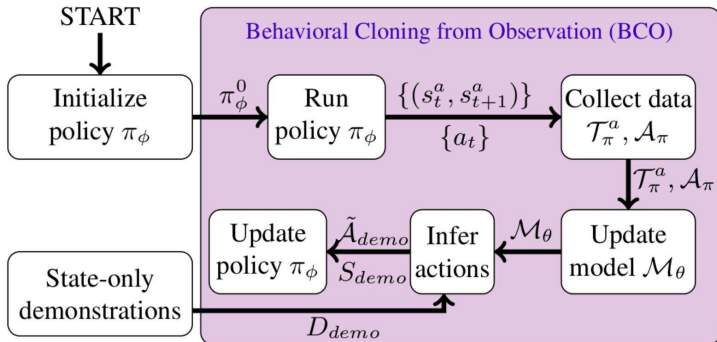
- Inverse dynamics model is learned using a random policy.

Solution: BCO(α)

- Update the model with the learned policy.
- Parameter α controls tradeoff between performance and environment interactions
 - ▶ $\alpha = 0$: no post-demonstration interaction.
 - ▶ Increasing α : increasing the number of interactions allowed at each iteration.

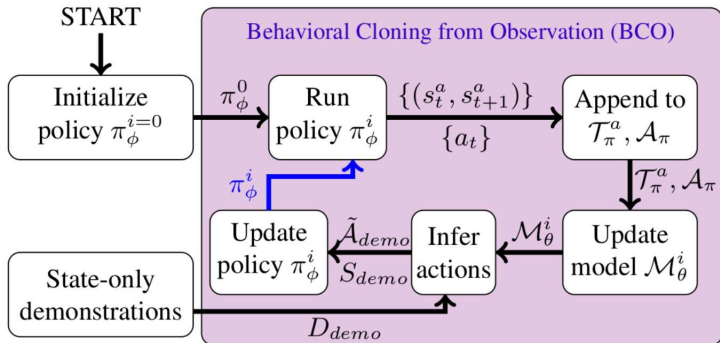
Behavioral Cloning from Observation (BCO(α))

Algorithm:



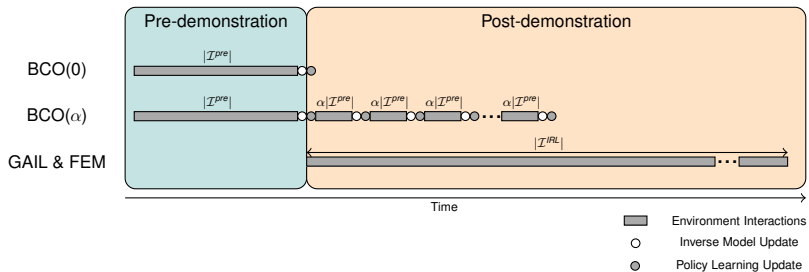
Behavioral Cloning from Observation (BCO(α))

Algorithm:



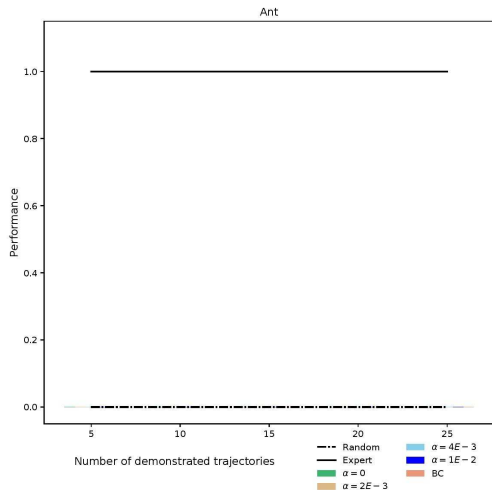
Behavioral Cloning from Observation (BCO(α))

Interaction time:



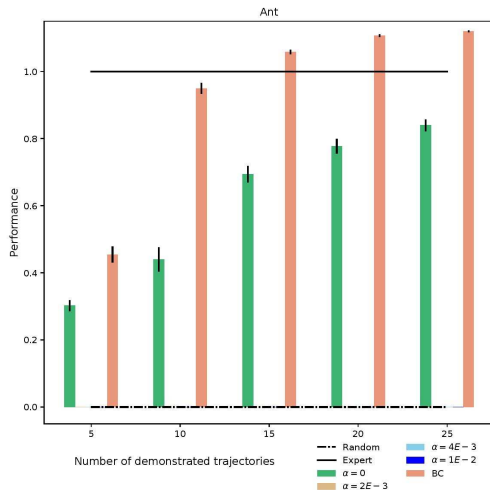
Behavioral Cloning from Observation (BCO(α))

Effect of varying α on BCO(α):



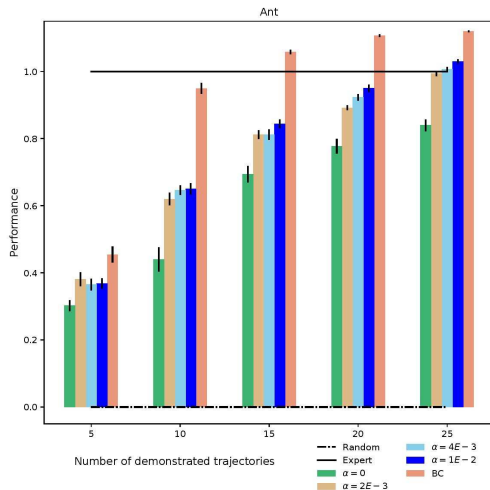
Behavioral Cloning from Observation (BCO(α))

Effect of varying α on BCO(α):



Behavioral Cloning from Observation (BCO(α))

Effect of varying α on BCO(α):

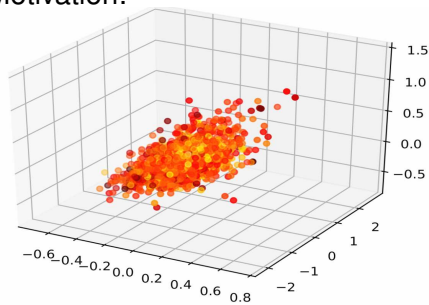


Efficient Robot Skill Learning

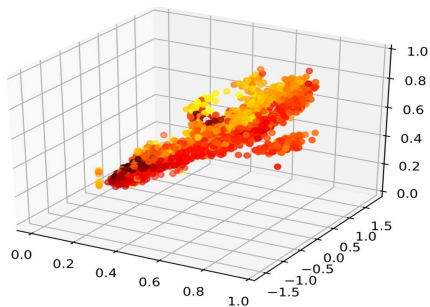
- Motivation: RoboCup
- Sim2Real: Grounded Simulation Learning
- Imitation Learning from Observation:
 - ▶ Model-based approach: BCO
 - ▶ **Model-free approach: GAlfO**

Gen. Adversarial Imitation from Observation (GAIfo)

Motivation:



(a) Random Policy



(b) Demonstration

Figure: State transition distribution in Hopper domain.

Gen. Adversarial Imitation from Observation (GAIfo)

Algorithm:

