

Understanding Adversarial Imitation Learning in Small Sample Regime: A Stage-coupled Analysis

Tian Xu

Nanjing University, School of Artificial Intelligence

Based on the paper: Tian Xu*, Ziniu Li*, Yang Yu, Zhi-Quan Luo. Understanding AIL in Small Sample Regime: A Stage-coupled Analysis. TPAMI.



Tian Xu
NJU



Ziniu Li
CUHKSZ



Yang Yu
NJU



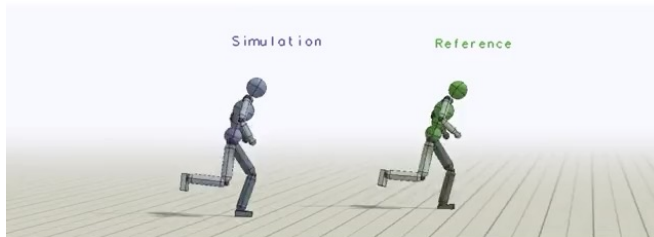
Zhi-Quan Luo
CUHKSZ

What is Imitation Learning

Imitation Learning (IL) aims to learn effective policies by **mimicking expert demonstrations**

Imitate to run

Humanoid: Run



Policy trained to imitate a running clip.

[<https://gfycat.com/bowedinexperiencedhogget>]

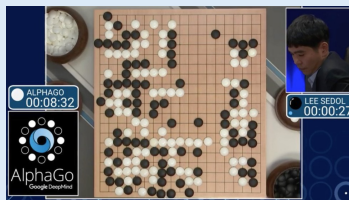
Imitate to take the banana



[<https://www.comp.nus.edu.sg/~rishav1/blog/2018/Understanding-GAIL/>]

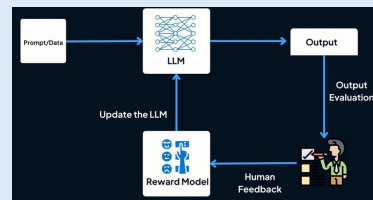
Applications of IL

DeepMind AlphaGo



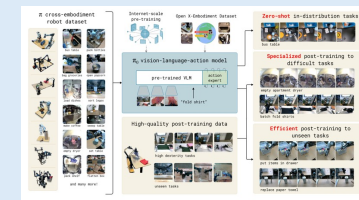
Imitate professional player for initial learning

OpenAI ChatGPT



Imitate internet corpora and high-quality human responses

Physical Intelligence π^0



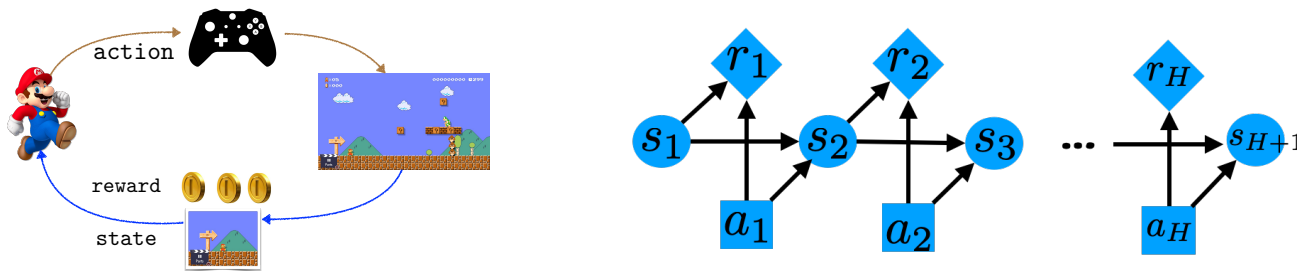
Imitate human demonstrations

IL is a general paradigm underpins multiple breakthrough AI systems

IL Set-up: Markov Decision Process

Definition (Finite-horizon Markov Decision Process). A finite-horizon Markov Decision Process is defined by a tuple $\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, r, H, s_1)$, where:

- \mathcal{S} is the state space;
- \mathcal{A} is the action space;
- $P = \{P_h\}_{h=1}^H$, where $P_h : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ is the transition probability at step h ;
- $r = \{r_h\}_{h=1}^H$, where $r_h : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ is the reward function at step h ;
- H is the decision horizon;
- s_1 is the initial state.



IL Set-up: Markov Decision Process

For policy $\pi = \{\pi_h\}_{h=1}^H$, $\pi_h : \mathcal{S} \rightarrow \Delta(\mathcal{A})$

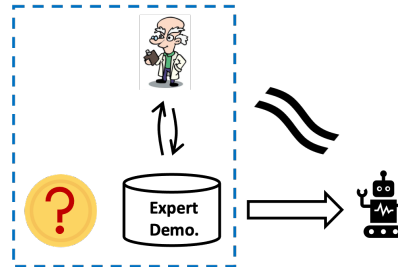
Policy Value V^π measures the quality of a policy:

$$V^\pi = \mathbb{E} \left[\sum_{h=1}^H r_h(s_h, a_h) \mid a_h \sim \pi_h(\cdot \mid s_h), s_{h+1} \sim P_h(\cdot \mid s_h, a_h), \forall h \in [H] \right].$$

State-action distribution $\{d_h^\pi(\cdot, \cdot)\}_{h=1}^H$:

$$d_h^\pi(s, a) = \mathbb{P}(s_h = s, a_h = a; \pi)$$

IL Set-up



IL Set-up

The learner has **no knowledge** of the reward function r , but has access to expert demonstrations \mathcal{D}^E . The expert demonstrations consist of N expert trajectories collected by the deterministic expert policy π^E .

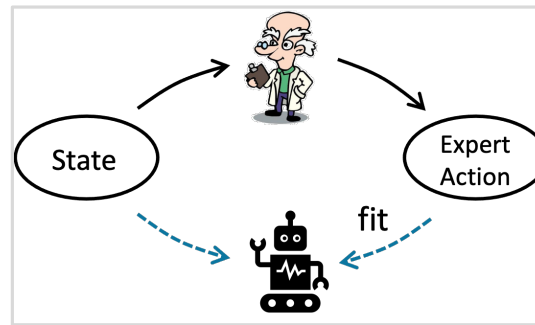
$$\mathcal{D}^E = \{\text{tr}^1, \dots, \text{tr}^N\}, \quad \text{tr}^i = \{(s_1, a_1), \dots, (s_H, a_H)\} \sim \pi^E$$

IL Target

The learner aims to output a policy $\hat{\pi}$ such that the policy value gap with π^E is small:

$$\text{Imitation Gap: } V^{\pi^E} - V^{\hat{\pi}}$$

Behavioral Cloning (BC)



Objective

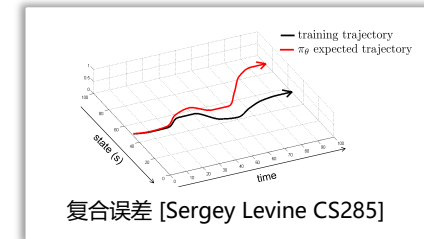
$$\max_{\pi} \mathbb{E}_{\tau \sim \pi^E} \left[\sum_{h=1}^H \log (\pi_h (a_h | s_h)) \right]$$

Feature

- Easy to implement
- **Offline learning:** purely trained on offline demonstrations

BC Suffers from Compounding Errors

- Due to its offline learning manner, BC cannot recover expert actions on states **out of demonstrations**
- Single-step error **accumulates over the decision horizon**, resulting in huge compounding errors



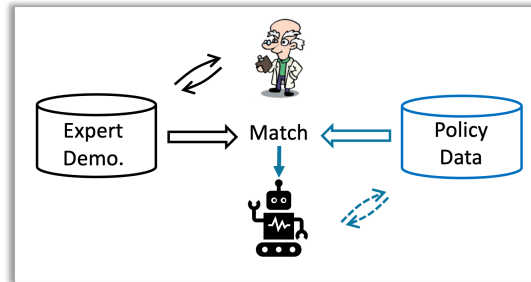
- Theory: [Ross and Bagnell, 2012] proved that BC has an imitation gap bound with **quadratic dependence on horizon**

$$V^{\pi^E} - V^{\pi^{BC}} \leq \mathcal{O}(H^2 \varepsilon) \quad \text{generalization error } \varepsilon = \frac{1}{H} \mathbb{E}_{\tau \sim \pi^E} \left[\sum_{h=1}^H \mathbb{E}_{a \sim \pi_h^{BC}(\cdot|s_h)} [\mathbb{I}\{a \neq \pi_h^E(s_h)\}] \right]$$

Stéphane Ross and Drew Bagnell. Efficient reductions for imitation learning. AISTATS 2010.

Key difference from supervised learning (**H=1**)

Adversarial Imitation Learning



“state-action distribution matching”

Objective

$$\min_{\pi} \sum_{h=1}^H D(d_h^{\pi}(\cdot, \cdot), d_h^{\pi^E}(\cdot, \cdot)) \quad , \quad D(\cdot, \cdot): \text{divergence measure}$$



$$\min_{\pi} \max_r \sum_{h=1}^H \underbrace{\mathbb{E}_{(s,a) \sim d_h^{\pi^E}} [r(s, a)]}_{V^{\pi^E, r}} - \sum_{h=1}^H \underbrace{\mathbb{E}_{(s,a) \sim d_h^{\pi}} [r(s, a)]}_{V^{\pi, r}}$$

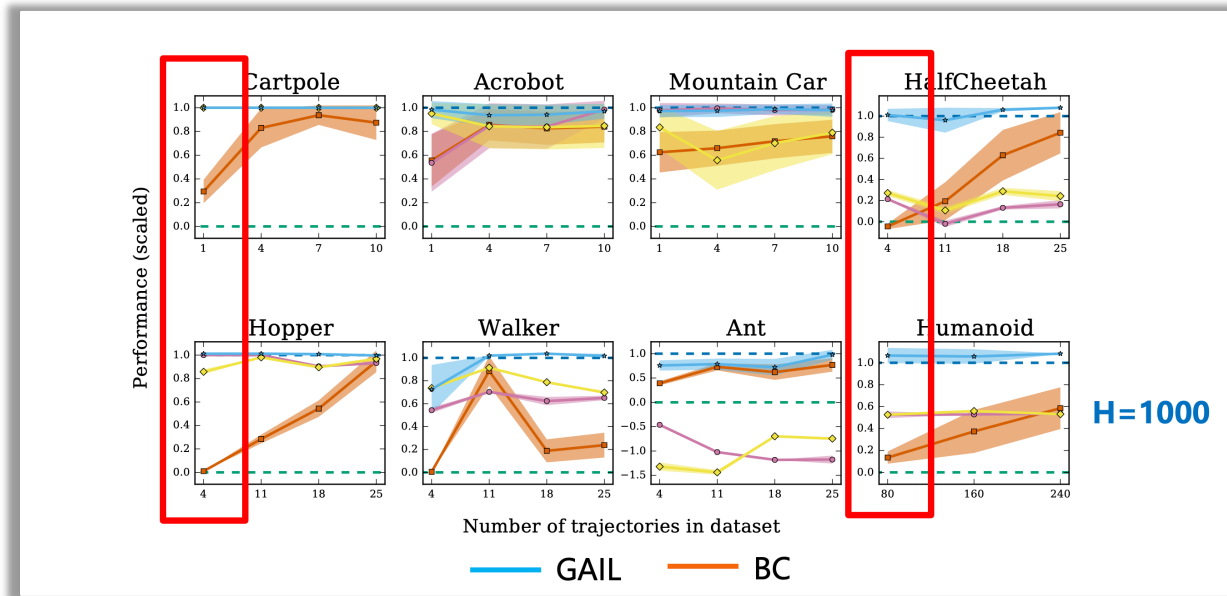
RL sub-problem

- inner: learn an **adversarial reward** to maximize the imitation gap
- outer: learn a **policy** to maximize the value w.r.t the adversarial reward

Feature

All involves solving RL problems, requiring **online interactions**

AIL Outperforms BC Significantly



Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning. NeurIPS 2016.

- With **limited demonstrations**, BC suffers from large performance gaps in **long-horizon tasks**
- AIL can achieve near-expert performance with **limited demonstrations** in **long-horizon tasks**

Understanding AIL in Small Sample Regime

**CoRL 2019
Best Paper**

**A Divergence Minimization Perspective
on Imitation Learning Methods**

Seyed Kamyar Seyed Ghasemipour
University of Toronto, Vector Institute
kamyar@cs.toronto.edu

Richard Zemel
University of Toronto, Vector Institute
zemel@cs.toronto.edu

Shixiang Gu
Google Brain
shanegu@google.com

gence, whereas AIRL and GAIL use divergences that exhibit more mode-seeking behaviour. These observations allow us to generate the following two hypotheses about why IRL methods outperform BC, particularly in the low-data regime,

Hypothesis 1 *In common MDPs of interest, the reward function depends more on the state than the action. Hence encouraging policies to explicitly match expert state marginals is an important learning criterion.*

Hypothesis 2 *It is known that optimization using the forward KL divergence results in distributions with a mode-covering behaviour, whereas using the reverse KL results in mode-seeking behaviour [39]. In RL we care about the “quality of trajectories”, as measured by the likelihood under the expert distribution. Therefore, being mode-seeking is more beneficial than mode-covering, particularly in the low-data regime.*

Open Question: “Why AIL methods outperform BC, particularly in the **low-data regime**?”

- Hypothesis 1: Matching state marginals is an important learning criterion (**high level**)
- Hypothesis 2: Divergence Choice: Forward KL V.S. Reverse KL (**not tailored to IL**)

Target: Develop **theory** to provide a **fundamental and concrete** understanding


Existing Theoretical Analysis for AIL

[Xu et al., 2020, Rajaraman et al., 2020] proved that AIL can achieve an imitation gap bound with **linear horizon dependence**:

$$V^{\pi^E} - \mathbb{E} \left[V^{\pi^{\text{AIL}}} \right] \leq \mathcal{O} \left(H \sqrt{\frac{|\mathcal{S}|}{N}} \right)$$

Tian Xu, et al. Error bounds of imitating policies and environments. NeurIPS 2020.

Nived Rajaraman, et al. Toward the fundamental limits of imitation learning. NeurIPS 2020.

- Compared with the $\mathcal{O}(H^2\epsilon)$ bound in BC, AIL achieves an improved $\mathcal{O}(H\epsilon)$ bound, explaining why AIL performs good in long-horizon tasks
- In the small sample regime ($N \leq |\mathcal{S}|$), this bound is meaningless because the maximum gap is H , contradicting with good practical performance of AIL 

Main Contributions

- The **first low data regime** imitation gap bound for AIL

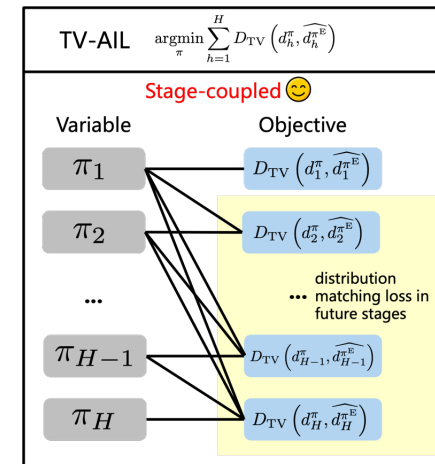
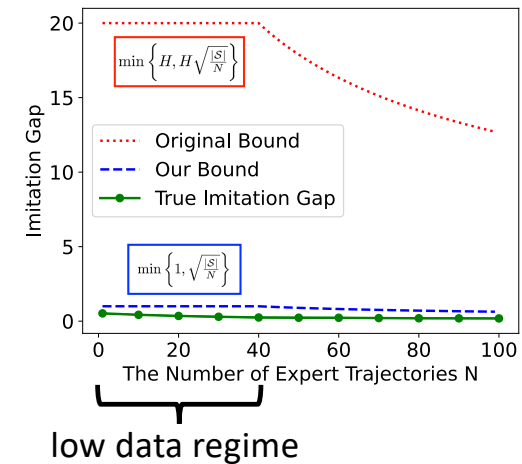
Theorem. Under mild assumptions, AIL can achieve an imitation gap bound:

$$V(\pi^E) - \mathbb{E}[V(\pi^{\text{AIL}})] \leq \mathcal{O}\left(\min\left\{1, \sqrt{\frac{|\mathcal{S}|}{N}}\right\}\right)$$

- A new **stage-coupled analysis** reveals why AIL has “**good generalization**” via the lens of “**coupled optimization**”

Key Mechanism

Future distribution matching loss guides AIL to recover expert actions on **preceding uncovered states**



Total-Variation Distance based AIL

Total-Variation Distance based AIL (TV-AIL)

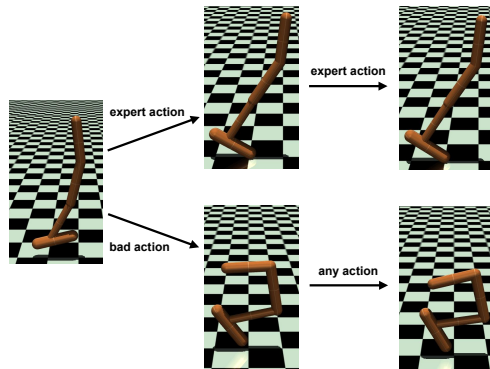
$$\pi^{\text{AIL}} \in \operatorname{argmin}_{\pi \in \Pi} \sum_{h=1}^H \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \left| d_h^\pi(s,a) - \widehat{d}_h^{\pi^{\text{E}}}(s,a) \right| \quad \widehat{d}_h^{\pi^{\text{E}}}: \text{empirical estimate of } d_h^{\pi^{\text{E}}} \text{ from demos.}$$

- **Linear programming formulation** (over d_h^π): the exact solution can be obtained in polynomial time
- **Minimax formulation**: standard gradient-descent-ascent yields an approximate optimal solution in polynomial time

$$\min_{\pi \in \Pi} \max_{c \in \mathcal{C}_{\text{TV}}} \sum_{h=1}^H \mathbb{E}_{(s,a) \sim d_h^\pi} [c_h(s,a)] - \mathbb{E}_{(s,a) \sim \widehat{d}_h^{\pi^{\text{E}}}} [c_h(s,a)],$$

$$\mathcal{C}_{\text{TV}} = \{c = (c_1, \dots, c_H) : c_h : \mathcal{S} \times \mathcal{A} \rightarrow [-1, 1], \forall h \in [H]\}$$

Reachable Bad Absorbing States



The robot takes an expert action and jumps forward, resulting in a positive reward

The robot takes an non-expert action, results in a fall and a zero reward on the **absorbing** state

Assumption (RBAS MDPs). For a tabular and episodic MDP and an expert policy, we define $\mathcal{S}^G = \{s : \exists h \in [H], d_h^{\pi^E}(s) > 0\}$ as expert states and $\mathcal{S}^B = \mathcal{S} \setminus \mathcal{S}^G$ as non-expert states. Assume that

- Non-expert states only have transitions to themselves: $\forall h \in [H], b \in \mathcal{S}^B, a \in \mathcal{A}, \sum_{b' \in \mathcal{S}^B} P_h(b'|b, a) = 1$.
- Expert states are reachable via expert actions and are non-reachable via non-expert actions: $\forall s, s' \in \mathcal{S}^G, h \in [H], P_h(s'|s, \pi_h^E(s)) > 0, \forall a \neq \pi_h^E(s), P_h(s'|s, a) = 0$

If taking a non-expert action, the agent transits into **absorbing non-expert states**

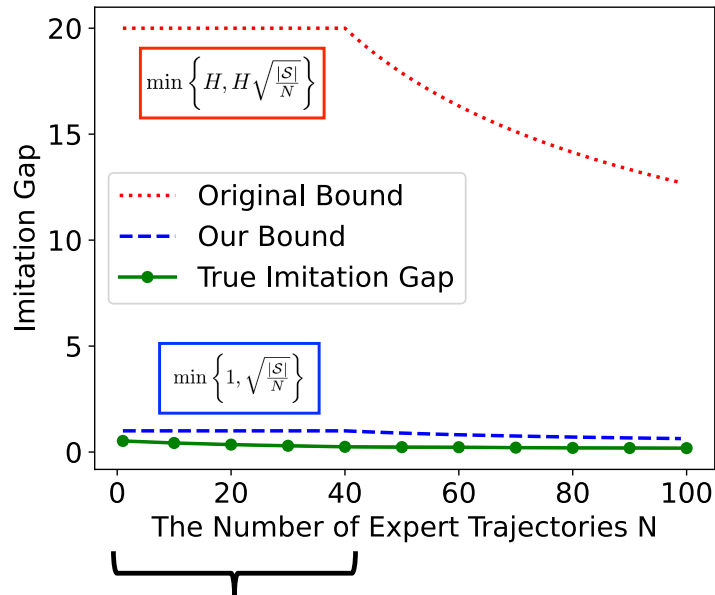
Low-Data-Regime Imitation Gap of TV-AIL

Theorem (Imitation Gap of TV-AIL). For any RBAS MDPs, suppose that π^{AIL} is the policy recovered by TV-AIL, we have that

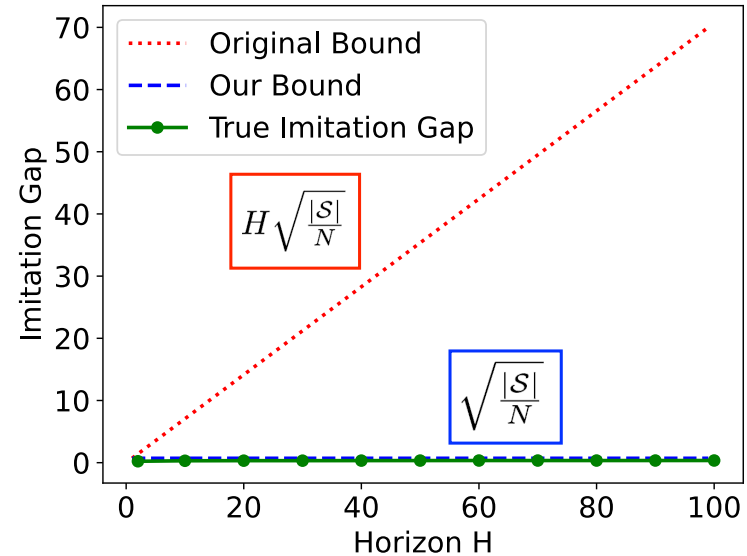
$$V(\pi^{\text{E}}) - \mathbb{E}[V(\pi^{\text{AIL}})] \leq \mathcal{O}\left(\min\left\{1, \sqrt{\frac{|\mathcal{S}|}{N}}\right\}\right)$$

- The **first low data regime bound** to our best knowledge
 - **Small sample regime** ($N \preceq |\mathcal{S}|$): the imitation gap of TV-AIL is at most 1, which is significantly smaller than H
 - **Large sample regime** ($N \succeq |\mathcal{S}|$): the imitation gap of TV-AIL is $\sqrt{|\mathcal{S}|/N}$, which diminishes to 0 as $N \rightarrow \infty$
- The bound is also **horizon-free** in the whole sample regime

New Theory Predicts Practice Well



low data regime



- **Our bound predicts the true imitation gap accurately** with varying numbers of expert trajectories and horizons
- The original bound provides a **loose prediction**

New Theory Predicts Practice Well

Performance of TV-AIL and BC on MuJoCo benchmarks

Varying N

		N = 1	N = 4	N = 7	N = 10
Hopper (scale = 3.2)	BC	784.97 \pm 28.09	887.04 \pm 31.26	666.44 \pm 106.58	460.72 \pm 74.95
	TV-AIL	10.38 \pm 11.25	1.81 \pm 2.88	-4.96 \pm 10.88	2.33 \pm 10.59
HalfCheetah (scale = 7.7)	BC	1058.48 \pm 8.27	1066.21 \pm 22.76	988.07 \pm 35.52	579.53 \pm 171.28
	TV-AIL	-22.45 \pm 101.65	-84.96 \pm 16.84	-78.69 \pm 6.98	-79.29 \pm 7.95
Walker2d (scale = 5.0)	BC	1002.13 \pm 9.68	939.27 \pm 19.10	528.91 \pm 181.46	222.98 \pm 52.97
	TV-AIL	12.89 \pm 19.61	9.04 \pm 17.99	-4.91 \pm 8.14	14.36 \pm 12.48

TV-AIL matches the expert performance **with limited demos. (N=1)** while BC suffers a large imitation gap, aligning with the **low-data property** of the new bound

Varying H

		H = 100	H = 500	H = 1000	H = 2000
Hopper (scale=3.2)	BC	0.80 \pm 1.72	178.56 \pm 79.27	784.97 \pm 28.09	1950.42 \pm 36.20
	TV-AIL	4.96 \pm 4.56	-1.73 \pm 6.77	10.38 \pm 11.25	-9.92 \pm 37.15
HalfCheetah (scale=7.7)	BC	56.44 \pm 9.23	491.91 \pm 21.01	1058.48 \pm 8.27	2198.61 \pm 12.93
	TV-AIL	8.49 \pm 8.26	-24.73 \pm 11.46	-22.45 \pm 101.65	-169.83 \pm 16.06
Walker2d (scale=5.0)	BC	10.25 \pm 7.89	413.02 \pm 13.65	1002.13 \pm 9.68	2158.84 \pm 2.05
	TV-AIL	-0.18 \pm 1.04	16.26 \pm 22.64	12.93 \pm 19.61	71.69 \pm 66.30

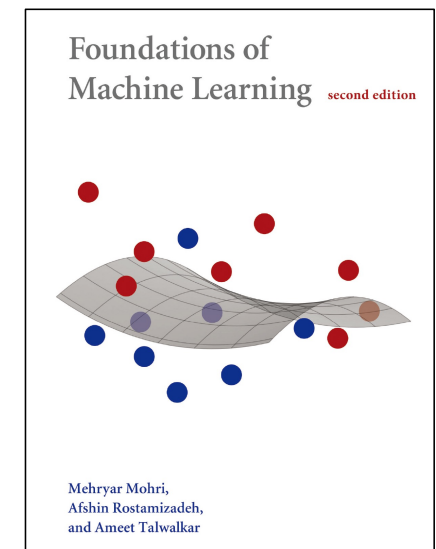
When increasing H, TV-AIL **maintains small imitation gaps** while BC's imitation gaps blow up, aligning with the **horizon-free property** of the new bound

An Illustration: Classical Machine Learning Theory

Classical machine learning theory often provides such a generalization error bound

$$R(h) \leq \hat{R}(h) + \tilde{O} \left(\sqrt{\frac{\log |\mathcal{H}|}{N}} \right)$$

- $R(h) = \mathbb{E}_{(x,y) \sim P} [\mathbb{I}\{h(x) \neq y\}]$ is the expected risk
- $\hat{R}(h) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}\{h(x_i) \neq y_i\}$ is the empirical risk
- $\log |\mathcal{H}|$ is the complexity of the hypothesis class



This bound is meaningless **in the small sample regime** $N \leq \log |\mathcal{H}|$ since the generalization error is 1

Classical Reduction-and-Estimation Analysis

Existing AIL analysis follows classical ml theory and decomposes the imitation gap into **statistical error** and **optimization error**

$$\begin{aligned} & |V(\pi^E) - V(\pi^{\text{AIL}})| \\ & \leq \dots \\ & \leq \underbrace{\sum_{h=1}^H \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \left| d_h^{\pi^E}(s,a) - \widehat{d}_h^{\pi^E}(s,a) \right|}_{\text{statistical error}} + \underbrace{\sum_{h=1}^H \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \left| \widehat{d}_h^{\pi^E}(s,a) - d_h^{\pi^{\text{AIL}}}(s,a) \right|}_{\text{optimization error}} \\ & \leq 2 \underbrace{\sum_{h=1}^H \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \left| \widehat{d}_h^{\pi^E}(s,a) - d_h^{\pi^E}(s,a) \right|}_{\text{statistical error}} \\ & \leq \mathcal{O} \left(H \sqrt{\frac{|\mathcal{S}|}{N}} \right) \end{aligned}$$

Classical Reduction-and-Estimation Analysis

Claim: The reduction to statistical error is very **loose** for ALL in the **low data regime**

TABLE 5: Imitation gap and estimation error of TV-AIL on RBAS MDPs with $N = 1$.

	$H = 100$	$H = 500$	$H = 1000$	$H = 2000$
Imitation Gap	$0.69_{\pm 0.00}$	$0.70_{\pm 0.00}$	$0.71_{\pm 0.00}$	$0.71_{\pm 0.00}$
Estimation Error	$189.47_{\pm 0.00}$	$947.37_{\pm 0.00}$	$1894.74_{\pm 0.00}$	$3789.47_{\pm 0.00}$

In the low data regime, despite **a large estimation error**, TV-AIL can generalize well and still achieve **a small imitation gap**

Sharp Characterization of TV-AIL



Target: Establish a sharp characterization of π^{AIL} instead of reducing to estimation error coarsely

$$\pi^{\text{AIL}} \in \operatorname{argmin}_{\pi \in \Pi} \sum_{h=1}^H \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \left| d_h^\pi(s,a) - \widehat{d}_h^{\pi^{\text{E}}}(s,a) \right| \quad \widehat{d}_h^{\pi^{\text{E}}}: \text{empirical estimate of } d_h^{\pi^{\text{E}}} \text{ from demos.}$$



Difficulty: This is a **non-smooth and non-convex** problem, making classical optimality conditions unapplicable

Proposition. There exist tabular and episodic MDPs such that the objective of TV-AIL is non-convex.

Stage-Coupled Analysis

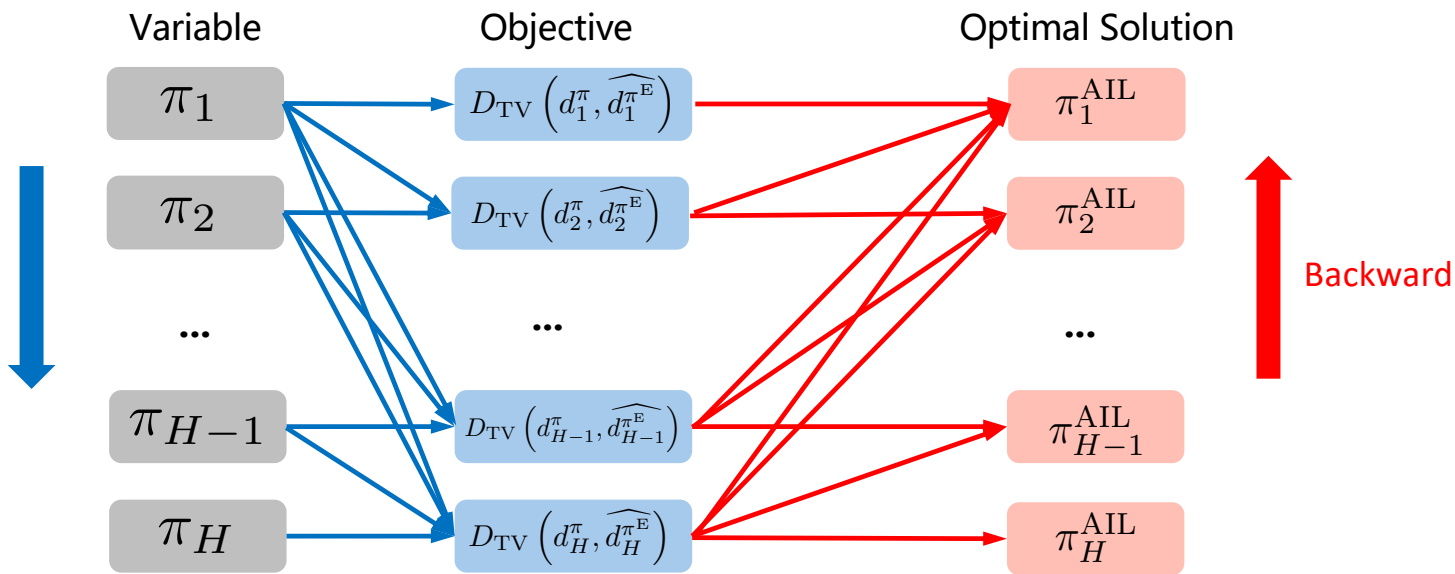
$$\pi^{\text{AIL}} \in \operatorname{argmin}_{\pi} \sum_{h=1}^H D_{\text{TV}} \left(d_h^{\pi}, \widehat{d}_h^{\pi^{\text{E}}} \right)$$



Forward Stage-coupled Structure

Backward Dynamic Programming Analysis

Forward



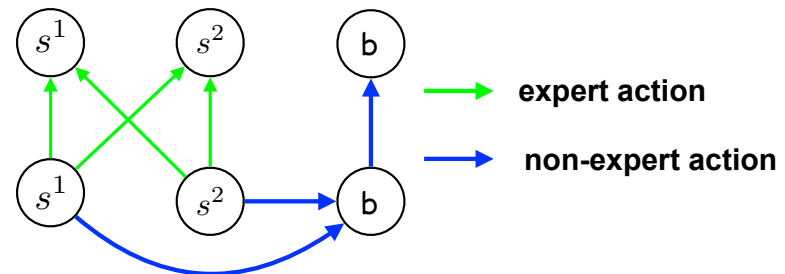
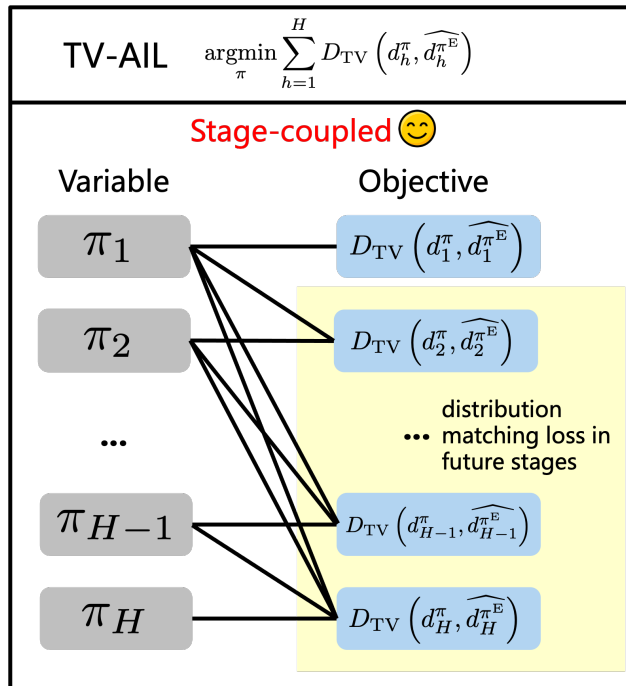
Optimality Condition of TV-AIL

Proposition. For any RBAS MDP, suppose that π^{AIL} is the policy recovered by TV-AIL. When $N \geq 1$, the following optimality condition holds almost surely:

$$\pi_h^{\text{AIL}}(\pi_h^{\text{E}}(s)|s) = 1, \forall s \in \mathcal{S}^G, h \in [H - 1].$$

- In the first $H-1$ stages, AIL can recover expert actions **even on states uncovered by demonstrations**, implying its **generalization ability beyond demonstrations**
- With additional analysis on stage H , we can obtain the final bound $\mathcal{O}(\min\{1, \sqrt{|\mathcal{S}|/N}\})$

Mechanism Under the Generalization of TV-AIL

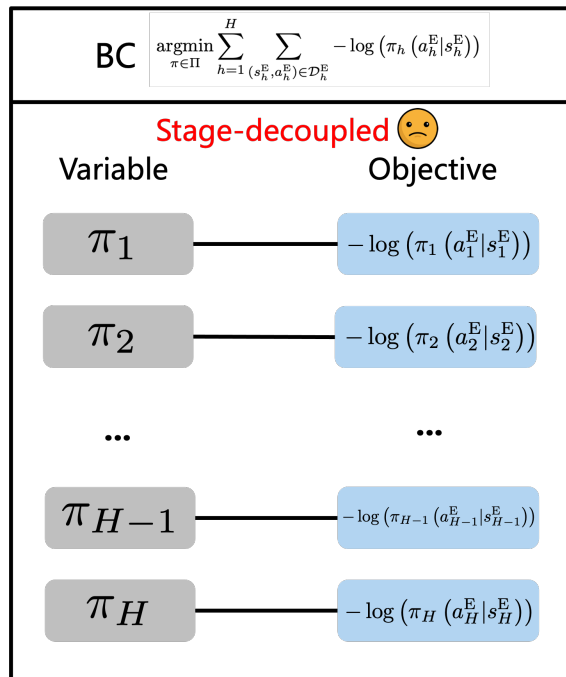


- demonstrations = $\{s^1 \rightarrow s^1\}$
- uncovered state s^2

$$\pi_1^{\text{AIL}}(a | s^2) = 1$$

- Key insight: Due to the **coupling structure**, **future distribution loss** provably guides AIL to recover expert actions on **preceding uncovered states**
- **The RL procedure** for multi-step distribution matching is the key for the generalization in AIL

Stage-decoupled Objective in BC



- The BC objective is **stage-decoupled**: each policy π_h is optimized solely against expert data at step h , with no signal from other time steps
- Thus BC cannot generalize beyond demonstrations

Extension for Approximate Solutions

Exact solution:

$$\pi^{\text{AIL}} \in \operatorname{argmin}_{\pi \in \Pi} \sum_{h=1}^H \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \left| d_h^\pi(s,a) - \widehat{d}_h^{\pi^{\text{E}}}(s,a) \right|$$



Approximate solution:

Definition (ε -optimal solution). A policy $\bar{\pi}$ is an ε -optimal solution for TV-AIL if

$$\sum_{h=1}^H \left\| d_h^{\bar{\pi}} - \widehat{d}_h^{\pi^{\text{E}}} \right\|_1 \leq \min_{\pi \in \Pi} \sum_{h=1}^H \left\| d_h^\pi - \widehat{d}_h^{\pi^{\text{E}}} \right\|_1 + \varepsilon.$$

Extension for Approximate Solutions

Theorem. For any RBAS MDP, the candidate policy set is defined as $\Pi^{\text{OPT}} = \{\pi \in \Pi : \forall h \in [H], \exists s \in \mathcal{S}^G, \pi_h(\pi_h^E(s)|s) > 0\}$. Suppose that $\bar{\pi} \in \Pi^{\text{OPT}}$ is an ε -optimal solution of TV-AIL. Then we have

$$V(\pi^E) - \mathbb{E}[V(\bar{\pi})] \lesssim \min \left\{ 1 + \frac{8\varepsilon}{c(\bar{\pi})}, \sqrt{\frac{|\mathcal{S}|}{N}} + \frac{8\varepsilon}{c(\bar{\pi})} \right\},$$

where $c(\bar{\pi}) > 0$ is defined as $c(\bar{\pi}) := \min_{1 \leq \ell < h \leq H, s, s' \in \mathcal{S}^G} \{\mathbb{P}^{\bar{\pi}}(s_h = s \mid s_\ell = s', a_\ell = \pi_h^E(s))\}$. Here $\mathbb{P}^{\bar{\pi}}(s_h = s \mid s_\ell = s', a_\ell = \pi_h^E(s))$ is the visitation probability of s in time step h by starting from $s', \pi_h^E(s)$ in time step ℓ , which is jointly determined by the transition function and policy $\bar{\pi}$.

- The low data regime imitation gap bound **still holds for approximate solutions** with $8\varepsilon/c(\bar{\pi})$ being the induced error
- We apply the dynamic programming analysis to establish the **growth condition of TV-AIL** $\text{dist}(\bar{\pi}, \pi^{\text{AIL}}) \leq f(\bar{\pi}) - f(\pi^{\text{AIL}})$, which is highly technical and of independent interests

Imitation Gap Bound for Any MDP



Q: Can TV-AIL achieve the desired low data regime bound in **all MDPs**?

Proposition. Suppose that π^{AIL} is the policy output by TV-AIL, there exists an MDP such that

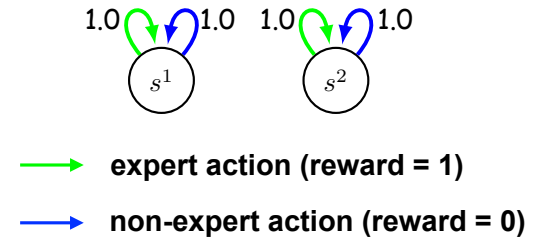
$$V(\pi^{\text{E}}) - \mathbb{E}[V(\pi^{\text{AIL}})] \geq \Omega\left(\min\left\{H, H\sqrt{\frac{|\mathcal{S}|}{N}}\right\}\right)$$

A: In the worst case, TV-AIL must suffer an imitation gap of $\min\{H, H\sqrt{|\mathcal{S}|/N}\}$, which equals to H in the low data regime

Hard Instance for AIL

Assumption. For a tabular and episodic MDP and an expert policy, we assume that

- Each state is absorbing and each action has the same transitions, i.e., $\forall (s, a) \in \mathcal{S} \times \mathcal{A}, h \in [H]$, we have $P_h(s|s, a) = 1$.
- For any state, a^1 is the expert action with a reward 1 and the others are non-expert actions with a reward 0.



- Why this instance is hard for AIL?

1. Each action has the same absorbing transitions

2. Key insight: state-action distributions becomes **decoupled across stages**

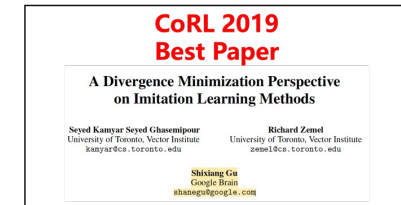
$$\forall (s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H], d_h^\pi(s, a) = \rho(s)\pi_h(a|s)$$

3. Due to this decoupling, TV-AIL cannot leverage multi-step distribution matching to generalize beyond demonstrations

- Remark: This instance is very artificial and rarely occur in practice where AIL exhibits excellent performance

Summary

Why AIL outperforms BC in the low-data regime?



- The **first** low-data regime imitation gap bound

Theorem. Under mild assumptions, AIL can achieve an imitation gap bound:

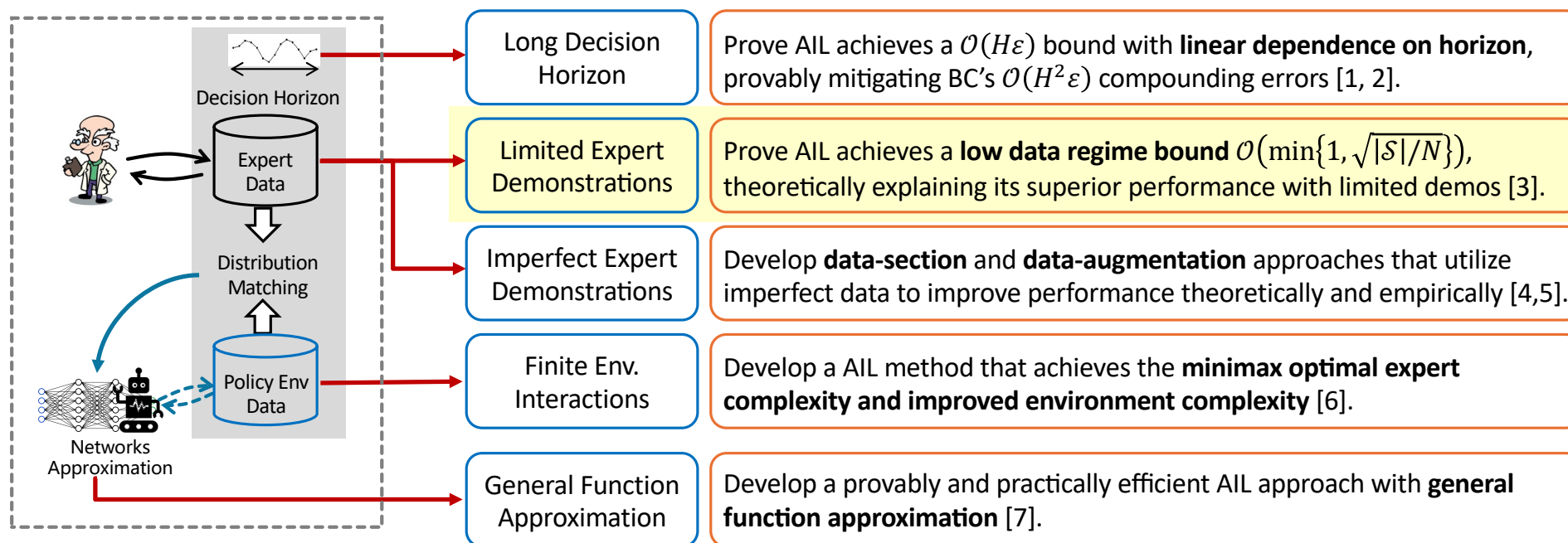
$$V(\pi^E) - \mathbb{E}[V(\pi^{\text{AIL}})] \leq \mathcal{O}\left(\min\left\{1, \sqrt{\frac{|S|}{N}}\right\}\right)$$

- A stage-coupled analysis

Reveal why AIL has **good “generalization”** through the lens of **coupled “optimization”**

- New insights
 - Algorithmic mechanism: **Future** distribution loss provably guides AIL to recover expert actions on **uncovered** states
 - Task understanding: **The coupling in state-action distributions** is the key for generalization

Our Exploration on Imitation Learning Theory



[1] Tian Xu, Ziniu Li, Yang Yu. Error bounds of imitating policies and environments for reinforcement learning. TPAMI.

[2] Yi-Chen Li*, Tian Xu*, Yang Yu*, et al. Generalist reward models: found inside large language models. arXiv. 2025.

[3] Tian Xu*, Ziniu Li*, Yang Yu, Zhi-Quan Luo. Understanding AIL in Small Sample Regime: A Stage-coupled Analysis. TPAMI.

[4] Ziniu Li*, Tian Xu*, Yang Yu, Zhi-Quan Luo. Imitation Learning from Imperfection: Theoretical Justifications and Algorithms. NeurIPS 2023.

[5] Zhilong Zhang*, Tian Xu*, Xinghao Du*, et al. Improving Reward Model Generalization from Adversarial Process Enhanced Preferences. ICML 2025.

[6] Tian Xu*, Ziniu Li*, Yang Yu, Zhi-Quan Luo. Provably Efficient Adversarial Imitation Learning with Unknown Transitions. UAI 2023.

[7] Tian Xu*, Zhilong Zhang*, et al. Adversarial Imitation Learning with General Function Approximation: Theoretical Analysis and Practical Algorithms. TPAMI.