



TEXAS A&M UNIVERSITY

Engineering

Learning from Demonstrations: Applications to Autonomous UAV Landing

Student: Prabhasa Kalkur
Advisor: Dr. Dileep Kalathil

Oct 05, 2020

What is imitation learning?

Learning to imitate from expert behavior

Sample-efficient learning: learn behavior from as little expert data as possible



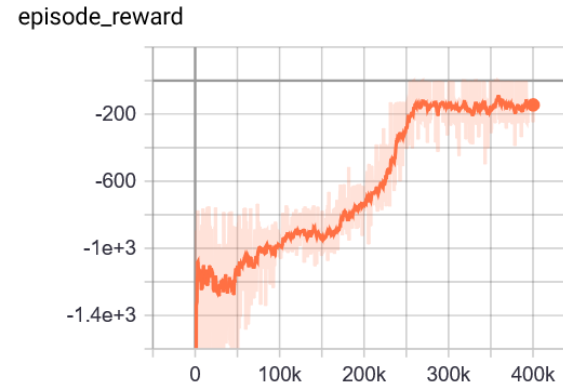
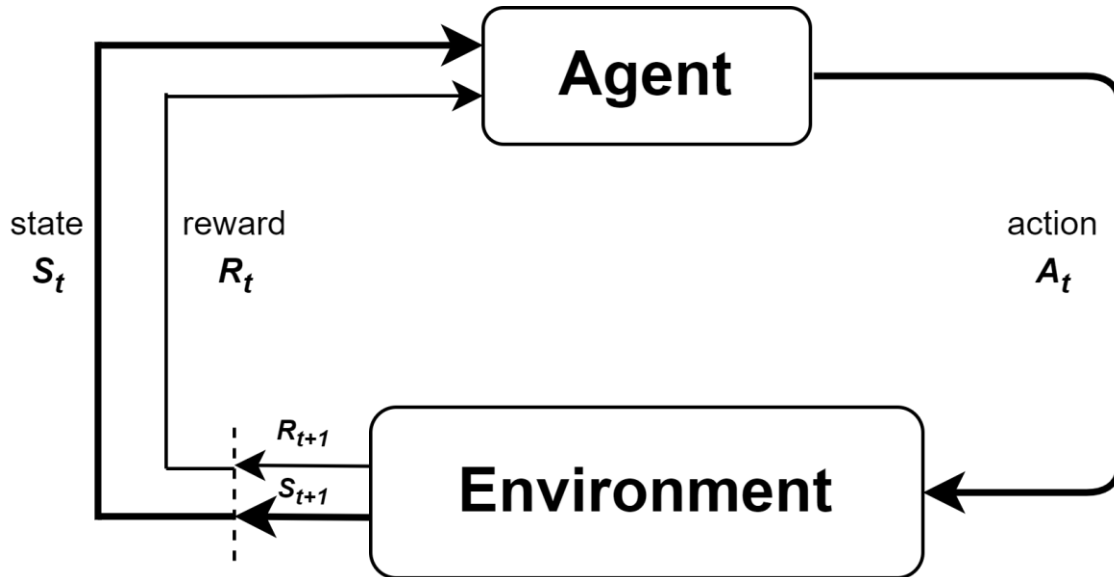


What is the presentation about?

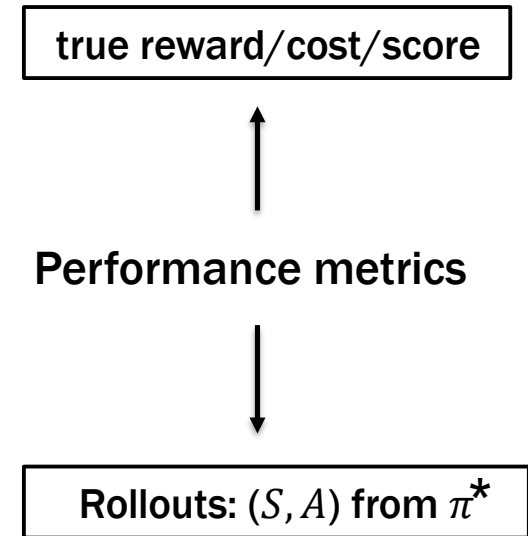
- Motivate the need for **sample-efficient** methods for learning behavior
- Pick Generative Adversarial Imitation Learning (**GAIL**) as our algorithm
- Apply GAIL to the **sparsely-rewarded** task of landing a drone (simulation)

Reinforcement Learning

- $s, s' \in S, a \in A$. Consider tuple $[S, A, P(s'|s, a), R(s, a), \gamma, H]$, define a policy (model) $\pi : S \rightarrow A$
 - Reinforcement Learning (RL): find an optimal π^* that maximizes $\sum_{t=0}^{\infty} \gamma^t R_t$



100 episodes of policy:
95/100 successful
Reward (mean, std): (-175, 50)





Organization of the talk

1. Need for sample-efficiency
2. Introduction to Imitation Learning
3. Application: Autonomous UAV Landing
4. Conclusions and Future Work

Sections



1. Need for sample-efficiency
2. Introduction to Imitation Learning
3. Application: Autonomous UAV Landing
4. Conclusions and Future Work

Why study imitation learning?

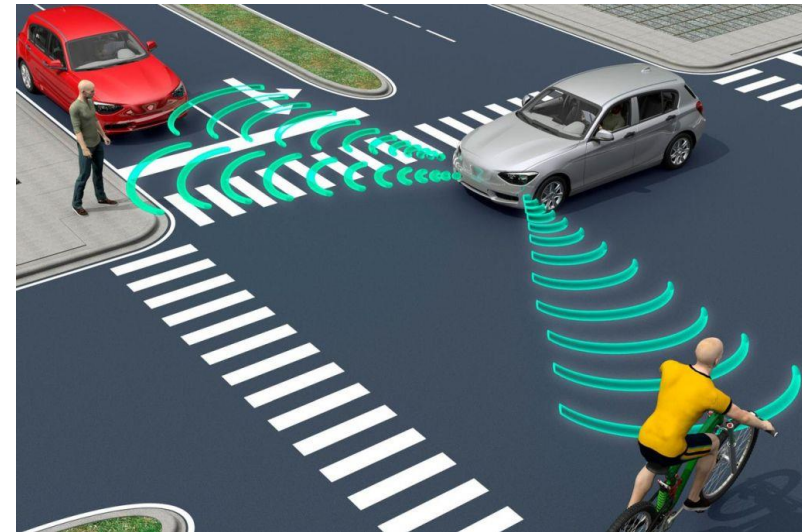
1. Rewards obvious in computer games: maximize score
 - Not so obvious in real-world scenarios: use a proxy instead

reward



Mnih et al. '15

VS



Why study imitation learning?

2. Can be easier to **demonstrate** desired behavior

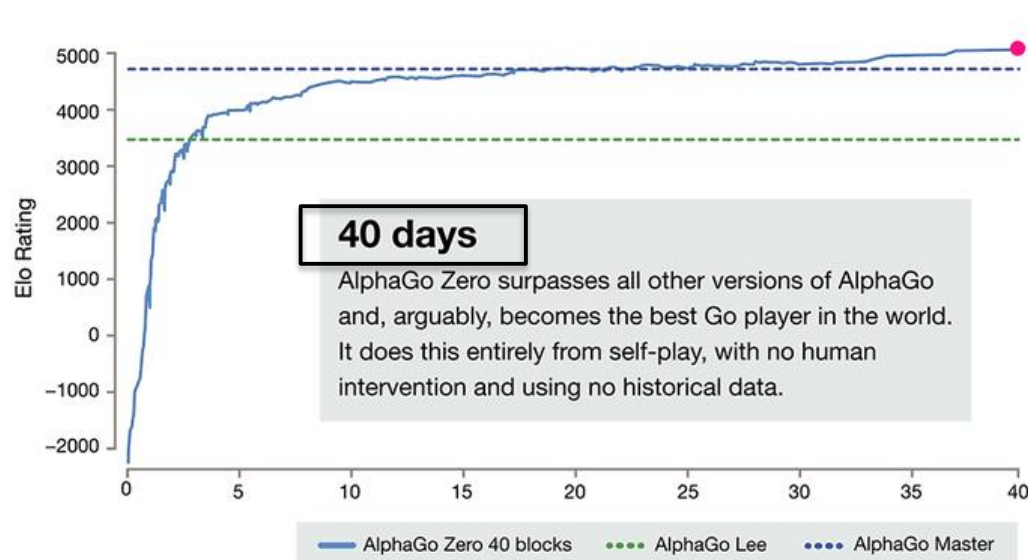


Levine et al. "Learning hand-eye coordination for robotic grasping with deep learning and large-scale data collection."

Why study imitation learning?

3. Modern Deep-RL requires exponentially increasing number of samples: **sample-inefficient**

- Challenging for the AI community to reproduce SOTA results



Go: AlphaGo Zero

	OPENAI 1V1 BOT	OPENAI FIVE
CPUs	60,000 CPU cores on Azure	128,000 <u>preemptible</u> CPU cores on GCP
GPUs	256 K80 GPUs on Azure	256 P100 GPUs on GCP
Experience collected	~300 years per day	~180 years per day (~900 years per day counting each hero separately)
Size of observation	~3.3 kB	~36.8 kB
Observations per second of gameplay	10	7.5
Batch size	8,388,608 observations	1,048,576 observations
Batches per minute	~20	~60

Dota 2: OpenAI Five

Why study imitation learning?

3. Modern Deep-RL requires exponentially increasing number of samples
 - **Not practical, especially when env samples are expensive, and compute is limited**
 - One approach: use sample-efficient methods like Imitation Learning

Many competitions trying to promote compute and sample-efficient learning:

- **NeurIPS 2019: Game of Drones**
- **NeurIPS 2019 & 2020: MineRL Challenge**

Why study imitation learning?

4. How humans and animals fundamentally learn behavior



Picture credits: Sapana

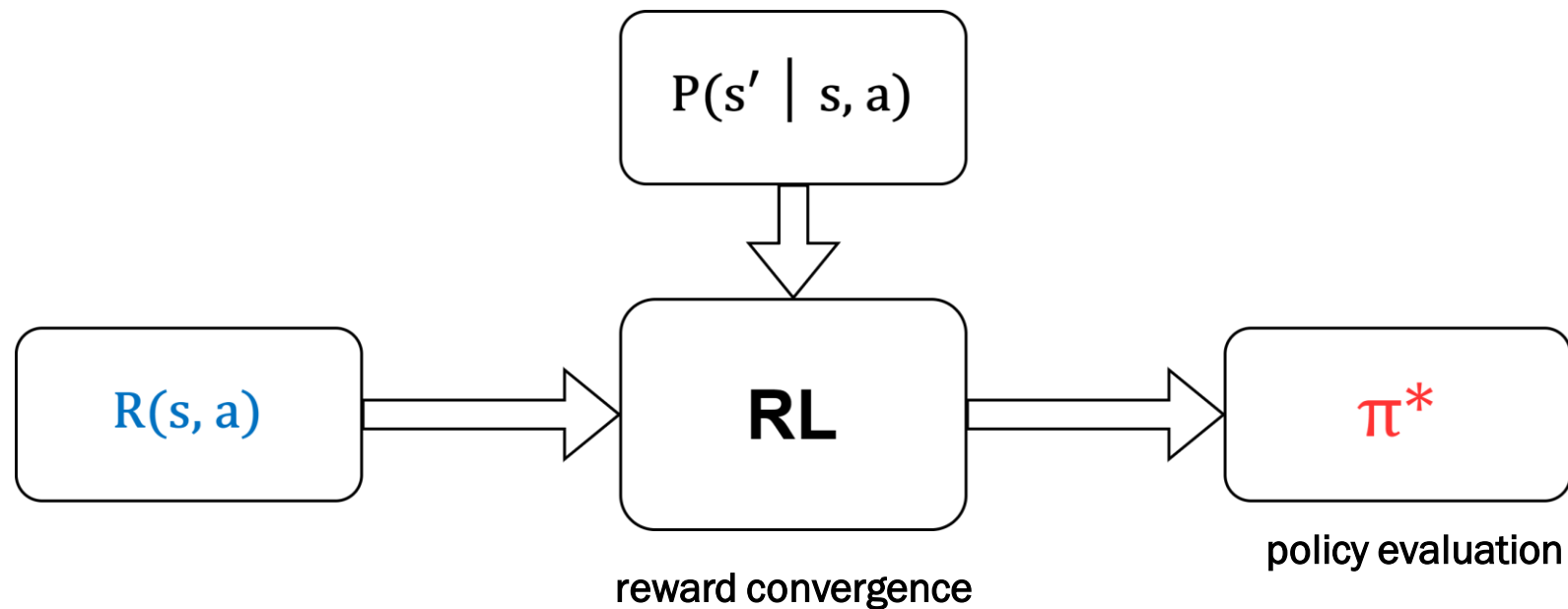


Sections

1. Need for sample-efficiency
2. Introduction to Imitation Learning
3. Application: Autonomous UAV Landing
4. Conclusions and Future Work

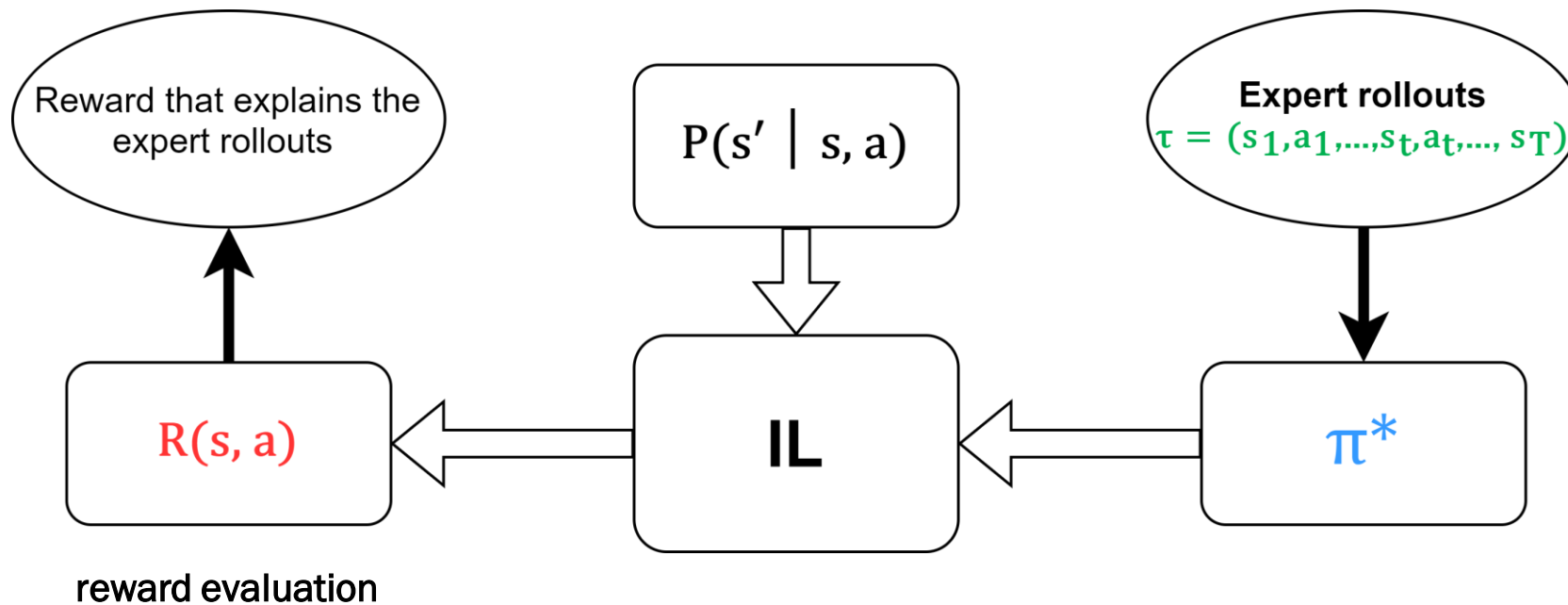
RL algorithms

- $s, s' \in S, a \in A$. For MDP $[S, A, P(s'|s, a), R(s, a), \gamma]$, define a policy $\pi : S \rightarrow A$
 - **Goal:** find an optimal π^* that maximizes $\sum_{t=0}^{\infty} \gamma^t R_t$
 - **Metric:** (i) Reward convergence, (ii) Policy evaluation (testing)



IL algorithms

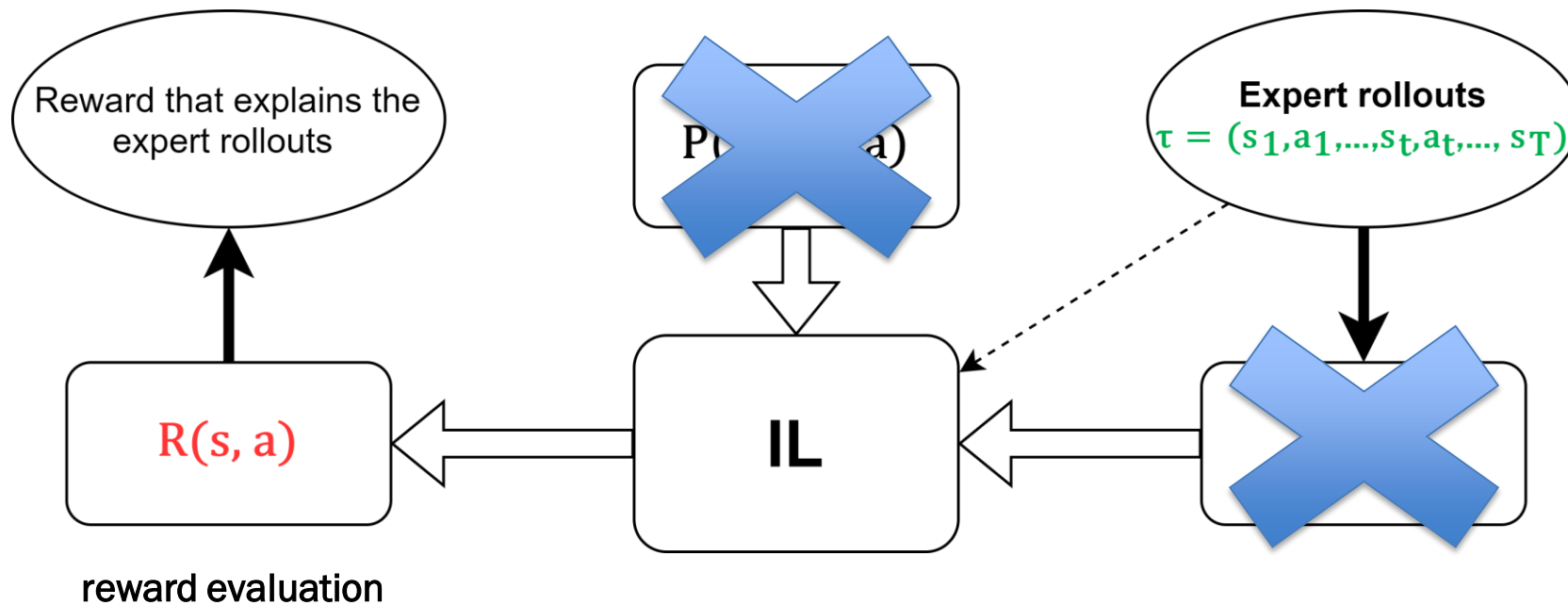
- $s, s' \in S, a \in A$. For MDP $[\mathcal{S}, \mathcal{A}, P(s'|s, a), R(s, a), \gamma]$, define a policy $\pi : S \rightarrow A$
 - **Goal:** given $\tau = (s_0, a_0, s_1, a_1, \dots, s_t, a_t, \dots, s_T)$ generated from a π^* , extract its $R(s, a)$
 - **Metric:** Reward evaluation (?)

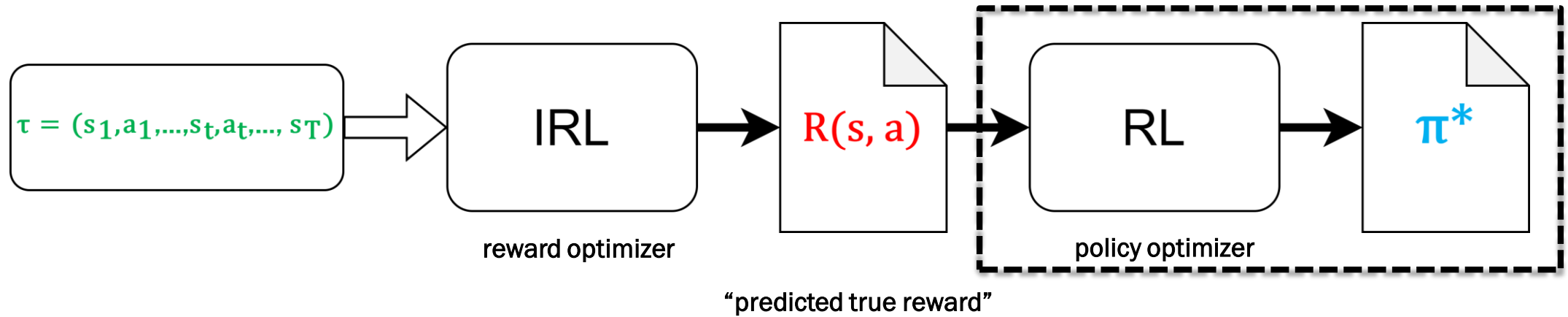
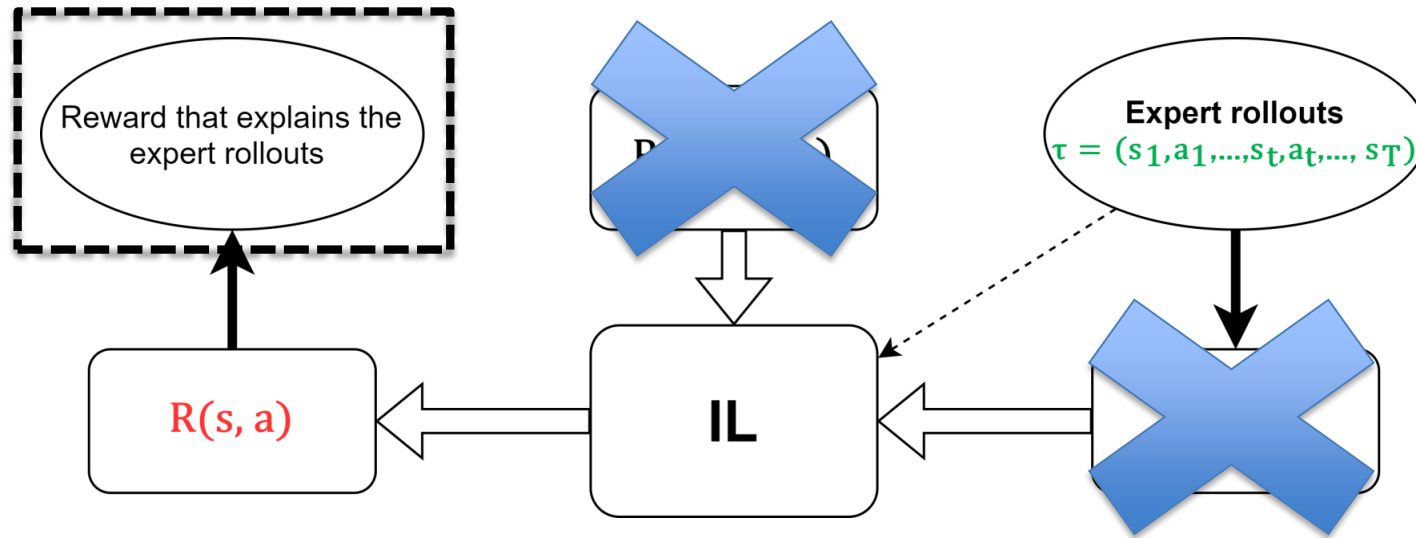


Flowchart credits: Sapana

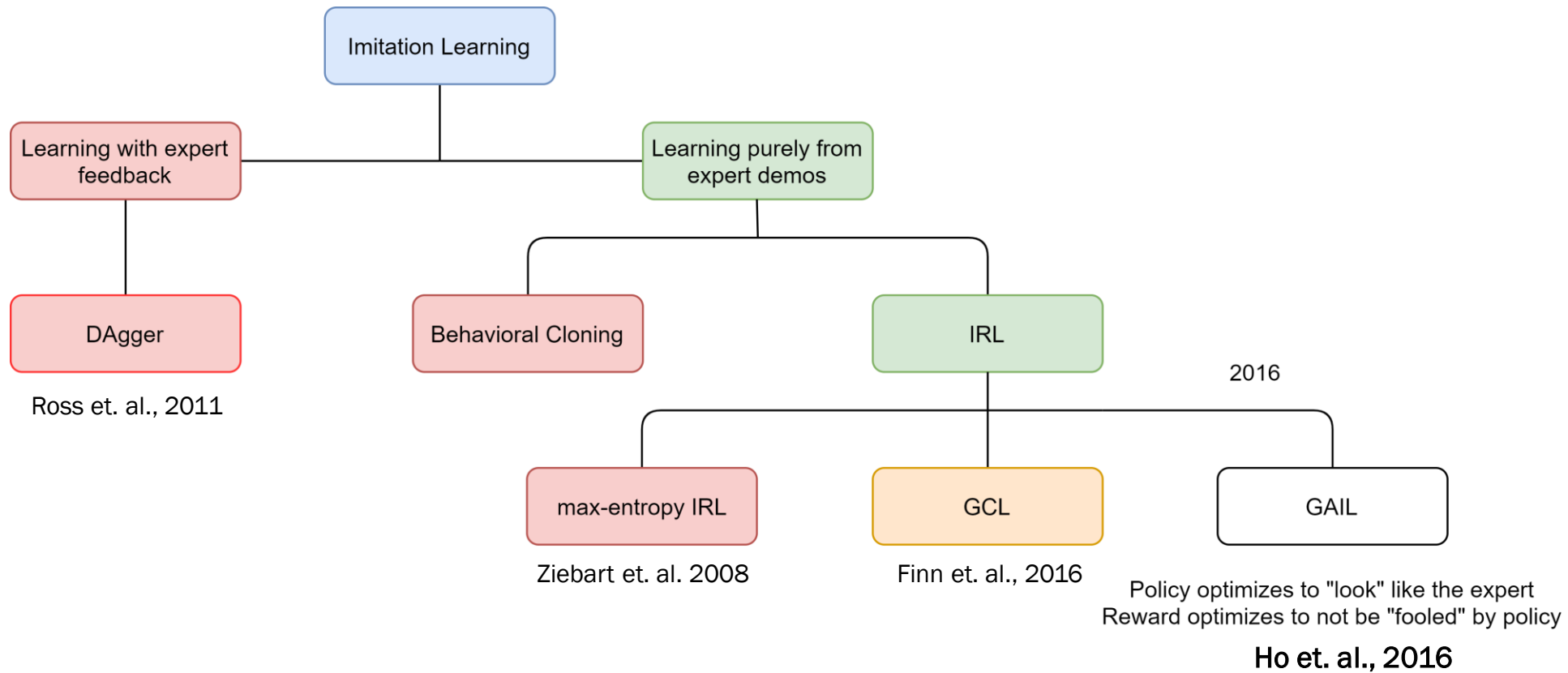
IL algorithms

- $s, s' \in S, a \in A$. For MDP $[\mathcal{S}, \mathcal{A}, P(s'|s, a), R(s, a), \gamma]$, define a policy $\pi : S \rightarrow A$
 - **Goal:** given $\tau = (s_0, a_0, s_1, a_1, \dots, s_t, a_t, \dots, s_T)$ generated from a π^* , extract its $R(s, a)$
 - **Metric:** Reward evaluation (?)





Imitation Learning approaches



Generative Adversarial Imitation Learning (GAIL) is the **SOTA IL algorithm**

Some questions...

1. How does imitation accuracy scale with problem dimensionality and demo data?
2. How 'smooth' are the learned policies compared to the expert policy?
3. Can behaviors with sparse rewards be learned? At what cost?
4. Can GAIL imitate suboptimal experts? At what cost?
5. Can GAIL generalize?

Let us learn how to imitate a simple control task: balance an inverted pendulum!



Problem setup

Train RL -> rollout **expert** -> Train GAIL -> **policy** evaluation (test)

Goal: GAIL should be able to 'imitate' expert (optimal/suboptimal?)

Discuss: imitation accuracy, sample efficiency, effect of reward quality on learning

- **Expert trajectories / rollout / demonstrations:** sample demos [5, 10, 20]
- **Policy evaluation / rollout / testing:** Check policy performance for 100 episodes
- **Task solved each episode:** True reward for 100 consecutive episodes during training

Tools

- **RL library:** Stable Baselines 2.10
- **Framework:** TensorFlow 1.14
- **Hyperparameters (HPs):** RL Baselines Zoo, etc.
- **Performance metrics (learned reward vs episodes, test scores):** Tensorboard 1.14, W&B 0.10

RL/IL Algorithms

- **SAC** – Soft Actor-Critic (optimal experts)
- **TRPO** – Trust Region Policy Optimization (policy optimizer for GAIL)
- **BC** - Behavioral Cloning* (comparison with GAIL)

*with policy: “MlpPolicy” [100, 100], optimizer: Adam, batch size: 256, train-val: 70-30



OpenAI Gym and MuJoCo

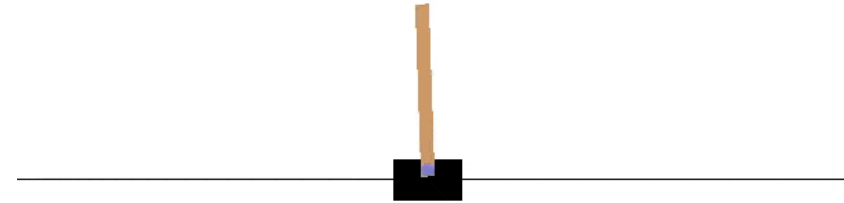
- **Gym:** “Toolkit for developing and comparing reinforcement learning algorithms”
- Platform for teaching agents to perform simulated tasks **under a true reward**
- E.g. Atari games, Robotic manipulation, control tasks

- **MuJoCo:** “A physics engine that does very detailed, efficient simulations with contacts”
- E.g. Continuous control tasks like hopping, walking, or running

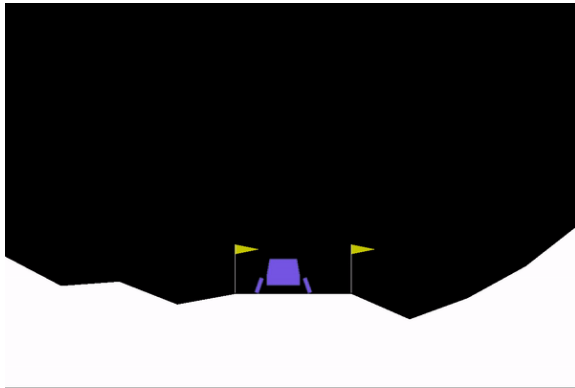
- **Why is this important?** Standard benchmark tasks for testing RL, IL algorithms



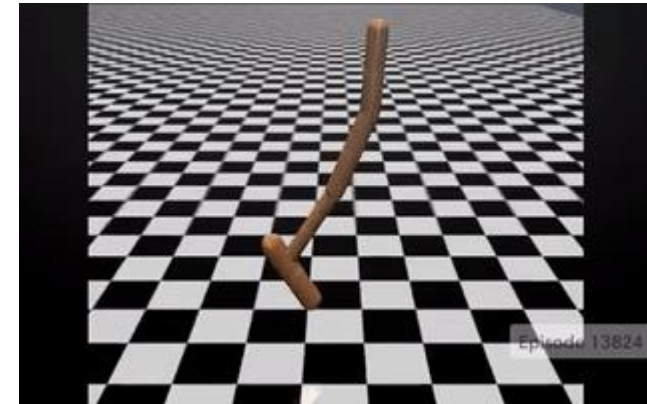
Pendulum-v0



CartPole-v1



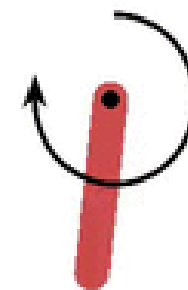
LunarLanderCts-v2



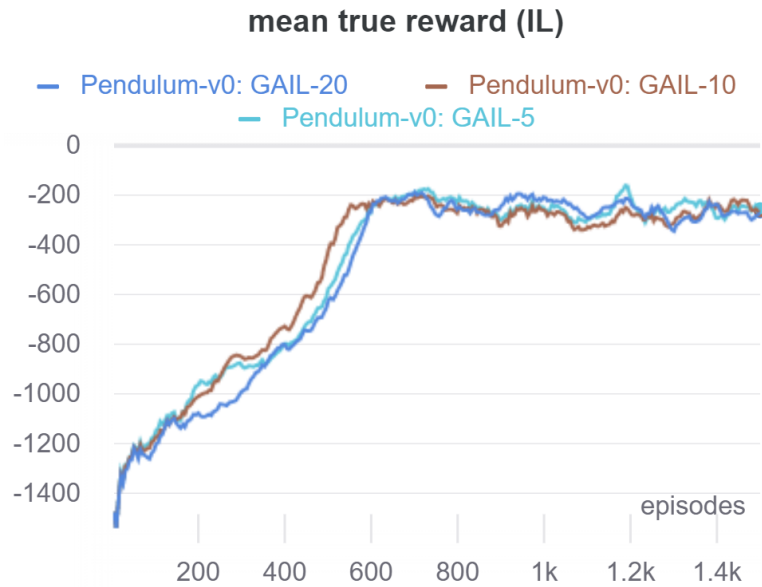
Hopper-v2

The Pendulum-v0 environment

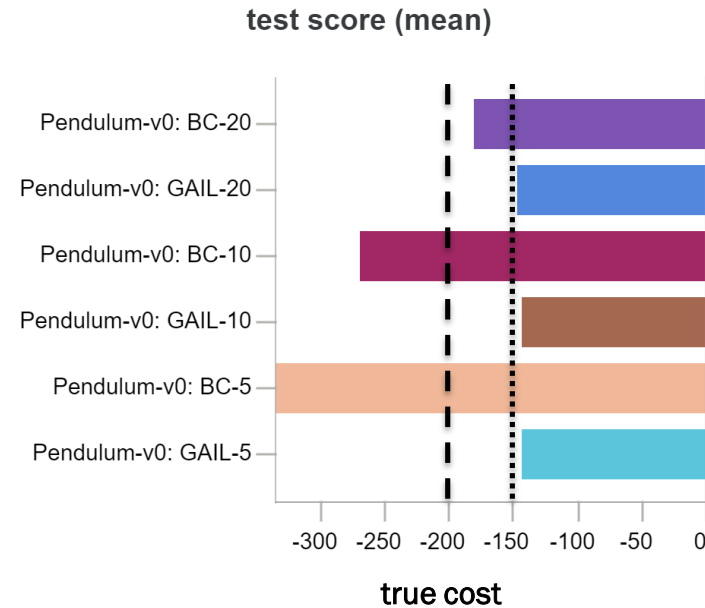
Properties	Description
State space (cts, dim = 3)	Cosine, sine of angle θ $[-1, 1]$, $\dot{\theta}_0$ $[-8, 8]$
Action space (cts, dim = 1)	Joint effort $[-2, 2]$
Reward	$-(\theta^2 + 0.1*\dot{\theta}_0^2 + 0.001*action^2)$, dense
Termination / Horizon	200 steps, finite
Solved / learned task	defined as -200 mean reward over 100 consecutive episodes of training
Expert Trajectories for IL	[5, 10, 20] with reward (mean, var): (-147, 84)



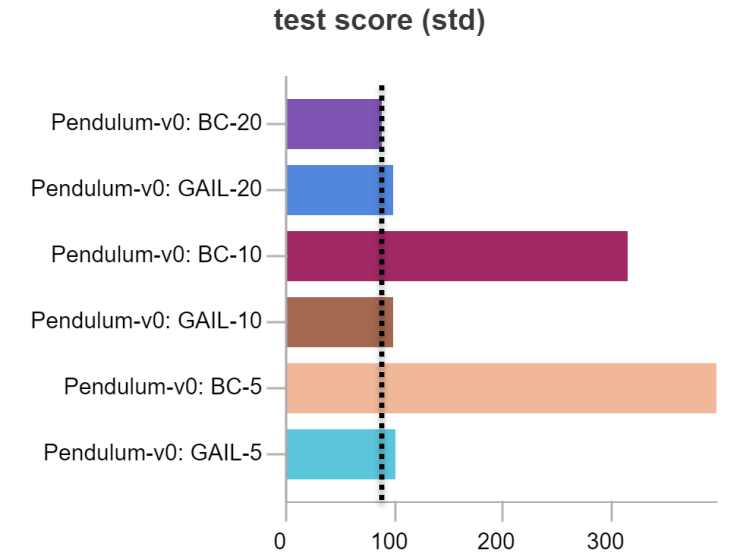
Pendulum-v0: GAIL and BC



reward convergence



policy evaluation

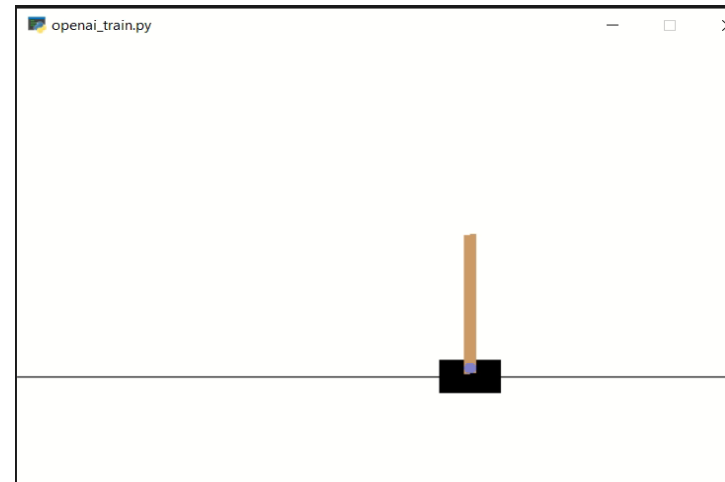


GAIL learns to achieve true cost **AND** imitate expert
GAIL score (mean, var) consistent over # demos – **sample-efficient**
BC improves over # demos, but only for optimal experts

Imitating suboptimal experts

Data	Optimal	Solved	Score	Mean length	Success
Expert	No	475	(402, 134)	402	63/100
GAIL policy	No	475	(410, 178)	411	59/100

- GAIL on suboptimal CartPole-v1 expert
- **Suboptimal experts can be imitated!**



Some questions...

1. How does imitation accuracy scale with dimensionality, demo data? **GAIL sample-efficient (low-dim)**
2. How 'smooth' are the learned policies compared to the expert policy? **Demo-dependent**
3. Can behaviors with sparse rewards be learned? At what cost?
4. Can GAIL imitate suboptimal experts? At what cost? **BC cannot. GAIL can, with the right HPs**
5. Can GAIL generalize?

Answered for a **low-dimensional, densely-rewarded, finite-horizon** control task. Let's try harder!

Sections

1. Need for sample-efficiency
2. Introduction to Imitation Learning
3. **Application: Autonomous UAV Landing**
4. Conclusions and Future Work

AirSim: Autonomous UAV Navigation and Landing

APPLICATION 1

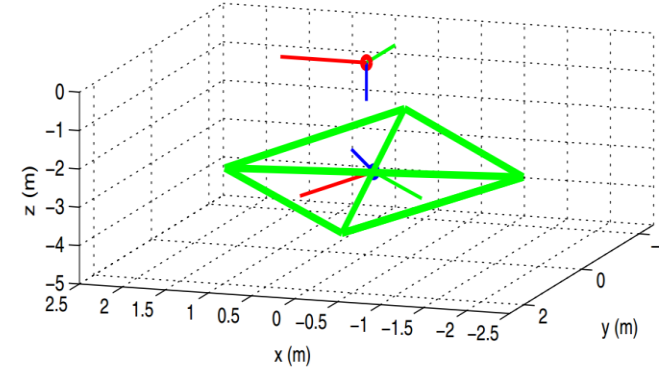
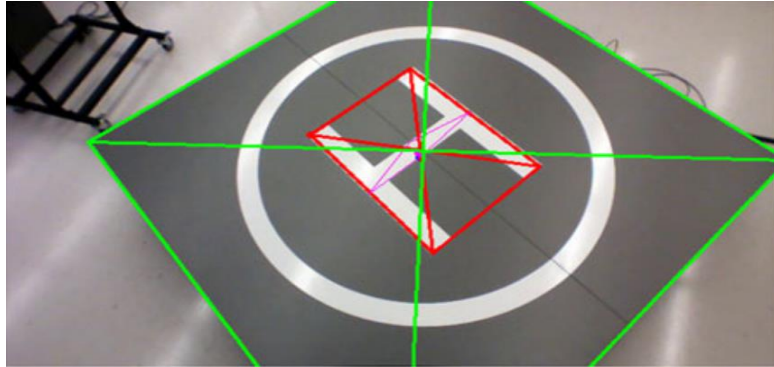
Landing on ships

Landing Zone	Ground	Ship
Space	Large	Limited
Motion	None	6 DOF
Visual References	More	Less
Alternate L/D Places	Many	Less
Weather	Affected	Extremely affected



Slide credits: Bochan

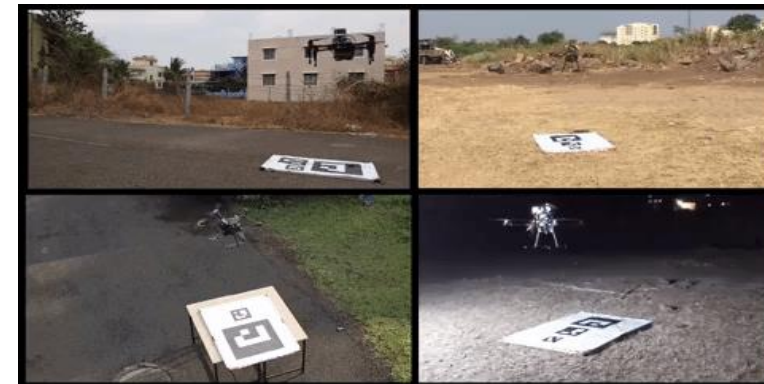
Common Approaches



ALL of them are looking at landing spot



Computer vision-based for autonomous landing
by G.Xu, Pattern Recogn.Lett., 2009



Flytdock by flytbase company, June 8, 2018

Slide credits: Bochan

How does a pilot approach the ship?



Slide credits: Bochan

CONTRIBUTIONS



Landing a UAV on a ship without looking at landing spot (simulation)



Refer to a visual cue for positioning, just like a pilot



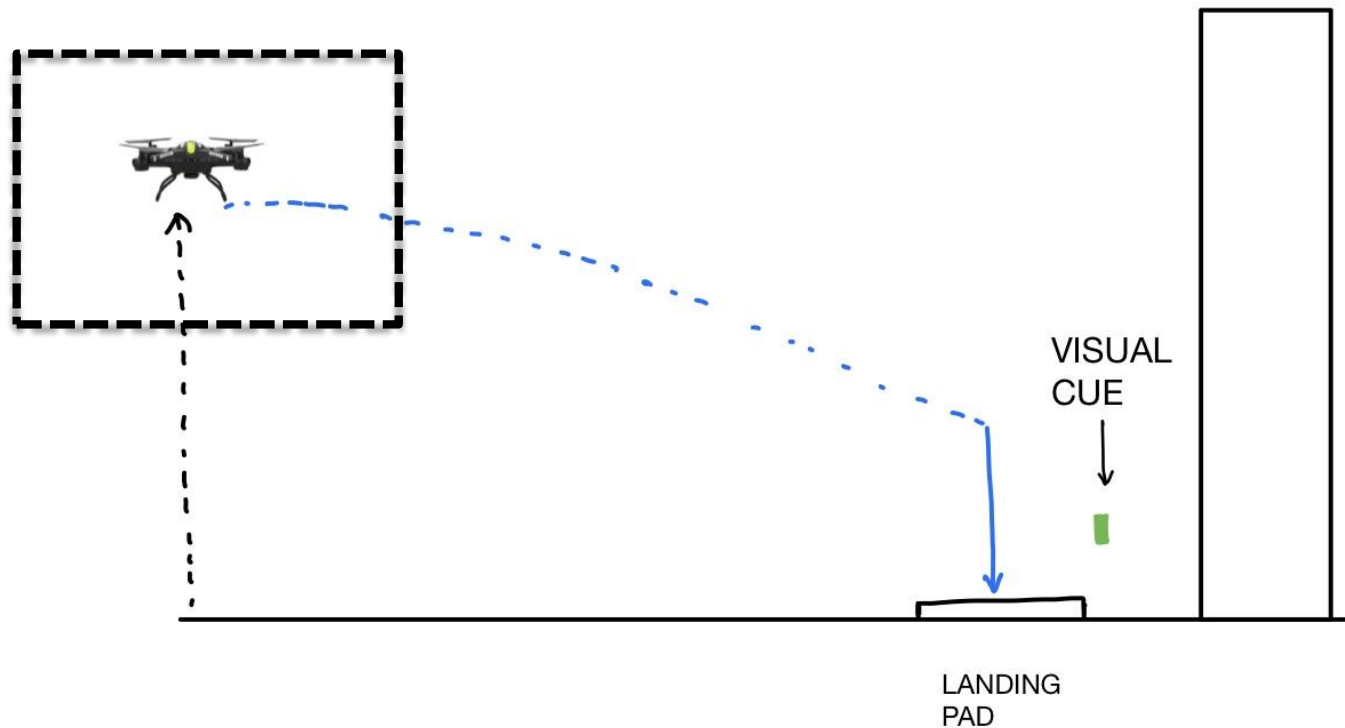
Bring pilot's intuition and flying skills using imitation learning (GAIL)

Environment simulator

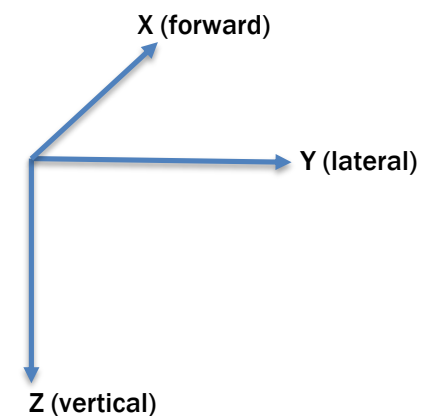
- Need a high-fidelity simulator environment for Unmanned Aerial Vehicles (UAVs)
- **Microsoft AirSim 2.0**
 - “An open source, cross platform simulator built on Unreal Engine”
 - Can integrate a flight controller for collecting demonstrations
 - Community support (NeurIPS 2019)
- Designed a custom ship deck
 - Landing pad, visual cue
 - Drone from AirSim



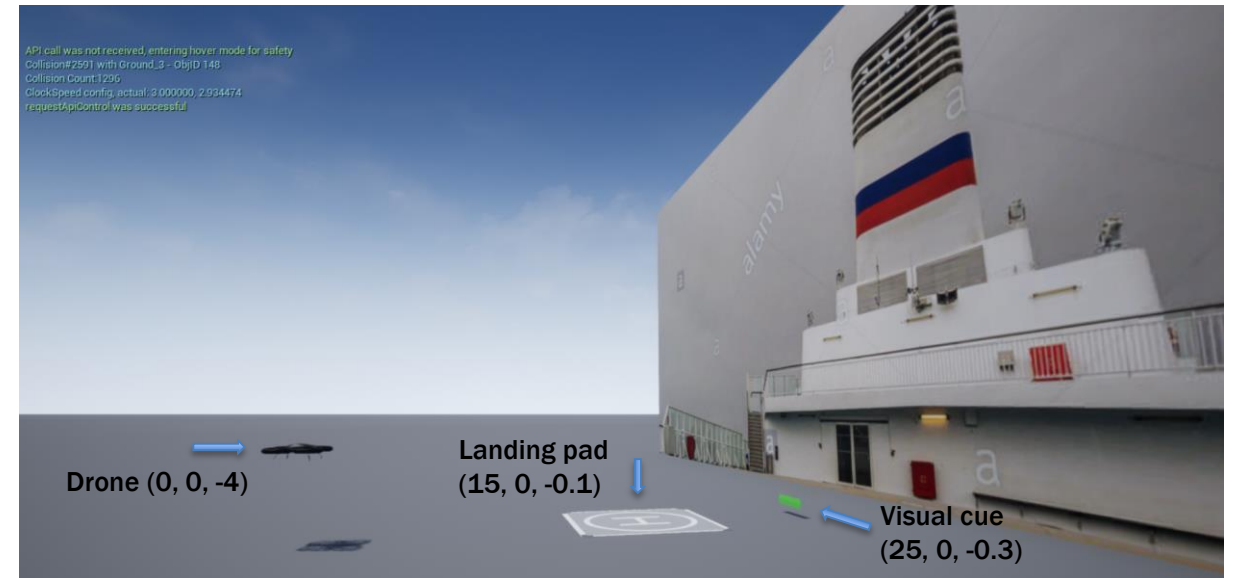
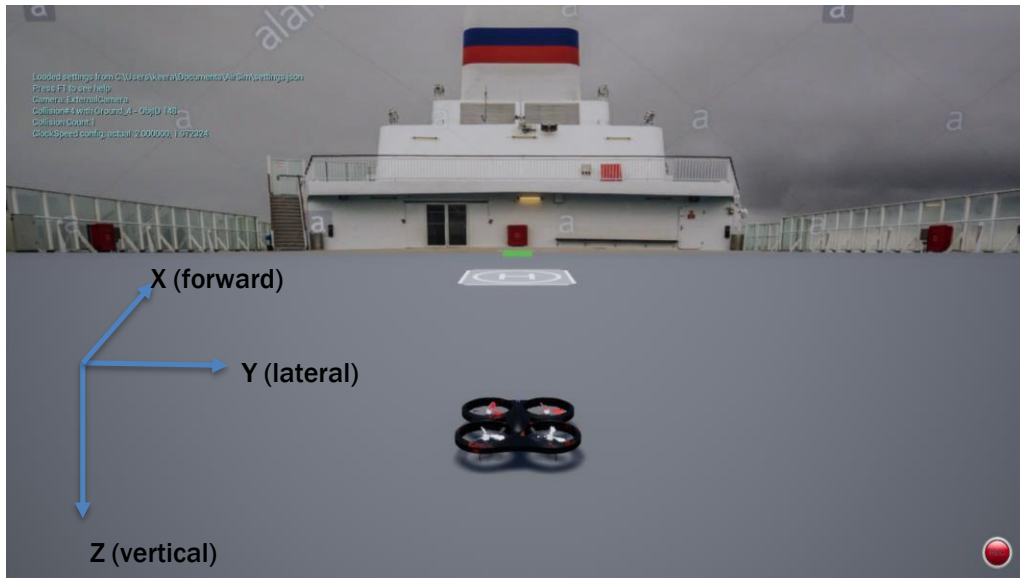
Model Concept



Component & Dimensions	Distance with respect to origin (world coordinate system)
Drone (1m x 1m)	At origin
Landing Pad (Centre) (4m X 4m)	X= 15 m, Y= 0, Z= 0
Visual Cue (Centre) (1m X 0.2m)	X= 25 m, Y= 0, Z= 0.3 m



AirSim environment: Front & Side view



The AirSim-v0 environment

Parameters	Details
State space (cts, dim = 6)	Drone position, velocity (X, Y, Z). Goal: 4x4 square around [15, 0, -0.1] Position: X [0, 17], Y [-2, 2], Z [-5, 0] – negative Z upwards Velocity: X [-1, 3], Y [-1, 1], Z [-4, 4]
Action space (cts, dim = 3)	[Pitch (rad), Roll (rad), Throttle (0, 1)]. Yaw zero. Negative pitch down
Termination / Horizon	Timeout (finite/infinite), out of bounds, below visual cue, crash, land

- Want to be able to **classify expert demos** as optimal/suboptimal. Assign a simple proxy reward
- Higher reward for getting closer to landing pad, penalty for termination without reaching goal

Generating human expert data

- Xbox controller:
 - Extremely sensitive
 - Cannot make custom calibration
- “Taranis x9d” flight controller:
 - Smoother data logging
 - Disabled yaw from the controller
- **Collected 120 demonstrations of landing UAV**
 - Started at random positions inside the box
 - Maneuver: different heights and at varying speeds
 - Collected (state, action, reward) pairs using AirSim APIs



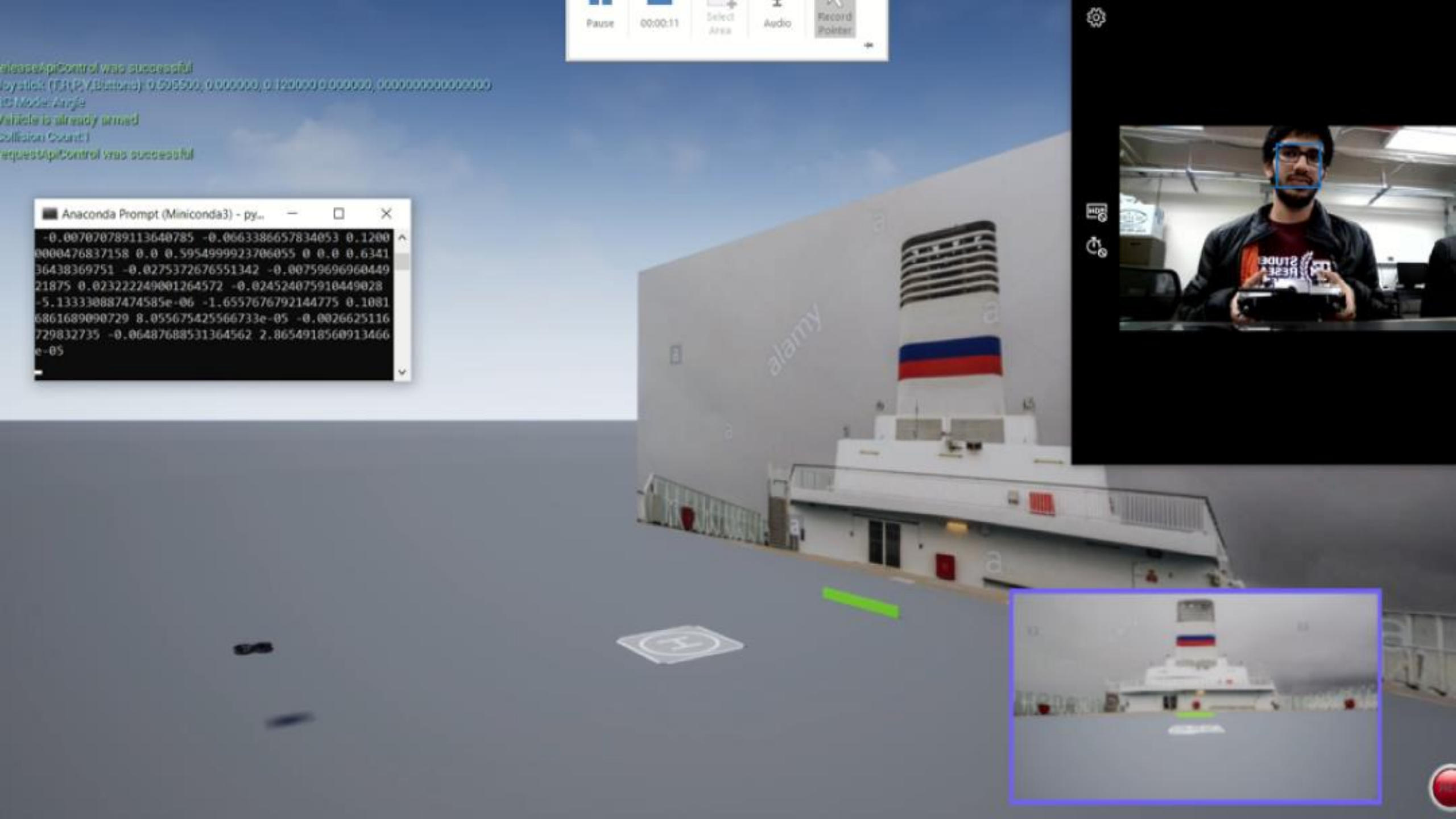
releaseApiControl was successful
Joystick (E31P/Bimond): 0.596516, 0.000000, 0.120000, 0.000000, 0.000000, 0.000000
3D Mode: Angle
Vehicle is already armed
Collision Count: 1
requestApiControl was successful

```
Anaconda Prompt (Miniconda3) - py...  
-0.007070789113640785 -0.0663386657834053 0.1200  
0000476837158 0.0 0.5954999923706055 0.0 0.6341  
36438369751 -0.0275372676551342 -0.00759696960449  
21875 0.023222249001264572 -0.024524075910449028  
-5.133330887474585e-06 -1.6557676792144775 0.1081  
6861689090729 8.055675425566733e-05 -0.0026625116  
729832735 -0.06487688531364562 2.8654918560913466  
e-05
```

Pause 00:00:11 Select Area Audio Record Pointer



HUD icons: a square with 'HUD', a refresh icon, and a camera icon.

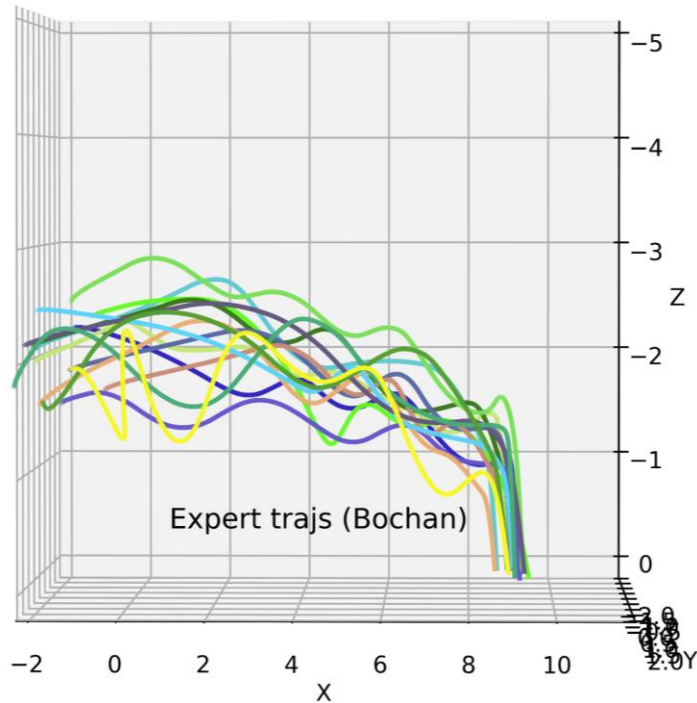


Human expert demos – stats

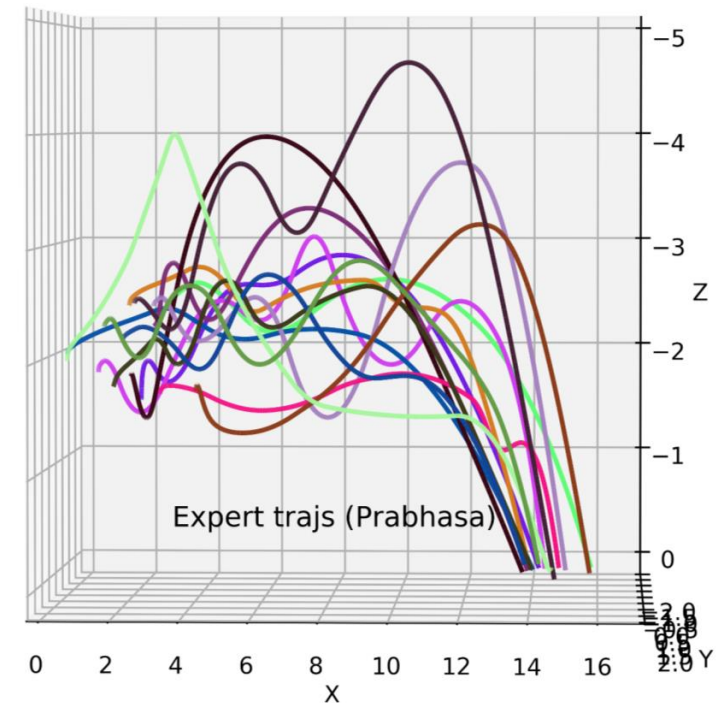
Expert	Optimal	Solved	Expert score (mean, std)	Expert length
Optimal	Yes	120/120	(1141, 27)	362
Suboptimal	No	132/140	(1116, 284)	307

- True reward (for proxy function): 1000
- Task: Train GAIL on expert samples [5, 10, 20, 50, 120] to learn behavior (optimal/suboptimal)
- Video of expert data collection: <https://youtu.be/e1no0lhzhQ4>

Expert Trajectories: humans



Optimal



Suboptimal

GAIL on suboptimal human expert

[20, 50, 120] experts, **finite horizon (400 steps, same as expert)**

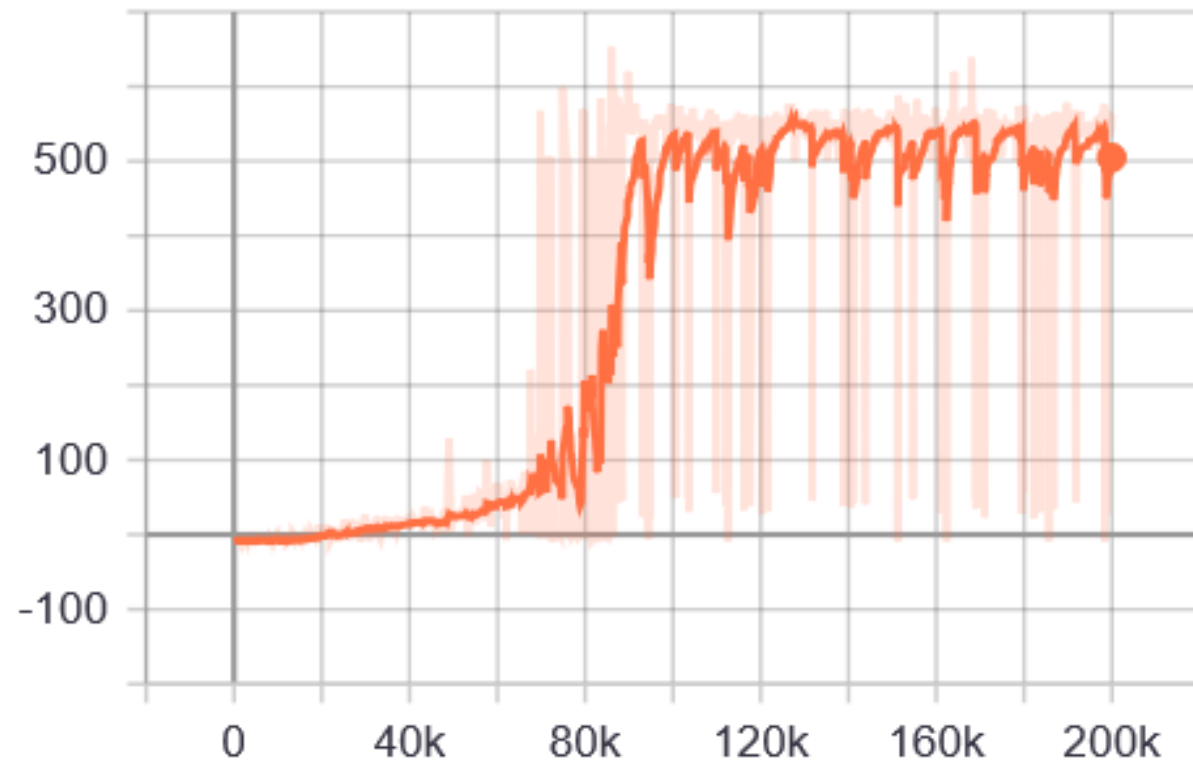
Learned model: <https://youtu.be/IUDpZna4uhk>

AirSim-v0 (IL) Suboptimal expert

GAIL hyperparameters

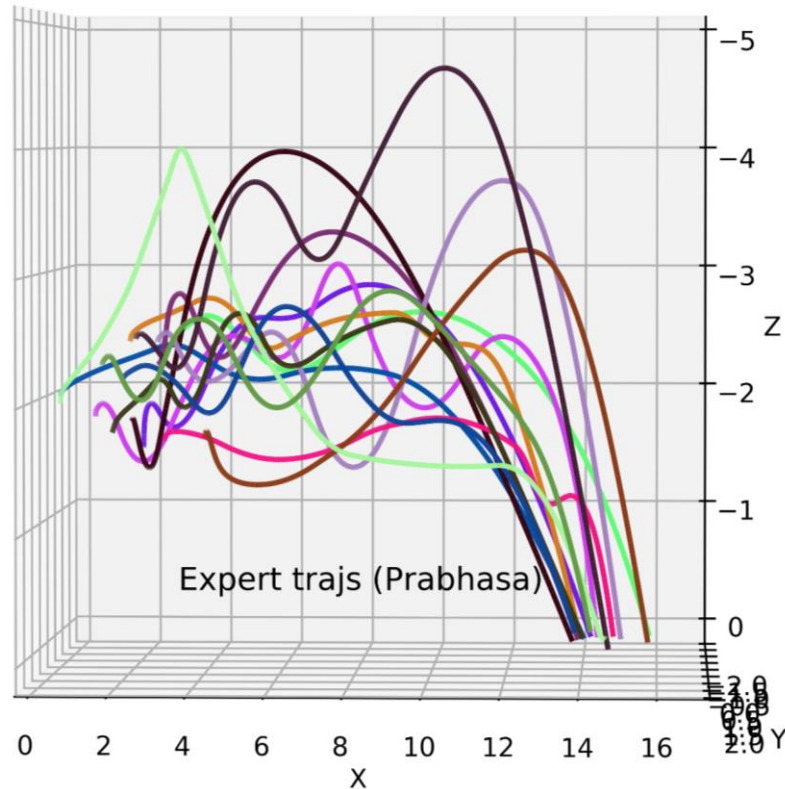
- n_timesteps: 2e5
- policy: 'MlpPolicy'
[128, 128]
- gamma: 0.99
- learning_rate: 3e-4
- timesteps_per_batch: 256
- buffer_size: 1e6

episode_reward

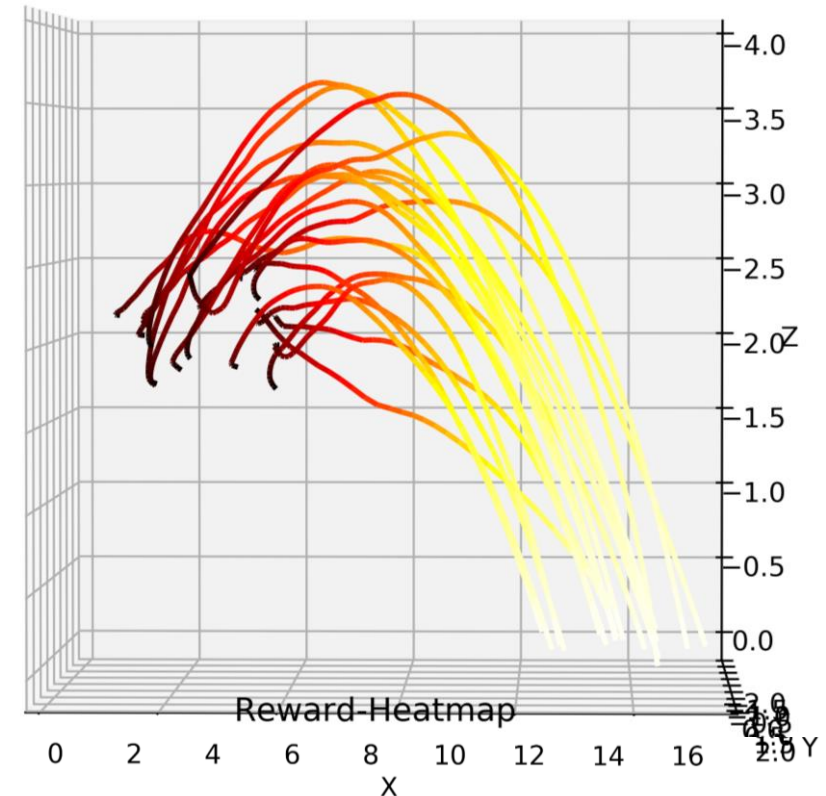


reward convergence

GAIL on suboptimal human expert



Suboptimal expert



GAIL-learned policy

GAIL on optimal human expert

[20, 50, 120] experts, **finite horizon (400 steps, same as expert)**

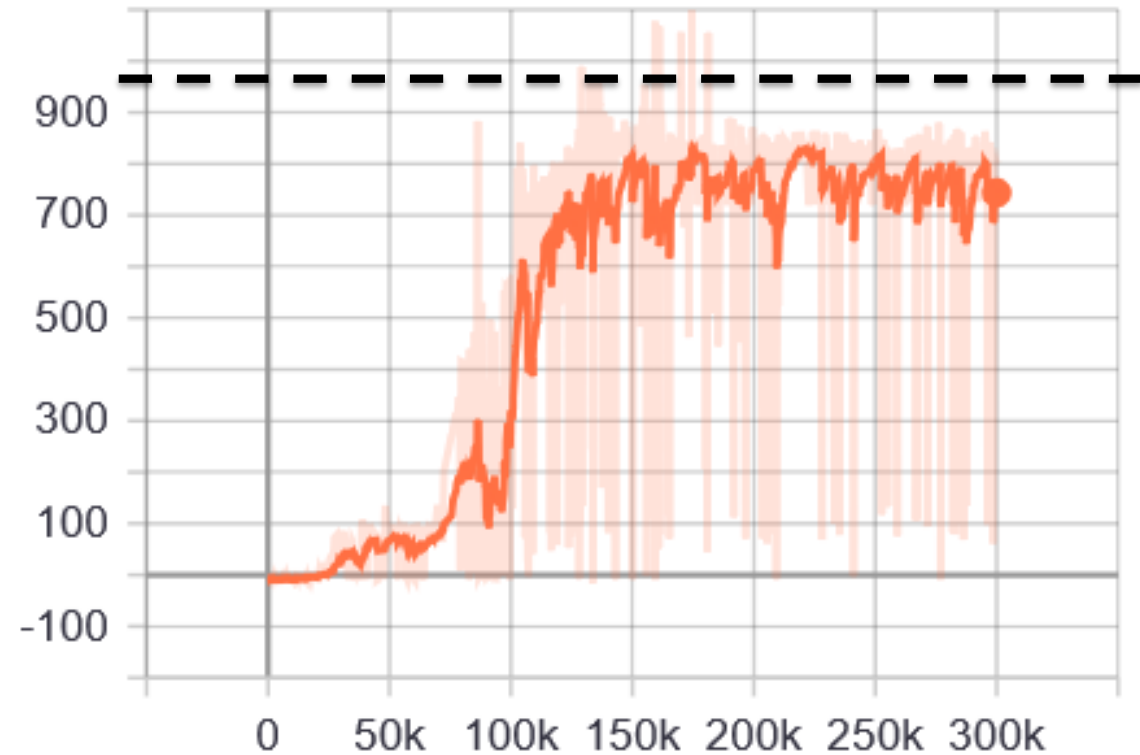
Learned model: <https://youtu.be/3ilW7Lzql2Y>

AirSim-v0 (IL) Optimal expert

GAIL hyperparameters

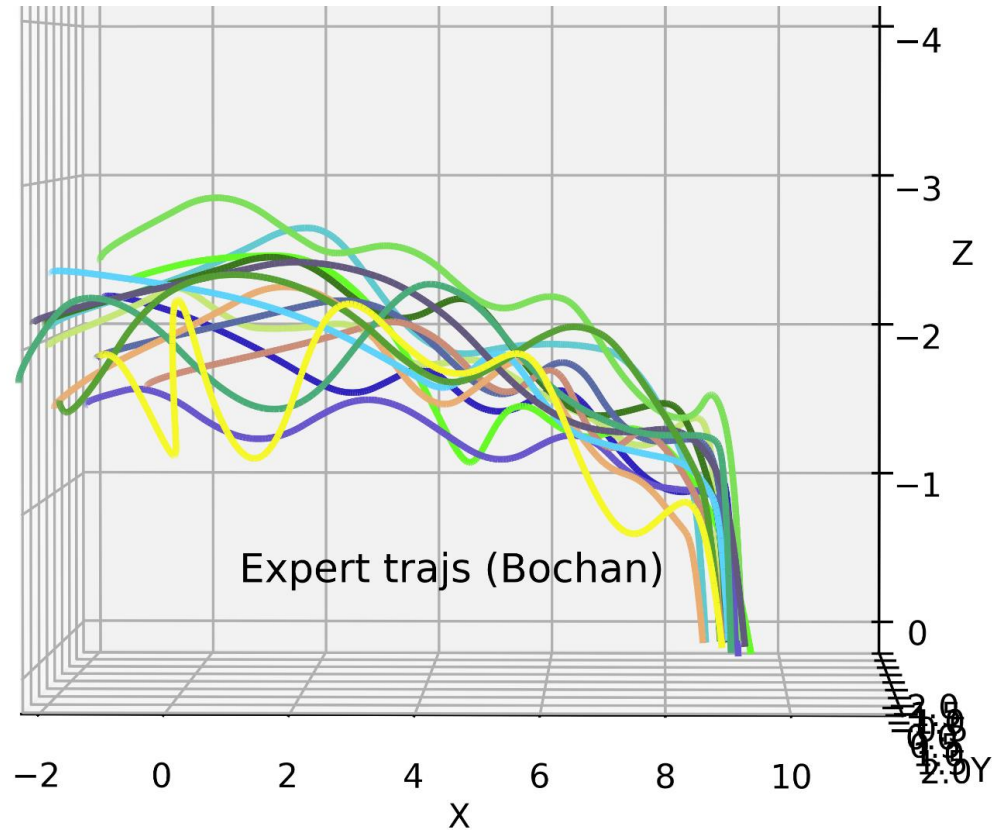
- n_timesteps: 3e5
- policy: 'MlpPolicy'
[128, 128]
- gamma: 0.99
- learning_rate: 3e-4
- timesteps_per_batch: 256
- buffer_size: 1e6

episode_reward

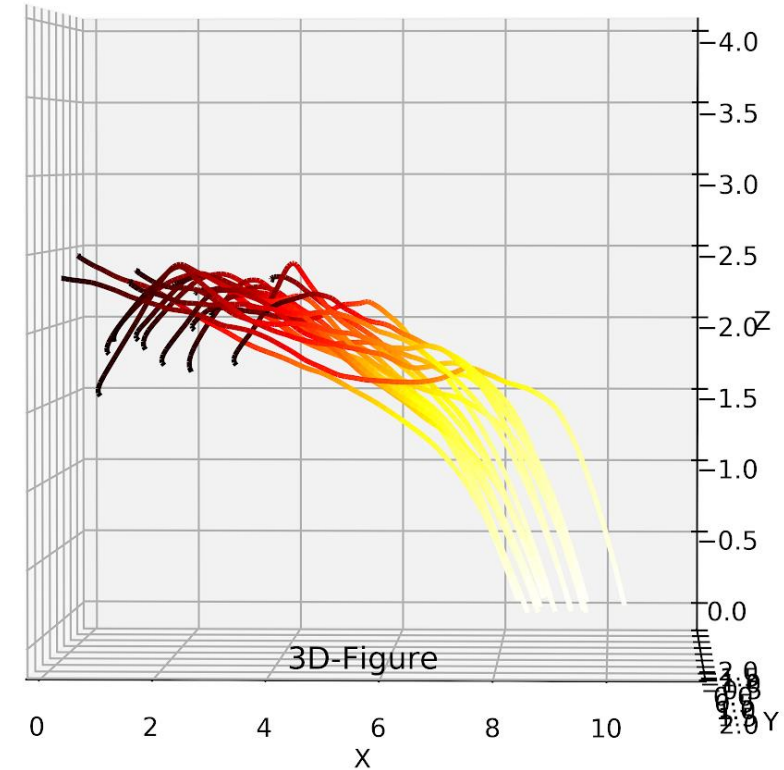


reward convergence

GAIL on optimal human expert



Optimal expert



GAIL-learned policy



Conclusions

- GAIL can imitate **navigation** (point A to point B)
- GAIL can learn suboptimal **landings**. HP-dependent
- Explanation: landings may be too 'non-smooth' for GAIL to learn
- Rendering of learned policy: [front-view](#)

- Can we perhaps **construct a proxy reward** that conveys smoother landings?
- Can RL algorithms learn smoother landings from this proxy?

EXPERT REWARD DESIGN

Can RL learn a 'smoother' landing than human expert?

Learned models:

Simple - <https://youtu.be/qJJ00OWfYcl>

Complex - <https://youtu.be/cFpFTDo-V7k>

Proxy reward function design

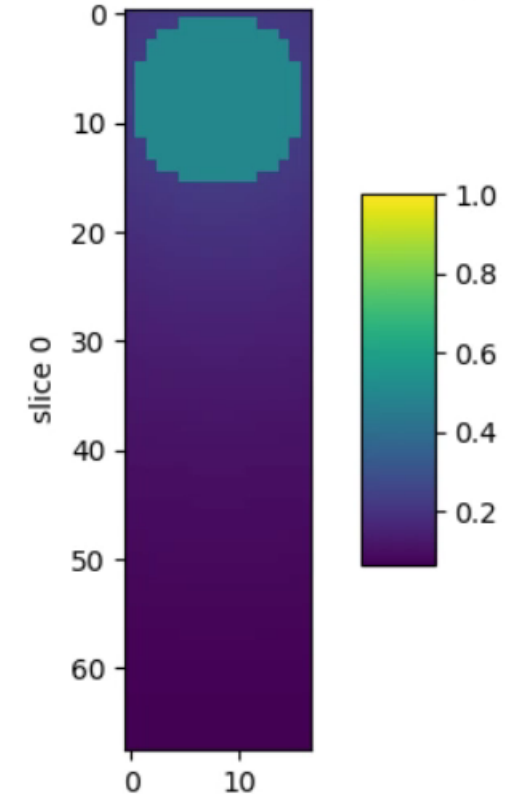
1. Simple reward (sparse):

- Increase reward as it gets close to landing pad ($1/x$)
- Large positive reward if it lands inside the landing pad (+1000)
- Other conditions: -10 (visual cue, out of bounds, crash, timeout)

2. Complex reward (sparse):

- Increase reward as it gets close to landing pad ($1/x$)
- **Scale** goal reward according to drone heading, speed (1250-750)
- Other conditions: -10 (visual cue, out of bounds, crash, timeout)

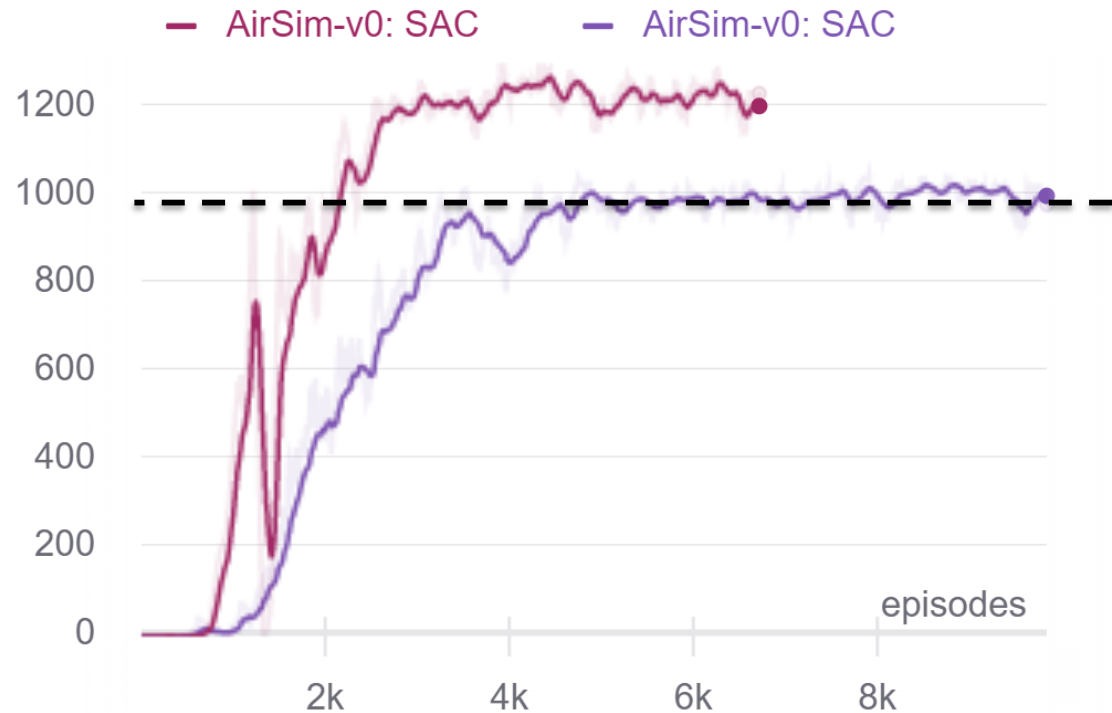
use scroll wheel to navigate images



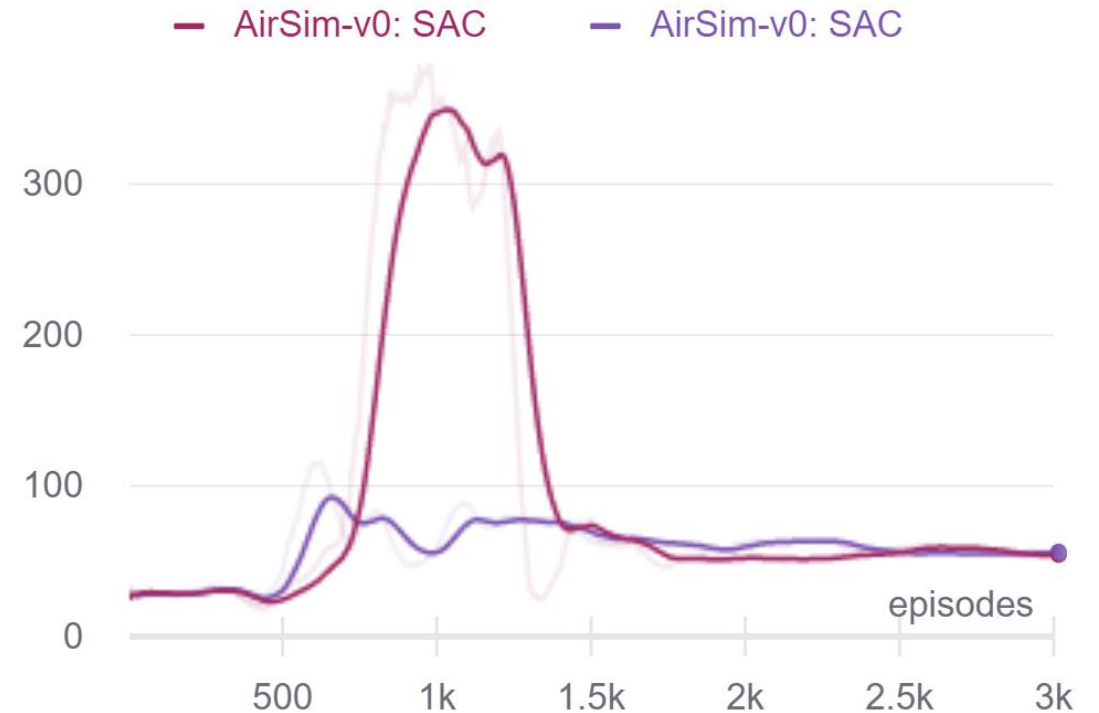
Reward Heatmap

SAC on AirSim-v0: Proxy rewards

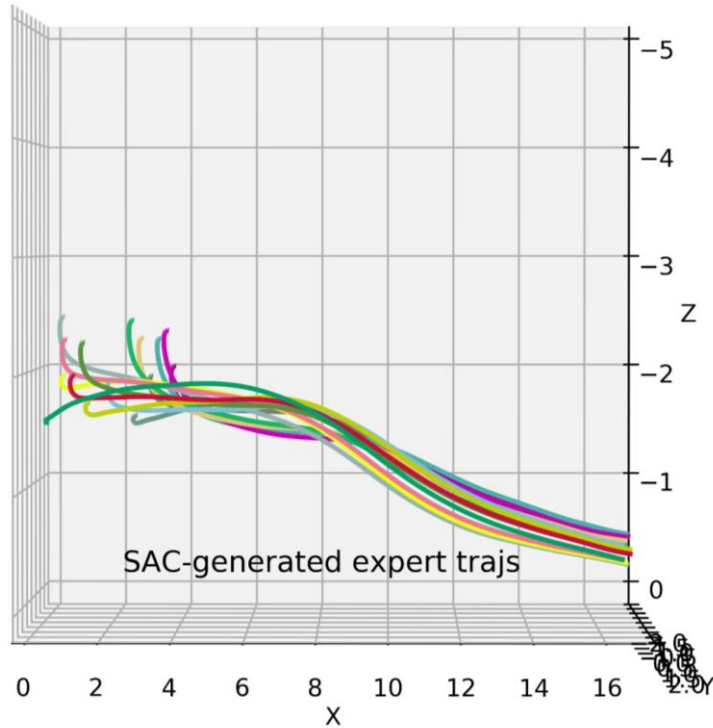
mean episode reward (RL)



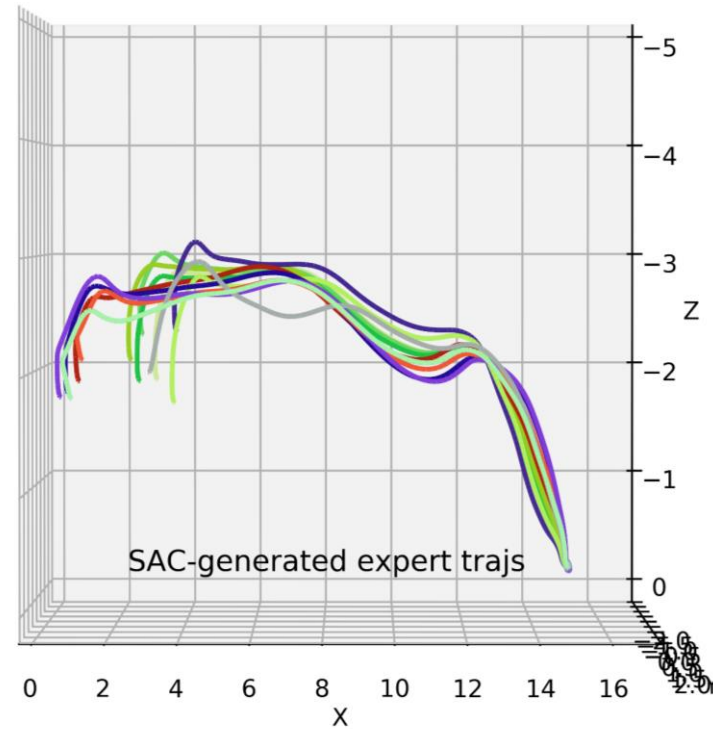
mean episode length



RL-generated expert demos

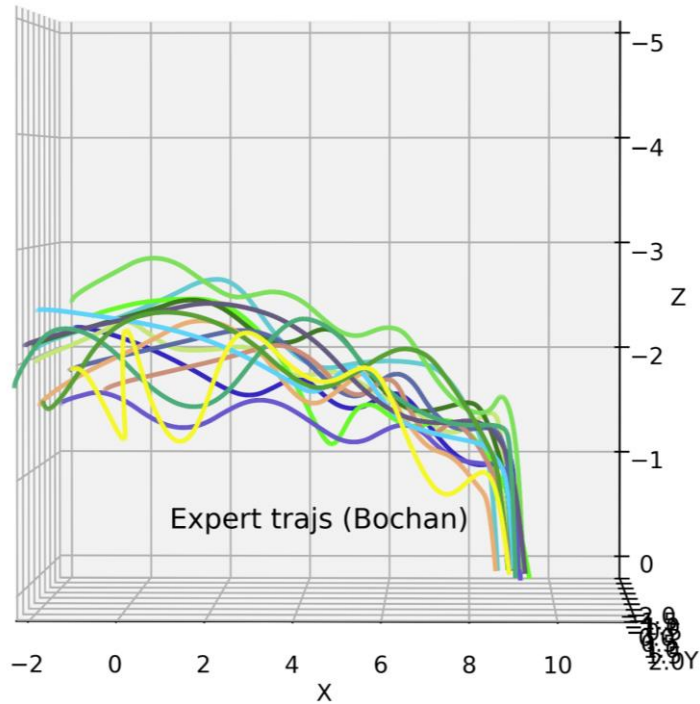


Proxy: simple

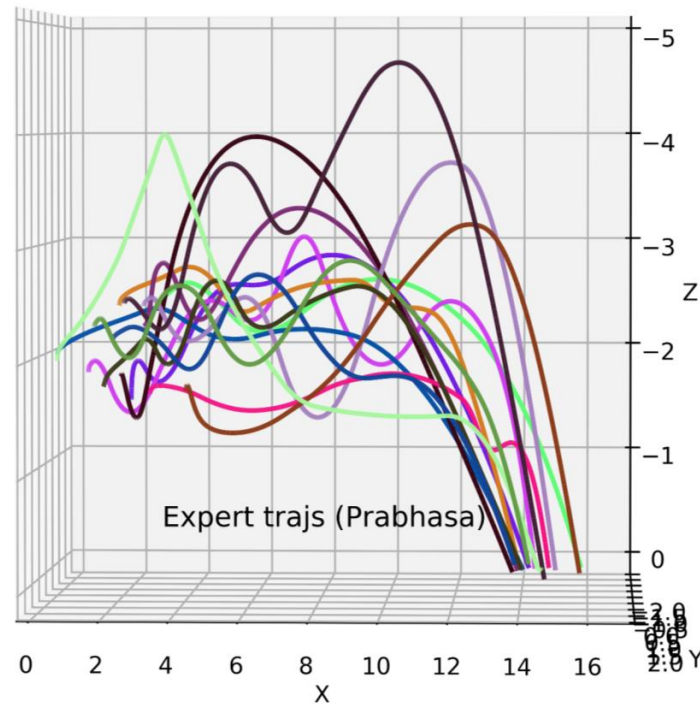


Proxy: complex

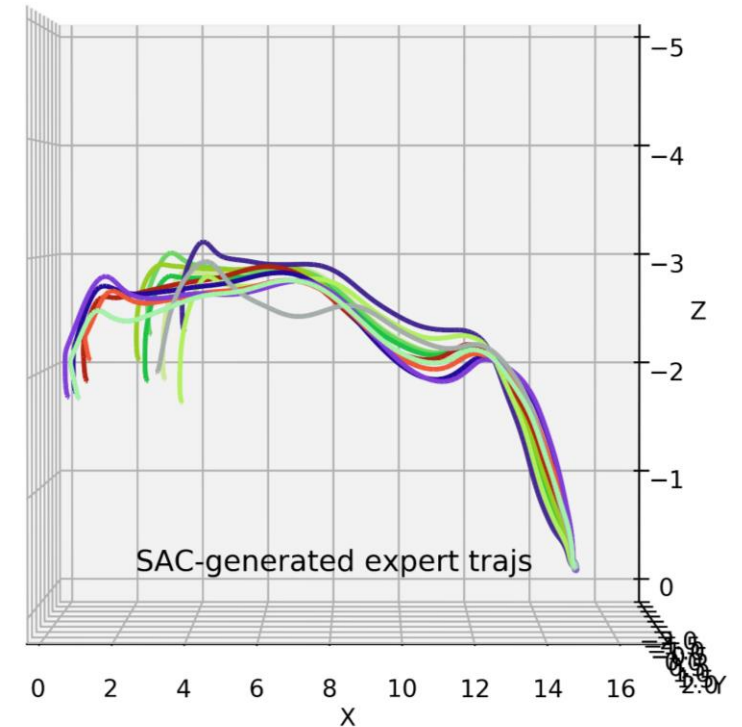
Expert demos: humans vs RL



Optimal



Suboptimal

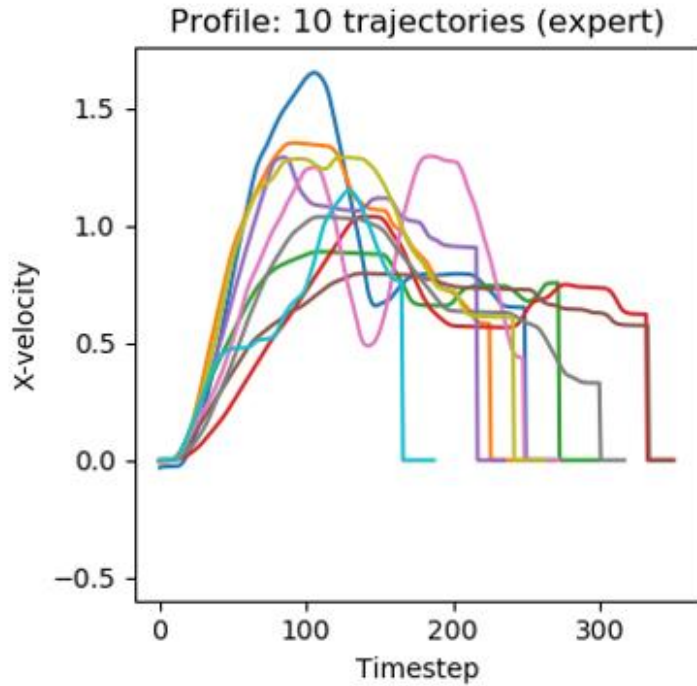


Proxy: complex

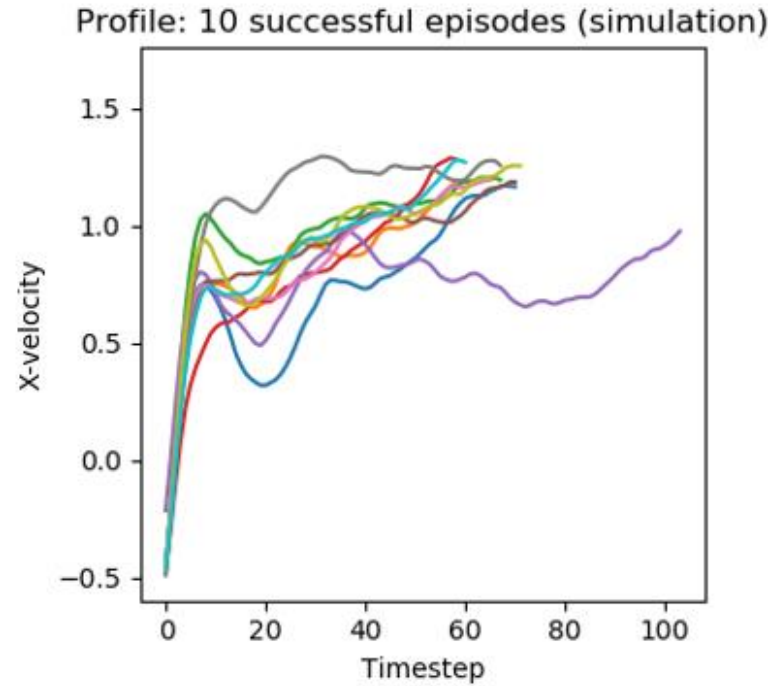
LANDING: HUMANS vs GAIL vs RL

Change in speed

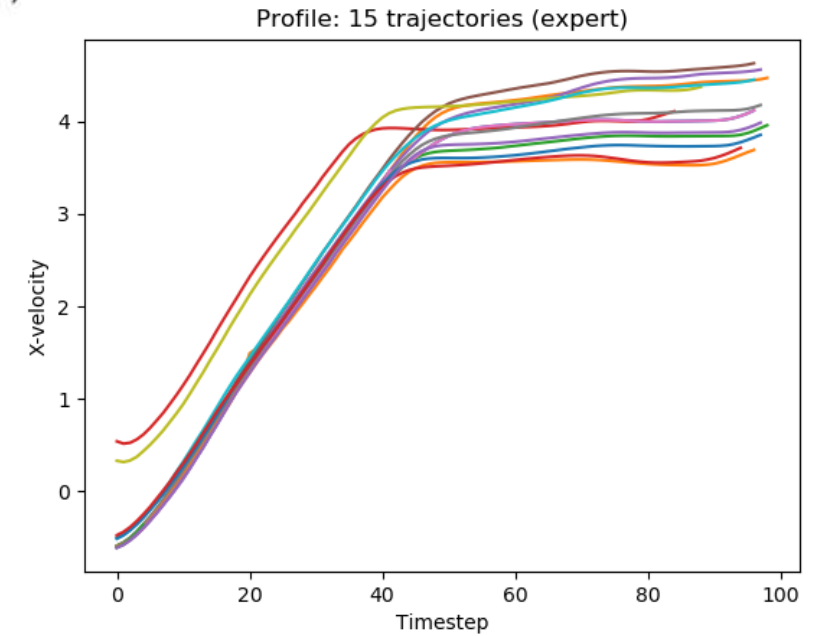
Forward speed



Human expert

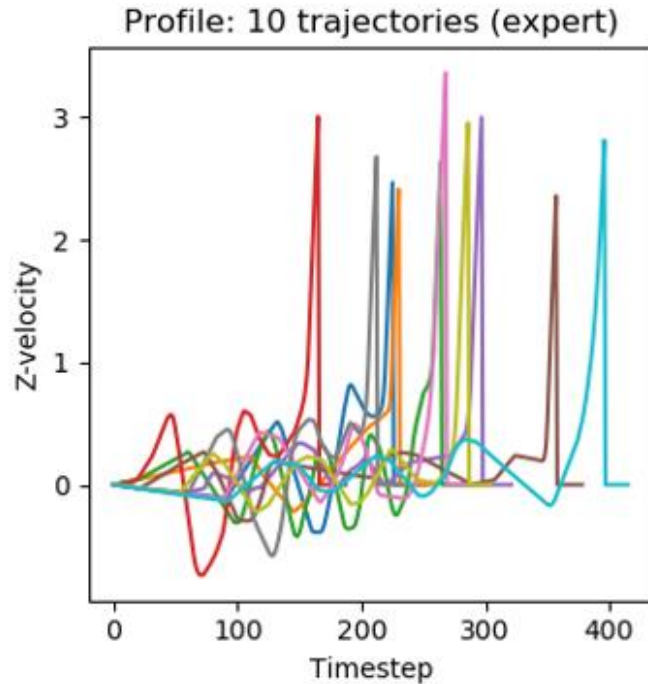


GAIL-learned policy

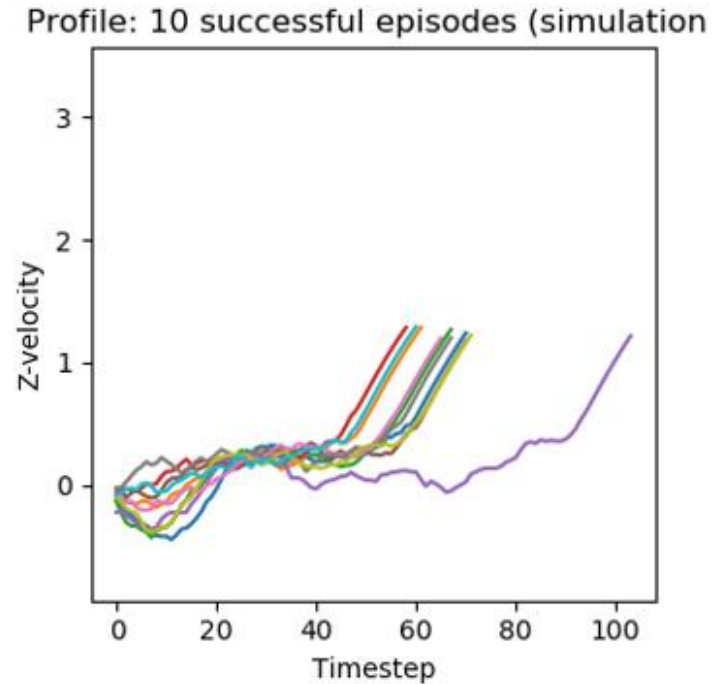


RL expert

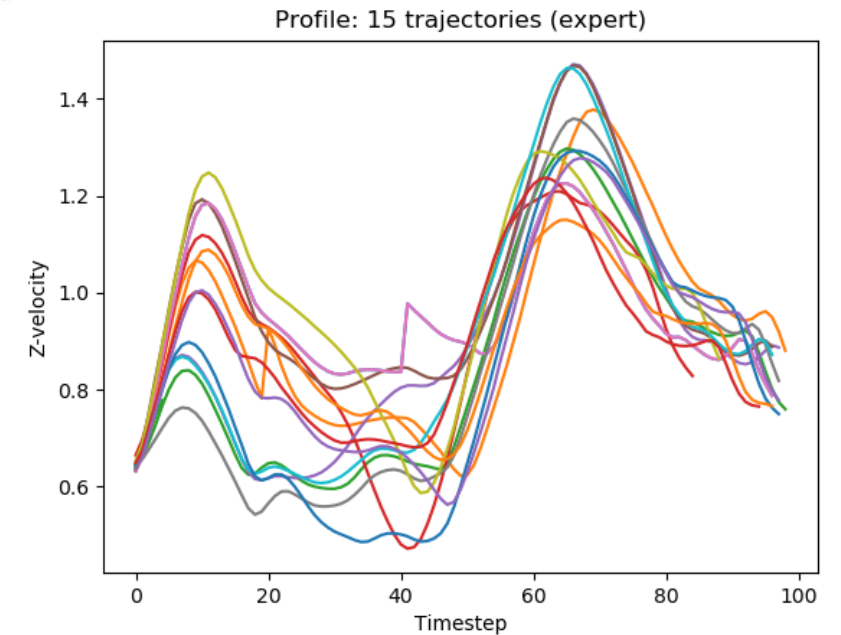
Descent speed



Human expert



GAIL-learned policy



RL expert

Summary: humans vs GAIL vs RL

Expert	Maneuver type demonstrated/learned	Optimal	Landed	Score (mean, std)	Expert length
Human: Optimal	Hard landing (pilot)	Yes	120/120	(1141, 27)	362
Human: Suboptimal	Large variance	No	132/140	(1116, 284)	307
GAIL: optimal (20 demos)	Navigation, *Landing	*Yes	84/100	(1048, 292)	80
GAIL: suboptimal (20 demos)	Navigation only	*No	12/100	(684, 580)	80
SAC: Simple	Shortest-path	Yes	99/100	(1106, 112)	99
SAC: Complex	Smooth landing	Yes	97/100	(1265, 225)	94

*Smooth landings crucial for perfect imitation with GAIL



Some questions...

1. How does imitation accuracy scale with dimensionality, demo data? **Sample-efficient**
2. How 'smooth' are the learned policies compared to the expert policy? **smooth if expert is smooth**
3. Can sparse rewards be learned? At what cost? **Yes, needs >20 demos, tuned HPs**
4. Can GAIL imitate suboptimal experts? **navigation easy, landing difficult. Tuned HPs**
5. Can GAIL generalize? Tried different bounding boxes for optimal expert, GAIL policy. Need tuned HPs

GAIL: Pros and Cons

- Pros
 - Can handle unknown dynamics
 - Can scale to large neural network reward functions
 - Can perform well on real-world tasks (**with an efficient policy optimizer**)
- Cons
 - **Adversarial optimization (GANs) hard to train!**
 - **Requires smooth experts for imitation**
 - First person demonstrations typically used (no “teaching” as such)

Chelsea Finn, RL Bootcamp, 2016

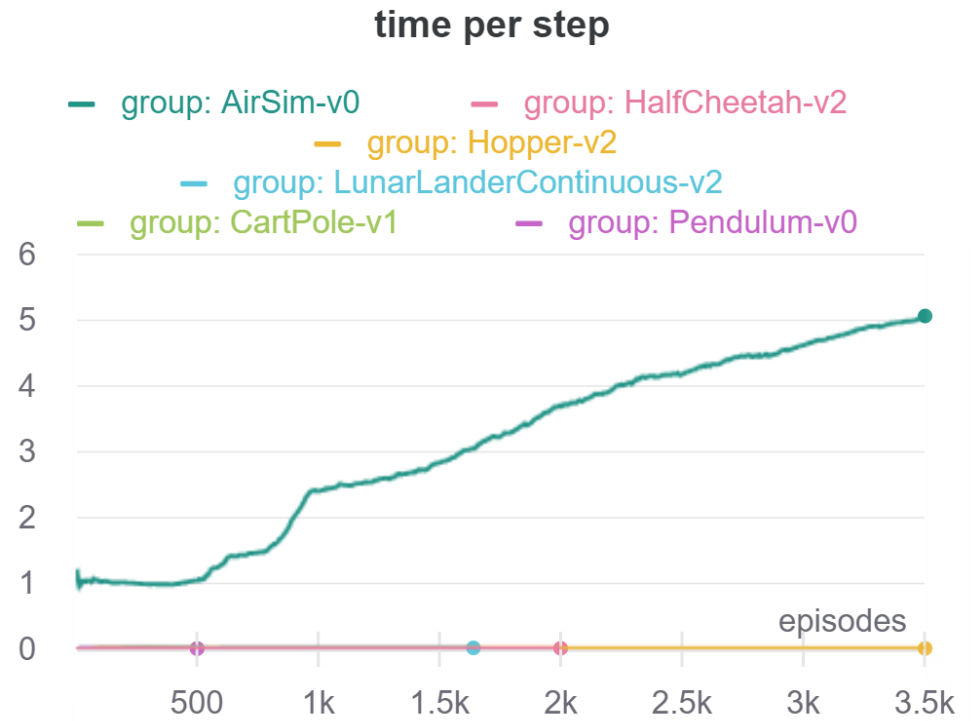
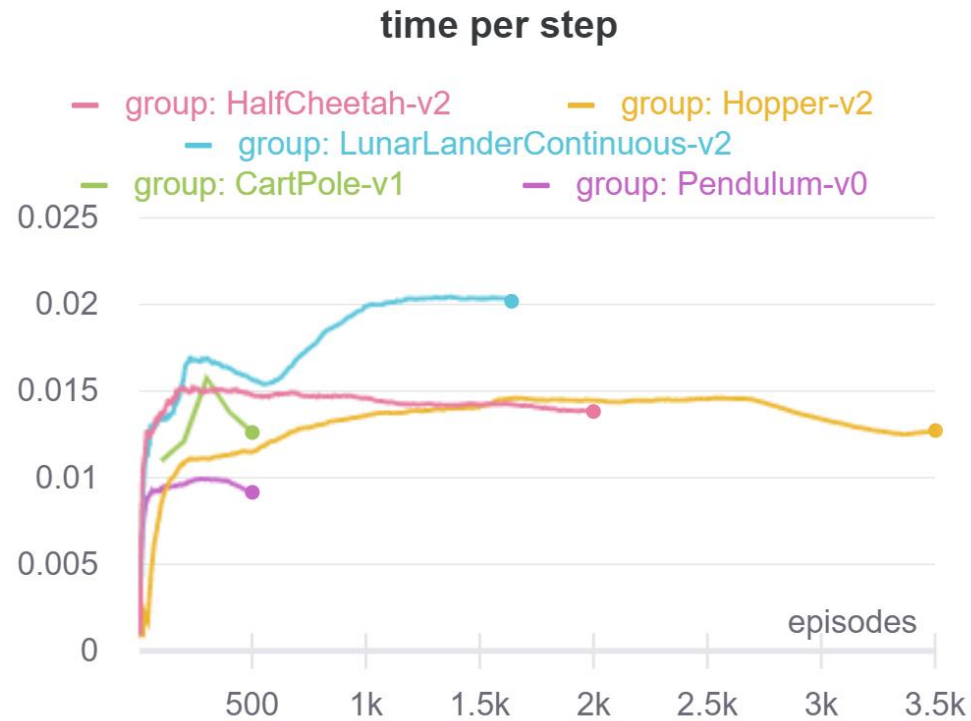
Summary

- **GAIL can imitate** AirSim-v0 human experts. Navigation easy, landing not-so-easy
- RL on proxy rewards can generate **smoother landings** – necessary for GAIL
 - reward ϵ space of cost functions explored
- **Expensive training time** limits number of experiments you can run
- **Lack of tuned HPs** affects imitation accuracy



ADDRESSING TIME BOTTLENECKS

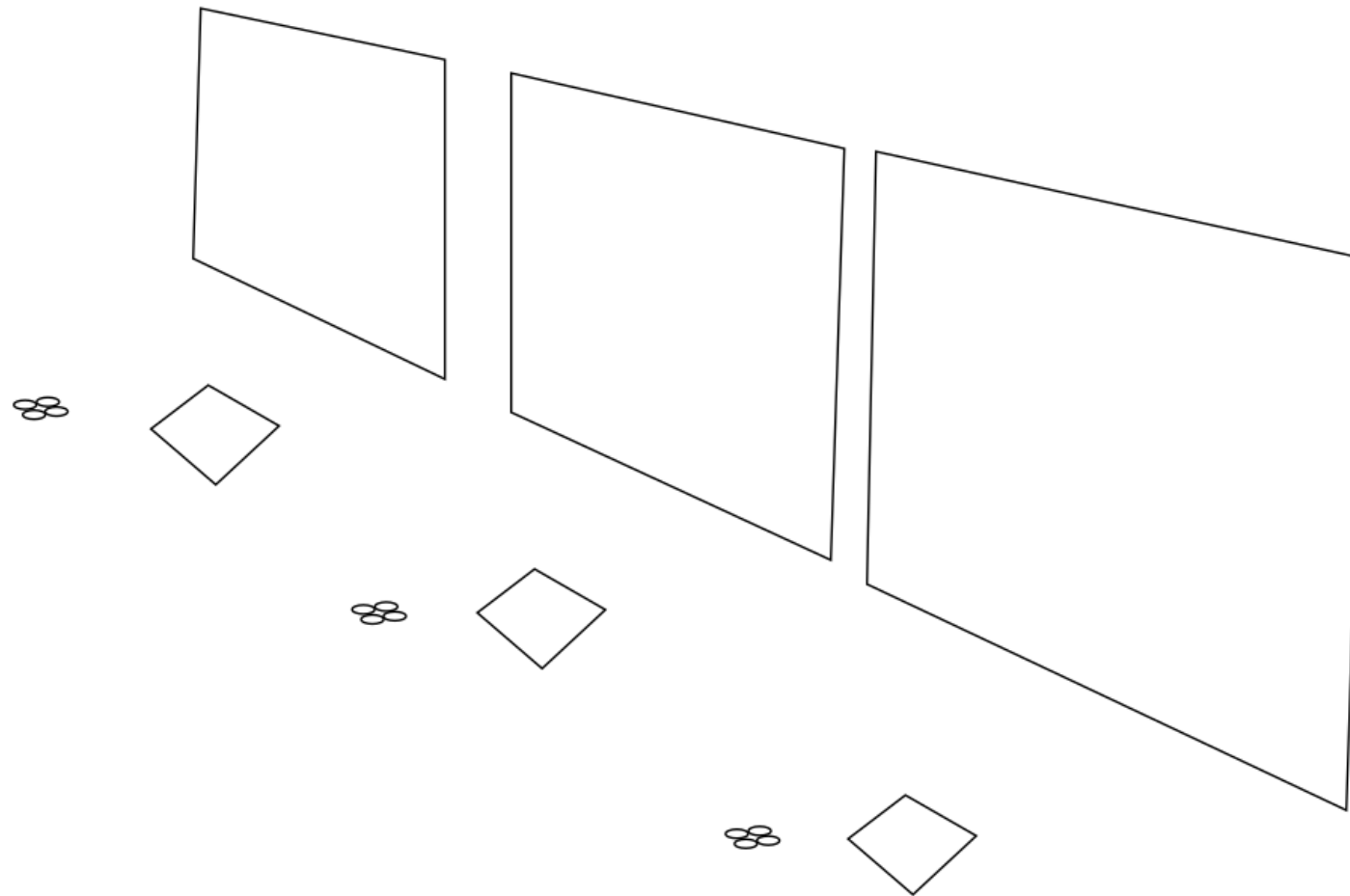
Environment samples are expensive!



Some numbers to crunch on...

Property	CartPole-v1	Hopper-v2	AirSim-v0
Dimension (state, action)	(4, 2)	(11, 3)	(6, 3)
Timesteps (GAIL)	3e5	1e6	1e6
Episode length (max)	500	1000	400 (human), 100 (RL)
Training time	20 minutes	2 hours	36 hours (at 4x)
Time/env interaction	4 ms	7.2 ms	129.6 ms
Time/episode	2 s	7.2 s	17.1 s (human), 4.3 s (RL)
Clock speed	Processor (4.8GHz)	Processor (4.8GHz)	4 x real time
Cumulative mean reward	Converged	Converged	Did not converge

Setting up multiple experiments



API call was not received, entering hover mode for safety
Collision#250 with Ground_4 - ObjID: 148
requestApiControl was successful
Collision Count:30
ClockSpeed config, actual: 4.000000, 3.989199





Sections

1. Need for sample-efficiency
2. Introduction to Imitation Learning
3. Application: Autonomous UAV Landing
4. Conclusions and Future Work

CONCLUSIONS

- Need sample-efficient learning for complex, long-horizon tasks
- IL (GAIL) is a sample-efficient approach to learn from demonstrations
- IL can be used to imitate (even suboptimal) experts from sparsely-rewarded environments
 - Requires smooth experts and careful HP tuning for perfect imitation
- Application: Designed a novel method of autonomous UAV landing (simulation)

Future Extensions: AirSim

- Learning reward functions with smoothness properties (e.g. WAIL)
- Collecting human expert data with smoother maneuvers, for better imitation
- Complex maneuvers. E.g. side-entry (yaw), landing on a moving platform, wind
- [Switch to a quadcopter for learning \(Parrot Anafi with Gazebo\)](#)
- Multi-agent, transfer learning, and meta-learning methods to learn behaviors that can be generalized to unknown environments (e.g. point-to-point navigation)





TEXAS A&M UNIVERSITY
Engineering

THANK YOU!