

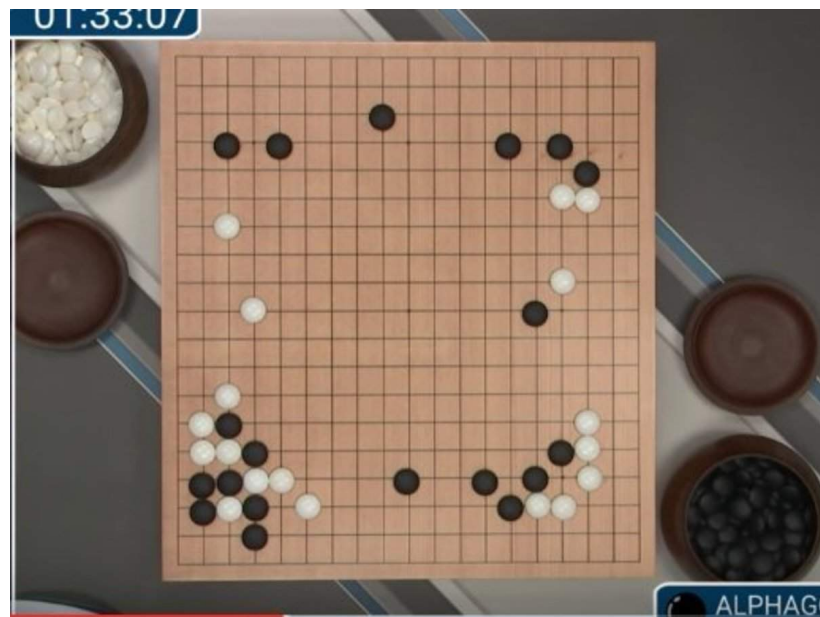


Deep Reinforcement Learning

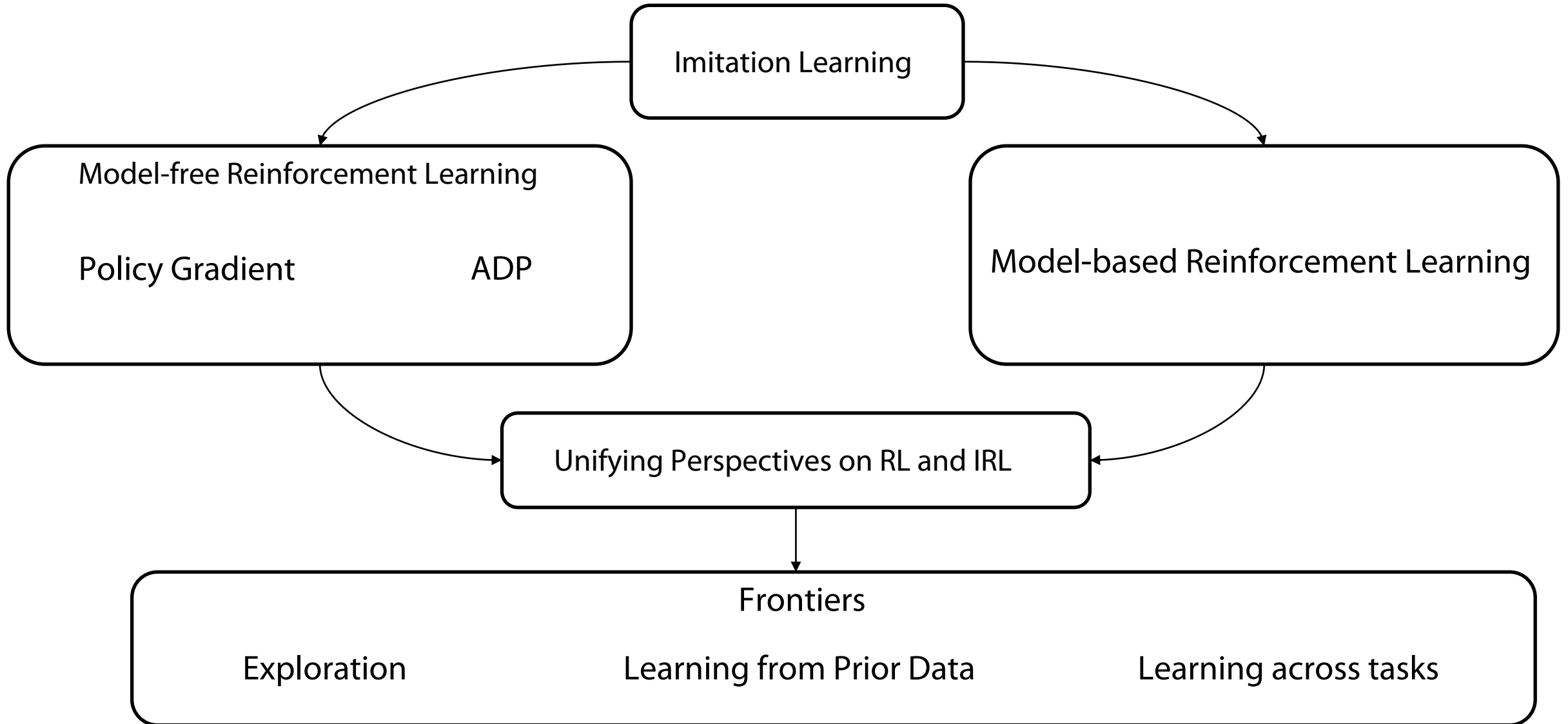
Spring 2026

Abhishek Gupta

TA: Mateo Guaman Castro



Class Structure



Lecture outline

Recap: Imitation Learning via Behavior Cloning



Challenges with DAgger



Modern Takes on Imitation Learning



Policy Gradient

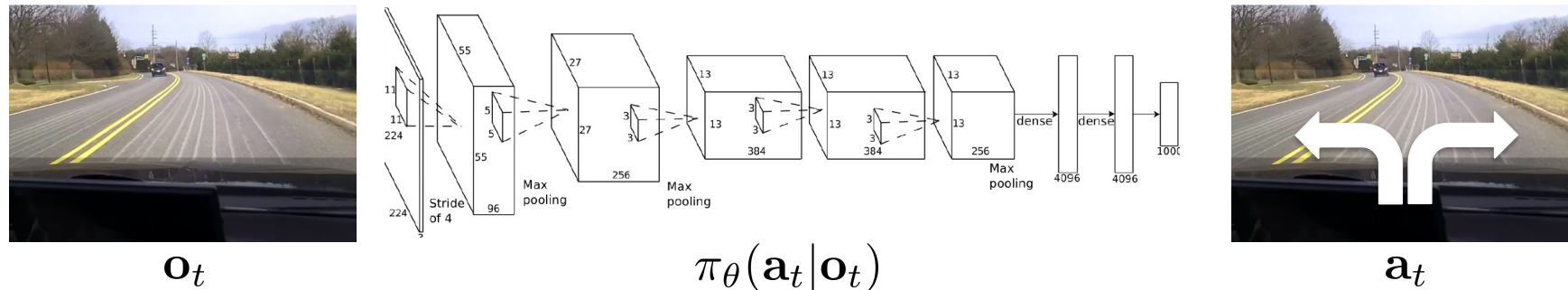
Idea 1: Imitation Learning via Supervised Learning

Given: Demonstrations of optimal behavior

$$\arg \max_{\theta} \mathbb{E}_{(s^*, a^*) \sim \mathcal{D}} [\log \pi_{\theta}(a^* | s^*)]$$

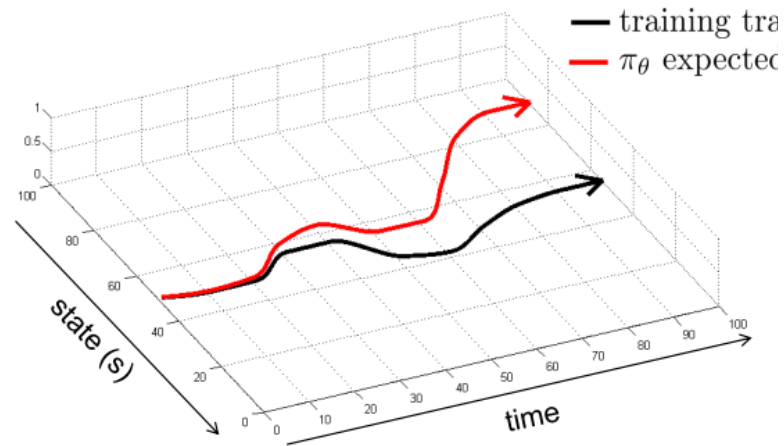
Goal: Train a policy to mimic the demonstrator

Idea: Treat imitation learning as a supervised learning problem! → Behavior Cloning



So does behavior cloning really work?

- Imitation Learning \neq Supervised Learning



Compounding error!

$$\arg \max_{\theta} \mathbb{E}_{(s^*, a^*) \sim \mathcal{D}} [\log \pi_{\theta}(a^* | s^*)]$$

$$\mathbb{E}_{(s, a) \sim \rho(\pi)} [\mathbf{1}(a = a^*)]$$

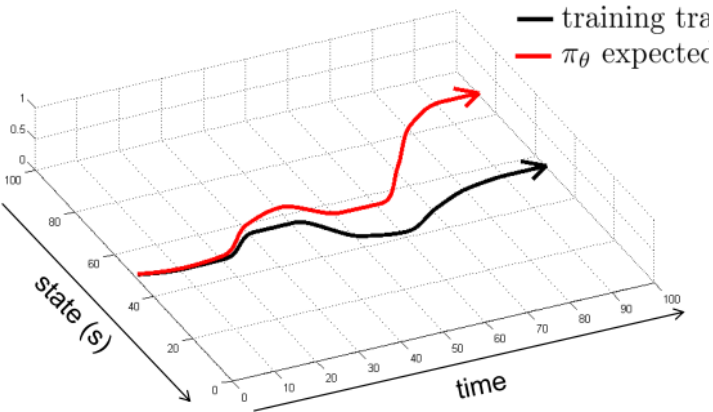


Not the same!

Let's try and understand where the problem lies?

Behavior cloning has challenges in both theory and practice

$$\sum_t \mathbb{E}_{(s_t, a_t) \sim p_{\pi_\theta}(s_t, a_t)} [c(s_t, a_t)] \leq O(\epsilon H^2)$$



Underfitting

$$\pi_\theta(a \neq \pi^*(s_t) | s_t) \leq \epsilon$$

Compounding error

$$\leq O(\epsilon H^2)$$

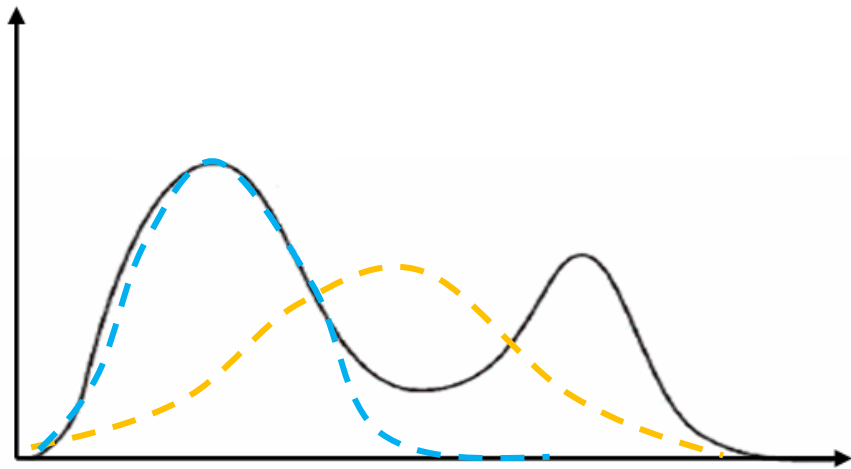
Effects of choice of f-divergence on behavior

Different divergences lead to different properties

$$\mathbb{E}_{s^* \sim p_{\pi_e}(\cdot)} [D_{\text{KL}}(\pi_e(\cdot|s^*) || \pi_\theta(\cdot|s^*))] \longrightarrow \mathbb{E}_{s^* \sim p_{\pi_e}(\cdot)} [D_f(\pi_e(\cdot|s^*), \pi_\theta(\cdot|s^*))]$$

Forward KL (behavior cloning)

More general class of divergences



$$D_f(p(x), q(x)) = \mathbb{E}_{q(x)} \left[f \left(\frac{p(x)}{q(x)} \right) \right]$$

- Forward KL (mode covering) $f(x) = x \log(x)$
- Reverse KL (mode seeking) $f(x) = -\log(x)$

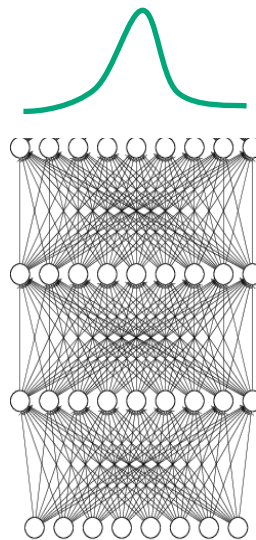
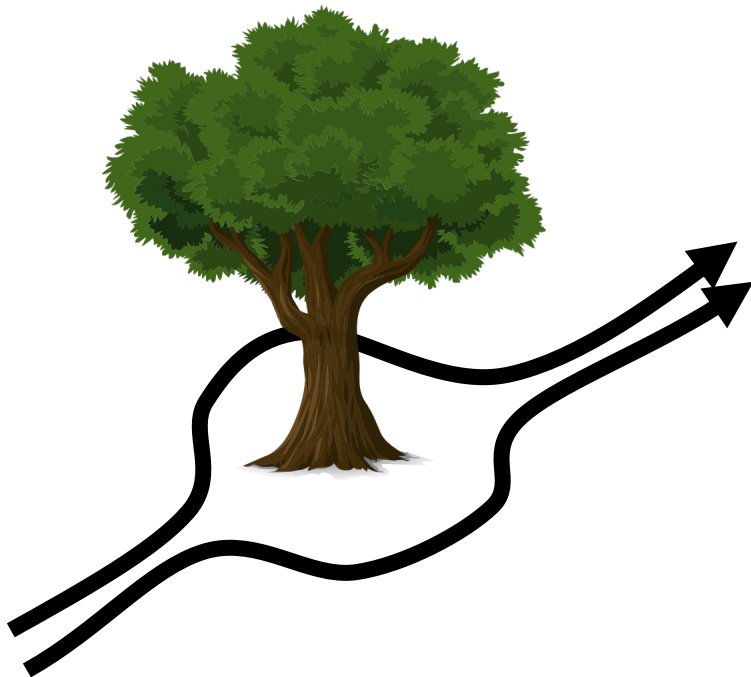
So how do we fix BC?

Use a different f-divergence!
(Change f)

or Use a richer distribution class!
(Change π_θ)

Using Richer Policy Distribution Classes

Multimodal behavior → use more **expressive** probability distributions, no mode averaging issues



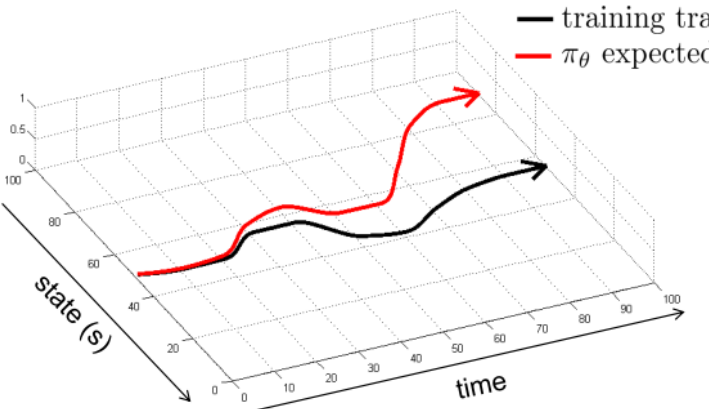
1. Output mixture of Gaussians
2. Latent variable models
3. Autoregressive discretization
4. Diffusion models
5. ...



Let's try and understand where the problem lies?

Behavior cloning has challenges in both theory and practice

$$\sum_t \mathbb{E}_{(s_t, a_t) \sim p_{\pi_\theta}(s_t, a_t)} [c(s_t, a_t)] \leq O(\epsilon H^2)$$



Underfitting

$$\pi_\theta(a \neq \pi^*(s_t) | s_t) \leq \epsilon$$

Compounding error

$$\leq O(\epsilon H^2)$$

Concrete Instantiation: DAgger

can we make $p_{\text{data}}(\mathbf{o}_t) = p_{\pi_\theta}(\mathbf{o}_t)$?

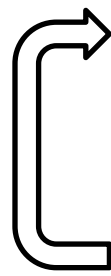
idea: instead of being clever about $p_{\pi_\theta}(\mathbf{o}_t)$, be clever about $p_{\text{data}}(\mathbf{o}_t)$!

DAgger: Dataset Aggregation

goal: collect training data from $p_{\pi_\theta}(\mathbf{o}_t)$ instead of $p_{\text{data}}(\mathbf{o}_t)$

how? just run $\pi_\theta(\mathbf{a}_t|\mathbf{o}_t)$

but need labels \mathbf{a}_t !

- 
1. train $\pi_\theta(\mathbf{a}_t|\mathbf{o}_t)$ from human data $\mathcal{D} = \{\mathbf{o}_1, \mathbf{a}_1, \dots, \mathbf{o}_N, \mathbf{a}_N\}$
 2. run $\pi_\theta(\mathbf{a}_t|\mathbf{o}_t)$ to get dataset $\mathcal{D}_\pi = \{\mathbf{o}_1, \dots, \mathbf{o}_M\}$
 3. Ask human to label \mathcal{D}_π with actions \mathbf{a}_t
 4. Aggregate: $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_\pi$

DAgger Example



Lecture outline

Recap: Imitation Learning via Behavior Cloning



Challenges with DAgger

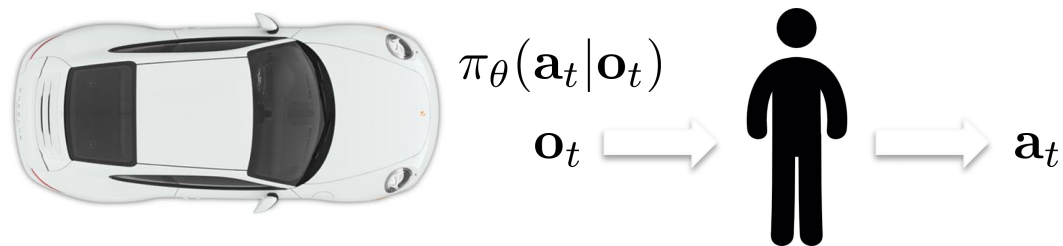
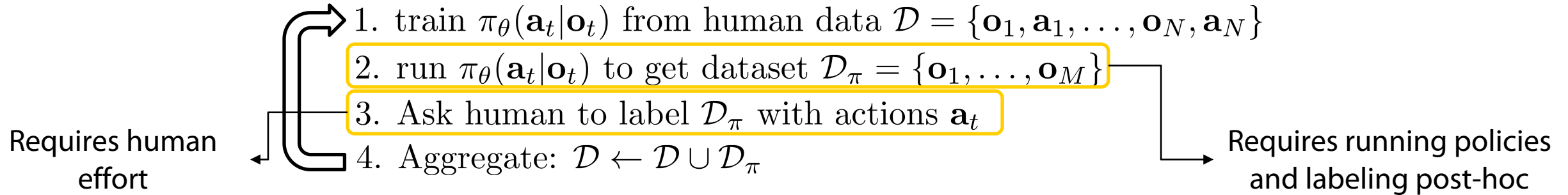


Modern Takes on Imitation Learning



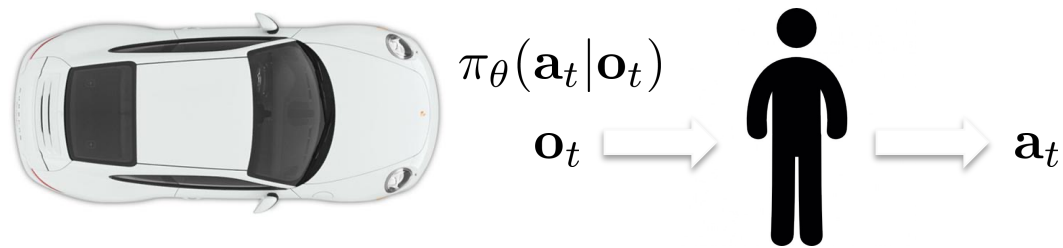
Policy Gradient

What's the problem?



How might we fix this?

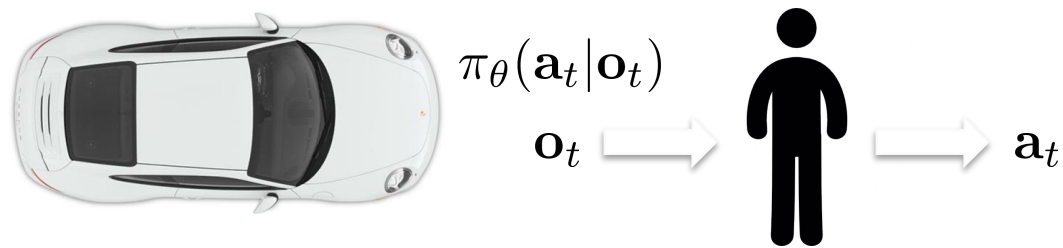
- Can we reduce human cost? →
1. train $\pi_\theta(\mathbf{a}_t|\mathbf{o}_t)$ from human data $\mathcal{D} = \{\mathbf{o}_1, \mathbf{a}_1, \dots, \mathbf{o}_N, \mathbf{a}_N\}$
 2. run $\pi_\theta(\mathbf{a}_t|\mathbf{o}_t)$ to get dataset $\mathcal{D}_\pi = \{\mathbf{o}_1, \dots, \mathbf{o}_M\}$ ← Can we avoid labeling post-hoc?
 3. Ask human to label \mathcal{D}_π with actions \mathbf{a}_t
 4. Aggregate: $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_\pi$



How might we fix this?

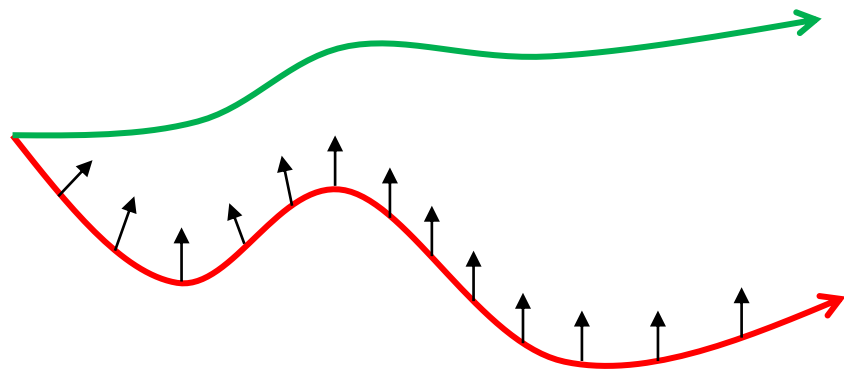
1. train $\pi_\theta(\mathbf{a}_t|\mathbf{o}_t)$ from human data $\mathcal{D} = \{\mathbf{o}_1, \mathbf{a}_1, \dots, \mathbf{o}_N, \mathbf{a}_N\}$
2. run $\pi_\theta(\mathbf{a}_t|\mathbf{o}_t)$ to get dataset $\mathcal{D}_\pi = \{\mathbf{o}_1, \dots, \mathbf{o}_M\}$
3. Ask human to label \mathcal{D}_π with actions \mathbf{a}_t
4. Aggregate: $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_\pi$

Can we avoid
labeling post-hoc?



Dagger is a bit impractical → HG-Dagger

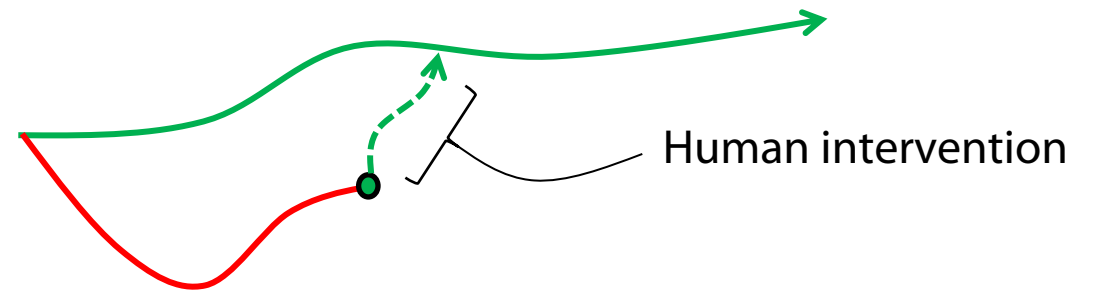
Dagger requires revisiting states and labeling post-hoc



$$s_t \sim d^{\pi_\theta}$$
$$a \sim \pi_h(a|s)$$

Can do “blended” control, but becomes hard to get precision

What if we allowed people to take control for longer?



$$s_t \sim d^{\pi_\theta} \quad \text{(point of intervention)}$$
$$s_{t+1:t+k} \sim d^{\pi_h} | s_t \quad \text{(human takeover)}$$
$$a \sim \pi_h(a|s) \quad \text{(human actions)}$$

Can we do this pre-emptively?

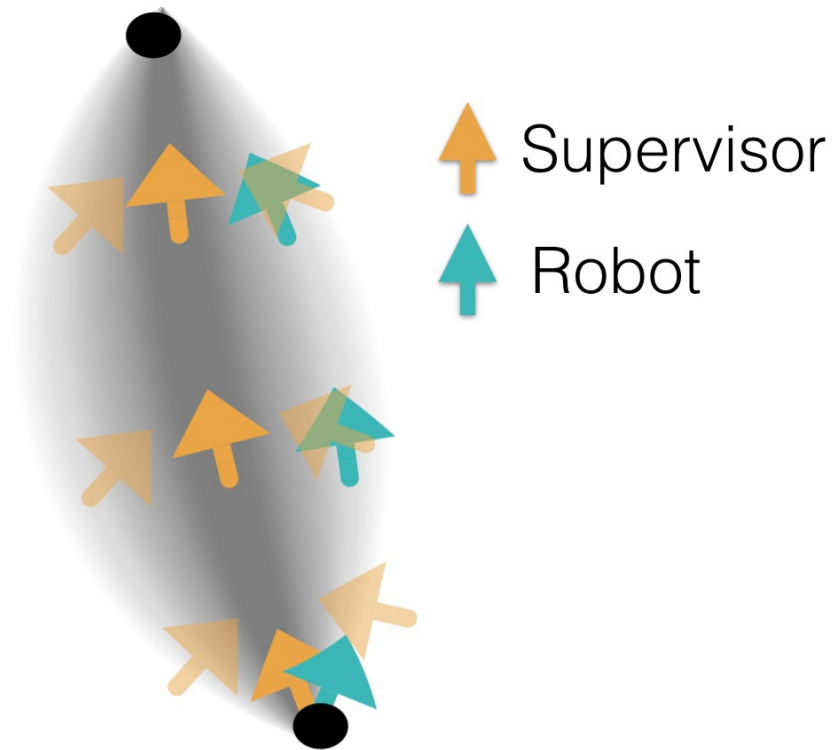
Noising the Data Collection Process

Key idea: force the human to correct for noise **during** training

Under noise during data collection

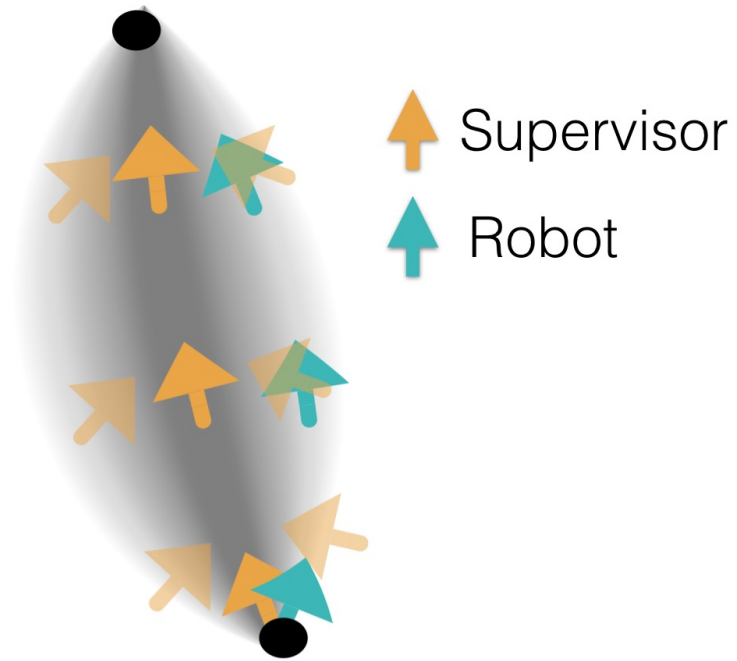
Maximize likelihood

$$\hat{\psi}_{k+1} = \underset{\psi}{\operatorname{argmin}} E_{p(\xi|\pi_{\theta^*}, \psi_k)} - \sum_{t=0}^{T-1} \log [\pi_{\theta^*}(\pi_{\hat{\theta}}(\mathbf{x}_t)|\mathbf{x}_t, \psi)]$$



Why might this not be enough?

Key idea: force the human to correct for noise during training



Noise Injection

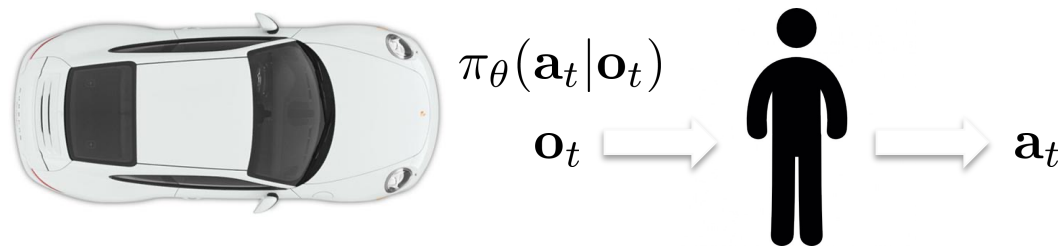


Assumes that the expert can actually perform behaviors under noise
→ Not always possible!

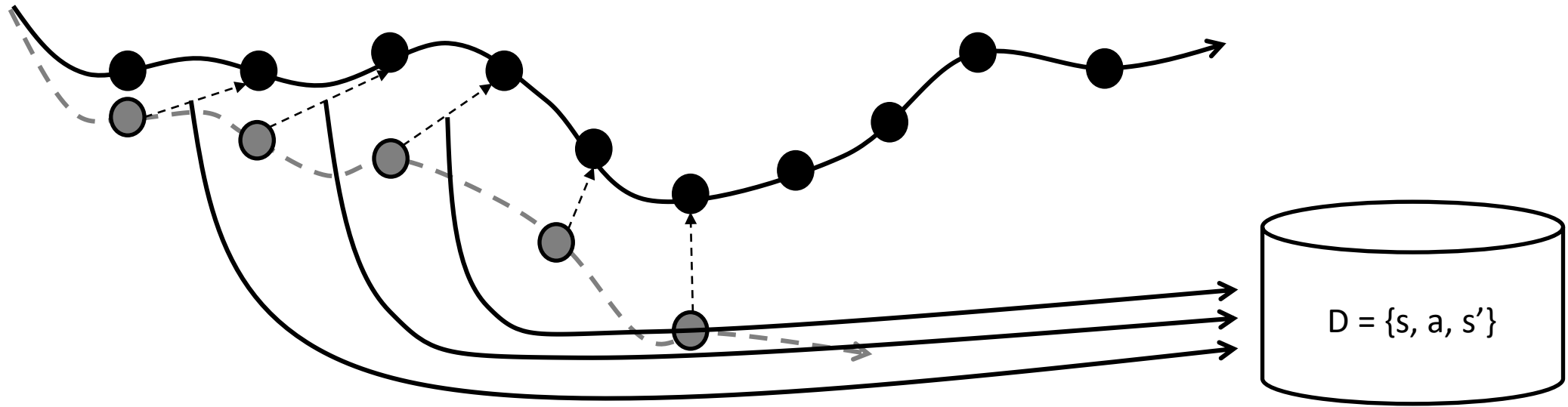
How might we fix this?

Can we reduce
human cost?

1. train $\pi_{\theta}(\mathbf{a}_t|\mathbf{o}_t)$ from human data $\mathcal{D} = \{\mathbf{o}_1, \mathbf{a}_1, \dots, \mathbf{o}_N, \mathbf{a}_N\}$
2. run $\pi_{\theta}(\mathbf{a}_t|\mathbf{o}_t)$ to get dataset $\mathcal{D}_{\pi} = \{\mathbf{o}_1, \dots, \mathbf{o}_M\}$
3. Ask human to label \mathcal{D}_{π} with actions \mathbf{a}_t
4. Aggregate: $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_{\pi}$



Can we avoid expensive online data collection/labeling?



Generate corrective labels
to dataset for imitation

How can we find corrective labels without an expensive human in the loop
and online data collection?



Abhay
Deshpande

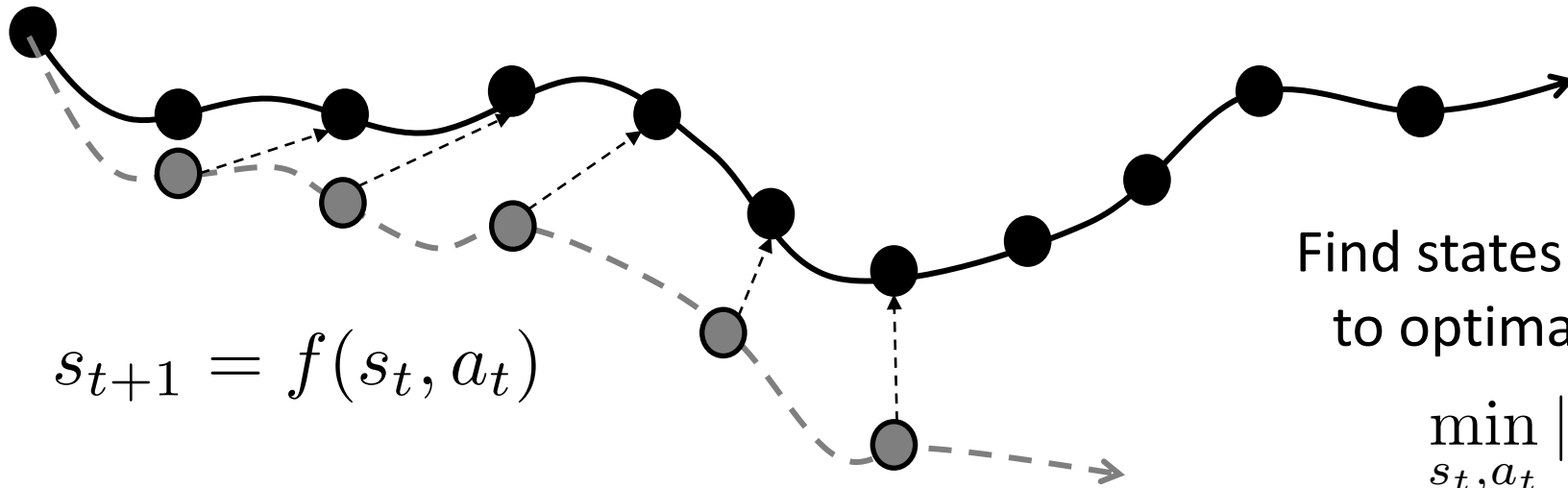


Yunchu
Zhang



Liyiming
Ke

Generating Corrective Labels From True Dynamics



Find states (s_t), actions (a_t) that lead back to optimal states under true dynamics

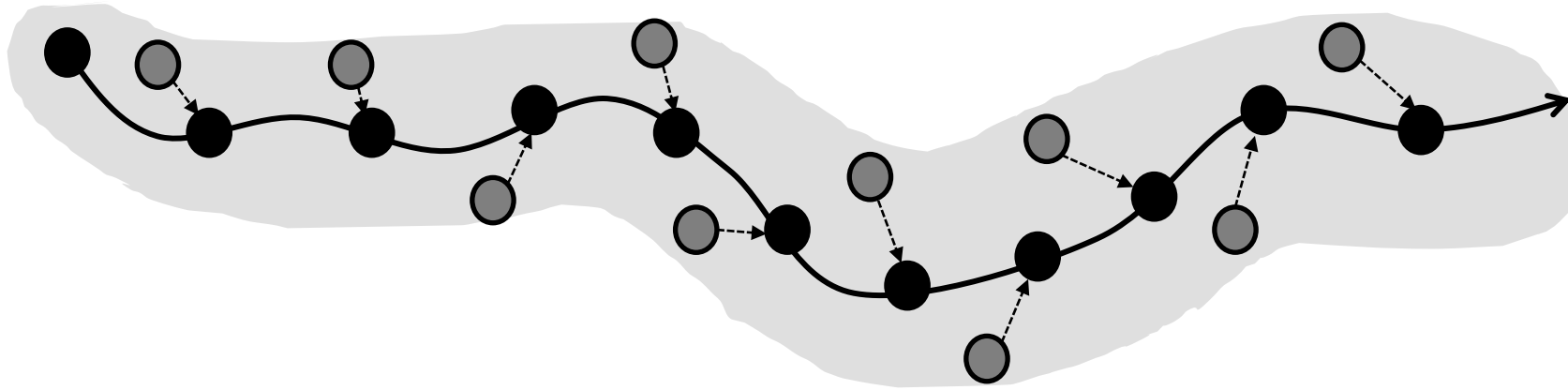
$$\min_{s_t, a_t} \|s_{t+1}^* - f(s_t, a_t)\| \leq \epsilon$$

Easy with known dynamics

Intuition: find labels to bring OOD states back in distribution

But models are unknown! ☹️

Generating Corrective Labels with Learned Dynamics



Ok models are unknown,
let's learn them!

$$\min_{\hat{f}} \mathbb{E}_{(s_t, a_t, s_{t+1}) \sim \mathcal{D}} \left[\|\hat{f}(s_t, a_t) - s_{t+1}\|_2 \right]$$

But learned dynamics \hat{f}_ϕ are not globally accurate?



Under approximately Lipschitz smooth models, trust models around training data

$$\|s_{t+1}^* - \hat{f}_\phi(s_t, a_t)\| \leq \epsilon$$

Find states (s_t), actions (a_t) that lead back to optimal states under ~~true~~ learned dynamics,
where learned dynamics can be trusted

$$\min_{s_t, a_t} \|s_{t+1}^* - \hat{f}_\phi(s_t, a_t)\| \leq \epsilon \longleftarrow \text{Corrective label}$$

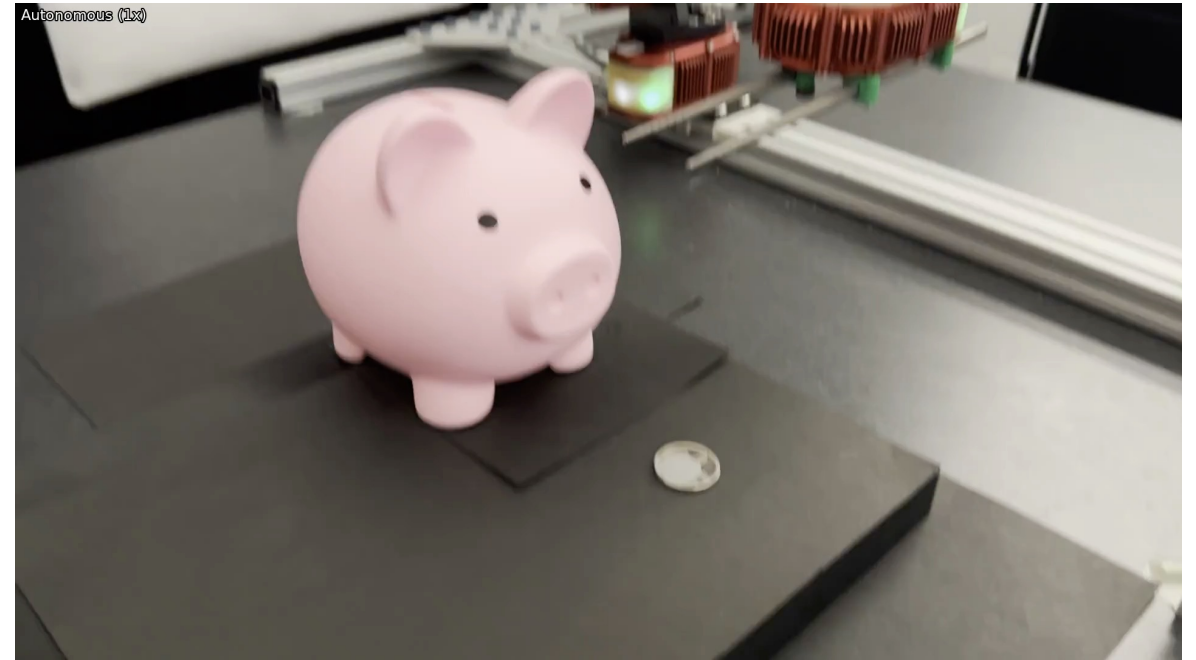
$$\text{s.t. } \|s_t^* - s_t\| \leq \epsilon_1, \|a_t^* - a_t\| \leq \epsilon_2 \longleftarrow \text{Close to data}$$

How well does generating corrective labels work?

With corrective labels

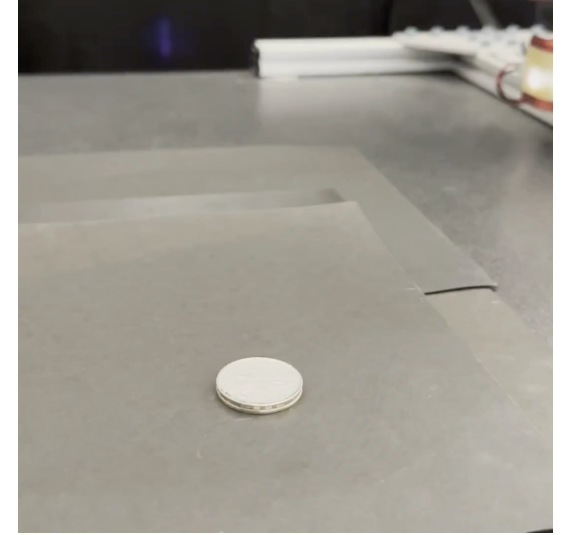
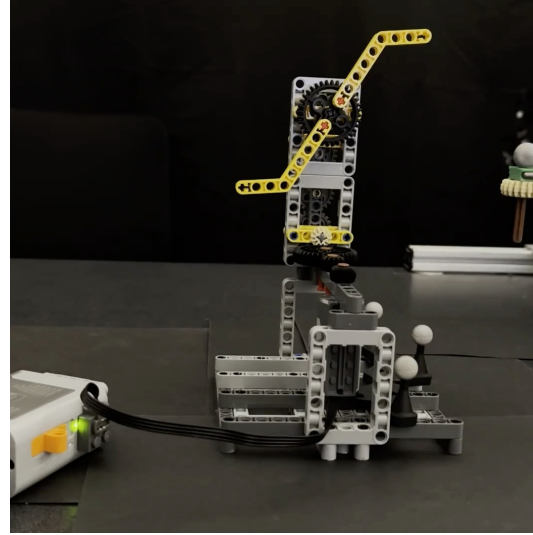
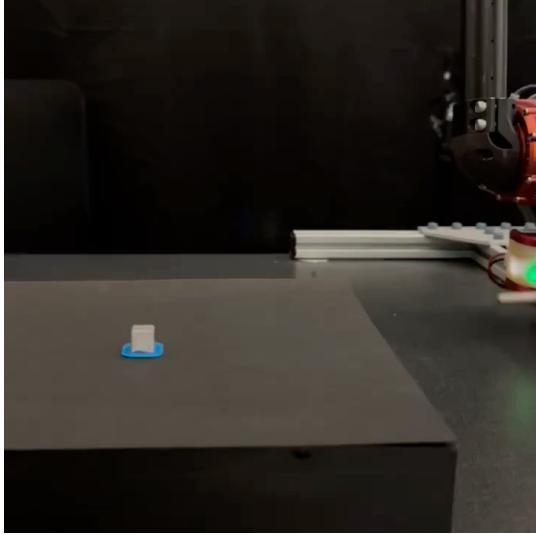


Without corrective labels

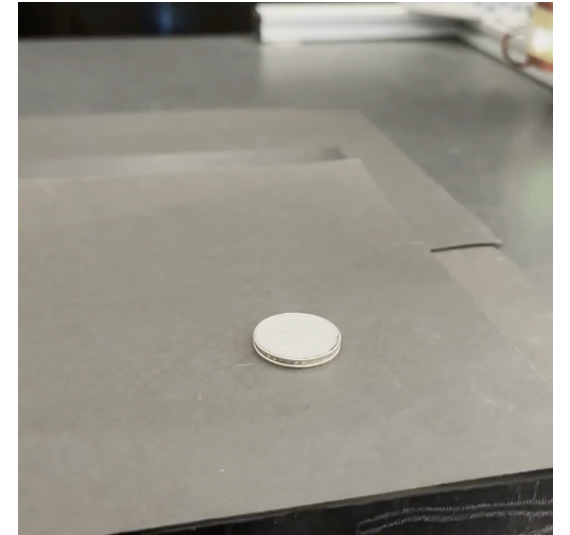
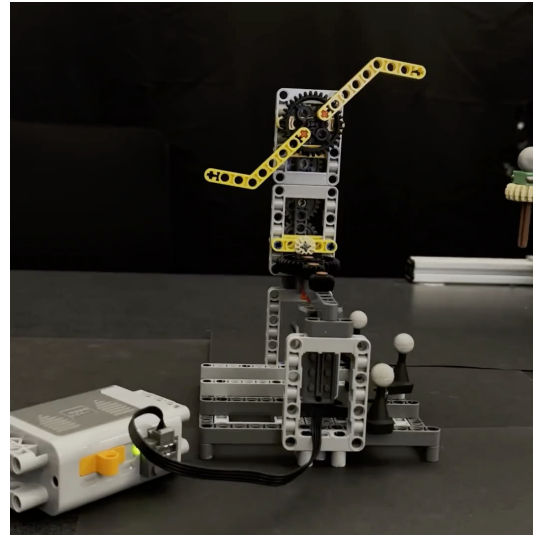
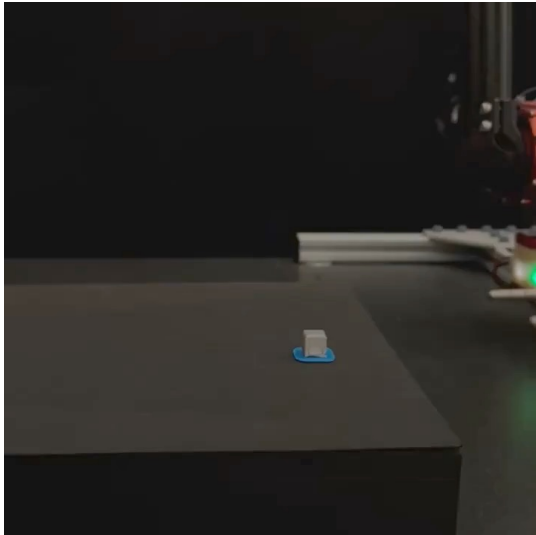


How well does generating corrective labels work?

With corrective labels



Without corrective labels



Lecture outline

Recap: Imitation Learning via Behavior Cloning



Challenges with DAgger



Modern Takes on Imitation Learning

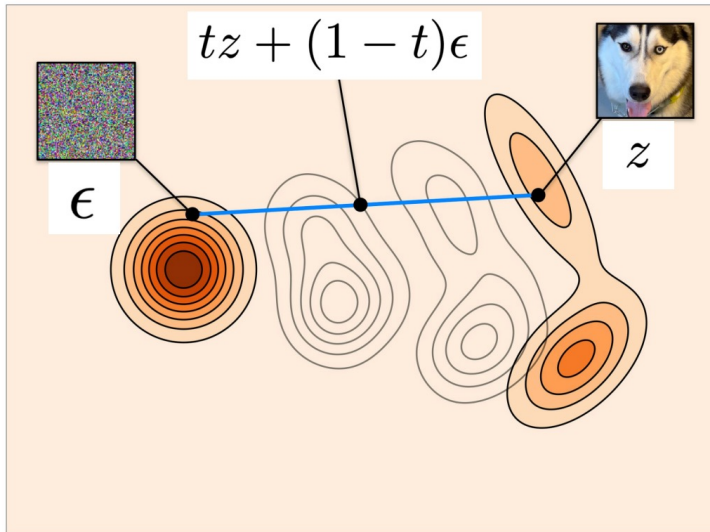


Policy Gradient

But BC works unreasonably well!

Is this by accident?

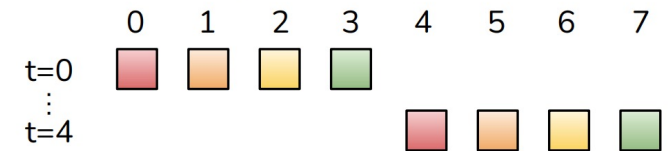
Generative Models



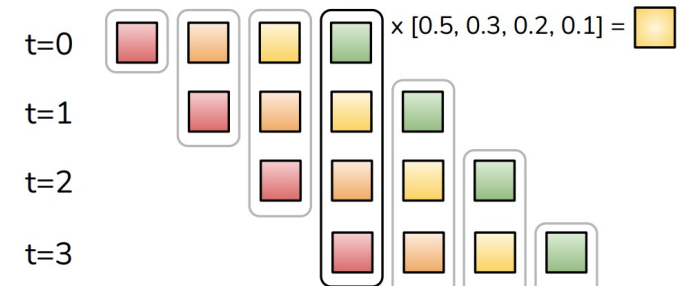
Flow matching, diffusion, etc

Action chunking

Action Chunking



Action Chunking + Temporal Ensemble

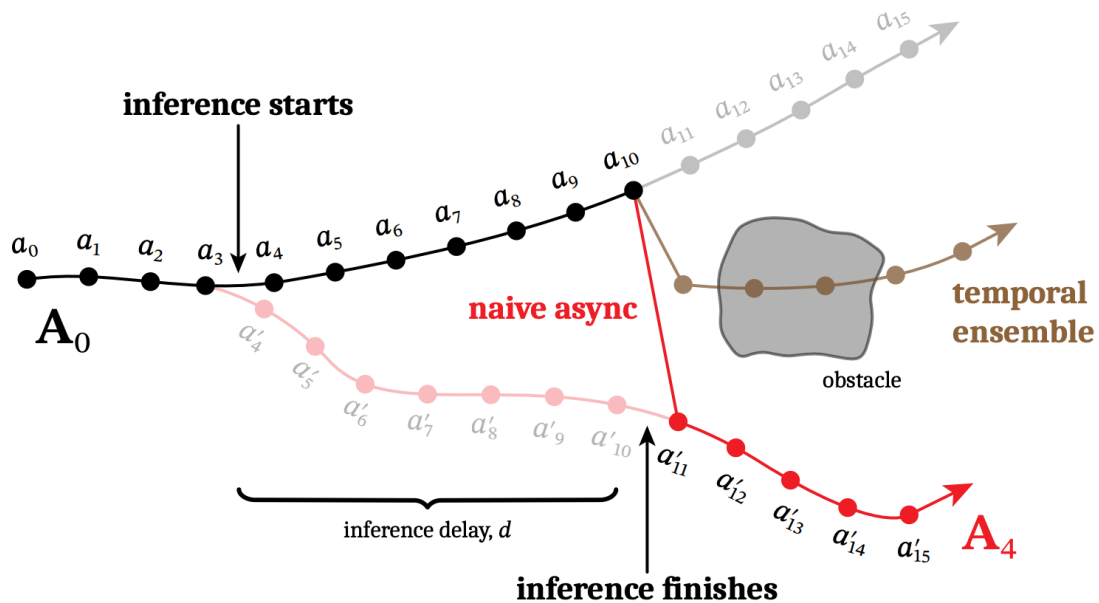


Short-horizon open-loop prediction

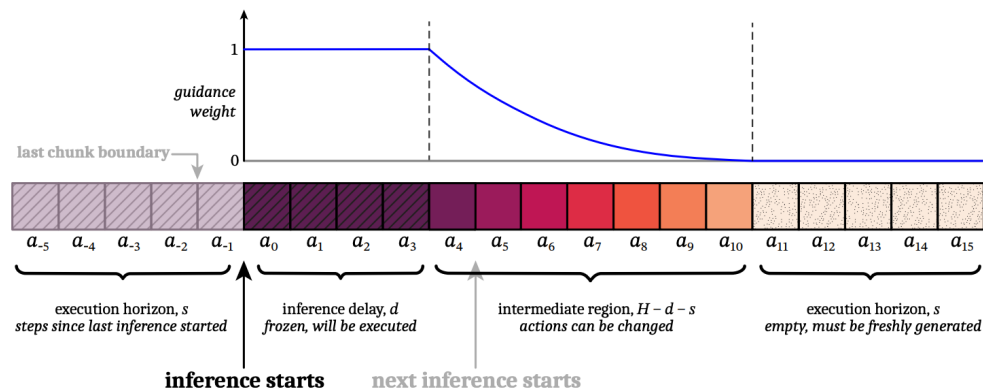
Both are non-trivially important!

“Real-Time” Action Chunking

Real challenge when trying to operate robotic policies quickly



Leads to jerky/suboptimal motion



Freeze k steps and inpaint the rest

Generative Modeling – why does it help?

Do generative models help because human data is multimodal?

Much Ado About Noising: Dispelling the Myths of Generative Robotic Control

Chaoyi Pan^{a,\$} Giri Anantharaman^a Nai-Chieh Huang^a Claire Jin^a Daniel Pfrommer^b
Chenyang Yuan^c Frank Permenter^c Guannan Qu^{a,†} Nicholas Boffi^{a,†} Guanya Shi^{a,†}
Max Simchowitz^{a,†}

Shows that the benefits of generative models is not only in multimodal settings

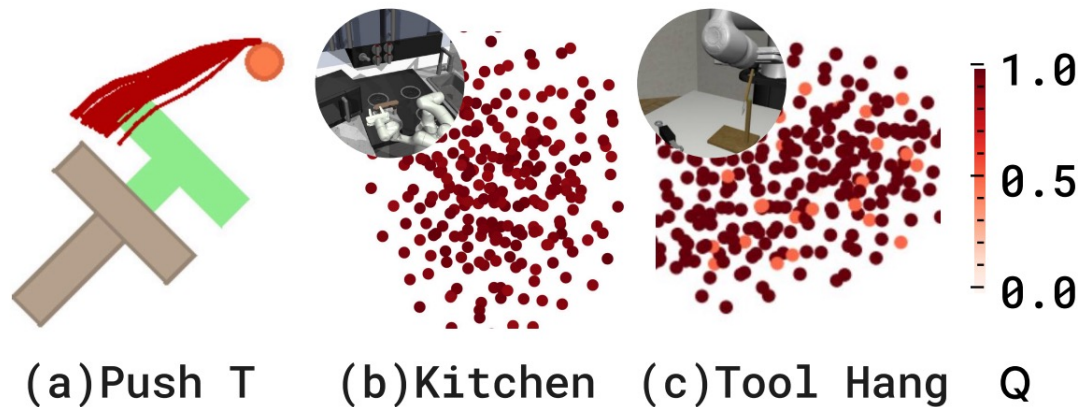
Shows that the benefits come from stochasticity injection + iteration

Generative Modeling – why does it help?

Do generative models help because human data is multimodal?

Not much multimodality in existing policies

Better even with deterministic experts



Task	$z \equiv 0$	$N(0, I)$	Mean z
Push-T	0.97	0.97	0.95
Kitchen	0.99	0.99	0.97
Tool-Hang	0.78	0.80	0.76

Dataset	Flow	Reg.
Original	0.78	0.58
Deterministic	0.72	0.64

Ok – so why do these flow models help?

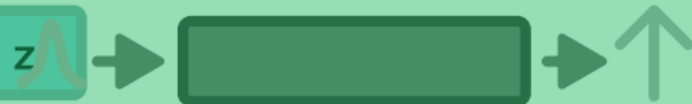
Generative Modeling – why does it help?

C1. Distributional Learning

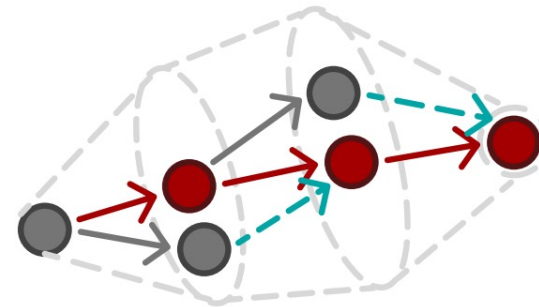


C2. Stochasticity Injection

Noise z

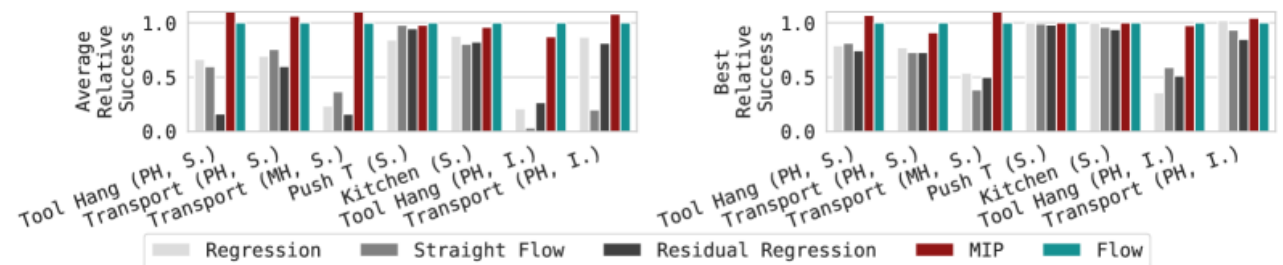


C3. Supervised Iterative Compute



- → ● Ideal GCP generation
- → ● GCP updates during training
- - - → Implicit correction due to stochasticity injection

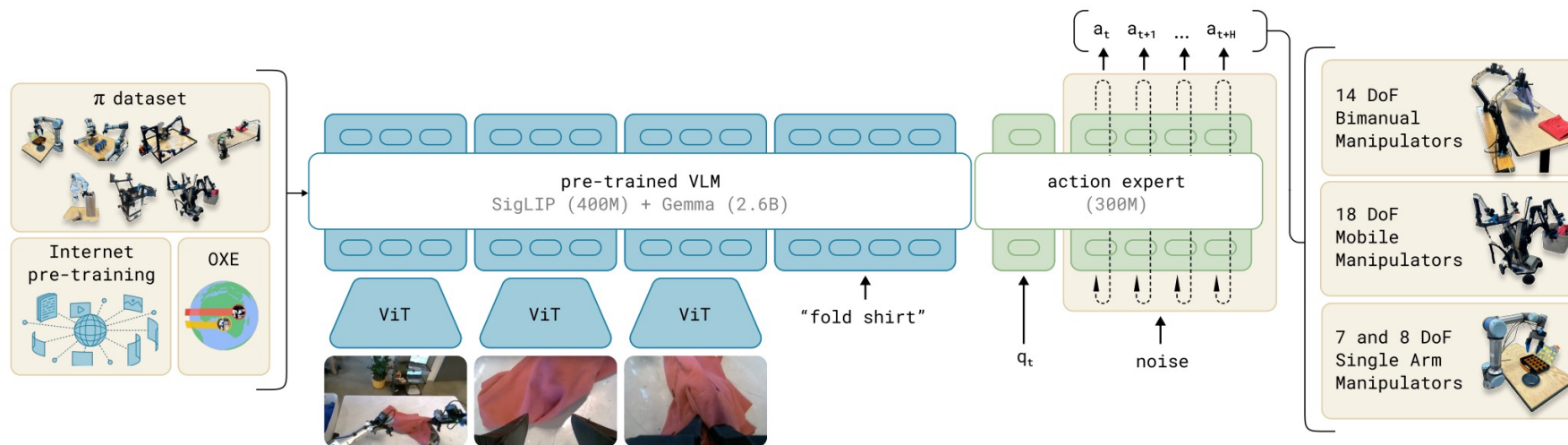
Kind of like DART but inside a single step of action generation



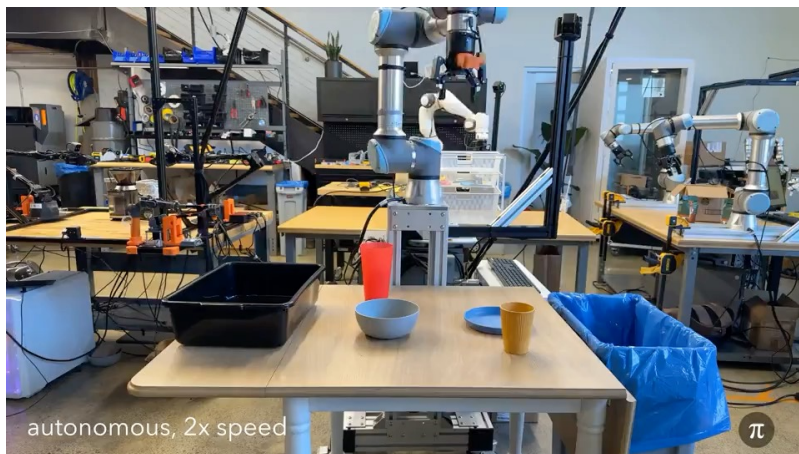
Train with stochasticity injected to “correct for errors” in generation

Minimal degradation when flow noise is removed

How is this being actualized today?



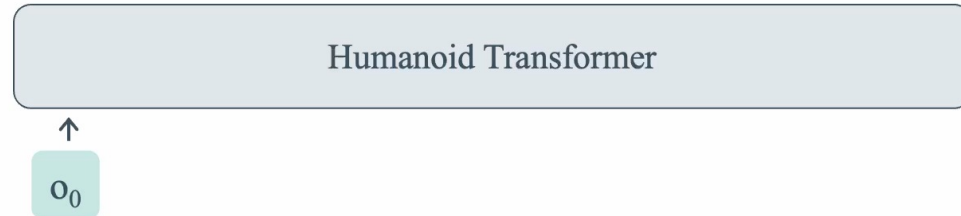
Vision-language-action policies



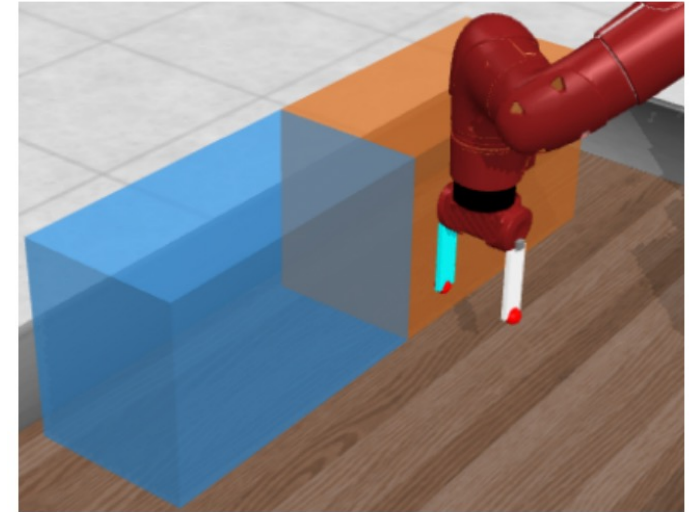
So does this solve all the issues in imitation?

Frontiers in Imitation Learning

Non-Markovian Demonstrators



Characterizing generalization

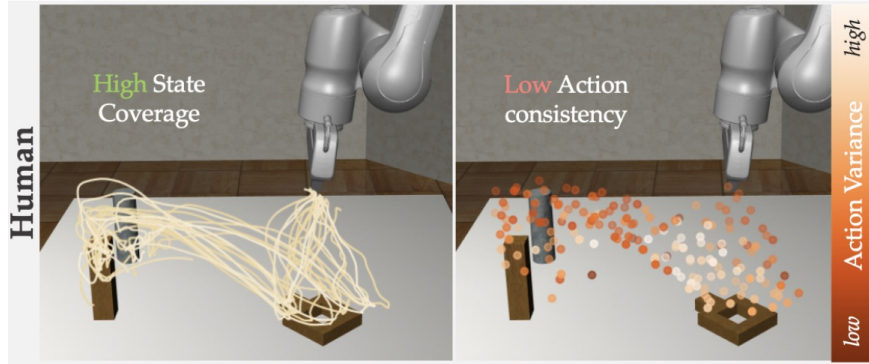


Action-Free Data

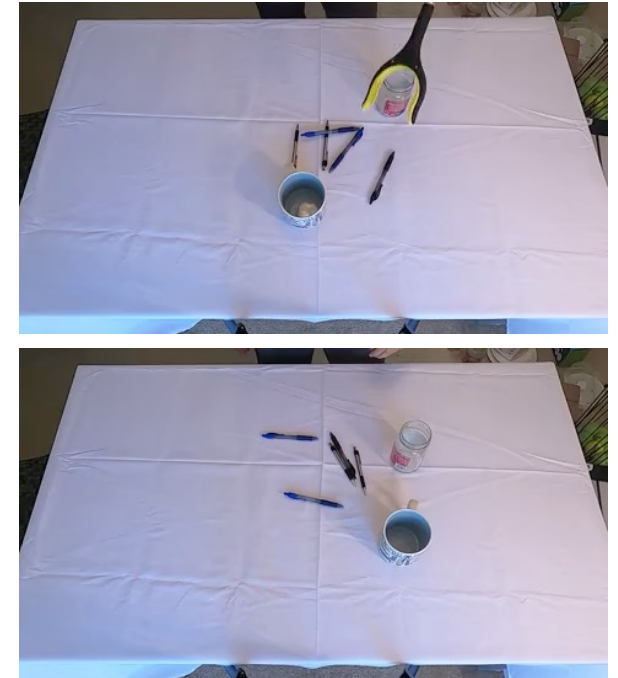


Frontiers in Imitation Learning

Data Curation and Quality



Embodiment Shift

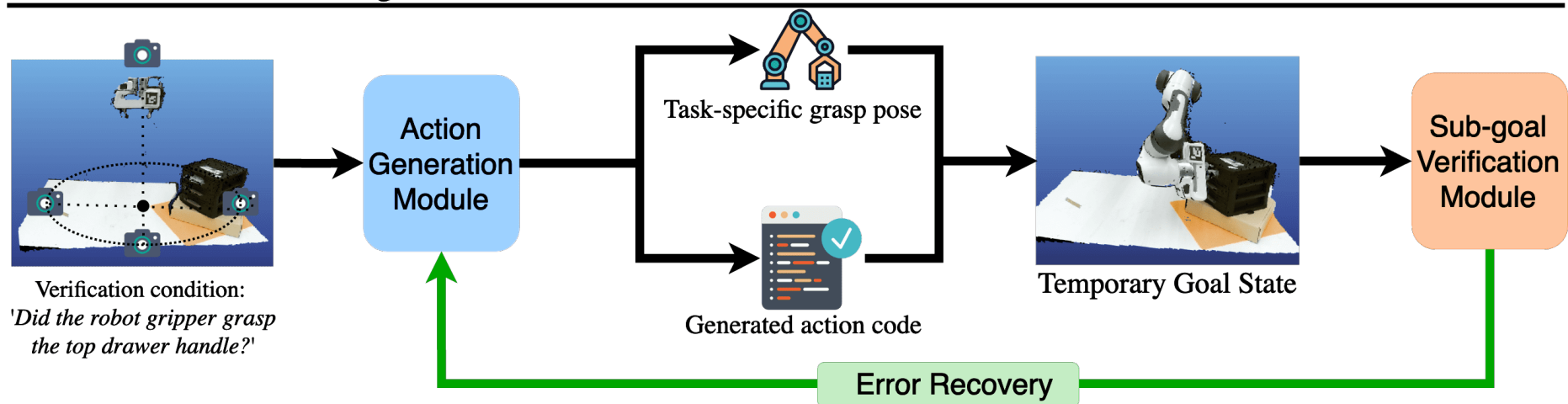
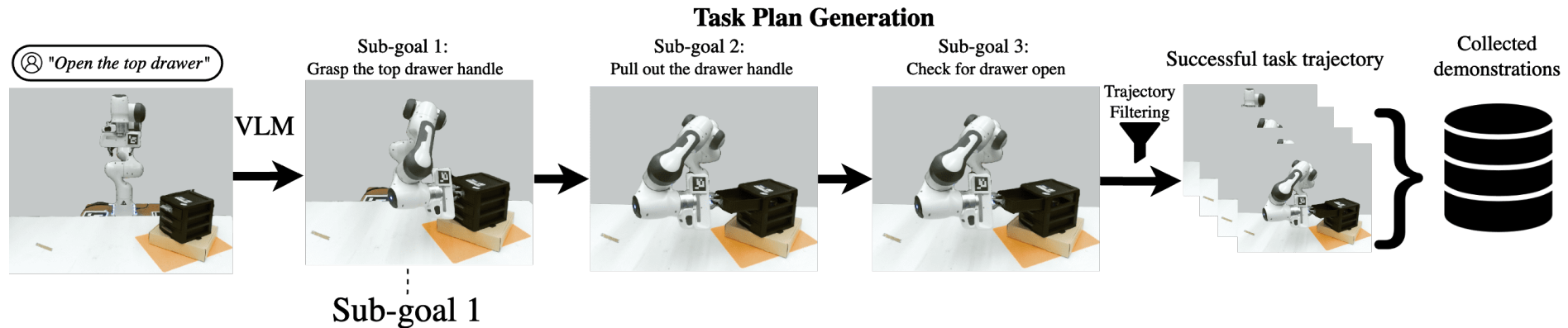


Teleoperation Interfaces



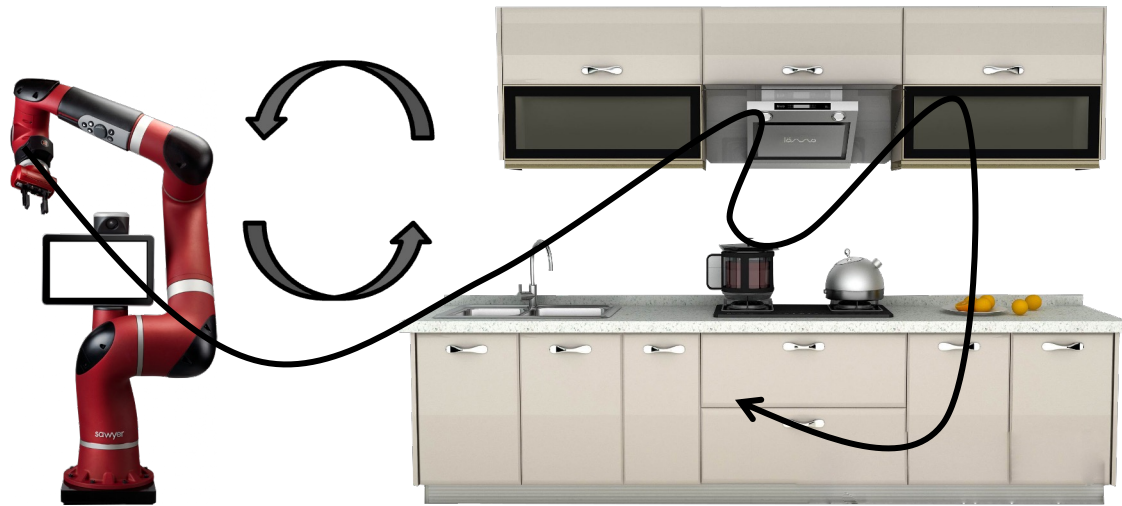
Frontiers in Imitation Learning

Learning how to retry and improve

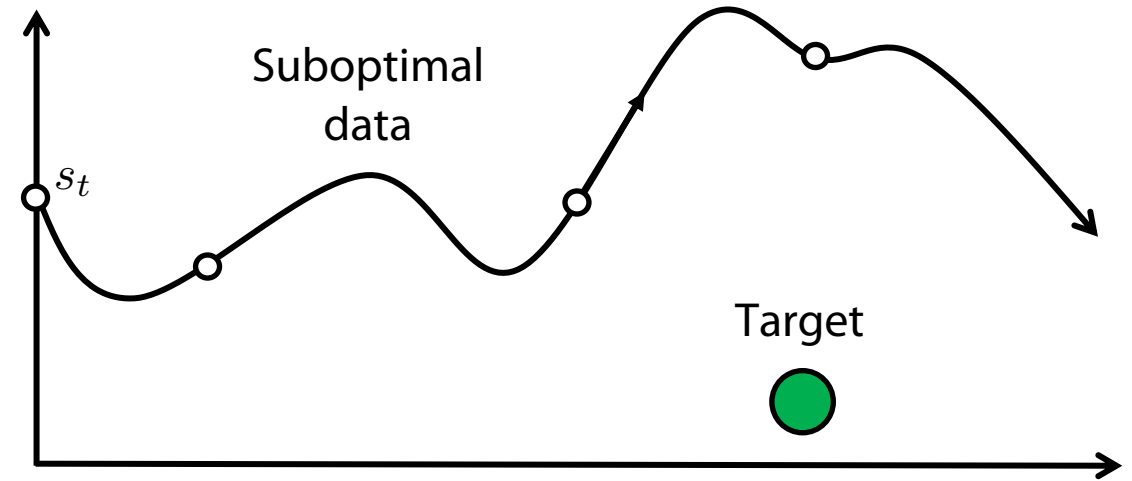


Duan et al

Accounting for Suboptimal Data

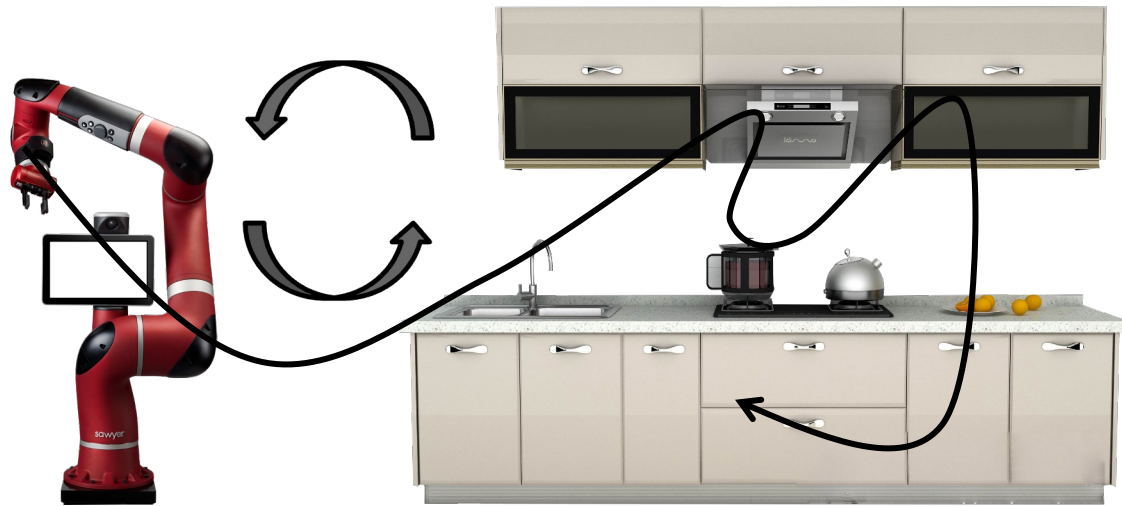


Random
Policy
 π_{θ}



How can we use this suboptimal data,
despite not reaching the target?

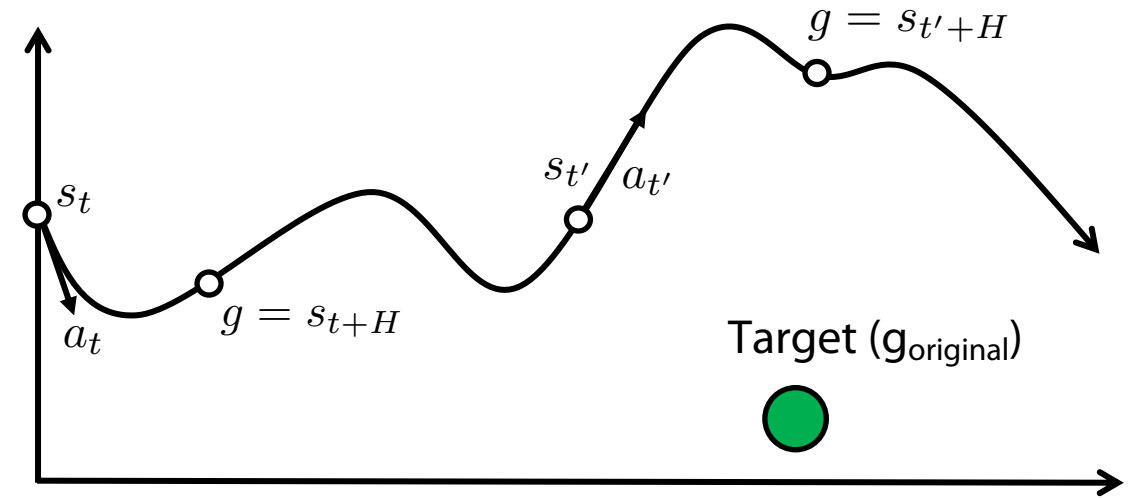
Hindsight relabeling for Imitation Learning



Key insight: maybe the data is not bad, it's just been labeled for the wrong problem!



Relabel the right goal in "hindsight"

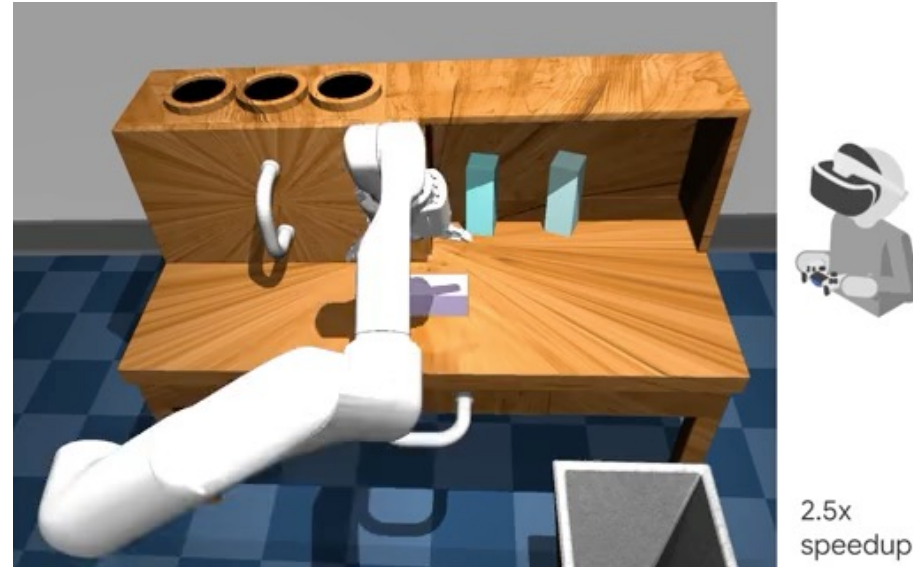


Learn a multi-goal policy $\pi_{\theta}(a|s, g)$

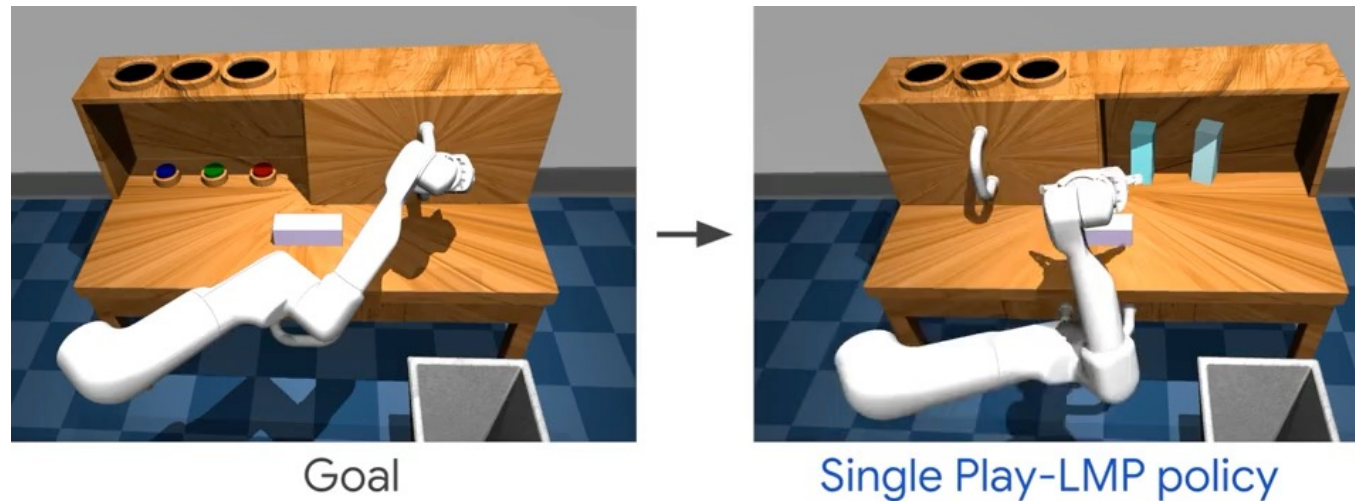


Treat reached states as **optimal** goals

What does this result in?



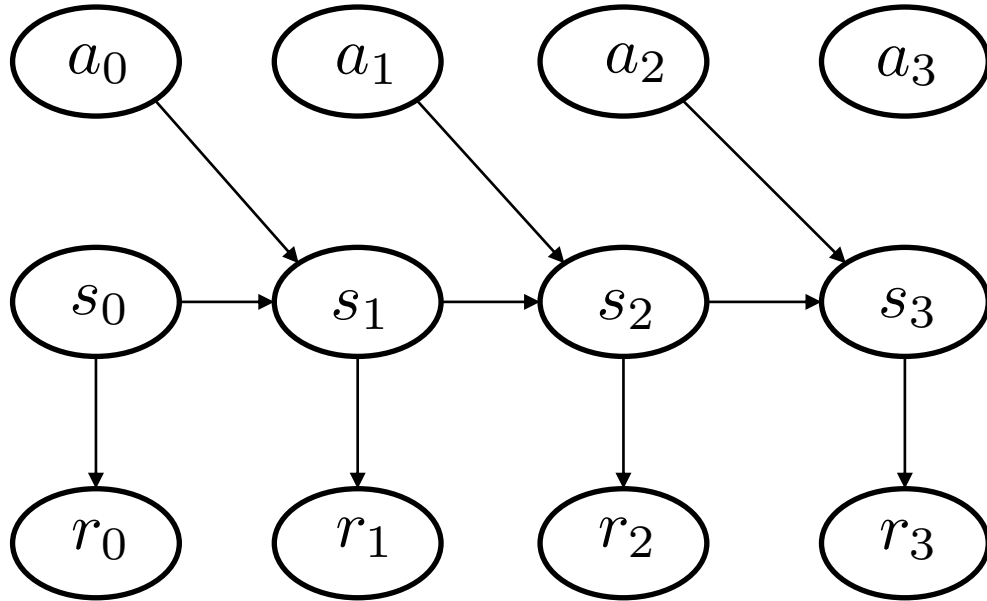
Undirected play data



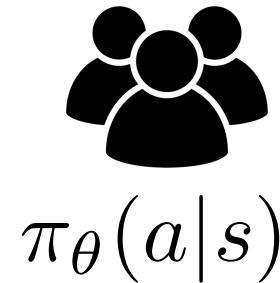
Goal-directed behavior

Dealing with non-Markovian demonstrators

Markov property $p(s_0, s_1, s_2, a_0, a_1, a_2) = p(s_0)p(a_0|s_0)p(s_1|s_0, a_0)p(a_1|s_1)p(s_2|s_1, a_1)p(a_2|s_2)$



Are human demonstrators Markovian?



If we see the same thing twice, we do the same thing twice, regardless of what happened before

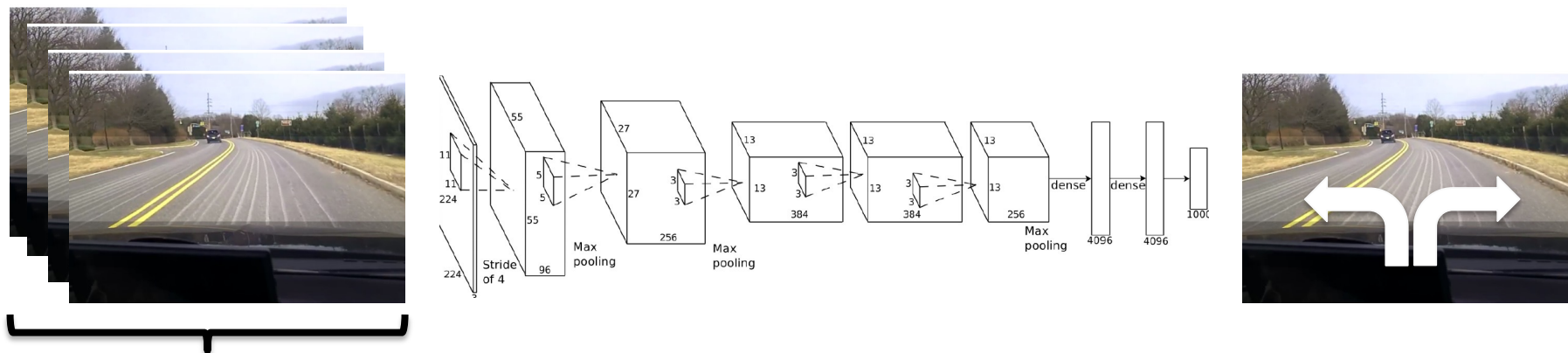
Not necessarily!

Humans often rely on history

Mixtures of Markovian humans may not be Markovian

How can we deal with non-Markovian demonstrators?

$$\text{Learn } \pi_{\theta}(a_t | s_t, s_{t-1}, \dots, s_0)$$

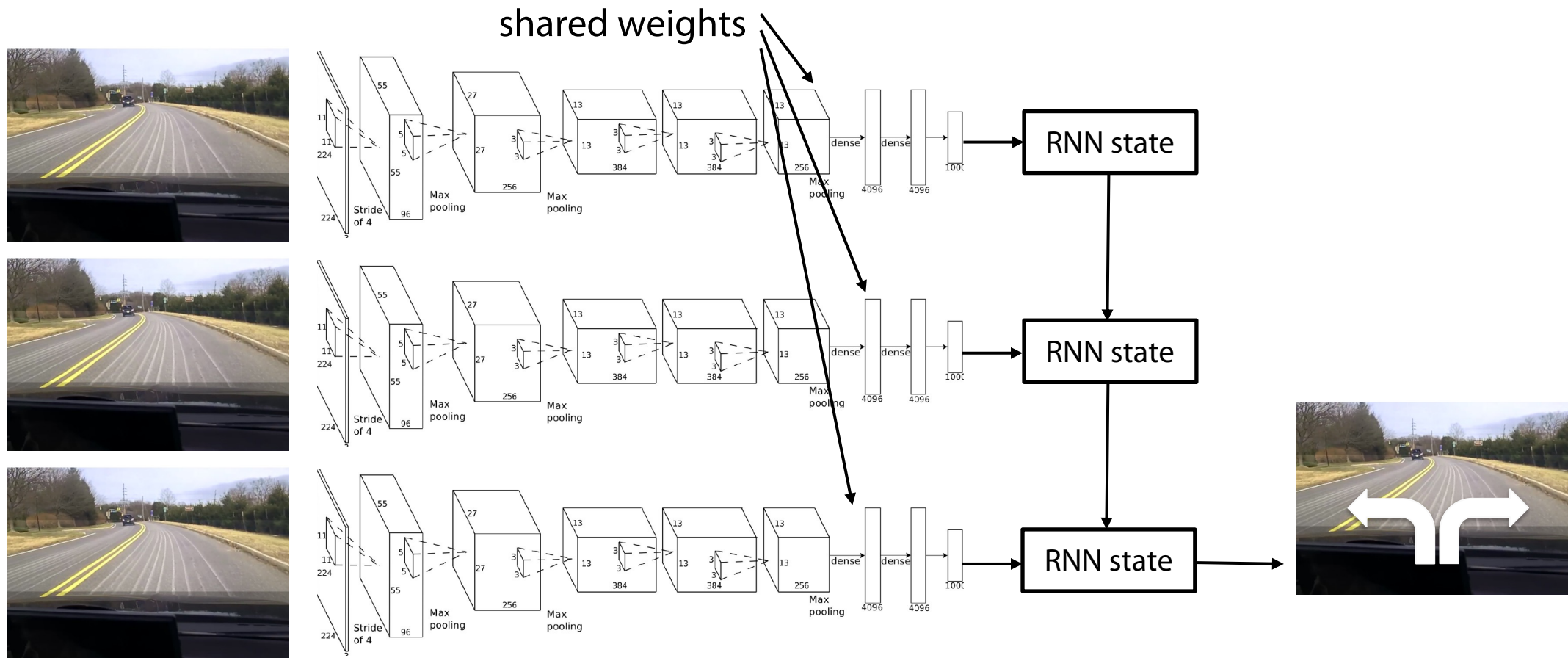


variable number of frames,
too many weights

Option 1: Stack all the past frames into a feedforward NN

How can we deal with non-Markovian demonstrators?

$$\text{Learn } \pi_{\theta}(a_t | s_t, s_{t-1}, \dots, s_0)$$



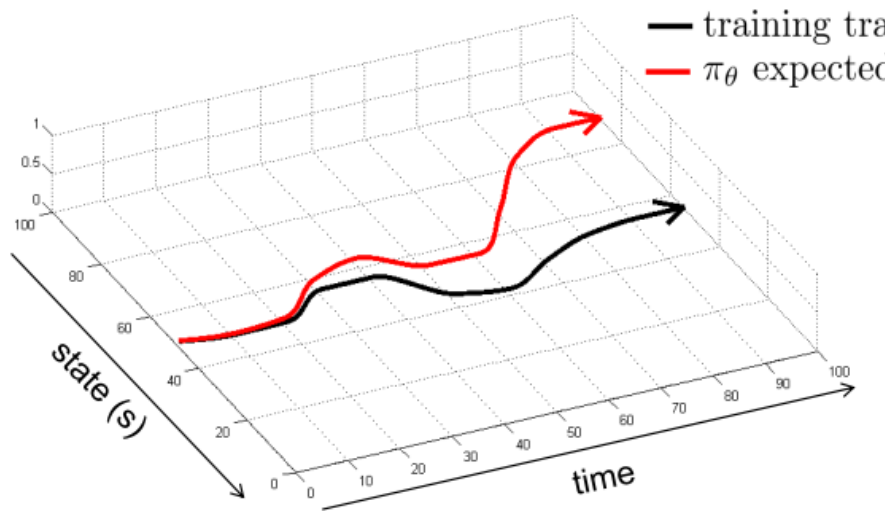
Option 2: Use a recurrent model (LSTM/transformer/RNN)

Credit: Sergey Levine

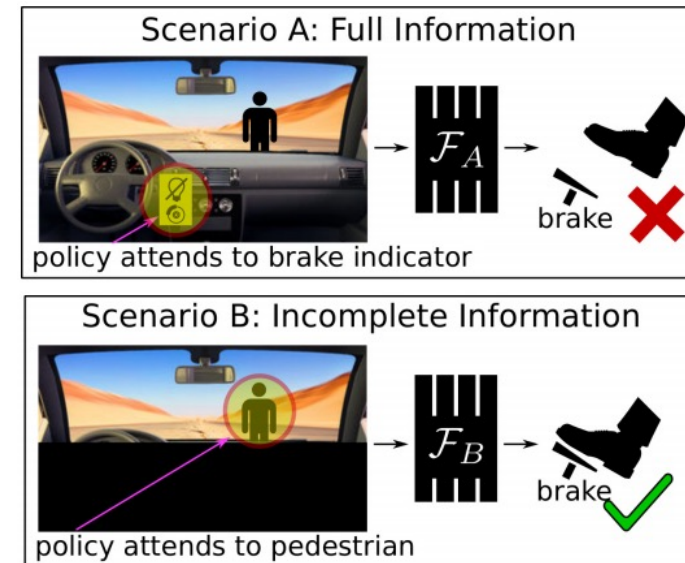
Why might this be challenging?

$$\text{Learn } \pi_{\theta}(a_t | s_t, s_{t-1}, \dots, s_0)$$

Easier to go OOD

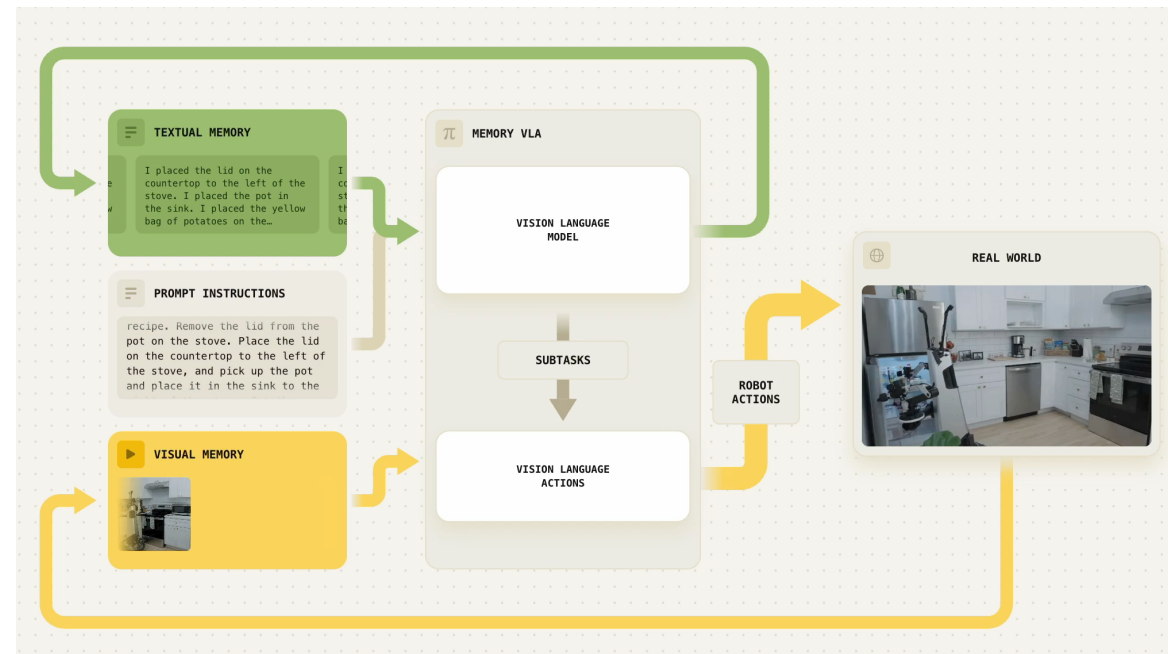


Learns spurious shortcut behaviors



Recent advances in memory for imitation learning

Use text + video encoder for long and short term memory



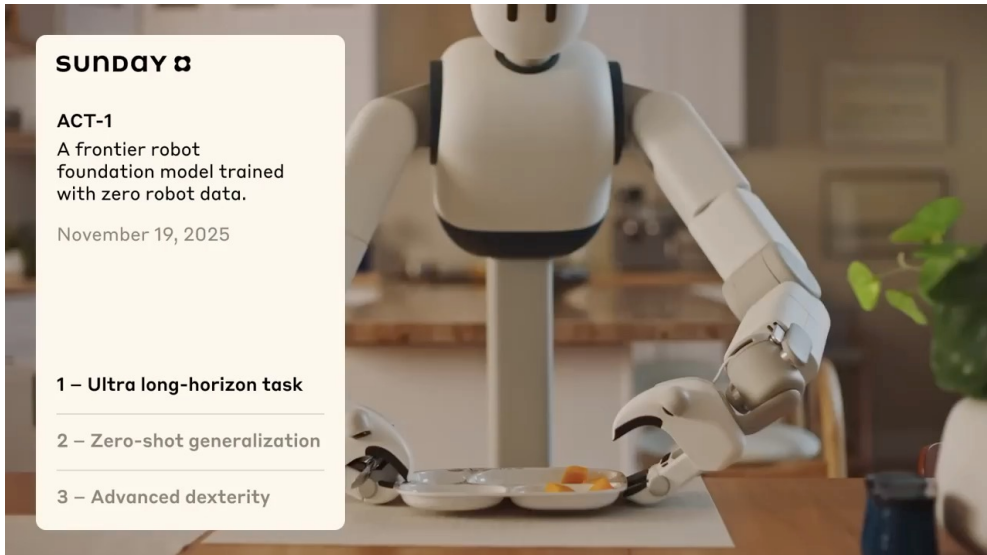
Still many open questions!

Some cool imitation videos

Generalist/Sunday Robots



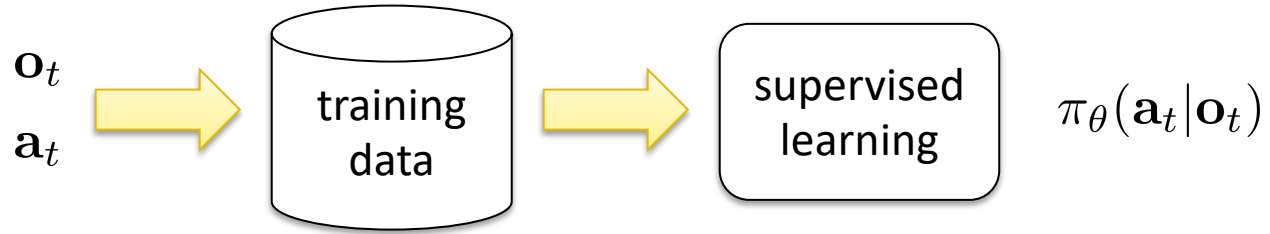
Sensorized people improve data throughput



Inheriting capabilities of vision-language models

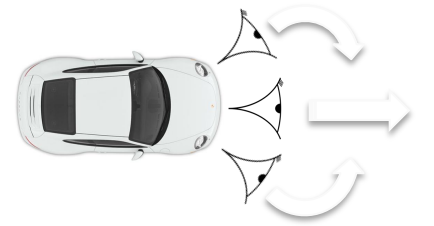


Perspectives on Imitation



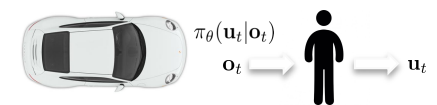
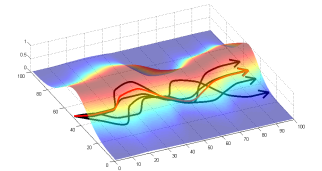
■ Pros:

- Easy to use, no additional infra
- Can sometimes be unreasonably effective



■ Cons:

- Challenges of compounding error, multimodality
- Doesn't really generalize
- Very expensive in terms of data collection!



Lecture outline

Recap: Imitation Learning via Behavior Cloning



Challenges with DAgger

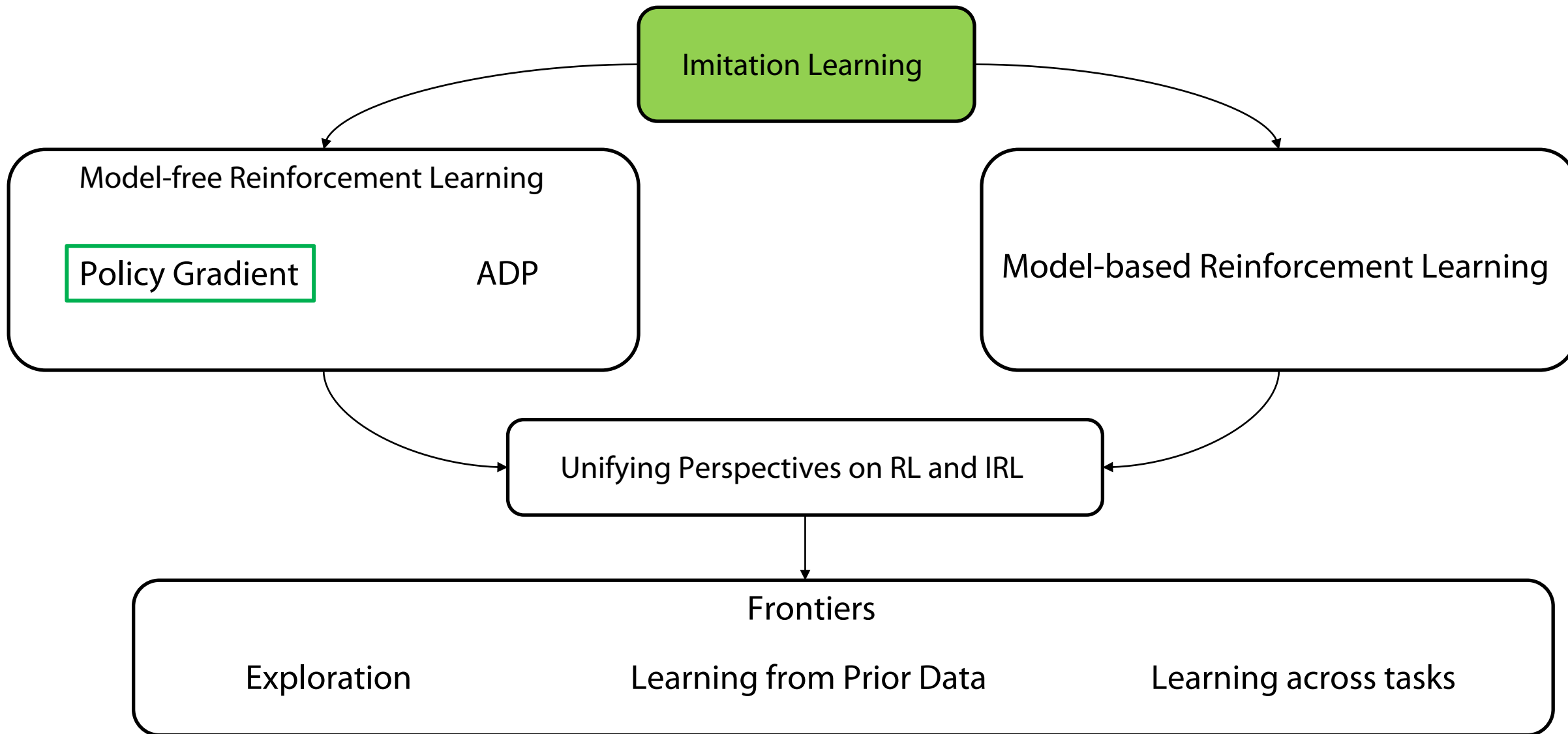


Modern Takes on Imitation Learning

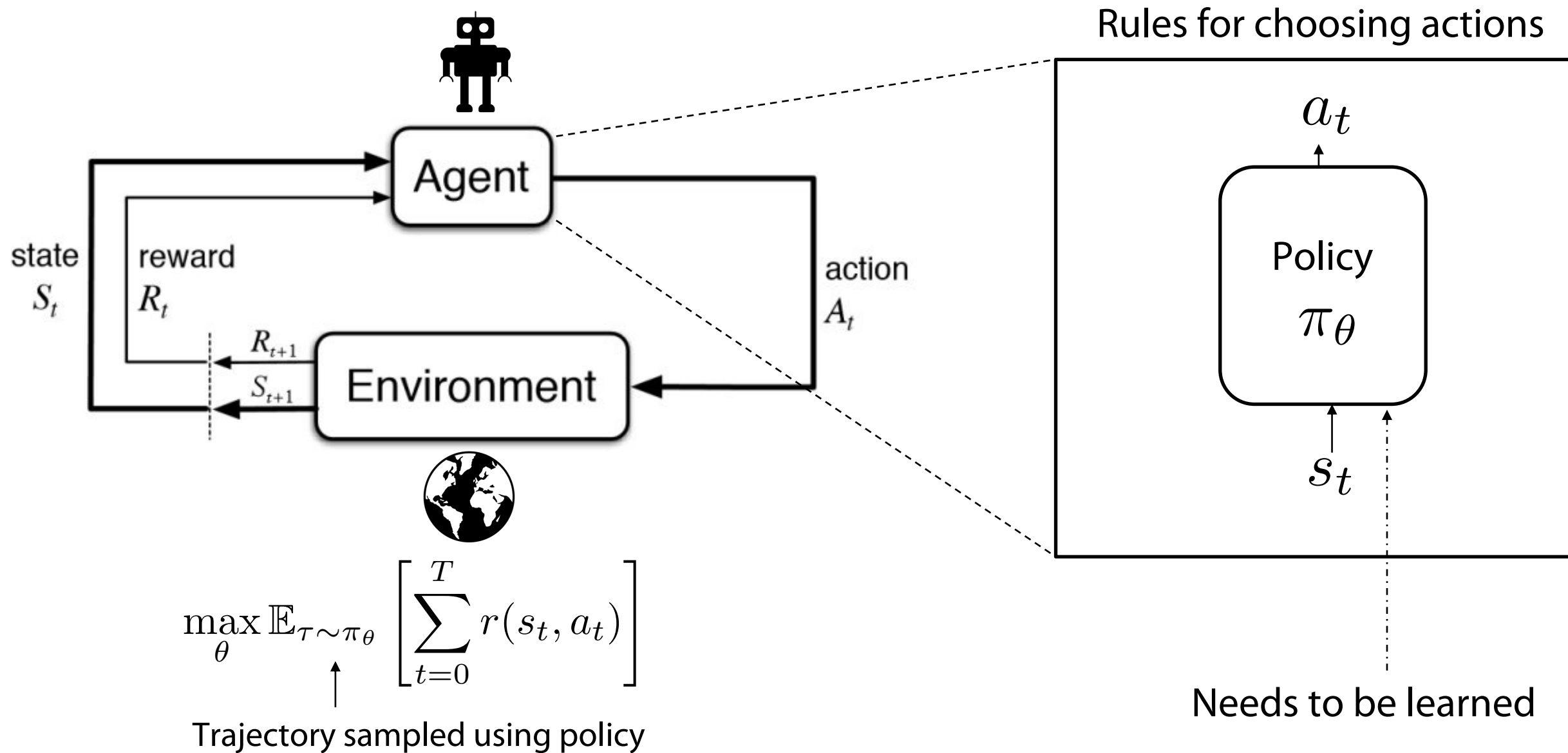


Policy Gradient

Class Structure

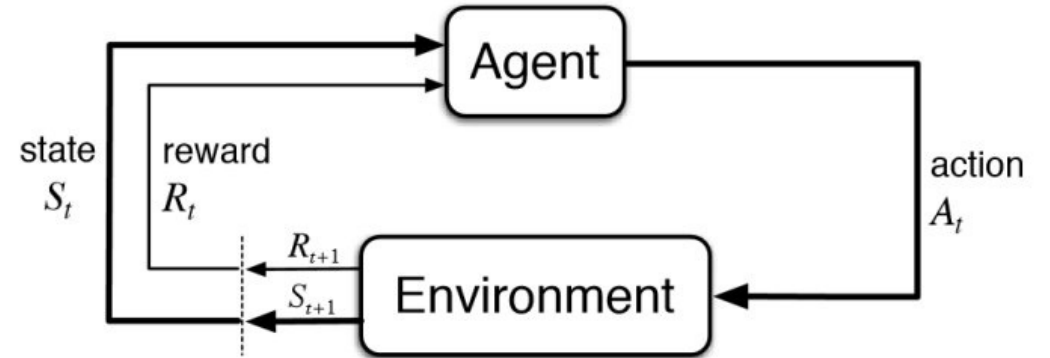


Objective of Reinforcement Learning



Finite horizon vs infinite horizon objective

$$\max_{\theta} \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^T r(s_t, a_t) \right]$$



Finite horizon

$$\mathbb{E}_{\pi_{\theta}^t} \left[\sum_{t=0}^T r(s_t, a_t) \right]$$

Time-dependent policy
(not stationary)

$$\mathbb{E}_{\pi_{\theta}^t} \left[\sum_{t=0}^T \gamma^t r(s_t, a_t) \right]$$

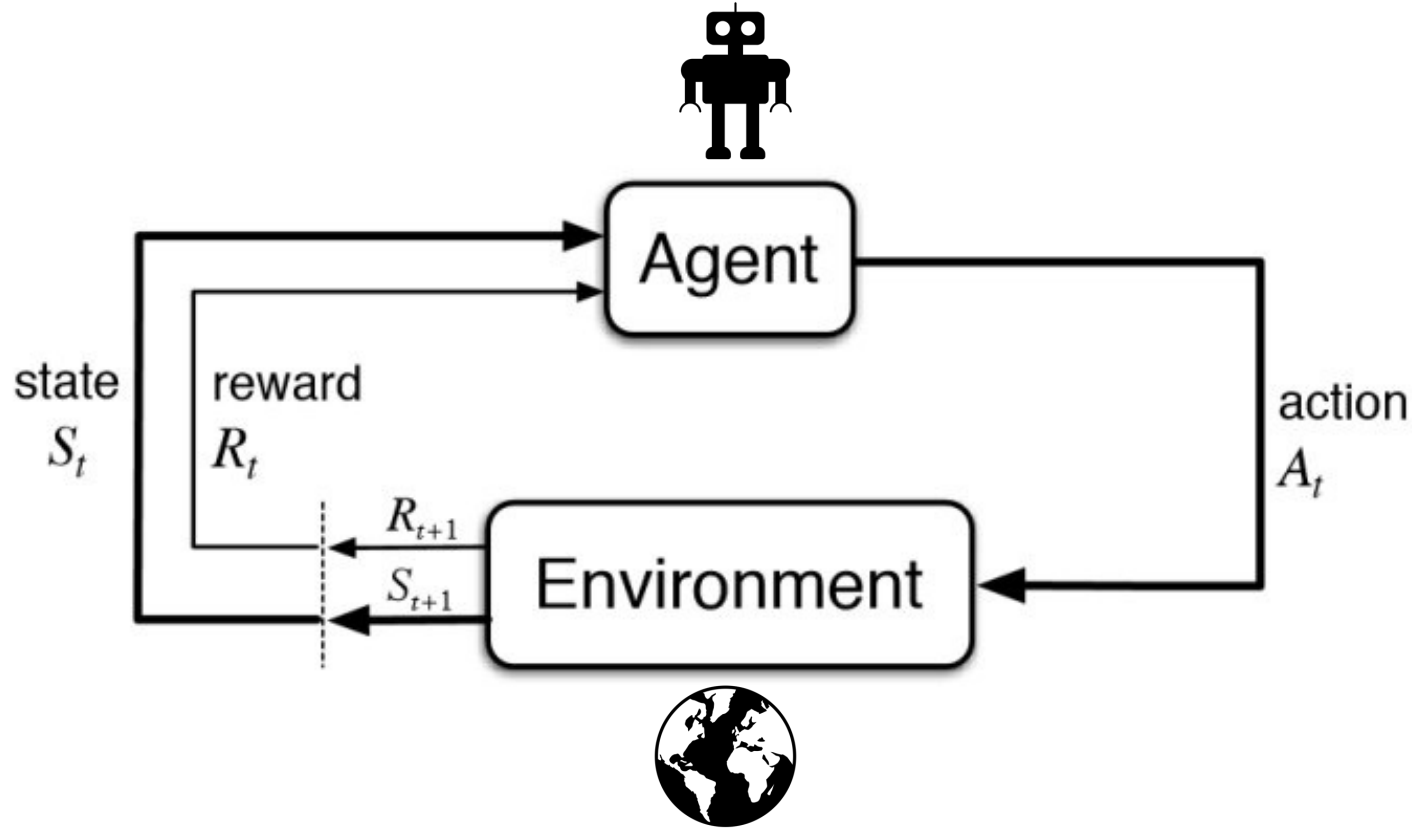
Infinite horizon discounted

$$\mathbb{E}_{\pi_{\theta}} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right]$$

Time-independent (stationary) policy
→ Need discount to prevent blow up

Lemma: there always exists a stationary optimal policy

Objective of Reinforcement Learning



Assumptions:

1. Rewards are additive
2. Dynamics can be sampled from, but functional form is unknown
3. Rewards are provided as every state is visited, functional form is unknown

$$\max_{\theta} \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^T r(s_t, a_t) \right]$$

↑
Trajectory sampled using policy

Connection to Optimal Control

Closely related: typically problem of finding control given a plant

$$\min_{x,u} \int_0^x L(t, x(t), u(t)).dx$$

w.r.t

$$x'(t) = f(x(t), u(t))$$

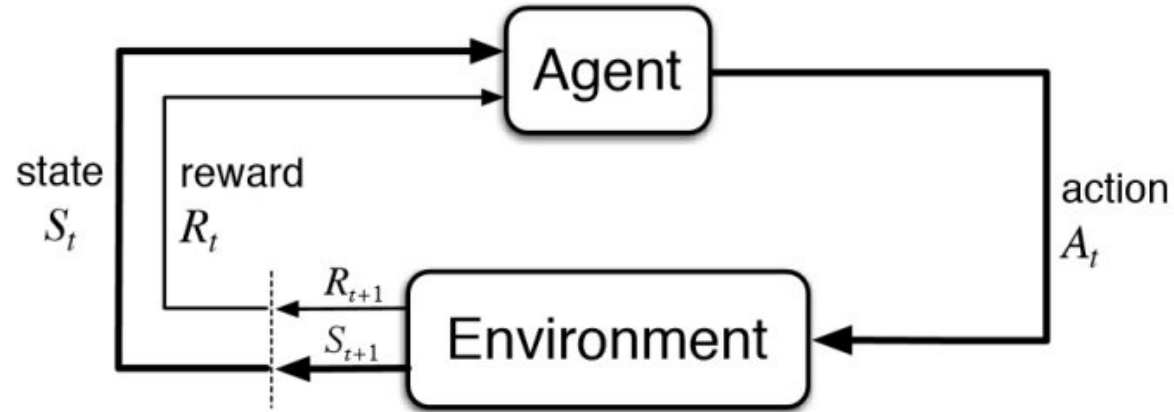


$$\max_{\theta} \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^T r(s_t, a_t) \right]$$

Main difference: model known vs unknown

Minor differences: Cost vs reward, discrete vs continuous time

How should we optimize this objective?



$$\max_{\theta} \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^T r(s_t, a_t) \right]$$

Gradient Ascent

Dynamic Programming

Model-Based Optimization

Each method has its own +/-

Lecture outline

Deriving the Policy Gradient

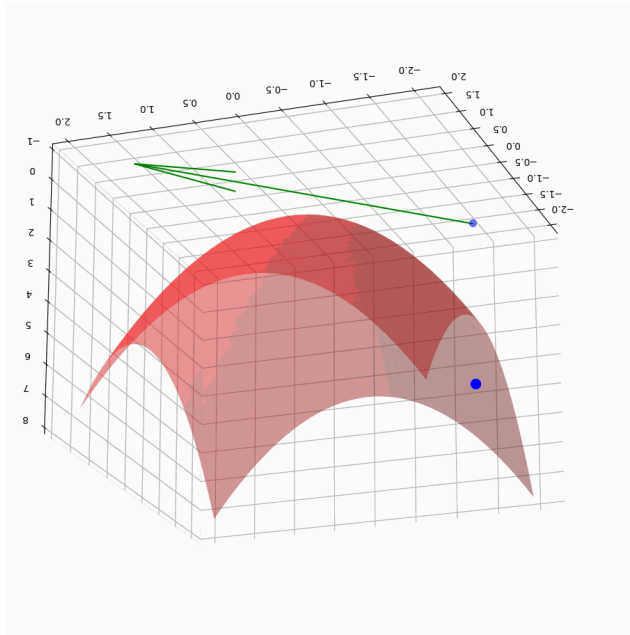
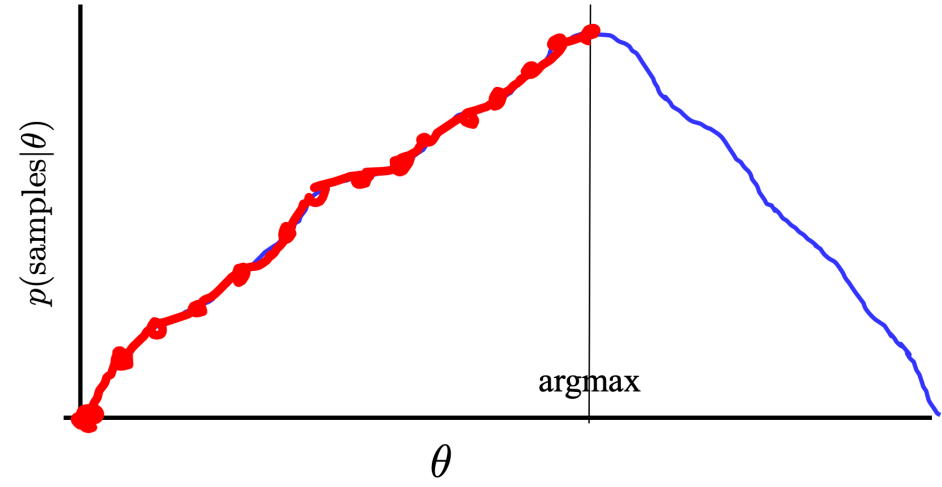
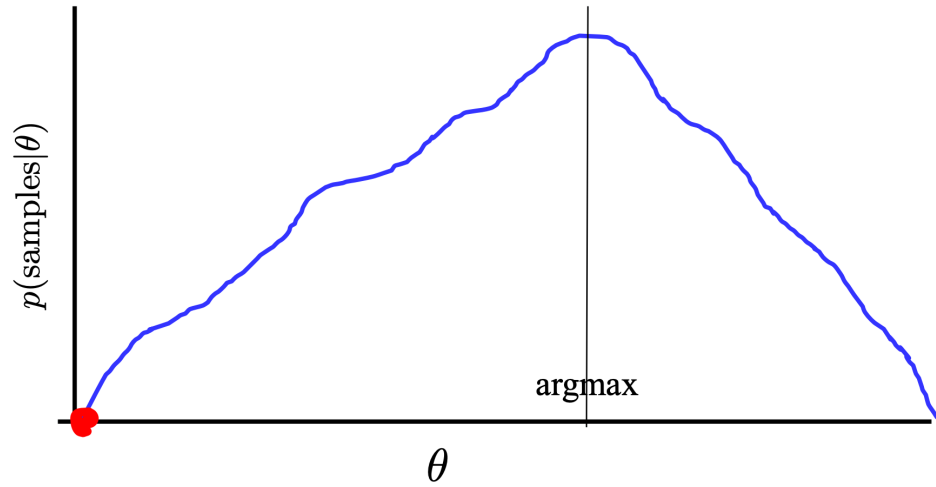


What makes the Policy Gradient Challenging? - Variance



What makes the Policy Gradient Challenging? – Covariant Parameterization

Gradient Ascent



Simple view – move the parameters in the direction of the gradient of the objective

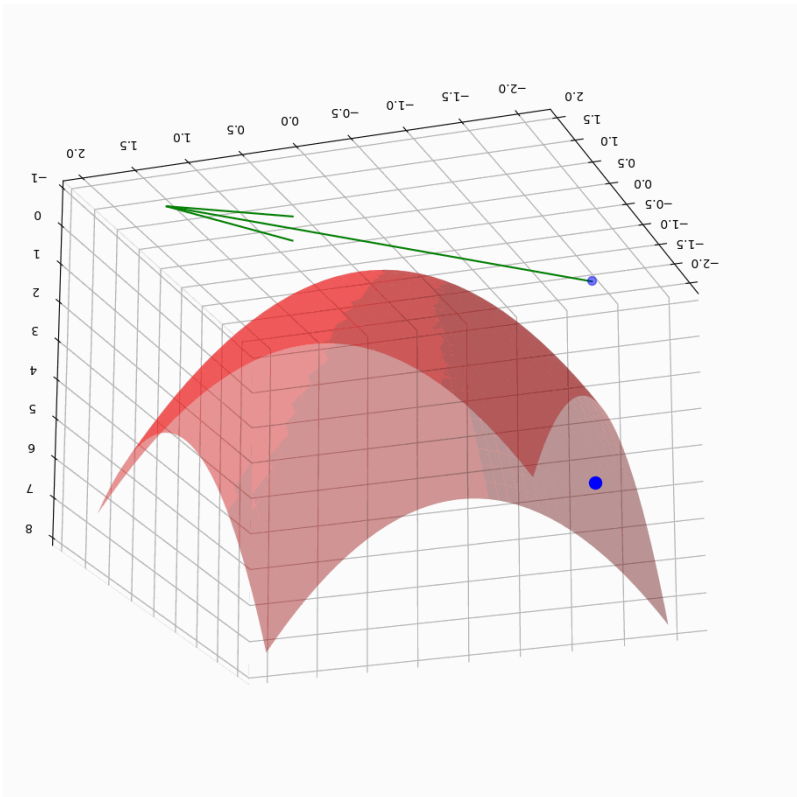
$$\theta_{i+1} = \theta_i + \alpha \nabla_{\theta} J(\theta) |_{\theta=\theta_i}$$

More later: can be derived as steepest ascent in Euclidean norm

Gradient Ascent for Supervised Learning

Recall our imitation learning objective $\arg \max_{\theta} \mathbb{E}_{(s^*, a^*) \sim \mathcal{D}} [\log \pi_{\theta}(a^* | s^*)]$

Let's apply gradient ascent



$$\nabla_{\theta} \mathbb{E}_{(s^*, a^*) \sim \mathcal{D}} [\log \pi_{\theta}(a^* | s^*)]$$

$$\nabla_{\theta} \int p(s^*, a^*) \log \pi_{\theta}(a^* | s^*) ds^* da^*$$

$$\int p(s^*, a^*) \nabla_{\theta} \log \pi_{\theta}(a^* | s^*) ds^* da^*$$

$$\mathbb{E}_{(s^*, a^*) \sim \mathcal{D}} [\nabla_{\theta} \log \pi_{\theta}(a^* | s^*)]$$

Compute gradient and average

Ok let's do gradient ascent for the RL objective

$$\max_{\theta} \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^T r(s_t, a_t) \right]$$
$$= \int p_{\theta}(\tau) R(\tau) d\tau$$



REINFORCE gradient descent (RL)

$$\nabla_{\theta} \mathbb{E}_{x \sim p_{\theta}(x)} [f(x)]$$

(Cannot simply compute average of expectation)

Standard gradient descent (supervised learning)

Gradient wrt expectation variable, not of integrand!

(Whiteboard)

$$\nabla_{\theta} \mathbb{E}_{x \sim g(x)} [f_{\theta}(x)]$$

(Gradient passes inside the expectation – compute gradient and average)

Taking the gradient of sum of rewards

$$\max_{\theta} \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^T r(s_t, a_t) \right]$$

Let's take the gradient of this objective

$$J(\theta) = \int p_{\theta}(\tau) R(\tau) d(\tau)$$

Let's think about this from the trajectory view

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} \int p_{\theta}(\tau) R(\tau) d(\tau)$$

We need to express this in a way that we can evaluate with expectations

$$= \int \nabla_{\theta} p_{\theta}(\tau) R(\tau) d(\tau) = \int \frac{p_{\theta}(\tau)}{p_{\theta}(\tau)} \nabla_{\theta} p_{\theta}(\tau) R(\tau) d(\tau)$$

$$= \int p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau) R(\tau) d(\tau) = \mathbb{E}_{p_{\theta}(\tau)} [\nabla_{\theta} \log p_{\theta}(\tau) R(\tau)]$$

$$\frac{d \log(x)}{d\theta} = \frac{d \log(x)}{dx} \frac{dx}{d\theta} = \frac{1}{x} \frac{dx}{d\theta}$$

Use chain rule

REINFORCE trick

Taking the gradient of return

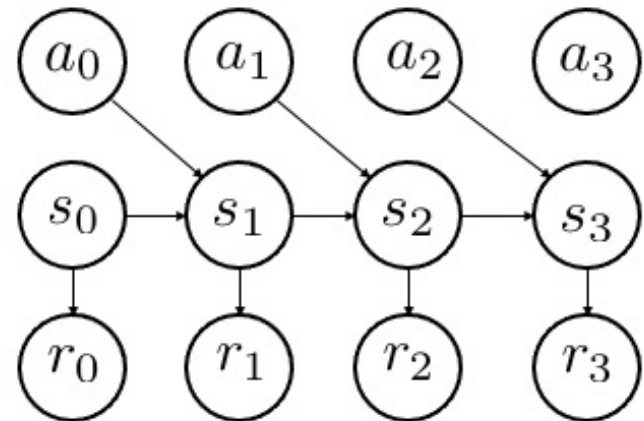
Initial State

Dynamics

Policy

$$p_{\theta}(\tau) = p(s_0) \prod_{t=0}^{T-1} p(s_{t+1} | s_t, a_t) \pi(a_t | s_t)$$

(Ancestral sampling)



$$\log p_{\theta}(\tau) = \log p(s_0) + \sum_{t=0}^{T-1} \log p(s_{t+1} | s_t, a_t) + \log \pi(a_t | s_t)$$

$$\nabla_{\theta} \log p_{\theta}(\tau) = \nabla_{\theta} \log p(s_0) + \sum_{t=0}^{T-1} \nabla_{\theta} \log p(s_{t+1} | s_t, a_t) + \nabla_{\theta} \log \pi(a_t | s_t)$$

$$\nabla_{\theta} \log p_{\theta}(\tau) = \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi(a_t | s_t)$$

Model Free!!

Taking the gradient of return

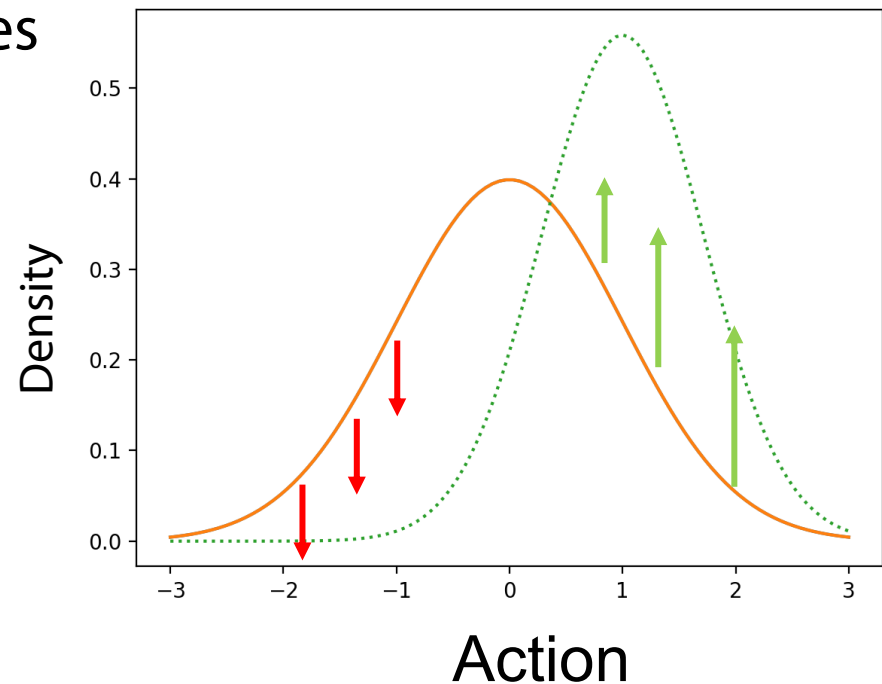
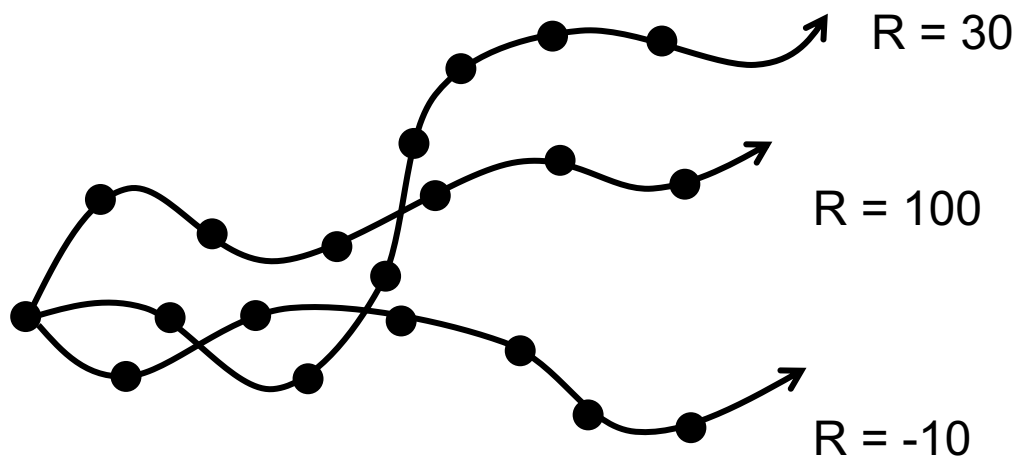
$$\begin{aligned} \nabla_{\theta} J(\theta) &= \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\nabla_{\theta} \log p_{\theta}(\tau) \sum_{t=0}^T r(s_t, a_t) \right] \\ \nabla_{\theta} J(\theta) &= \mathbb{E}_{\substack{s_0 \sim p(s_0) \\ s_{t+1} \sim p(s_{t+1} | s_t, a_t) \\ a_t \sim \pi(a_t | s_t)}} \left[\sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \sum_{t'=0}^T r(s_{t'}, a_{t'}) \right] \\ &\approx \frac{1}{N} \sum_{i=0}^N \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t^i | s_t^i) \sum_{t'=0}^T r(s_{t'}^i, a_{t'}^i) \quad (\text{approximating using samples}) \end{aligned}$$

(Monte-Carlo approximation)

What does this mean?

$$\nabla_{\theta} J(\theta) = \int p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau) d\tau \approx \frac{1}{N} \sum_{i=0}^N \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t^i | s_t^i) \sum_{t'=0}^T r(s_{t'}^i, a_{t'}^i)$$

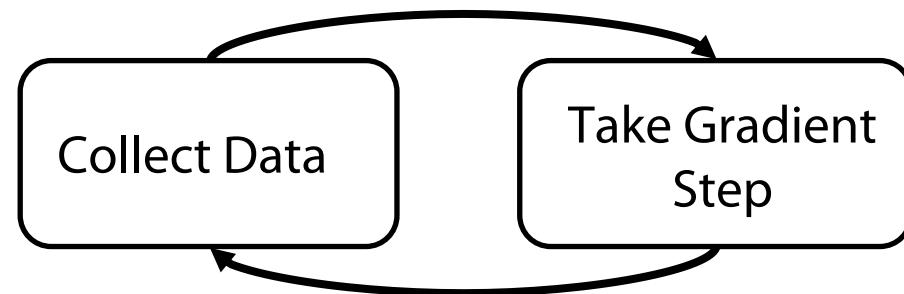
Increase the likelihood of actions in high return trajectories



Resulting Algorithm (REINFORCE)

$$\nabla_{\theta} J(\theta) = \int p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau) d\tau$$

$$\theta_{i+1} = \theta_i + \alpha \nabla_{\theta} J(\theta) |_{\theta=\theta_i}$$



REINFORCE algorithm:

On-policy



1. sample $\{\tau^i\}$ from $\pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)$ (run it on the robot)
2. $\nabla_{\theta} J(\theta) \approx \sum_i (\sum_t \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t^i | \mathbf{s}_t^i)) (\sum_t r(\mathbf{s}_t^i, \mathbf{a}_t^i))$
3. $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$

How is this related to supervised learning?

Reinforcement Learning

$$\nabla_{\theta} J(\theta) = \int p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau) d\tau$$

$$\approx \frac{1}{N} \sum_{i=0}^N \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t^i | s_t^i) \sum_{t'=0}^T r(s_{t'}^i, a_{t'}^i)$$

Supervised Learning

$$\max_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\log p_{\theta}(y|x)]$$

$$\approx \frac{1}{N} \sum_i \nabla_{\theta} \log p_{\theta}(y^i | x^i)$$

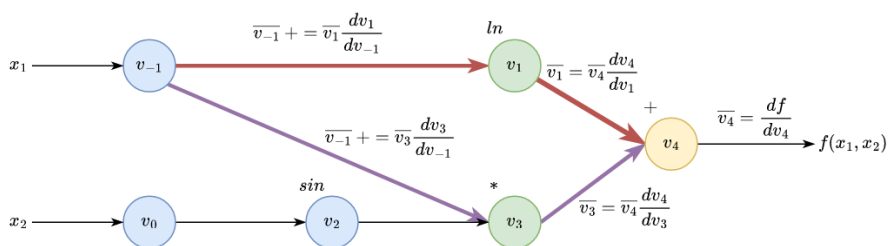
PG = select good data + increase likelihood of selected data

How do we implement this?

REINFORCE algorithm:

1. sample $\{\tau^i\}$ from $\pi_\theta(\mathbf{a}_t|\mathbf{s}_t)$ (run it on the robot)
2. $\nabla_\theta J(\theta) \approx \sum_i (\sum_t \nabla_\theta \log \pi_\theta(\mathbf{a}_t^i|\mathbf{s}_t^i)) (\sum_t r(\mathbf{s}_t^i, \mathbf{a}_t^i))$
3. $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$

$$\nabla_\theta J(\theta) = \int p_\theta(\tau) \nabla_\theta \log p_\theta(\tau) d\tau \approx \frac{1}{N} \sum_{i=0}^N \sum_{t=0}^T \nabla_\theta \log \pi_\theta(a_t^i | s_t^i) \sum_{t'=0}^T r(s_{t'}^i, a_{t'}^i)$$



Compute gradients with autodiff

Sum up rewards in a trajectory



How do we implement this?

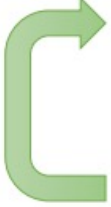
Maximum likelihood:

```
# Given:  
# actions - (N*T) x Da tensor of actions  
# states - (N*T) x Ds tensor of states  
# Build the graph:  
logits = policy.predictions(states) # This should return (N*T) x Da tensor of action logits  
negative_likelihoods = tf.nn.softmax_cross_entropy_with_logits(labels=actions, logits=logits)  
loss = tf.reduce_mean(negative_likelihoods)  
gradients = loss.gradients(loss, variables)
```

^Standard maximum likelihood training

How do we implement this?

REINFORCE algorithm:

- 
1. sample $\{\tau^i\}$ from $\pi_\theta(\mathbf{a}_t|\mathbf{s}_t)$ (run it on the robot)
 2. $\nabla_\theta J(\theta) \approx \sum_i (\sum_t \nabla_\theta \log \pi_\theta(\mathbf{a}_t^i|\mathbf{s}_t^i)) (\sum_t r(\mathbf{s}_t^i, \mathbf{a}_t^i))$
 3. $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$

Policy gradient:

```
# Given:
# actions - (N*T) x Da tensor of actions
# states - (N*T) x Ds tensor of states
# q_values - (N*T) x 1 tensor of estimated state-action values → Sum of rewards
# Build the graph:
logits = policy.predictions(states) # This should return (N*T) x Da tensor of action logits
negative_likelihoods = tf.nn.softmax_cross_entropy_with_logits(labels=actions, logits=logits)
weighted_negative_likelihoods = tf.multiply(negative_likelihoods, q_values)
loss = tf.reduce_mean(weighted_negative_likelihoods)
gradients = loss.gradients(loss, variables)
```

Formalizes the notion of trial and error

How do we implement this?

$$\nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

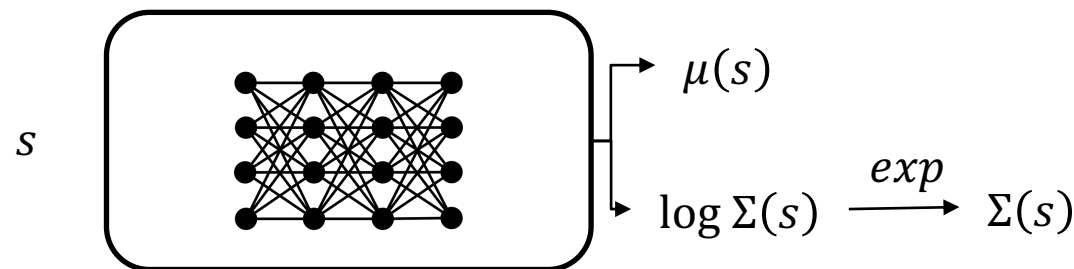
Let's try it for a Gaussian

$$\pi(\mathbf{a} | \mathbf{s})$$

$$= \pi(\mathbf{a} | \boldsymbol{\mu}_{\theta}(\mathbf{s}), \boldsymbol{\Sigma}_{\theta}(\mathbf{s}))$$

$$= \pi(\mathbf{a} | \boldsymbol{\mu}_{\theta}(\mathbf{s}), \boldsymbol{\Sigma}_{\theta}(\mathbf{s})) = \frac{1}{\sqrt{(2\pi)^k |\boldsymbol{\Sigma}_{\theta}(\mathbf{s})|}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_{\theta}(\mathbf{s}))^{\top} \boldsymbol{\Sigma}_{\theta}(\mathbf{s})^{-1}(\mathbf{x} - \boldsymbol{\mu}_{\theta}(\mathbf{s}))\right)$$

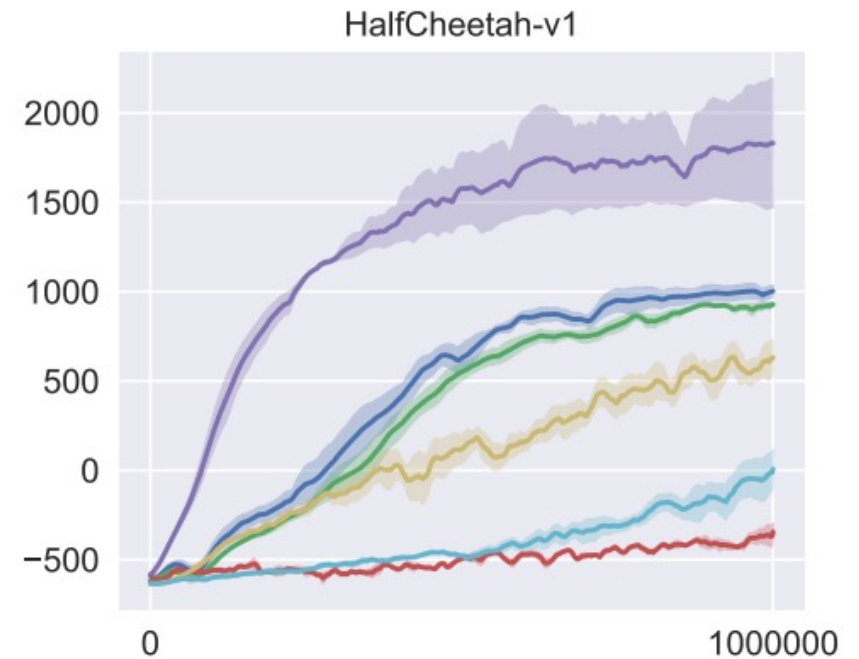
Similar for categorical or other distributions



Easier for distributions where likelihood can be expressed symbolically

Does this work?

Comparison of
RL algorithms
in Humanoid-v2
using CleanRL



Kind of?

Lecture outline

Recap: Imitation Learning via Behavior Cloning



Challenges with DAgger



Modern Takes on Imitation Learning



Policy Gradient

Fin.

