

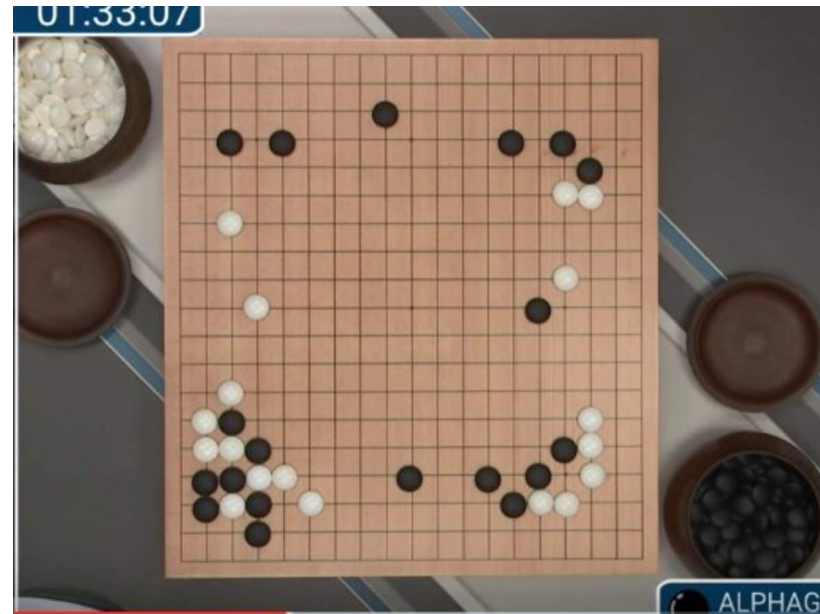


Reinforcement Learning

Autumn 2024

Abhishek Gupta

TA: Jacob Berg



Lecture outline

Recap: Imitation Learning + Why it is hard



Multimodality and Underfitting in Imitation



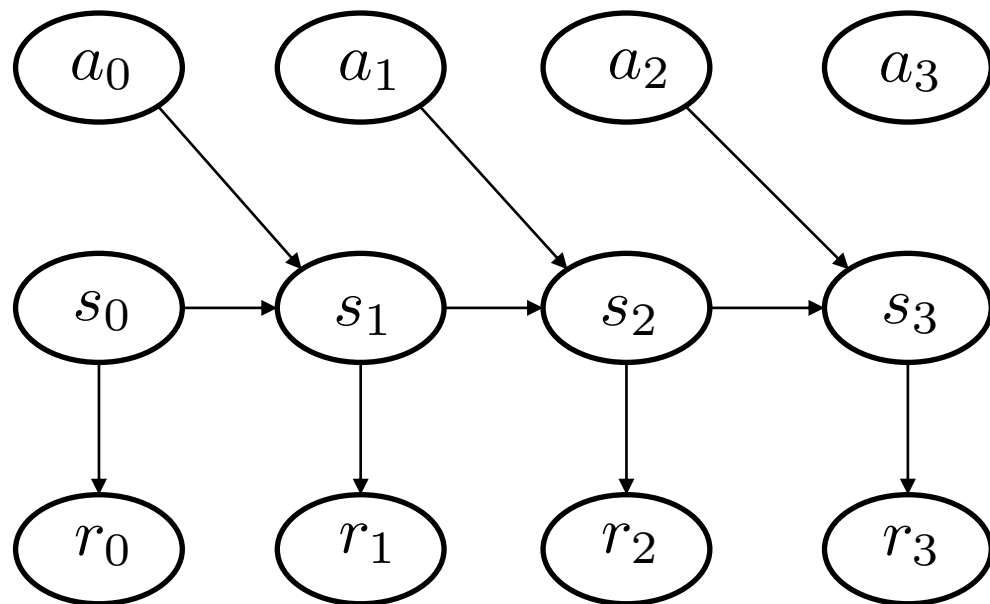
Compounding Error in Imitation



Frontiers in Imitation

Framework for RL - Markov Decision Process

Augment Markov chain with rewards and actions



States: \mathcal{S}

Initial state dist: $\rho_0(s)$

Actions: \mathcal{A}

Discount: γ

Rewards: \mathcal{R}

Transition Dynamics - $p(s_{t+1}|s_t, a_t)$

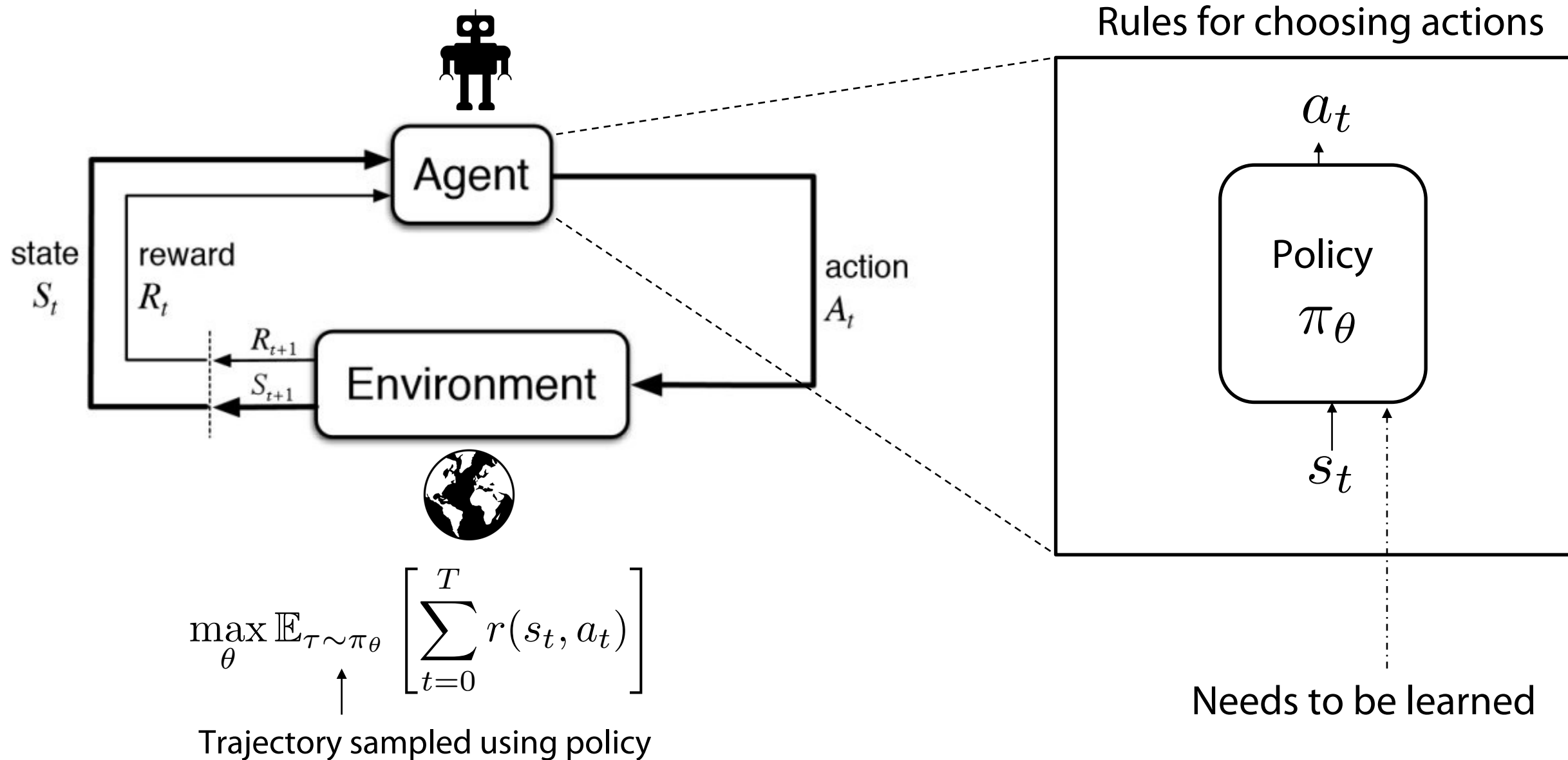
Markov property

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

Trajectory

$$\tau = (s_0, a_0, r_0, s_1, a_1, r_1, \dots, s_T, a_T, r_T)$$

Reinforcement Learning Formalism



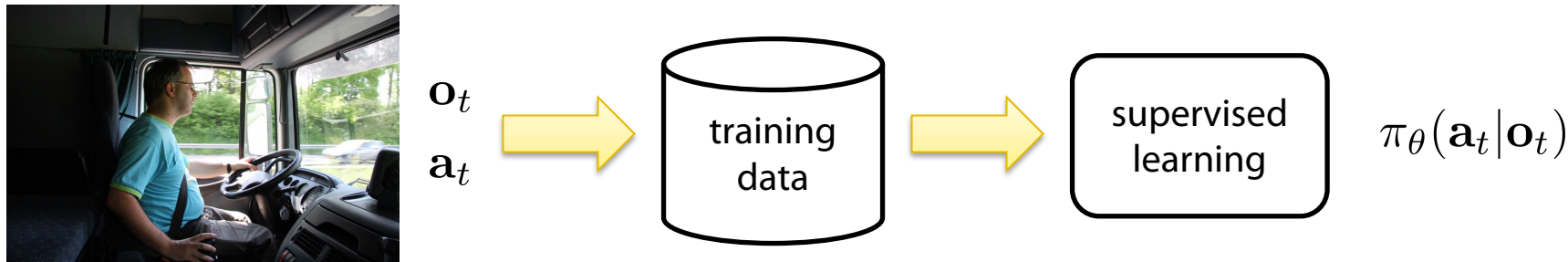
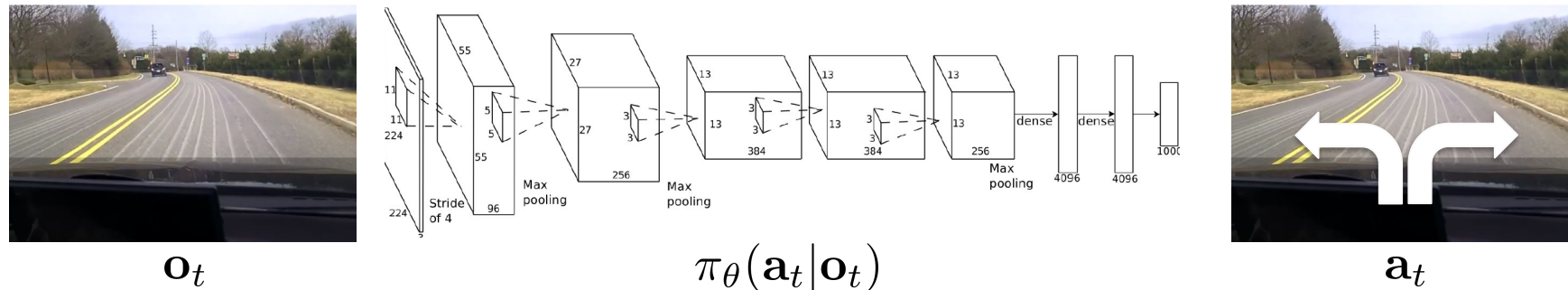
Idea 1: Imitation Learning via Supervised Learning

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! \rightarrow Behavior Cloning



Idea 1: Imitation Learning via Supervised Learning

Given: Demonstrations of optimal behavior

Goal: Train a policy to mimic the demonstrator

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

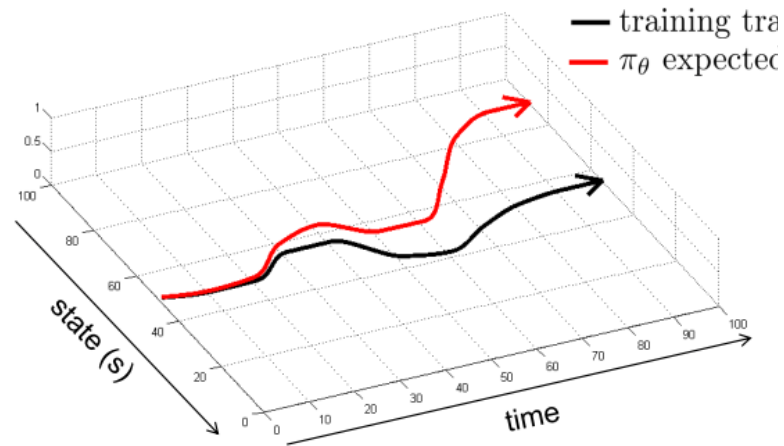
Discrete vs continuous

```
if isinstance(env.action_space, gym.spaces.Box):
    criterion = nn.MSELoss()
else:
    criterion = nn.CrossEntropyLoss()
# Extract initial policy
model = student.policy.to(device)
def train(model, device, train_loader, optimizer):
    model.train()
    for batch_idx, (data, target) in enumerate(train_loader):
        data, target = data.to(device), target.to(device)
        optimizer.zero_grad()
        if isinstance(env.action_space, gym.spaces.Box):
            if isinstance(student, (A2C, PPO)):
                action, _, _ = model(data)
            else:
                action = model(data)
            action_prediction = action.double()
        else:
            dist = model.get_distribution(data)
            action_prediction = dist.distribution.logits
            target = target.long()
        loss = criterion(action_prediction, target)
        loss.backward()
        optimizer.step()
```

Maximum likelihood

So does behavior cloning really work?

- Imitation Learning \neq Supervised Learning



Compounding error!

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

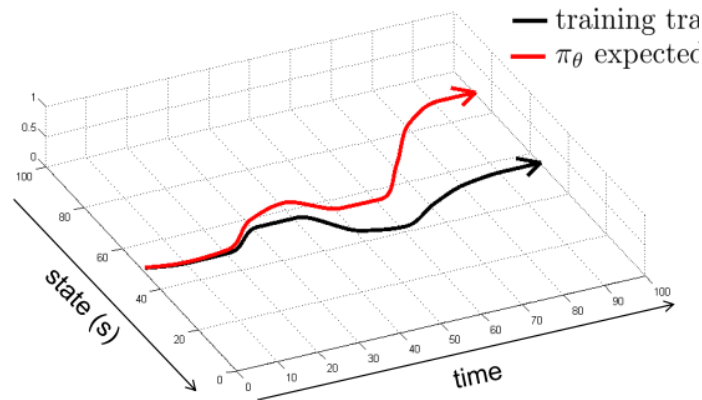
$$\mathbb{E}_{(s, a) \sim \rho(\pi)} [\mathbf{1}(a = a^*)]$$



Not the same!

How well does BC do?: Intuition

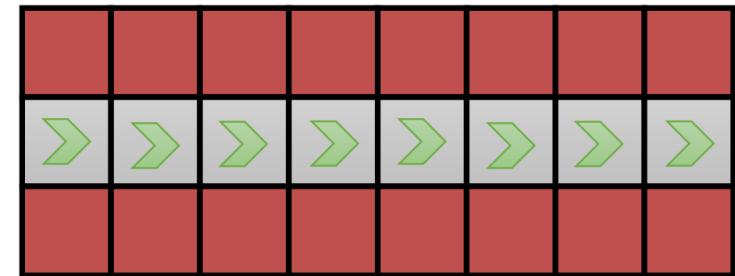
Behavior cloning has quadratically compounding error



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

Horizon H

If you fall off,
assume the worst



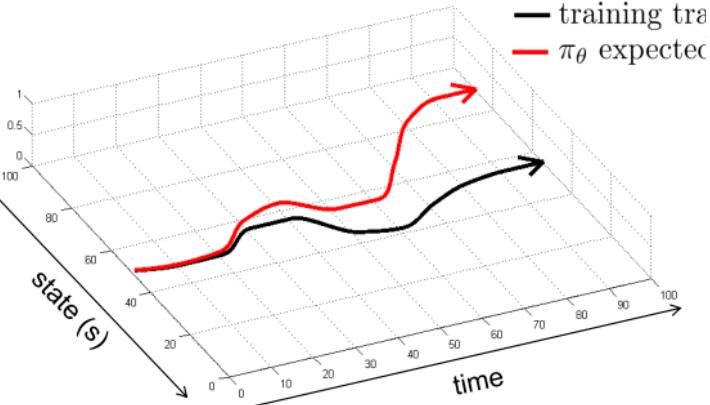
$$\underbrace{\mathbb{E} \left[\sum_t c(s_t, a_t) \right]}_{O(\epsilon H^2)} \leq \epsilon H + \dots + \dots$$

Union bound

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

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

Lecture outline

Recap: Imitation Learning + Why it is hard



Multimodality and Underfitting in Imitation



Compounding Error in Imitation

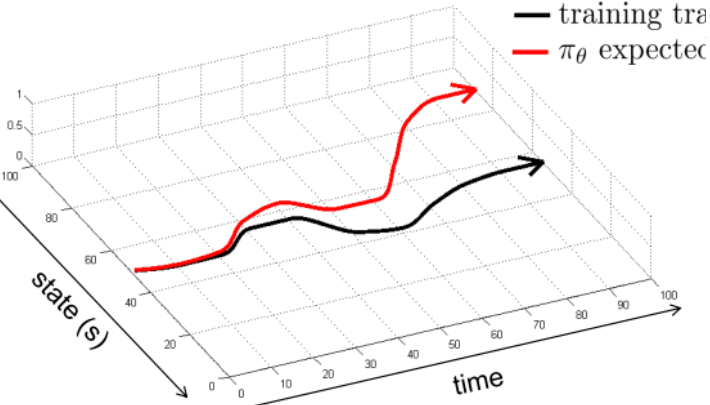


Frontiers in Imitation

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

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

But won't a bigger neural net just solve this?

- Behavior cloning can underfit the data

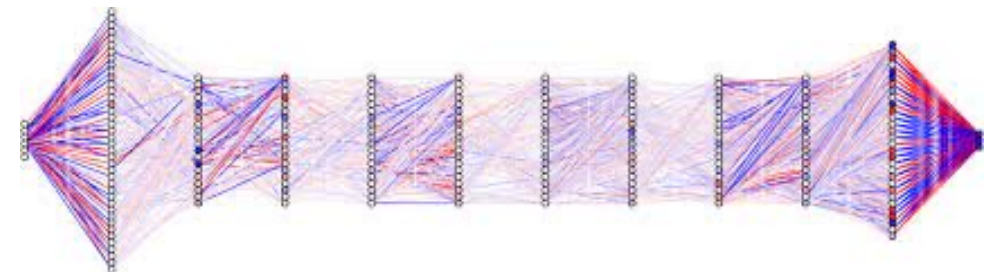
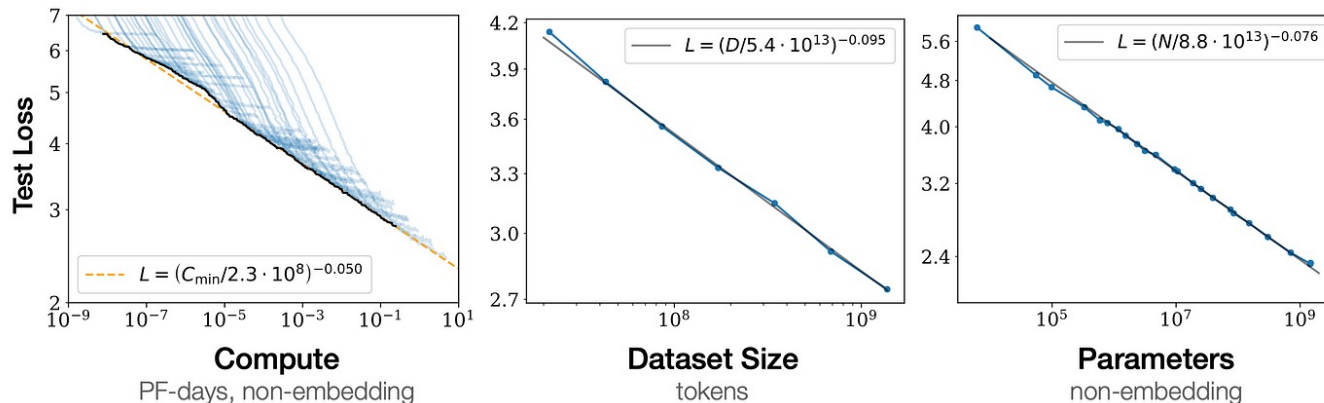
$$\sum_t \mathbb{E}_{(s_t, a_t) \sim p_{\pi_\theta}(s_t, a_t)} [c(s_t, a_t)] \leq O(\epsilon H^2)$$

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

for $s_t \sim p_{\text{train}}(s_t)$

May not be able to satisfy this

Q: won't a bigger model just solve the problem?

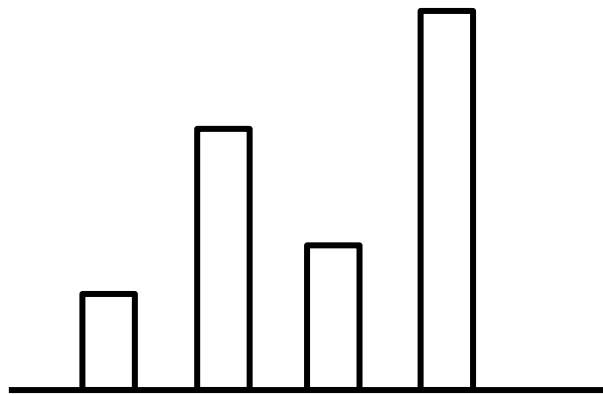


Kind of, but there's a fundamental problem!

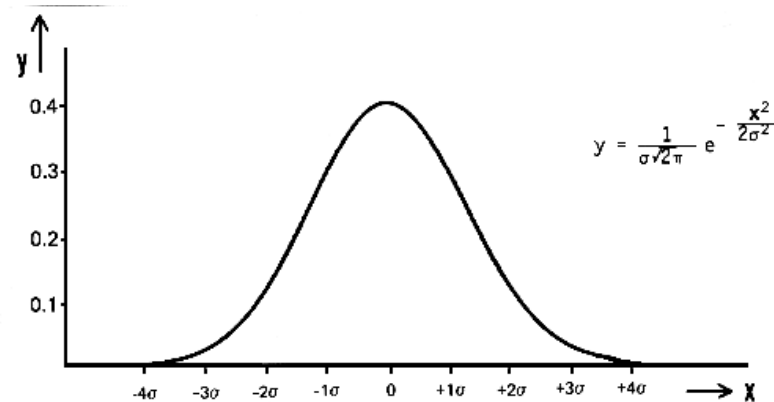
Distributional Expressivity

- Policy expressivity is a combination of expressivity of the function approximator and of the distribution family

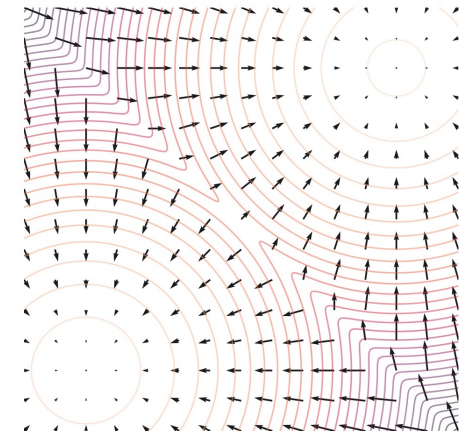
Categorical



Gaussian



Diffusion policy

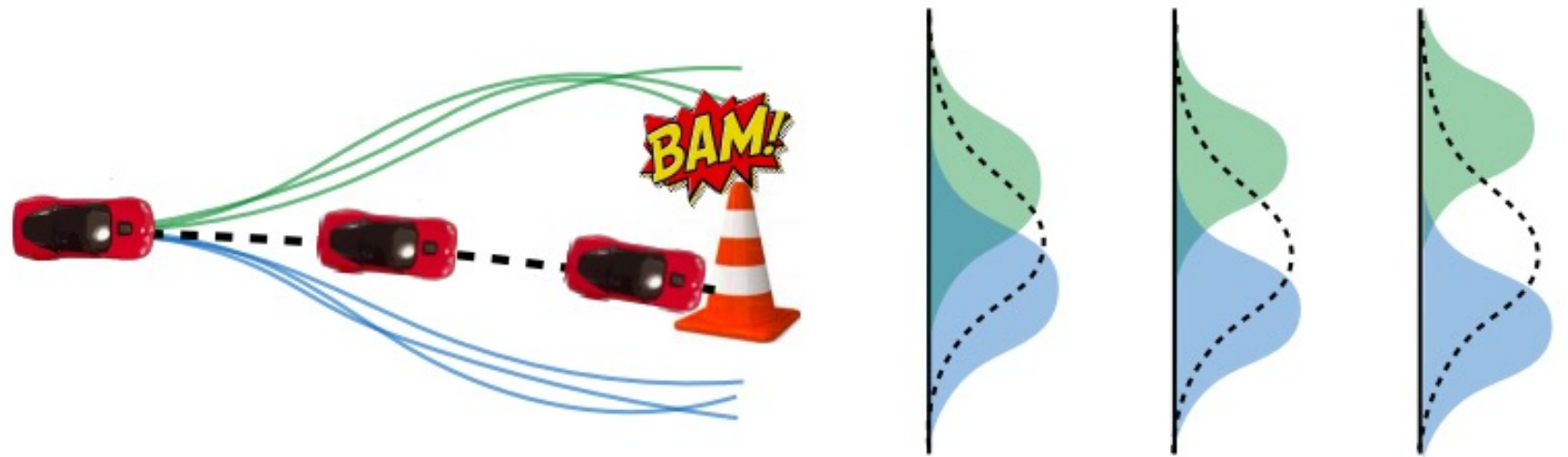
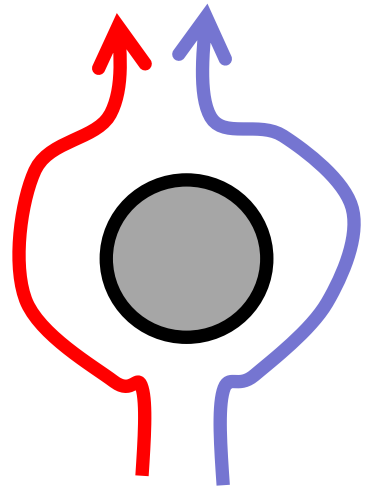


Tradeoff between expressivity and tractability

How does this reflect on imitation learning?

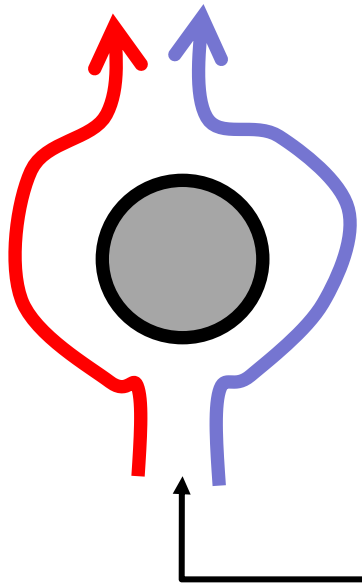
Let us consider a case with Gaussian policy

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



A combination of distributional expressivity and objective lead to mode averaging

Let's take a closer look at the objective



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

$$\max_{\theta} \mathbb{E}_{s^* \sim p_{\pi_e}(\cdot)} [\mathbb{E}_{a^* \sim \pi_e(\cdot | s^*)} [\log \pi_{\theta}(a^* | s^*) - \log \pi_e(a^* | s^*)]]$$

$$\min_{\theta} \mathbb{E}_{s^* \sim p_{\pi_e}(\cdot)} \left[\mathbb{E}_{a^* \sim \pi_e(\cdot | s^*)} \left[\log \frac{\pi_e(a^* | s^*)}{\pi_{\theta}(a^* | s^*)} \right] \right] = \mathbb{E}_{s^* \sim p_{\pi_e}(\cdot)} [D_{\text{KL}}(\pi_e(\cdot | s^*) || \pi_{\theta}(\cdot | s^*))]$$

Leads to mode averaging

Forward KL divergence

One instance of a broader class of divergences – f divergences $D_f(p(x), q(x)) = \mathbb{E}_{q(x)} \left[f \left(\frac{p(x)}{q(x)} \right) \right]$

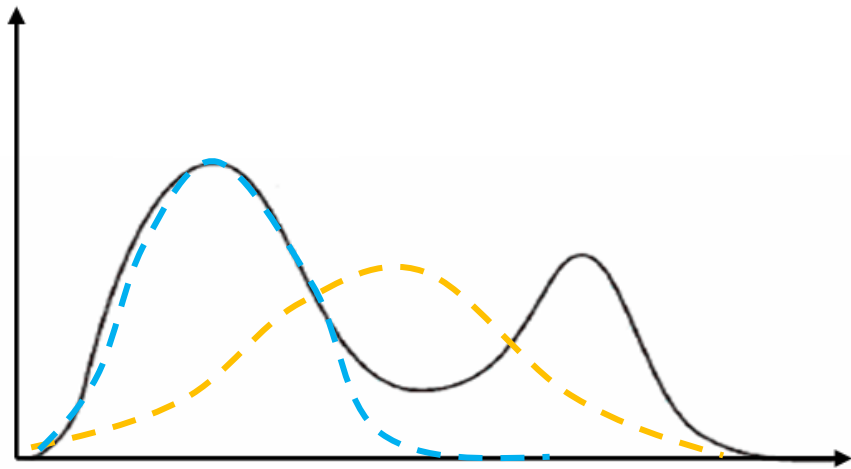
Effects of choice of f-divergence on behavior

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 alternative f-divergences: Reverse KL

- Reverse KL helps, is mode seeking $D_{\text{RKL}}(\pi_e(\cdot|s^*), \pi^\theta(\cdot|s^*)) = \mathbb{E}_{\pi^\theta(\cdot|s^*)} \left[\log \left(\frac{\pi^\theta(\cdot|s^*)}{\pi_e(\cdot|s^*)} \right) \right]$
- Challenge – requires known expert likelihood
- We need a sample based estimate!

Imitation Learning as f-Divergence Minimization

Liyiming Ke¹, Sanjiban Choudhury¹, Matt Barnes¹, Wen Sun², Gilwoo Lee¹,
and Siddhartha Srinivasa¹

Go read this!

$$\min_{\theta} \mathbb{E}_{\pi^\theta(\cdot|s^*)} \left[\log \left(\frac{\pi^\theta(\cdot|s^*)}{\pi_e(\cdot|s^*)} \right) \right] \longleftrightarrow \min_{\theta} \max_{\phi} \mathbb{E}_{a \sim \pi^\theta(\cdot|s^*)} [\phi(a)] - \mathbb{E}_{a \sim \pi_e(\cdot|s^*)} [f^*(\phi(a))]$$

(Intractable) (Tractable – GAN style optimization)

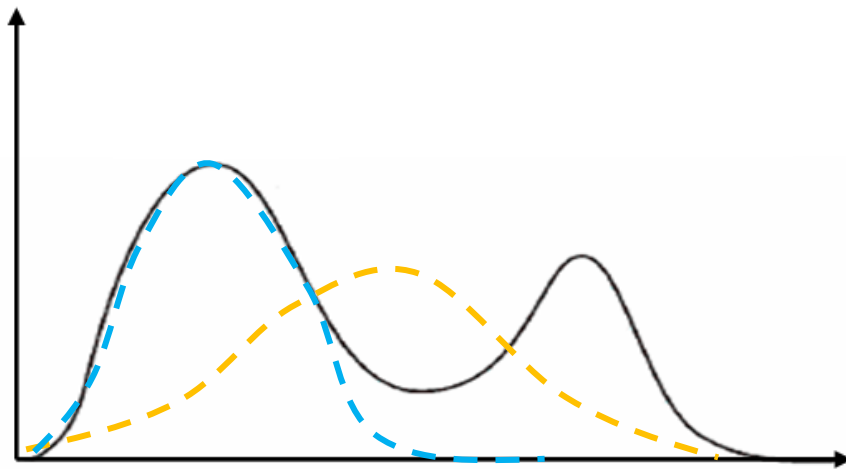
Effects of choice of f-divergence on behavior

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?

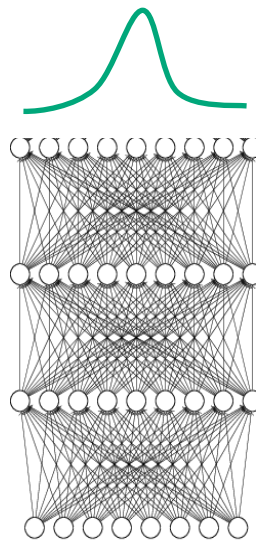
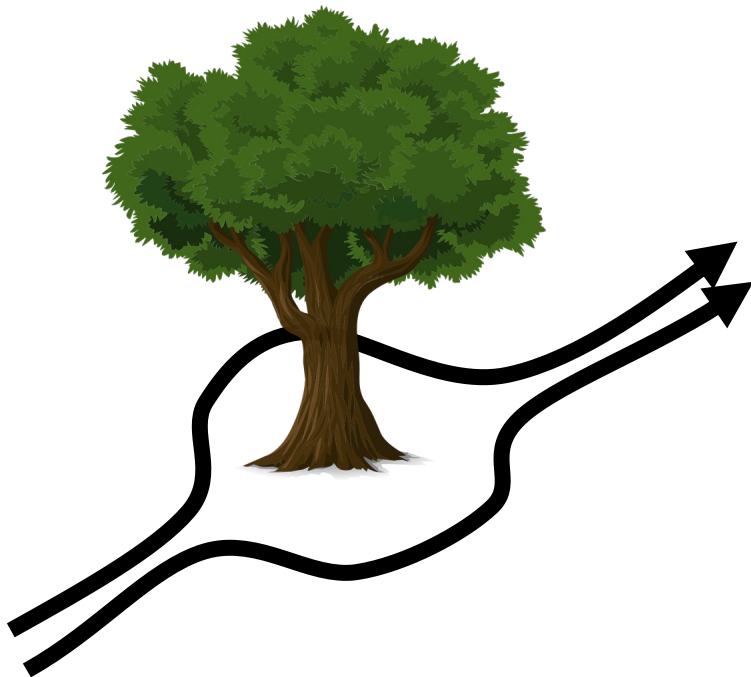
Use a different f-divergence!
(Change f)

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



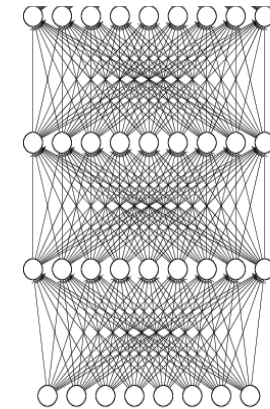
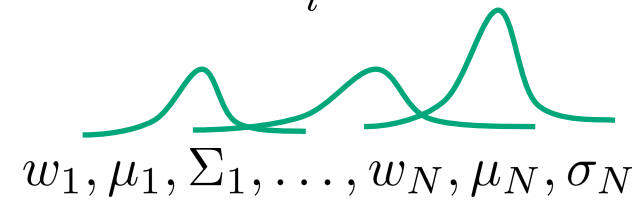
1. Output mixture of Gaussians
2. Latent variable models
3. Autoregressive discretization
4. Diffusion models
5. ...



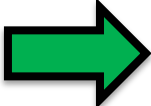
Why might we fail to fit the expert?

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

$$\pi(\mathbf{a}|\mathbf{o}) = \sum_i w_i \mathcal{N}(\mu_i, \Sigma_i)$$



Why might we fail to fit the expert?

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

Why does this work?

first step: $p(a_{t,0}|\mathbf{s}_t)$

second step: $p(a_{t,1}|\mathbf{s}_t, a_{t,0})$

third step: $p(a_{t,2}|\mathbf{s}_t, a_{t,0}, a_{t,1})$

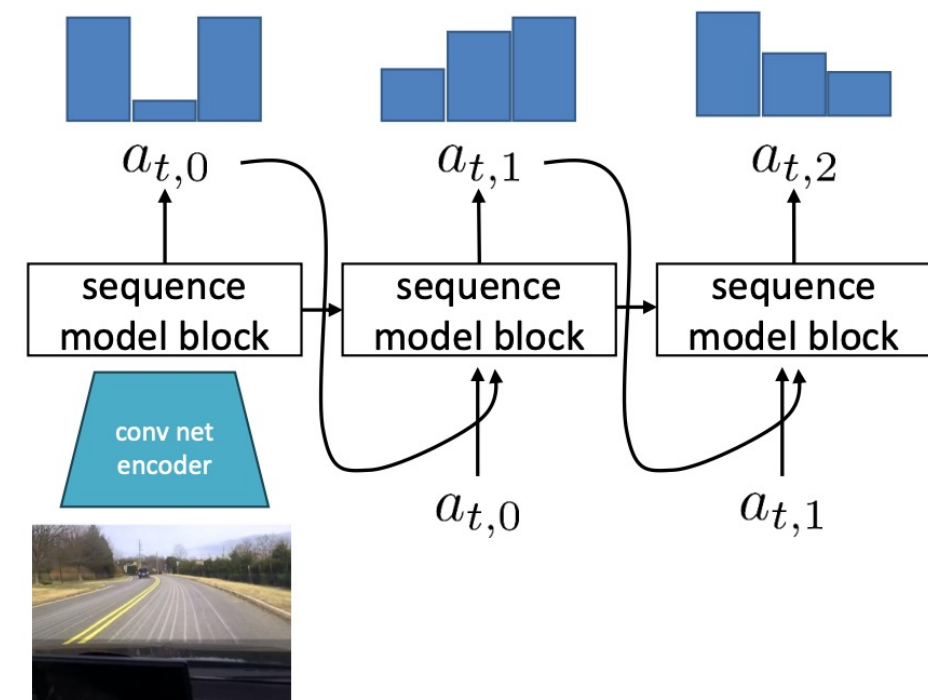
$$p(a_{t,2}|\mathbf{s}_t, a_{t,0}, a_{t,1})p(a_{t,1}|\mathbf{s}_t, a_{t,0})p(a_{t,0}|\mathbf{s}_t)$$

$$= p(a_{t,0}, a_{t,1}, a_{t,2}|\mathbf{s}_t)$$

$$= p(\mathbf{a}_t|\mathbf{s}_t)$$

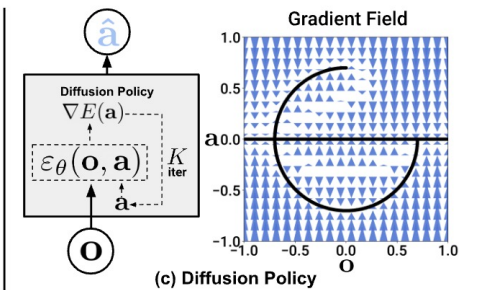
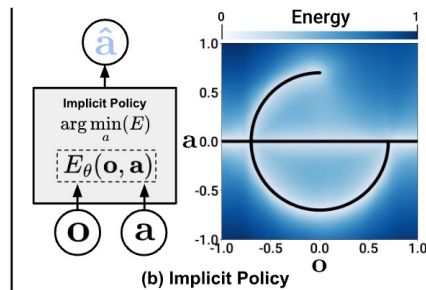
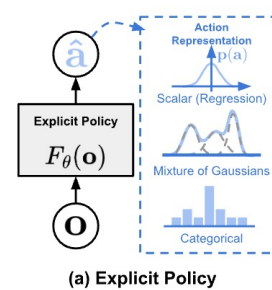
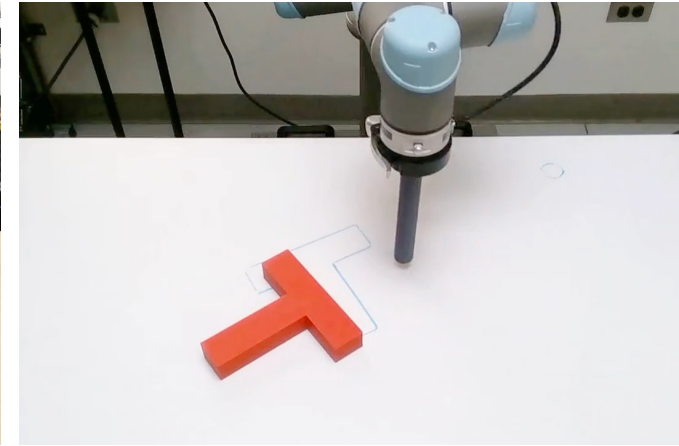
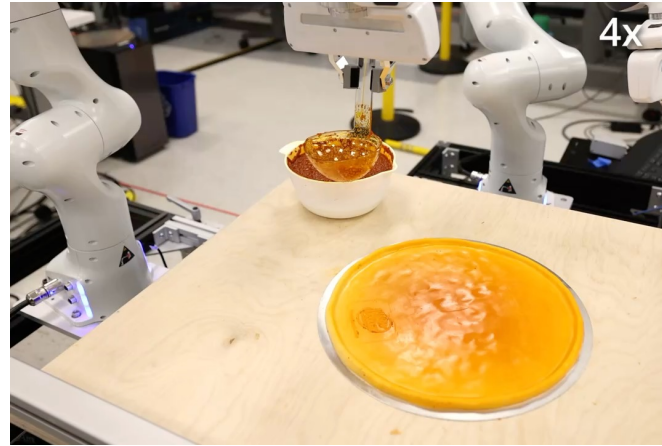
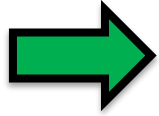
$$\mathbf{a}_t = \begin{pmatrix} 0.1 \\ 1.2 \\ -0.3 \end{pmatrix} \begin{matrix} a_{t,0} \\ a_{t,1} \\ a_{t,2} \end{matrix}$$

use LSTM or
Transformer



Why might we fail to fit the expert?

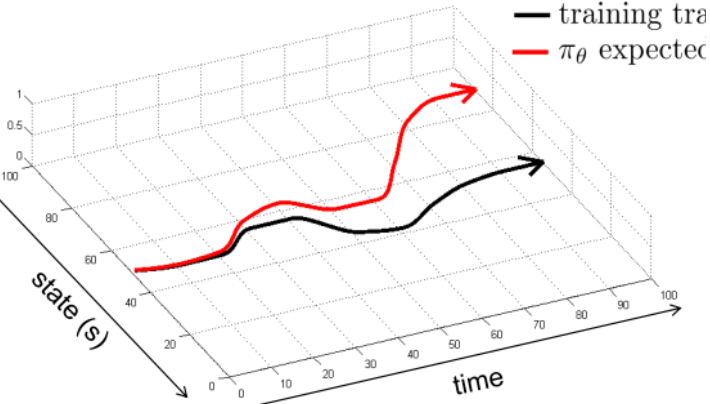
1. Output mixture of Gaussians
2. Latent variable models
3. Autoregressive discretization
4. Diffusion models
5. ...



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

Richer policy class Alternative Divergence

Lecture outline

Recap: Imitation Learning + Why it is hard



Multimodality and Underfitting in Imitation



Compounding Error in Imitation

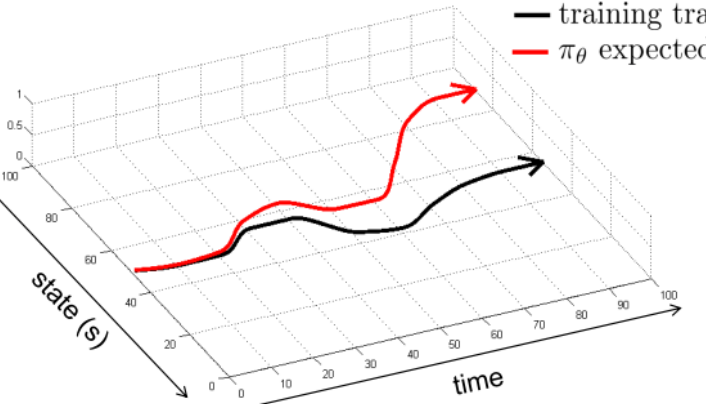


Frontiers in Imitation

Let's try and understand where the problem lies?

Behavior cloning has challenges in both theory and practice

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



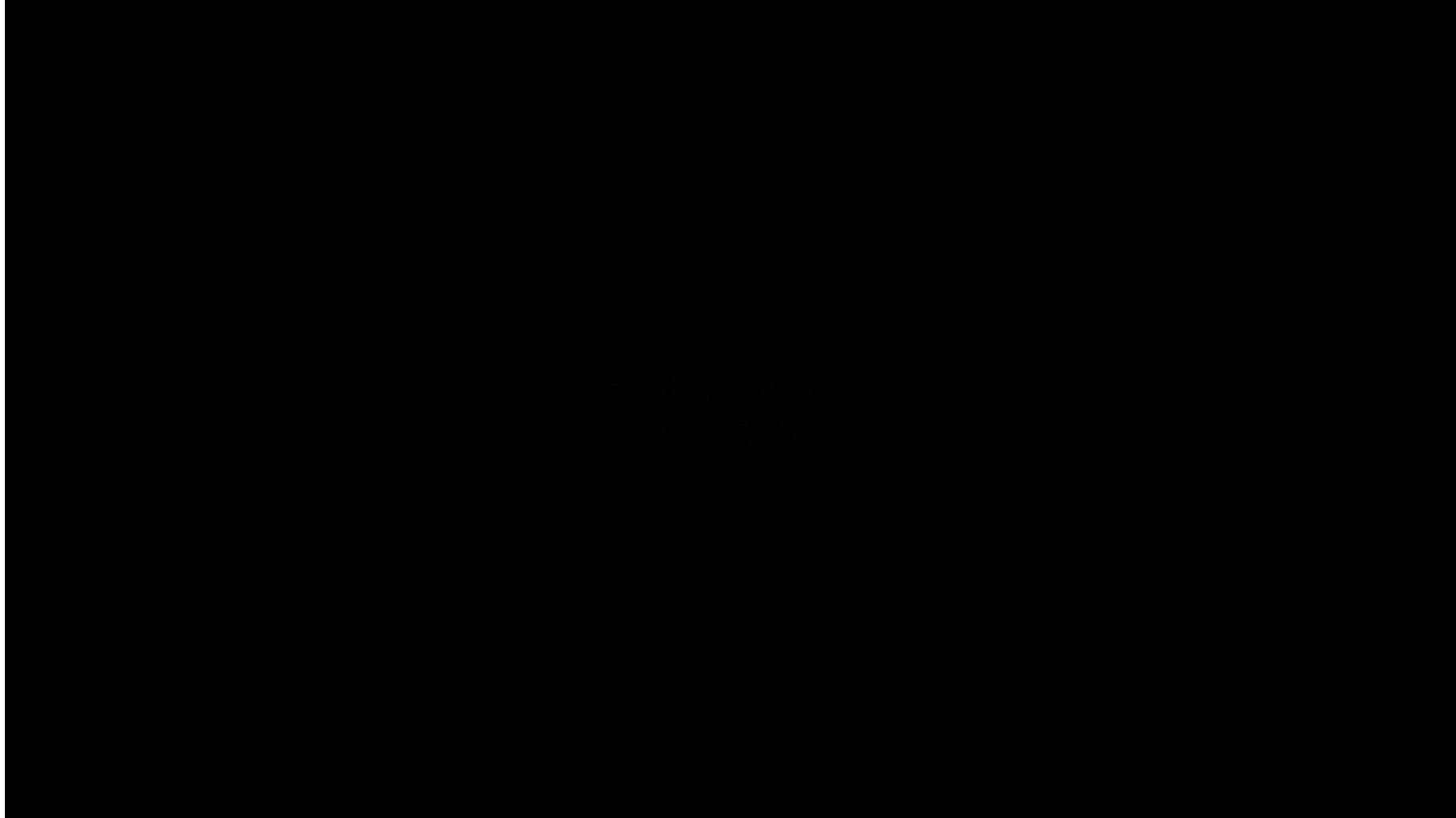
Underfitting

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

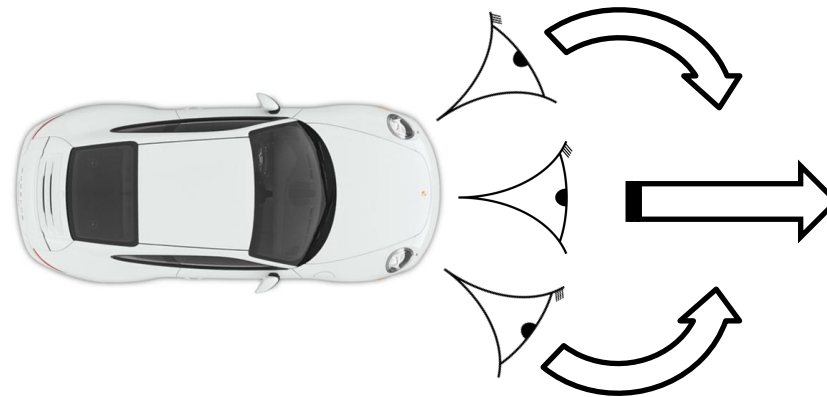
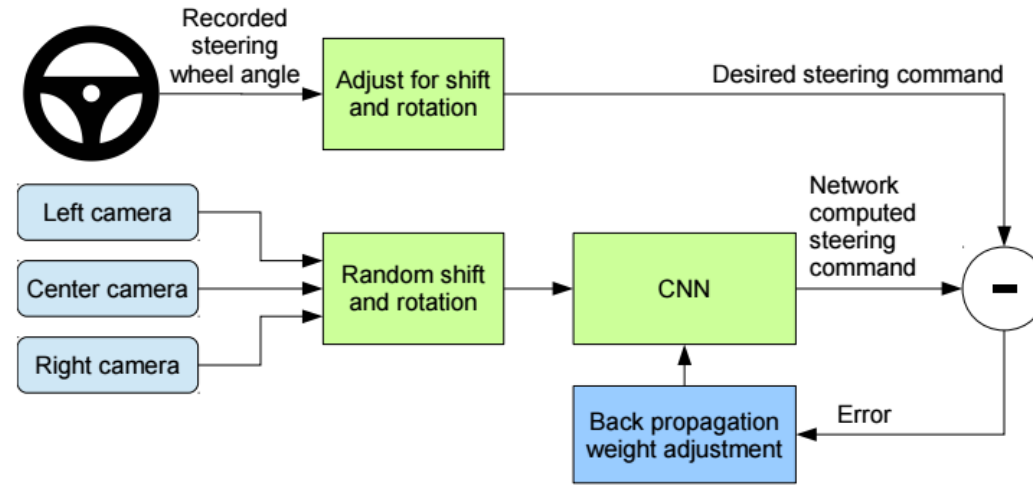
Compounding error

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

Can we avoid compounding error in special cases?

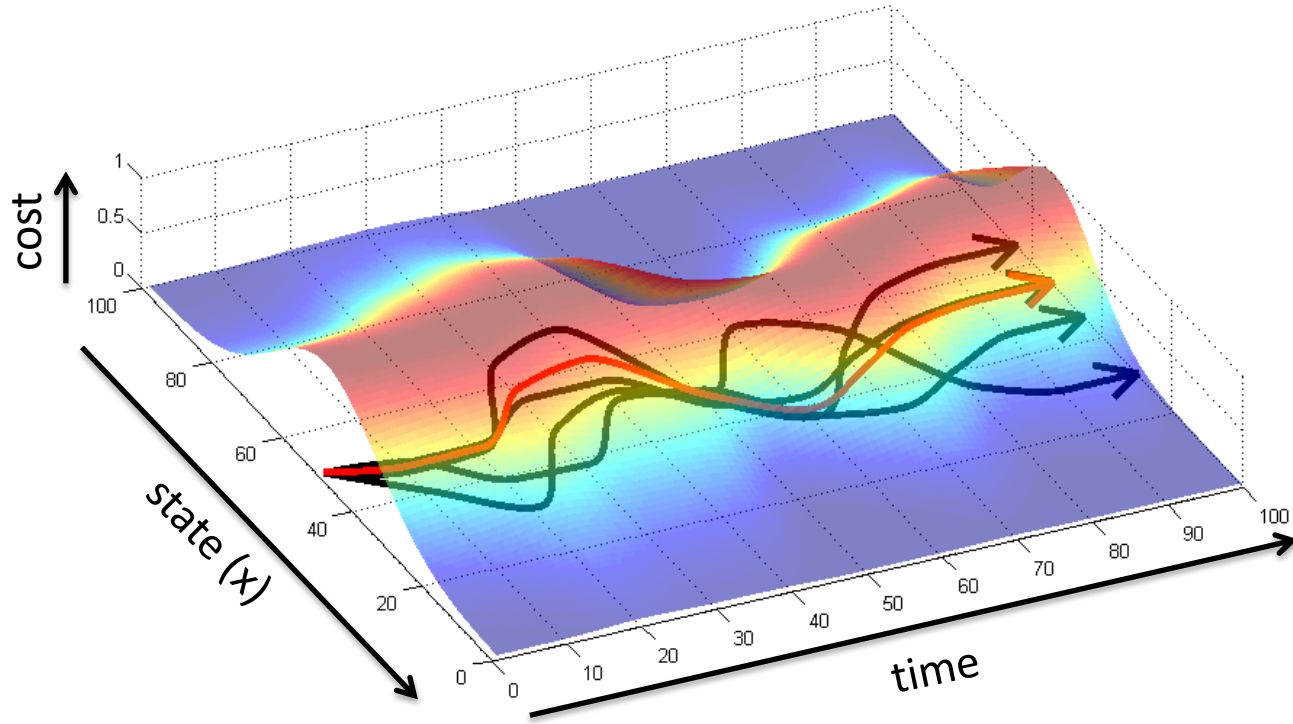


Why did that work?



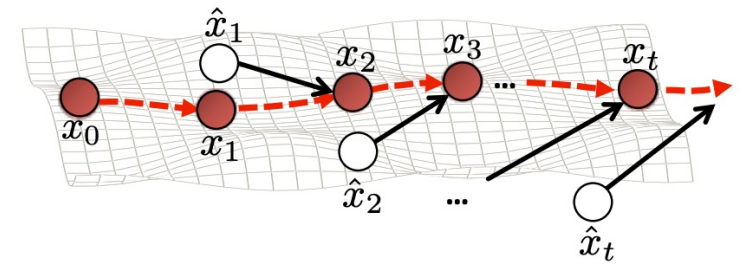
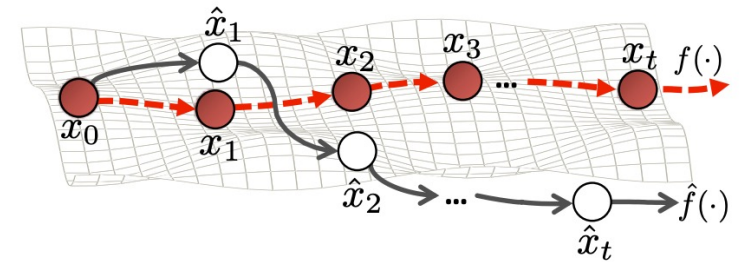
What is the general principle?

- training trajectory
- π_θ expected trajectory

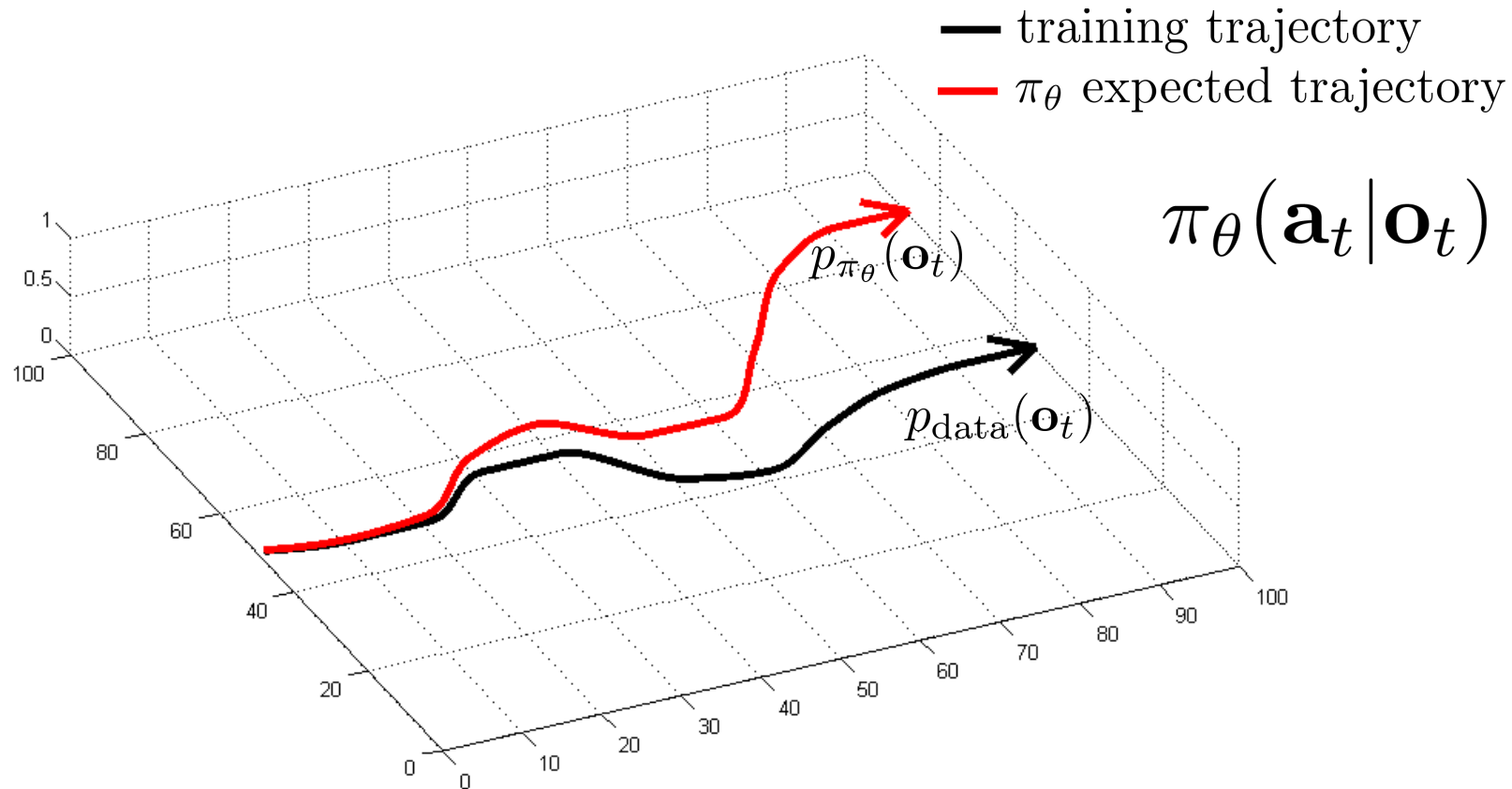


stability

Corrective labels that bring you back to the data



What might this mean mathematically?



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

Concrete Instantiation: DAgger

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

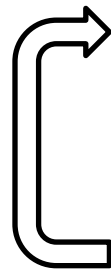
idea: instead of being clever about $p_{\pi_\theta}(\mathbf{o}_t)$, be clever about $p_{\text{data}}(\mathbf{o}_t)$!

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 !

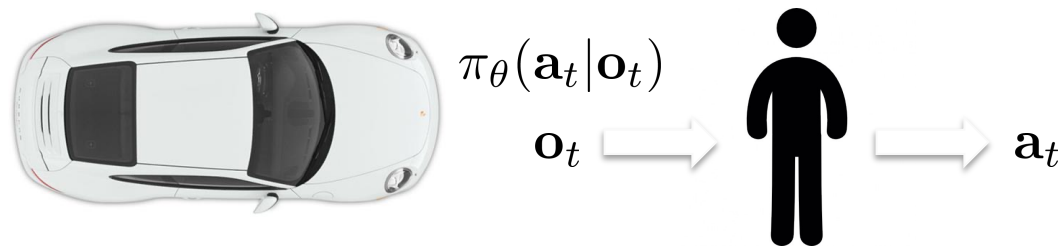
- 
1. train $\pi_\theta(\mathbf{a}_t|\mathbf{o}_t)$ from human data $\mathcal{D} = \{\mathbf{o}_1, \mathbf{a}_1, \dots, \mathbf{o}_N, \mathbf{a}_N\}$
 2. run $\pi_\theta(\mathbf{a}_t|\mathbf{o}_t)$ to get dataset $\mathcal{D}_\pi = \{\mathbf{o}_1, \dots, \mathbf{o}_M\}$
 3. Ask human to label \mathcal{D}_π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

Dagger Example



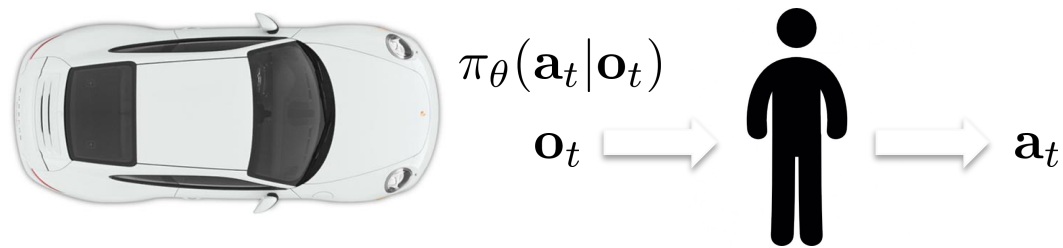
What's the problem?

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: $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_\pi$



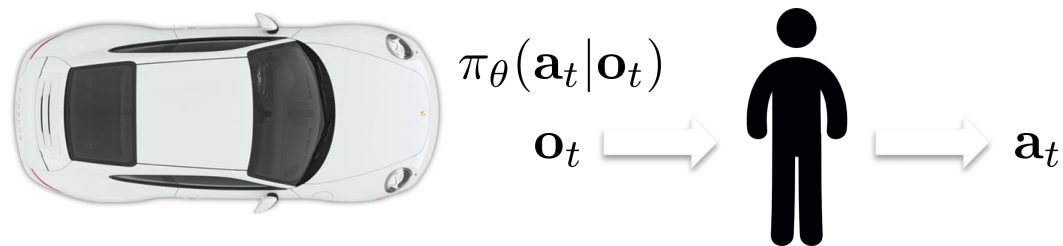
How might we fix this?

- "Generate" corrective labels automatically
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: $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_{\pi}$
- Do at data collection time



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 $\mathcal{D}_\pi = \{\mathbf{o}_1, \dots, \mathbf{o}_M\}$ ← Do at data collection time
3. Ask human to label \mathcal{D}_π with actions \mathbf{a}_t
4. Aggregate: $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_\pi$



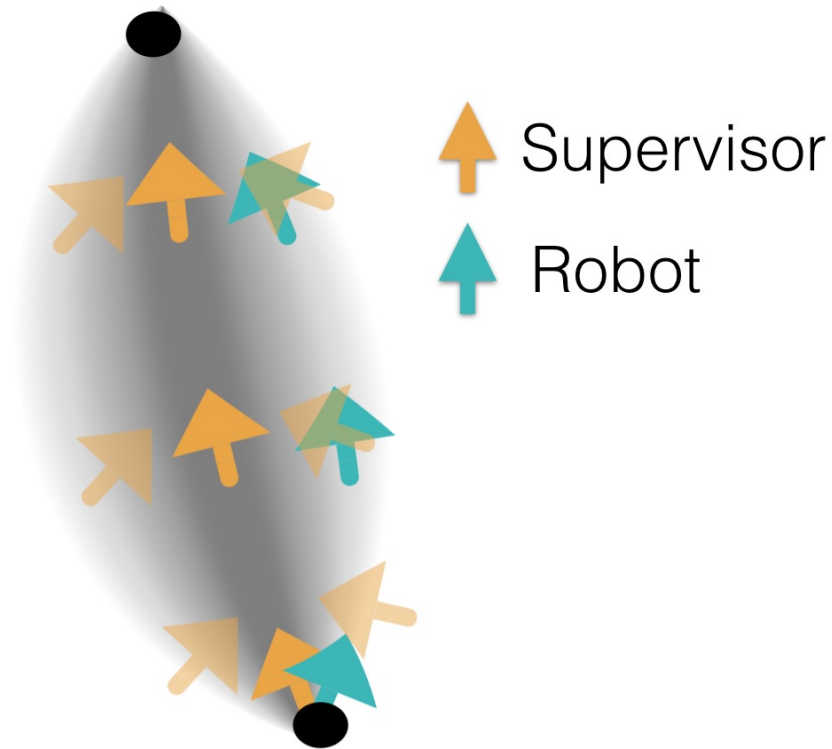
Noising the Data Collection Process

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

Under noise during data collection

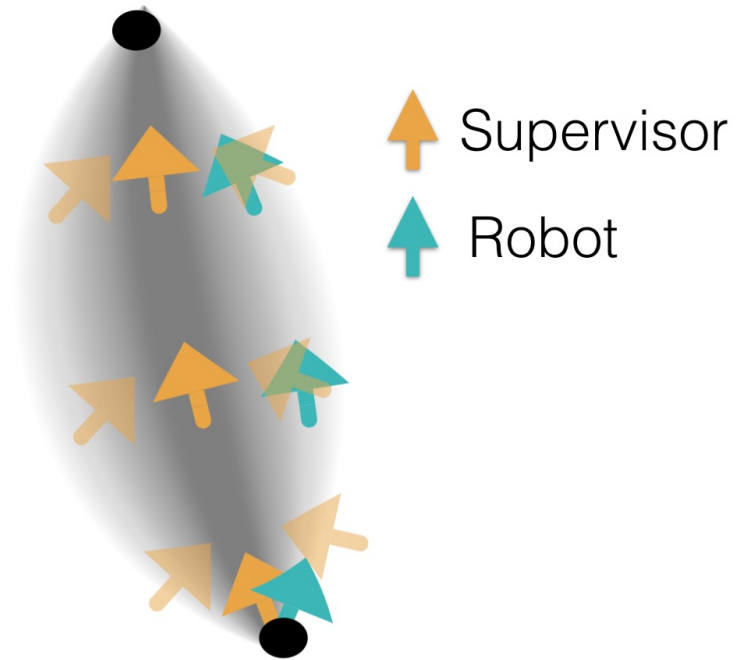
Maximize likelihood

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



Why might this not be enough?

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



Noise Injection

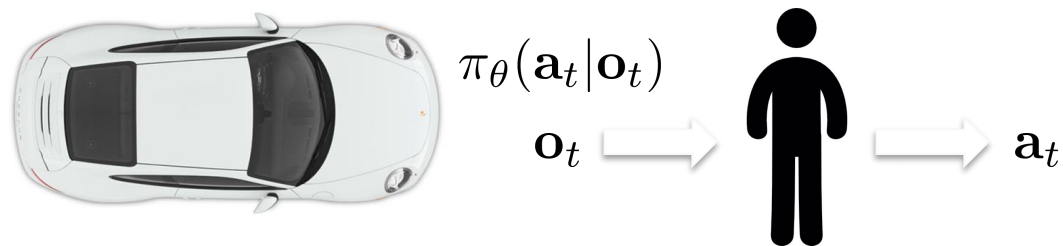


Assumes that the expert can actually perform behaviors under noise
→ Not always possible!

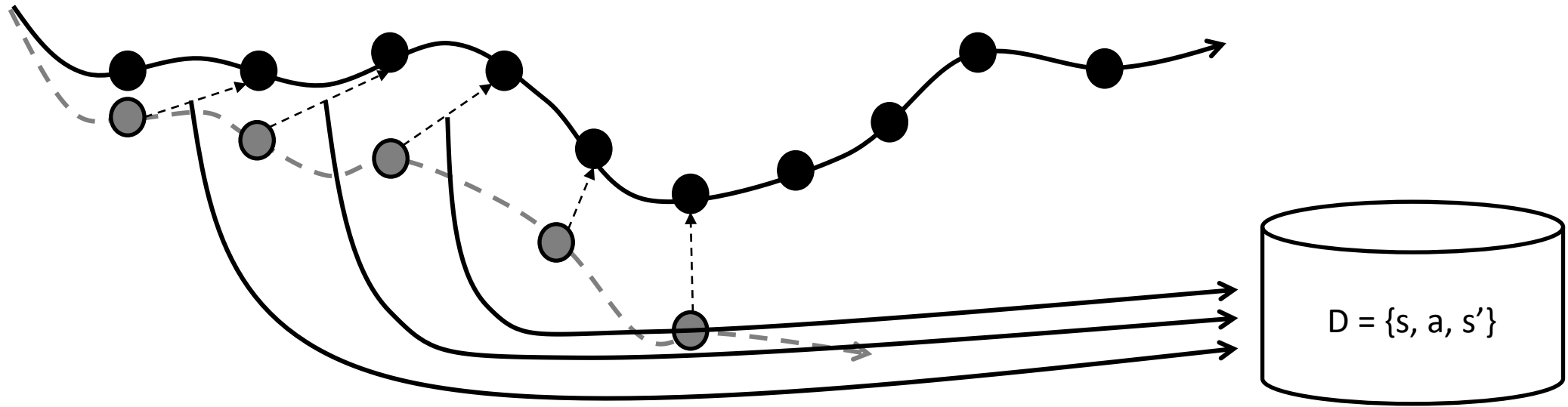
How might we fix this?

"Generate"
corrective labels
automatically

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

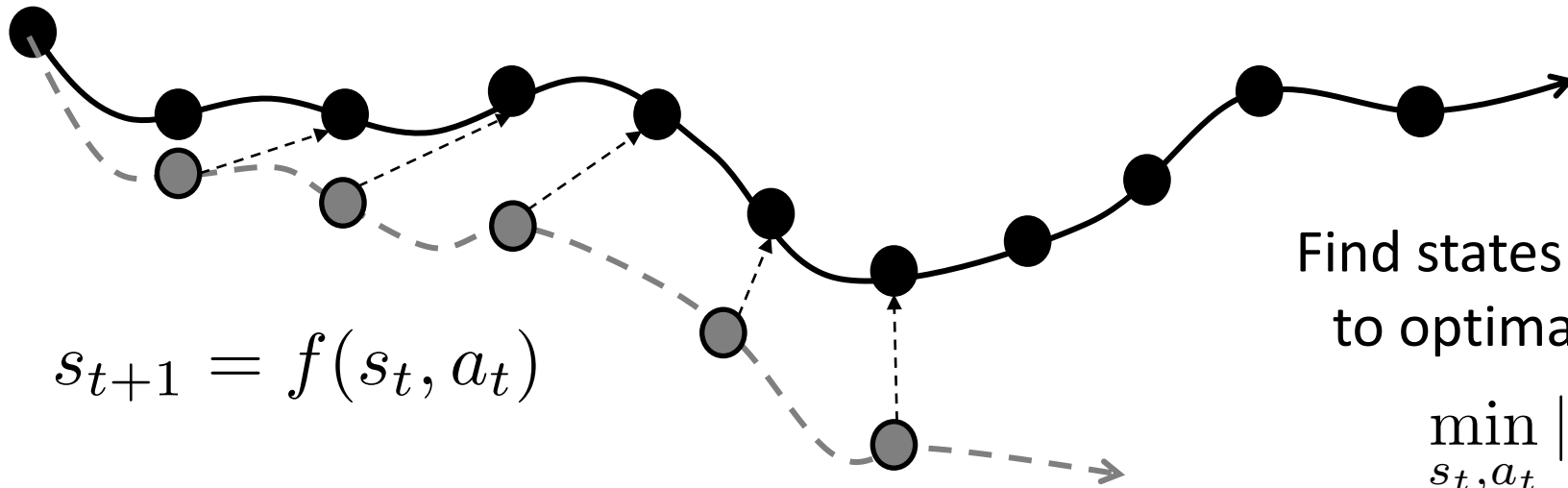


Yunchu
Zhang



Liyiming
Ke

Generating Corrective Labels From True Dynamics



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

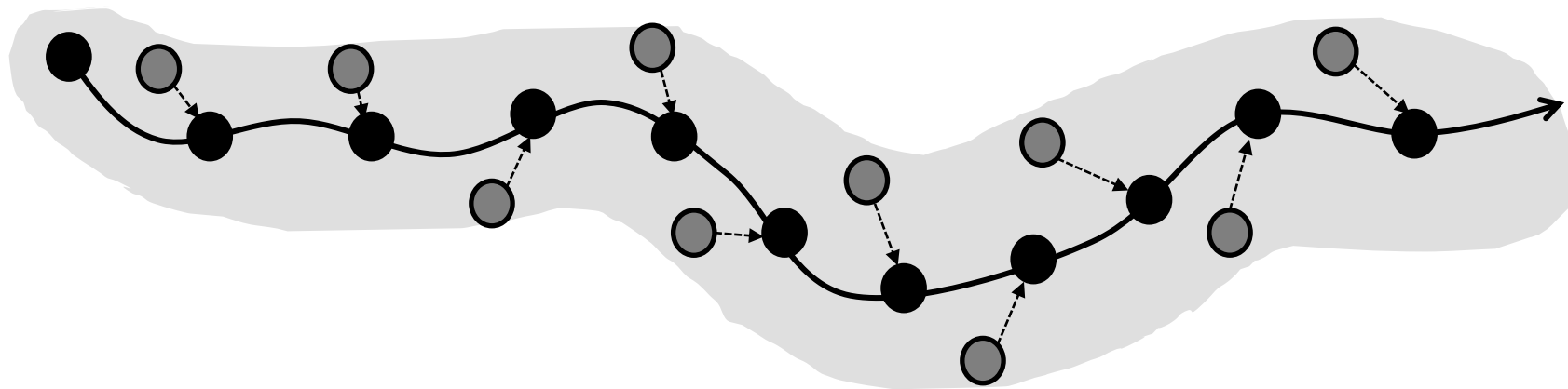
$$\min_{s_t, a_t} \|s_{t+1}^* - f(s_t, a_t)\| \leq \epsilon$$

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

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

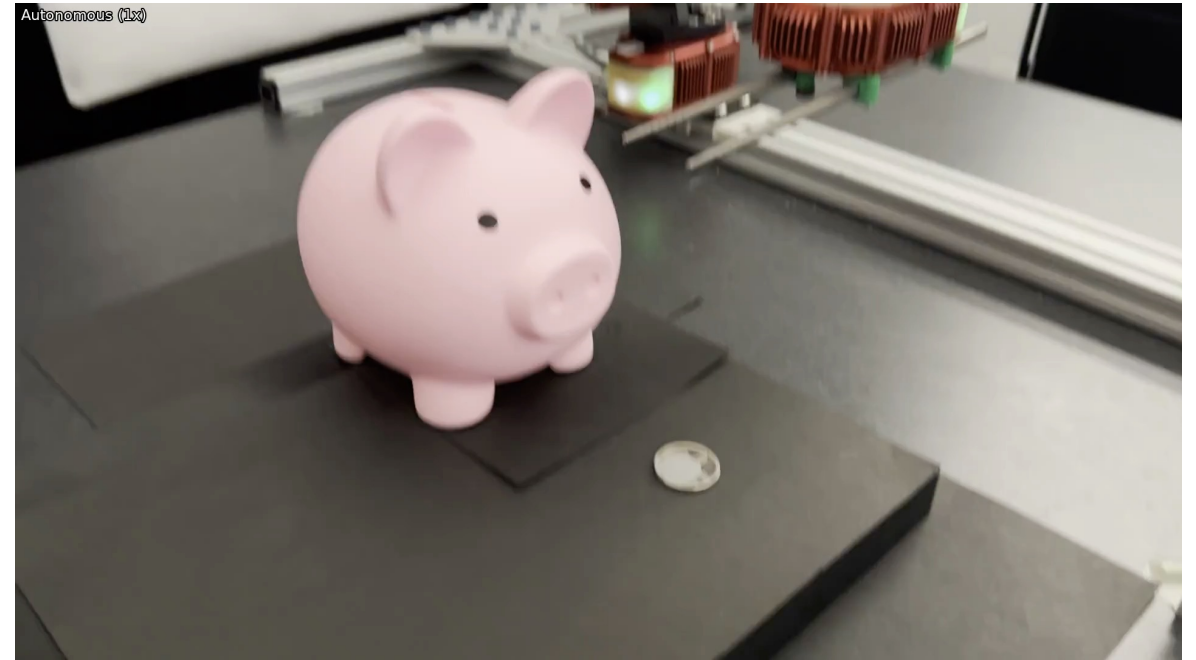
$$\text{s.t. } \|s_t^* - s_t\| \leq \epsilon_1, \|a_t^* - a_t\| \leq \epsilon_2 \longleftarrow \text{Close to data}$$

How well does generating corrective labels work?

With corrective labels

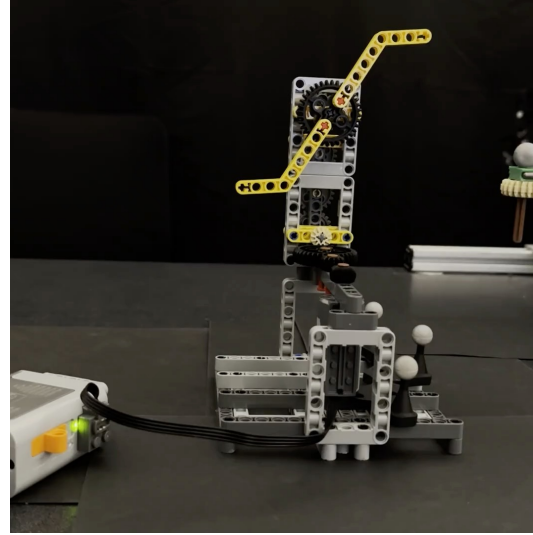
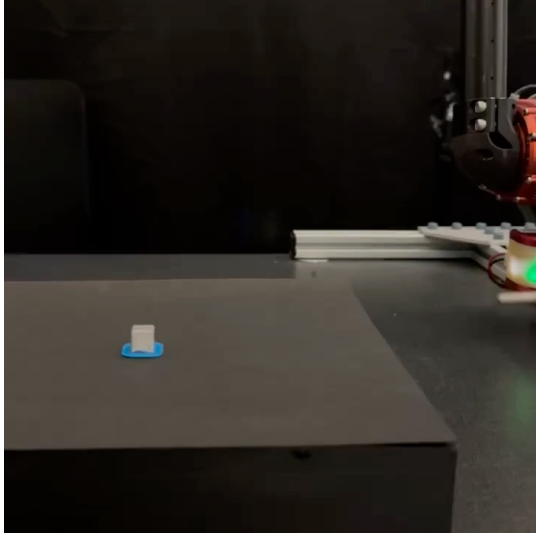


Without corrective labels

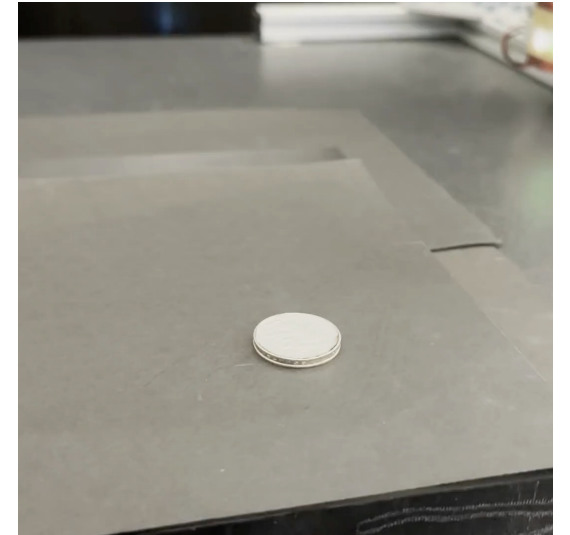
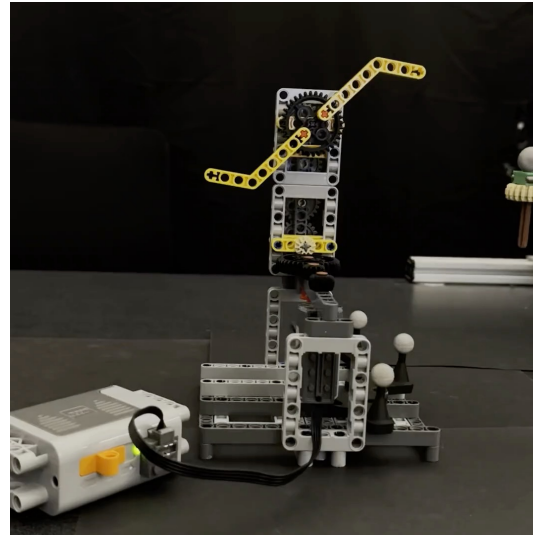
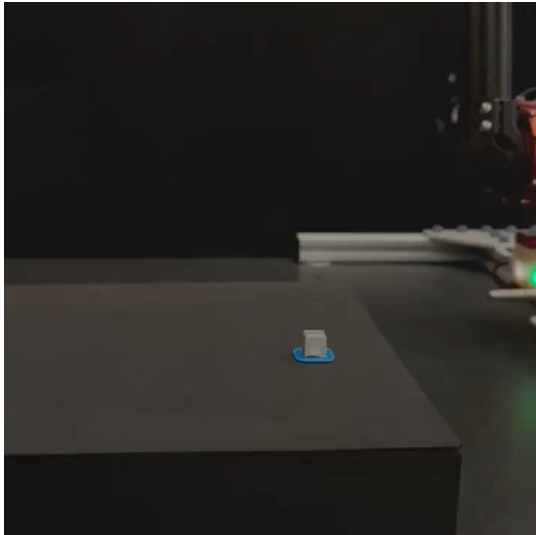


How well does generating corrective labels work?

With corrective labels



Without corrective labels



Lecture outline

Recap: Imitation Learning + Why it is hard



Multimodality and Underfitting in Imitation



Compounding Error in Imitation

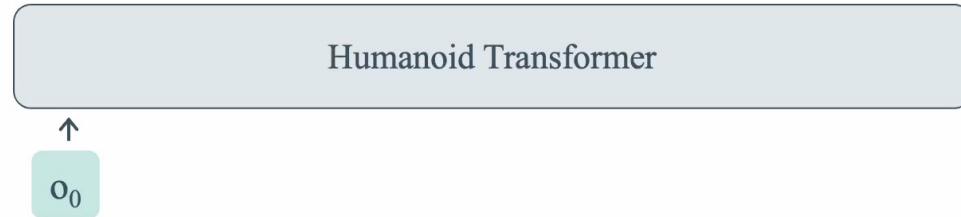


Frontiers in Imitation

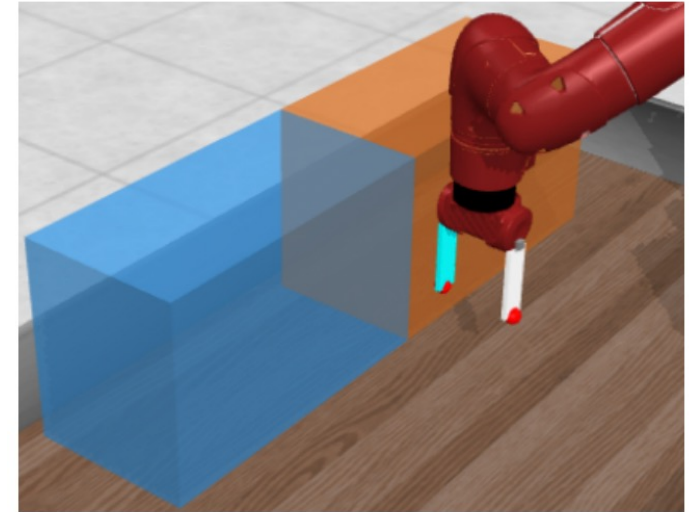
So does this solve all the issues in imitation?

Frontiers in Imitation Learning

Non-Markovian Demonstrators



Characterizing generalization

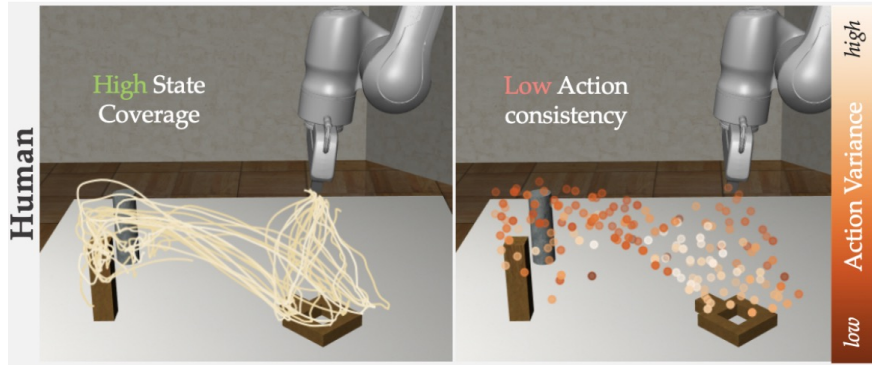


Action-Free Data

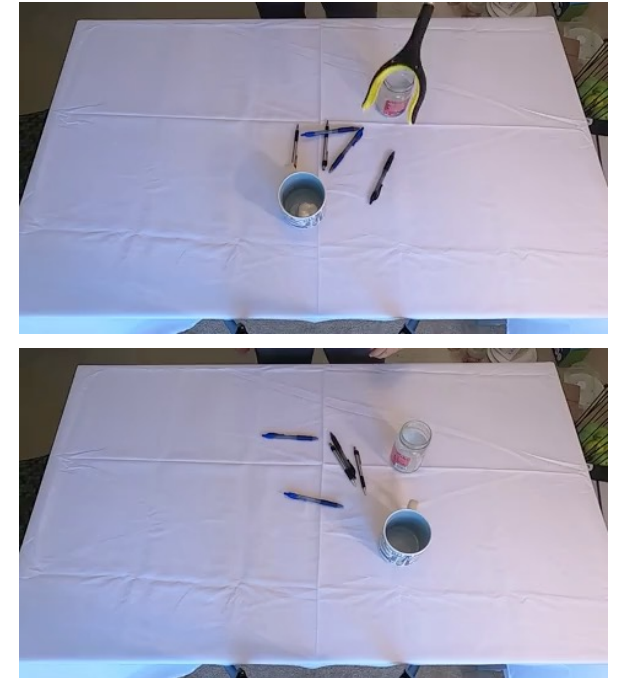


Frontiers in Imitation Learning

Data Curation and Quality



Embodiment Shift

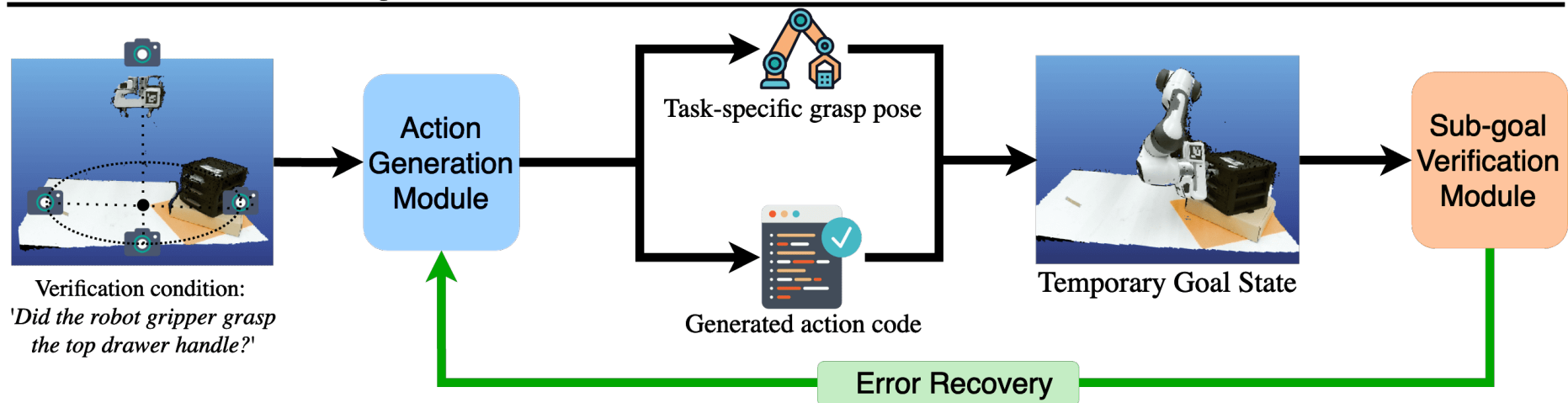
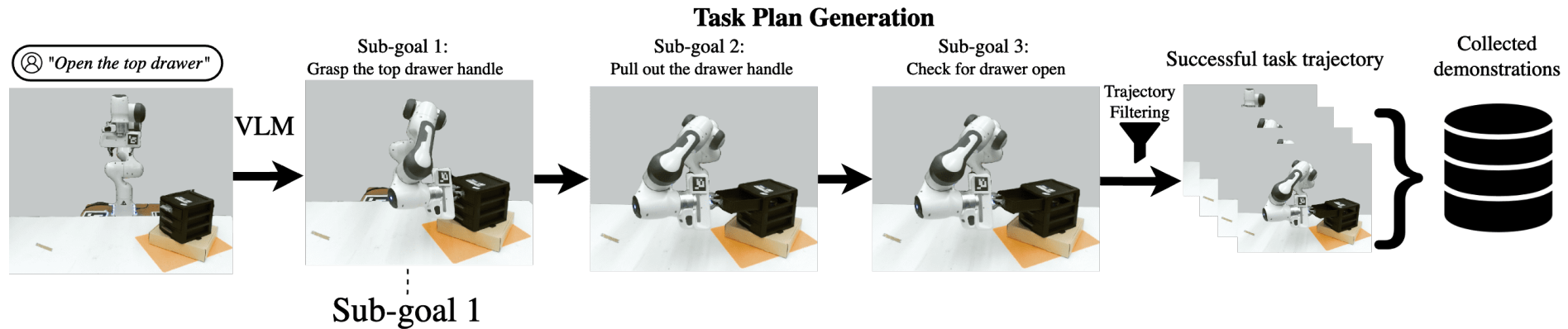


Teleoperation Interfaces



Frontiers in Imitation Learning

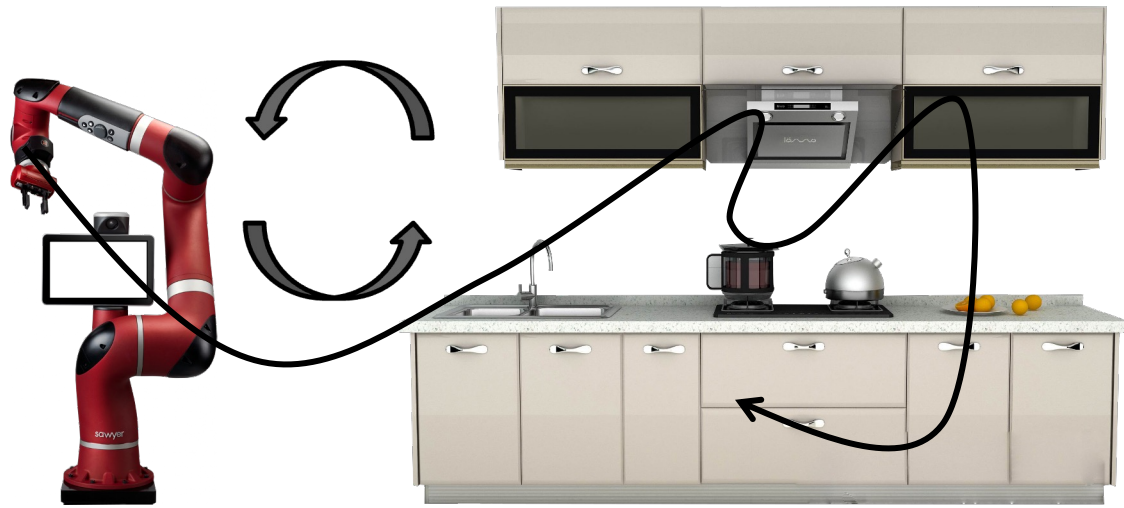
Learning how to retry and improve



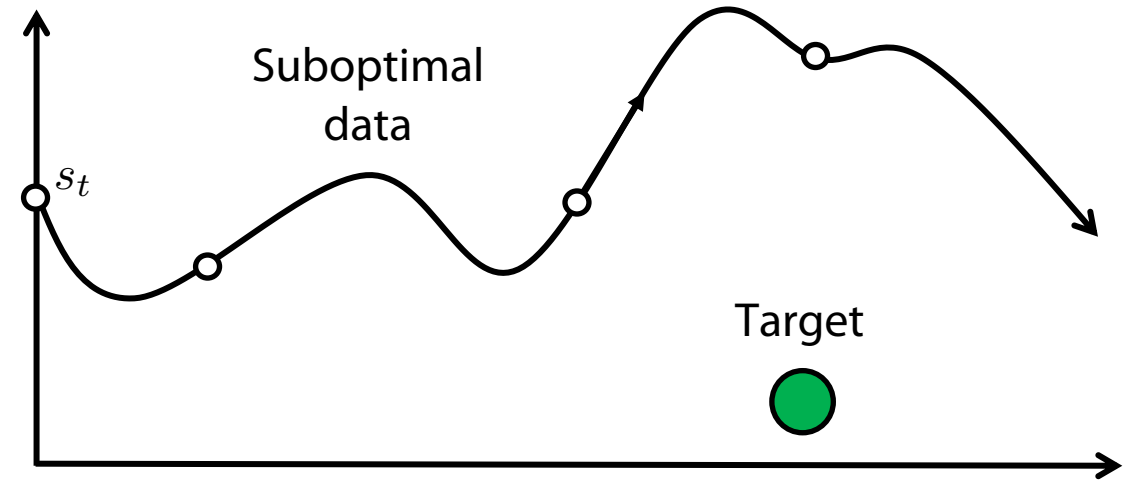
Duan et al

Let's dive into a few

Accounting for Suboptimal Data

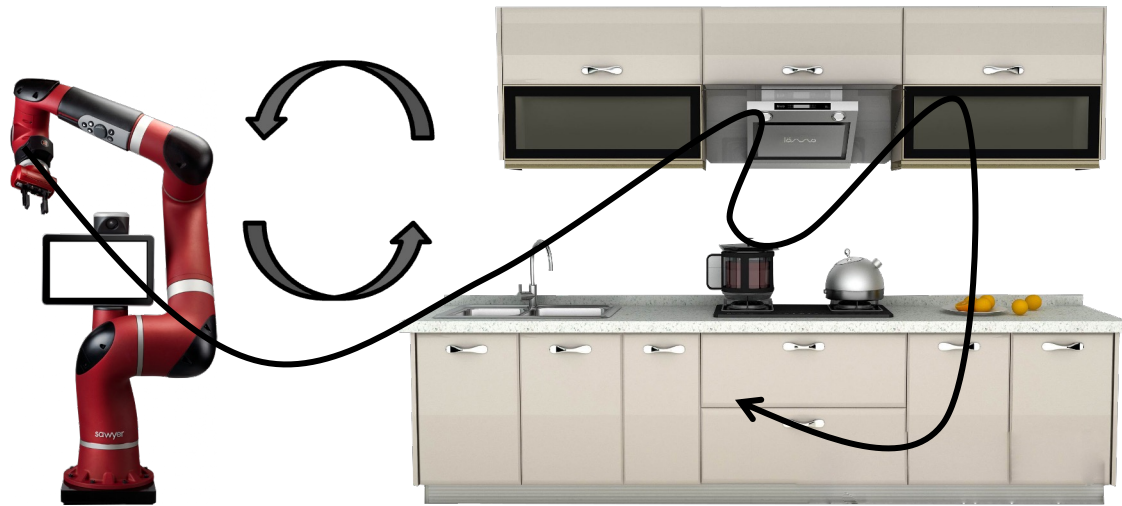


Random
Policy
 π_{θ}



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

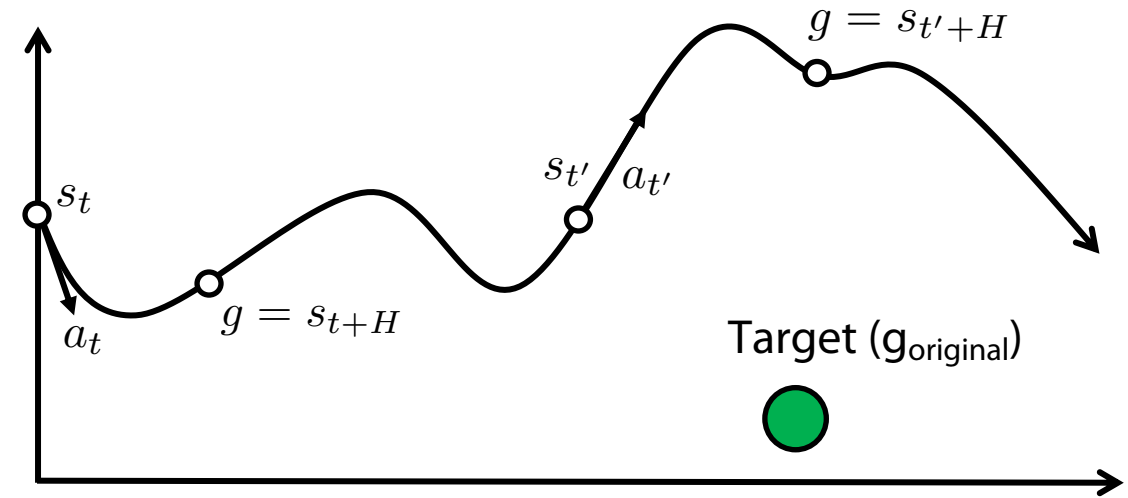
Hindsight relabeling for Imitation Learning



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



Relabel the right goal in "hindsight"

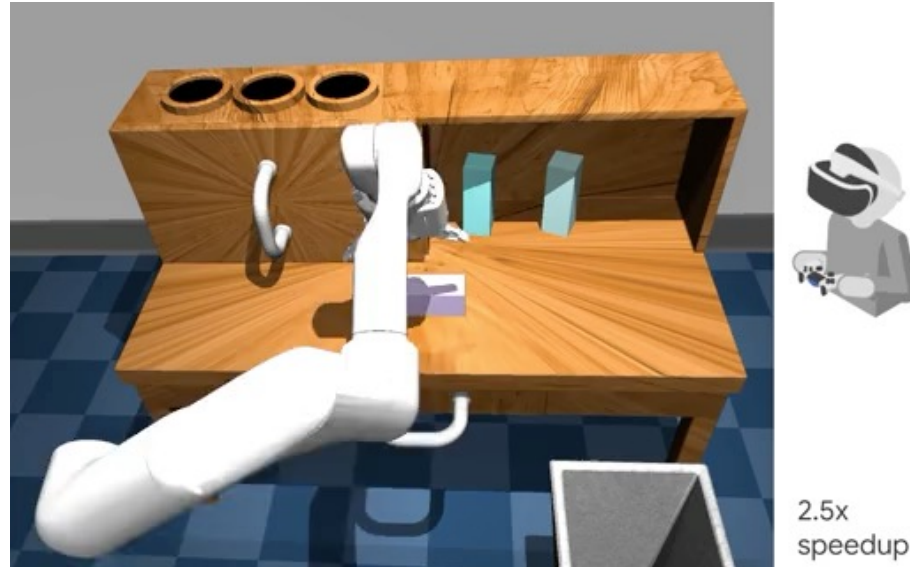


Learn a multi-goal policy $\pi_{\theta}(a|s, g)$



Treat reached states as **optimal** goals

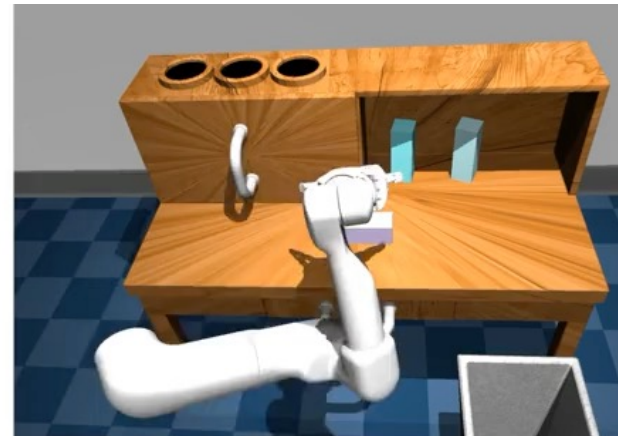
What does this result in?



Undirected play data



Goal

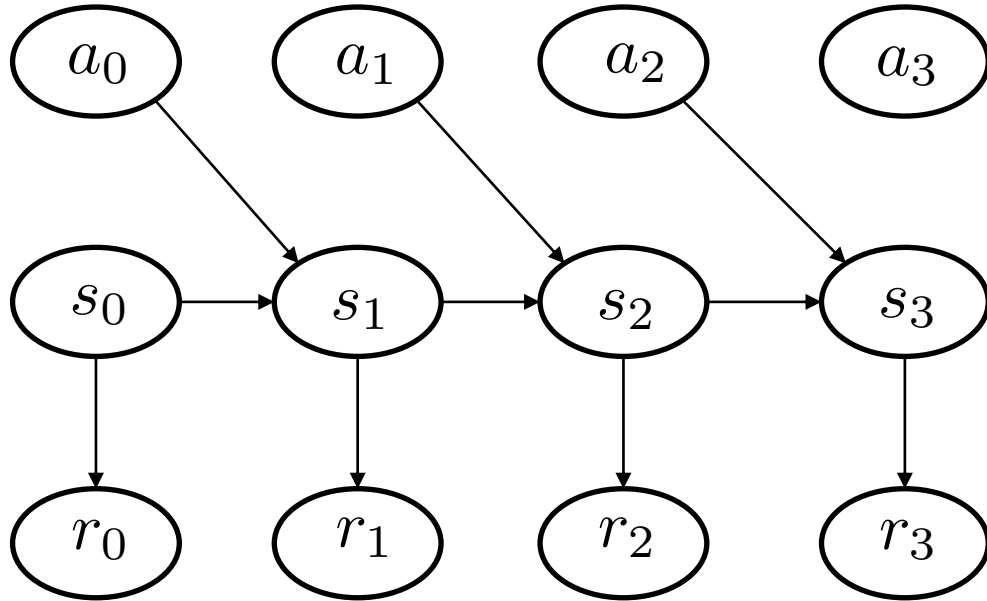


Single Play-LMP policy

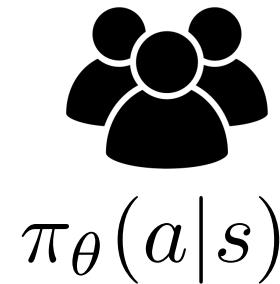
Goal-directed behavior

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

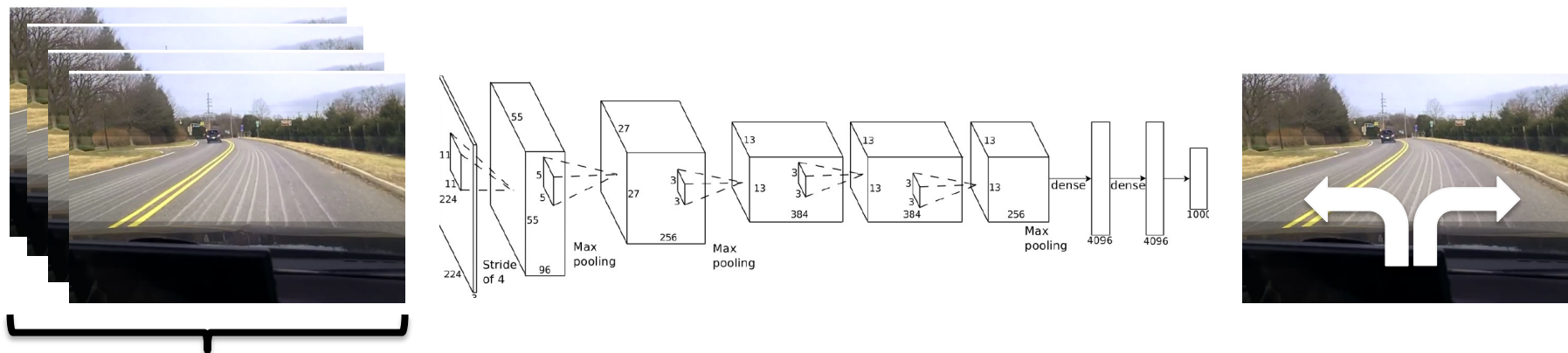
Not necessarily!

Humans often rely on history

Mixtures of Markovian humans may not be Markovian

How can we deal with non-Markovian demonstrators?

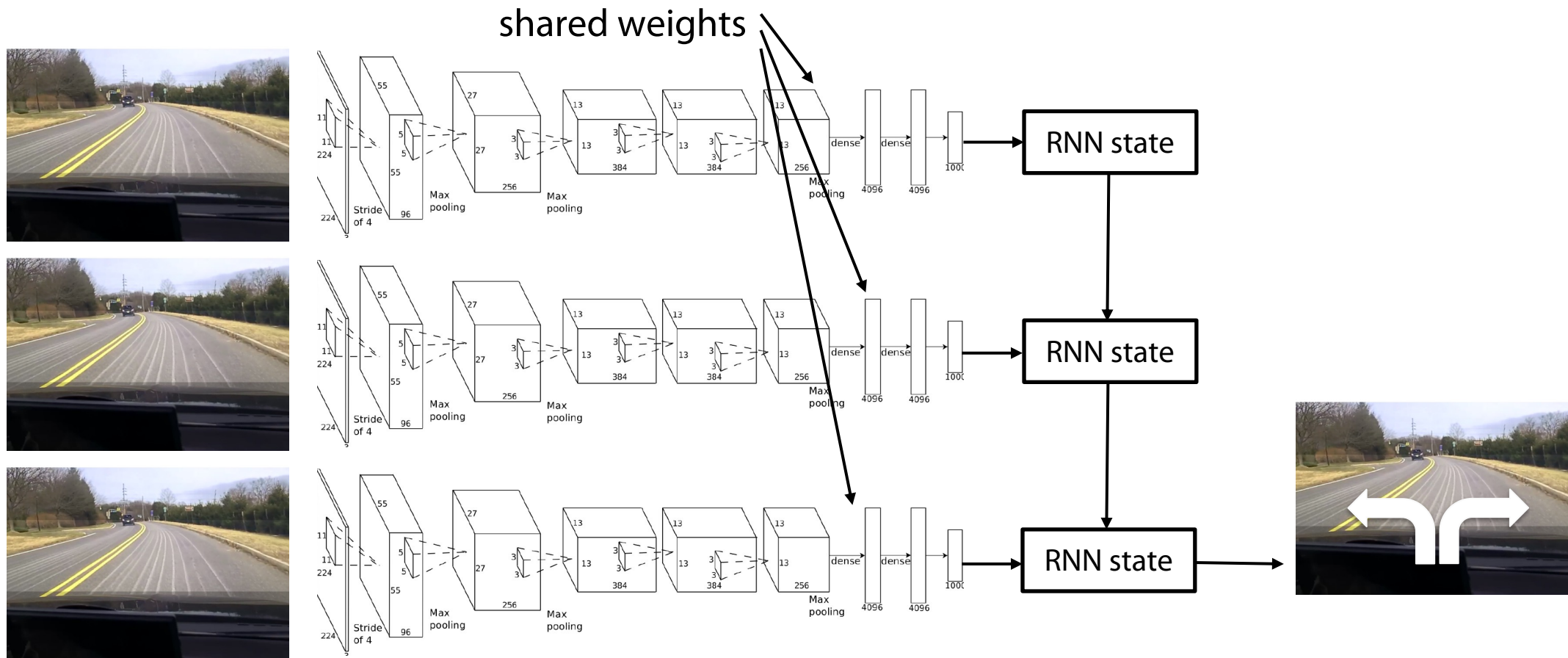
$$\text{Learn } \pi_{\theta}(a_t | s_t, s_{t-1}, \dots, s_0)$$



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

How can we deal with non-Markovian demonstrators?

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



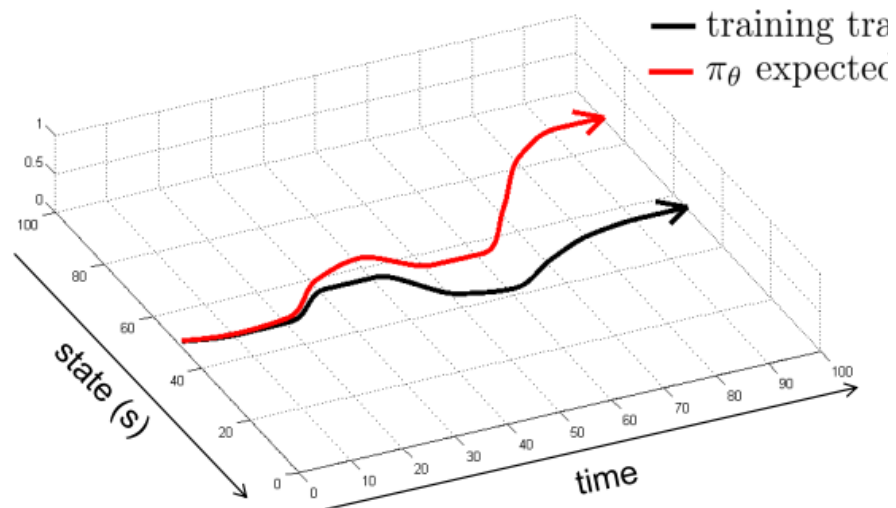
Option 2: Use a recurrent model (LSTM/transformer/RNN)

Credit: Sergey Levine

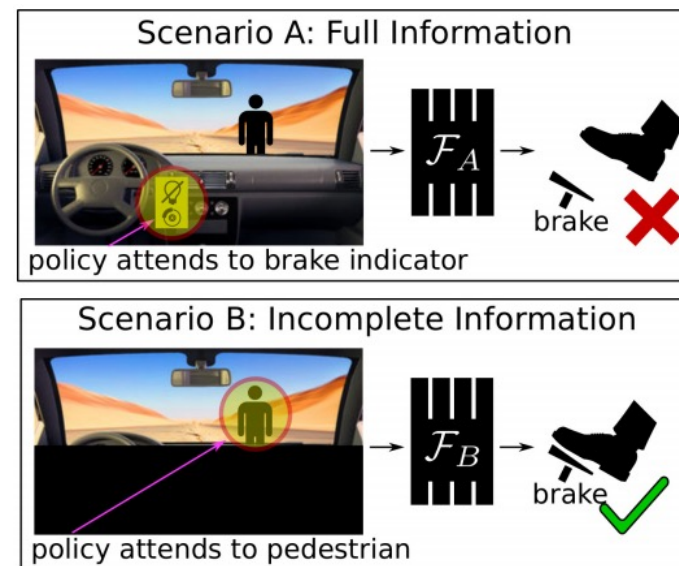
Why might this be challenging?

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

Easier to go OOD



Learns spurious shortcut behaviors



Some cool imitation videos

1x and tesla humanoid robots



● 1X END-TO-END AUTONOMY
UPDATE, JAN 2024

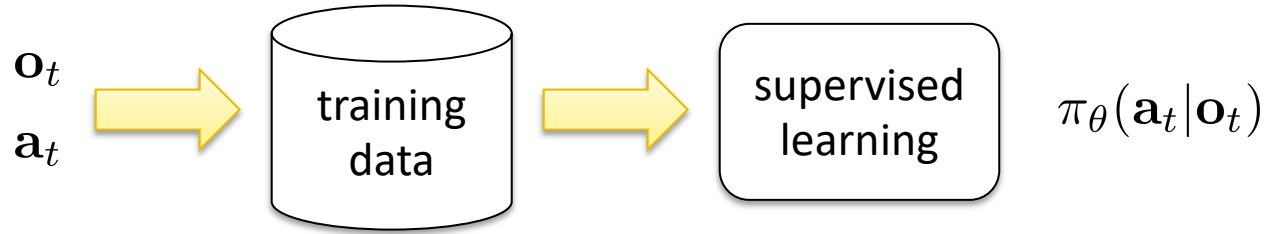
ALOHA and CherryBot Fine Manipulation



TRI Diffusion Policies

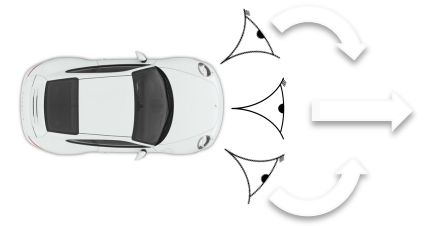


Perspectives on Imitation



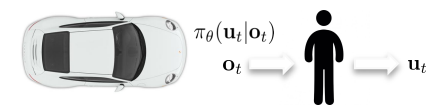
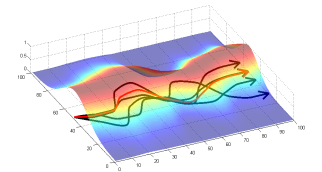
■ Pros:

- 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 + Why it is hard



Multimodality and Underfitting in Imitation



Compounding Error in Imitation



Frontiers in Imitation