

# Supervised Learning of Behaviors

CS 285: Deep Reinforcement Learning, Decision Making, and Control

Sergey Levine

# Class Notes

1. Homework 1 is out this evening
2. Remember to start forming final project groups
  - Final project assignment document is now out!
  - Proposal due Sep 25

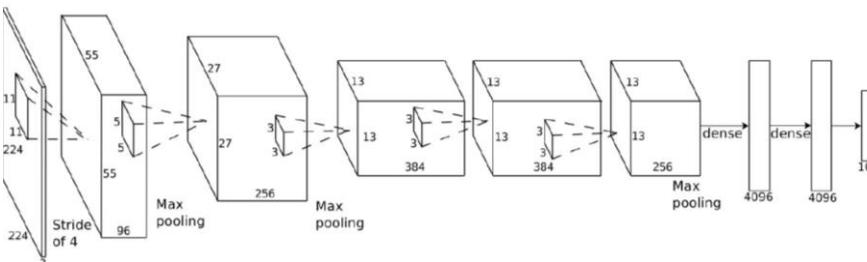
# Today's Lecture

1. Definition of sequential decision problems
  2. Imitation learning: supervised learning for decision making
    - a. Does direct imitation work?
    - b. How can we make it work more often?
  3. A little bit of theory
  4. Case studies of recent work in (deep) imitation learning
- Goals:
    - Understand definitions & notation
    - Understand basic imitation learning algorithms
    - Understand tools for theoretical analysis

# Terminology & notation



$\mathbf{o}_t$



$$\pi_{\theta}(\mathbf{a}|\mathbf{o}_t)$$



$\mathbf{a}_t$

$\mathbf{s}_t$  – state

$\mathbf{o}_t$  – observation

$\mathbf{a}_t$  – action

partially observed

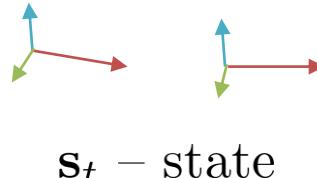
$\pi_{\theta}(\mathbf{a}_t|\mathbf{o}_t)$  – policy

$\pi_{\theta}(\mathbf{a}_t|\mathbf{s}_t)$  – policy (fully observed)

The important distinction between states and observations is that the states fully describe everything that's going on in the world, whereas the observations may contain a loss of meaning that you might have 2 observations that are perhaps indistinguishable to you but that actually correspond to different states.



$\mathbf{o}_t$  – observation

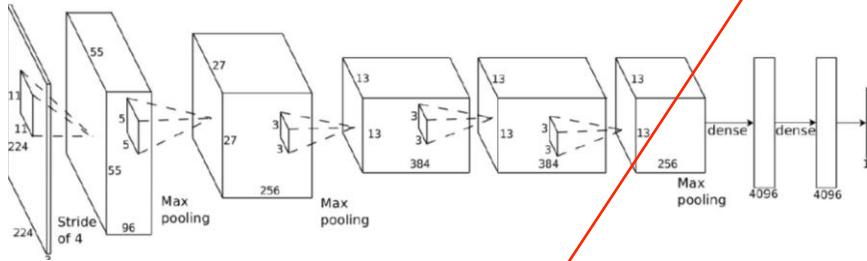


$\mathbf{s}_t$  – state

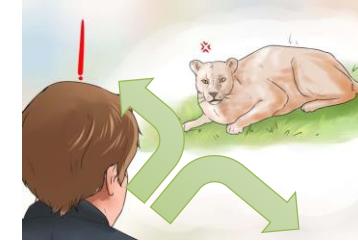
# Terminology & notation



$\mathbf{o}_t$



$\pi_\theta(\mathbf{a}_t | \mathbf{o}_t)$



$\mathbf{a}_t$

$s_t$  – state

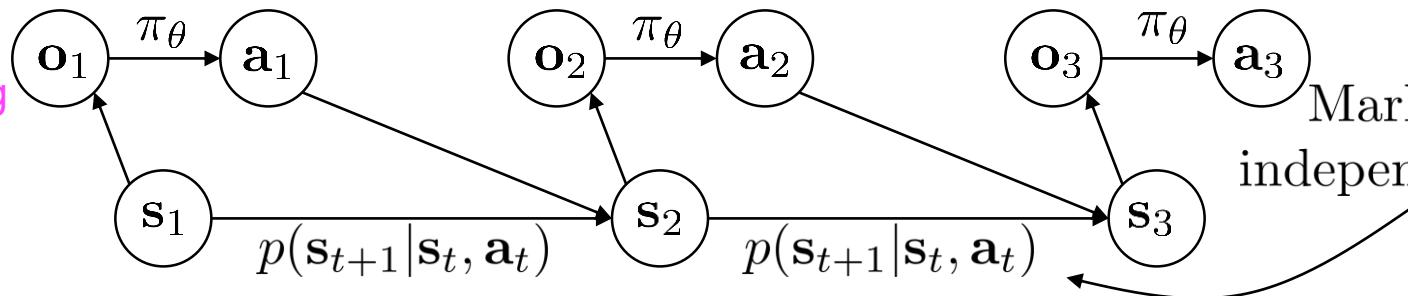
$\mathbf{o}_t$  – observation

$\mathbf{a}_t$  – action

$\pi_\theta(\mathbf{a}_t | \mathbf{o}_t)$  – policy

$\pi_\theta(\mathbf{a}_t | \mathbf{s}_t)$  – policy (fully observed)

formal mathematically  
distinction between states  
and observations: drawing  
a probabilistic graphical  
model that describes how  
all 3 of these things relate  
to 1 another.



observations do not form a Markov chain.

# Aside: notation

$s_t$  – state

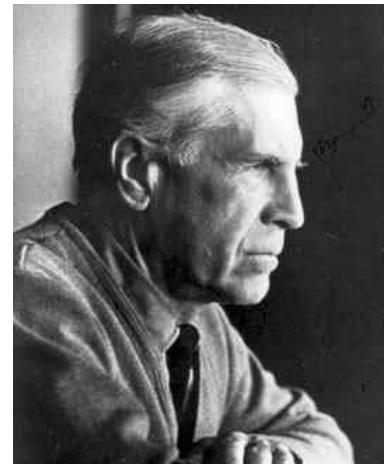
$a_t$  – action



Richard Bellman

$x_t$  – state

$u_t$  – action      управление

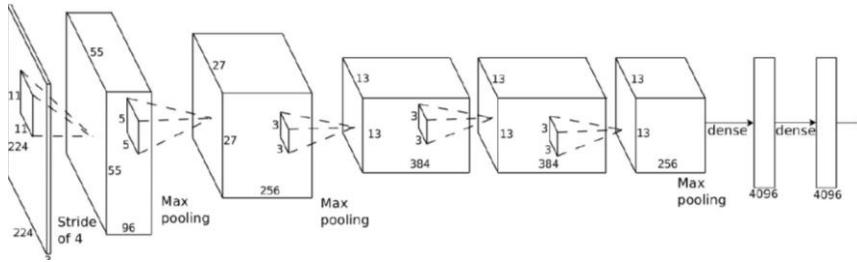


Lev Pontryagin

# Imitation Learning



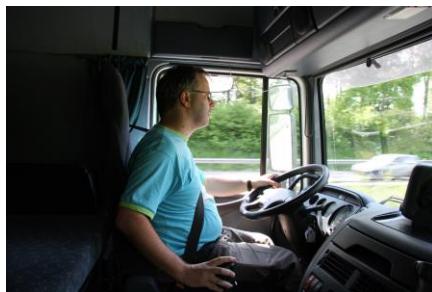
$\mathbf{o}_t$



$$\pi_{\theta}(\mathbf{a}_t | \mathbf{o}_t)$$



$\mathbf{a}_t$



$\mathbf{o}_t$   
 $\mathbf{a}_t$



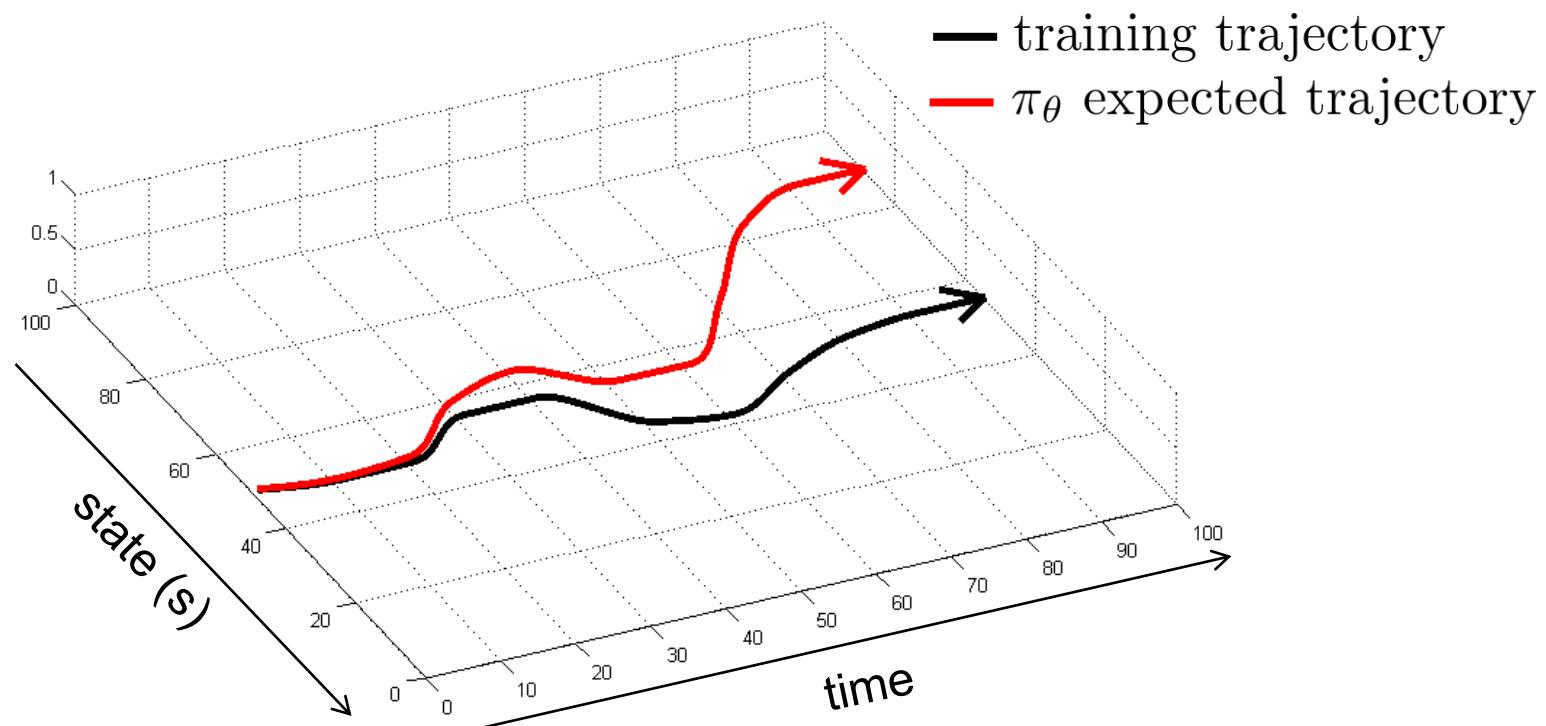
$$\pi_{\theta}(\mathbf{a}_t | \mathbf{o}_t)$$

## behavioral cloning

# Does it work?

# No!

small mistake makes big mistake.

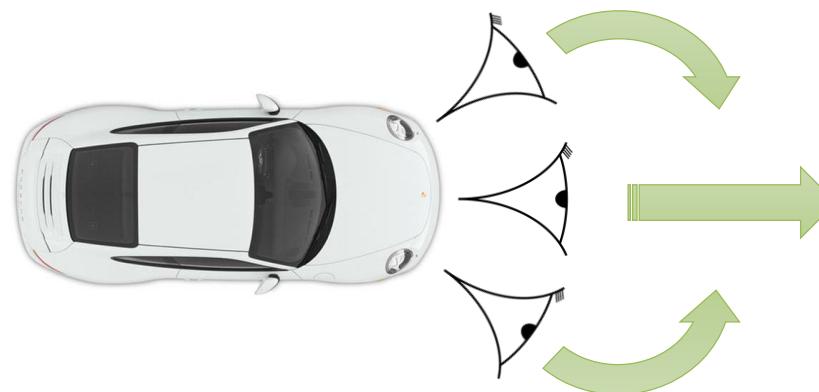
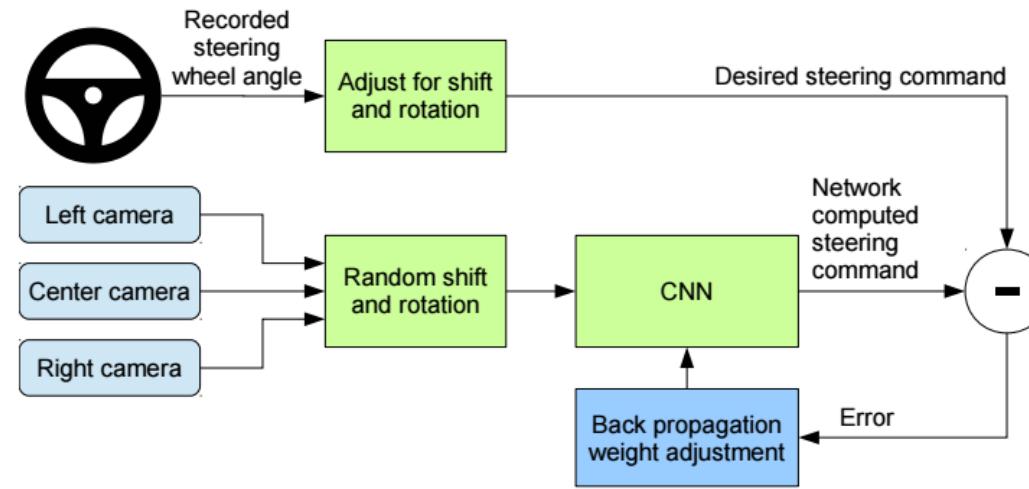


Does it work?

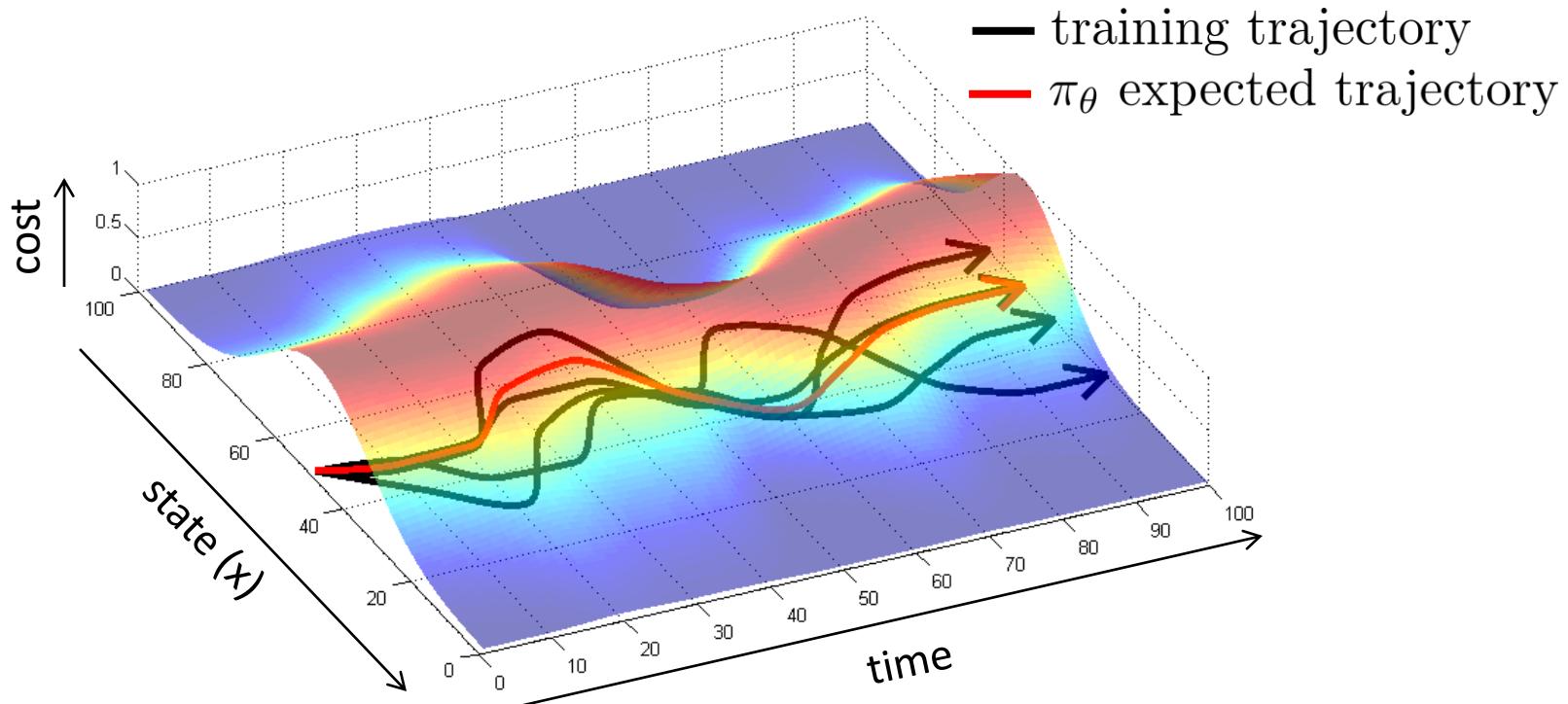
Yes!



# Why did that work?



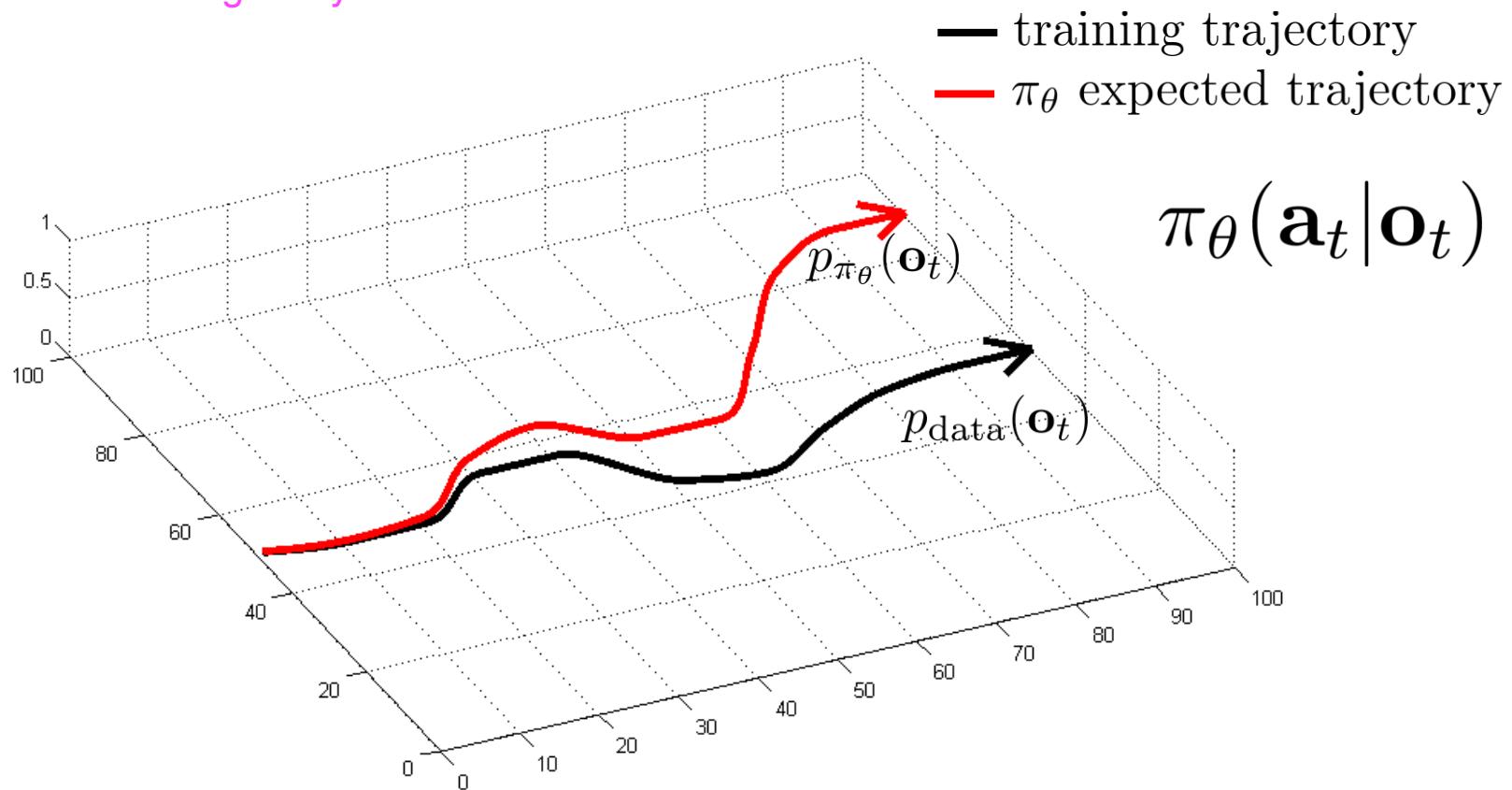
# Can we make it work more often?



stability  
(more on this later)

# Can we make it work more often?

What is the actual mathematical explanation for the problem of the behavior cloning really is.



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

# Can we make it work more often?

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

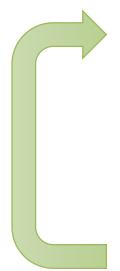
idea: instead of being clever about  $p_{\pi_\theta}(\mathbf{o}_t)$ , be clever about  $p_{\text{data}}(\mathbf{o}_t)$ !

## 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  $\mathbf{a}_t$ !

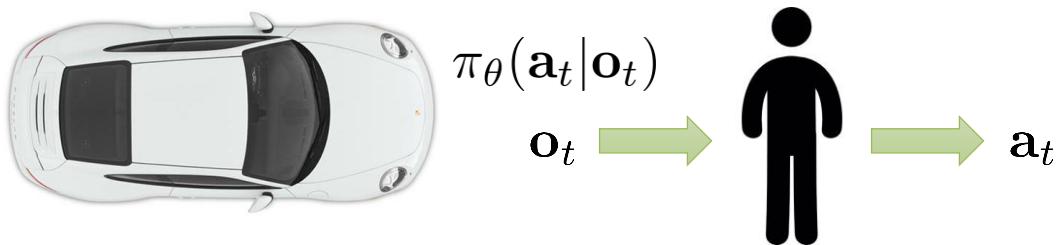
- 
1. train  $\pi_\theta(\mathbf{a}_t | \mathbf{o}_t)$  from human data  $\mathcal{D} = \{\mathbf{o}_1, \mathbf{a}_1, \dots, \mathbf{o}_N, \mathbf{a}_N\}$
  2. run  $\pi_\theta(\mathbf{a}_t | \mathbf{o}_t)$  to get dataset  $\mathcal{D}_\pi = \{\mathbf{o}_1, \dots, \mathbf{o}_M\}$
  3. Ask human to label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

# DAgger Example



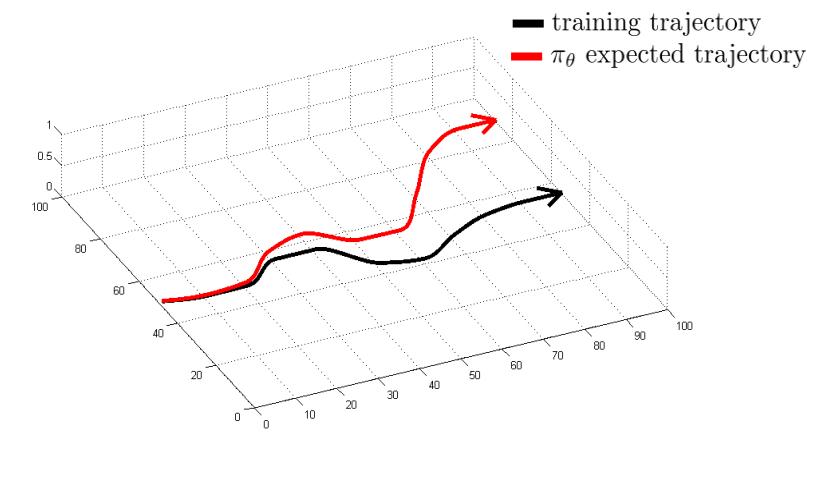
# What's the problem?

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  $\mathcal{D}_\pi$  with actions  $\mathbf{a}_t$
4. Aggregate:  $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_\pi$



# Can we make it work without more data?

- DAgger addresses the problem of distributional “drift”
- What if our model is so good that it doesn’t drift?
- Need to mimic expert behavior very accurately
- But don’t overfit!



Theoretically the answer is no and what does it mean is that you can construct a counterexample a pathological situation where it will not happen. But in real-world situations sometimes you can actually make it work.

So (above) let's discuss a few practical considerations that can help us make it work in practice even if there are no guarantees.

# Why might we fail to fit the expert?

- 
1. Non-Markovian behavior
  2. Multimodal behavior

$$\pi_{\theta}(\mathbf{a}_t | \mathbf{o}_t)$$


behavior depends only  
on current observation

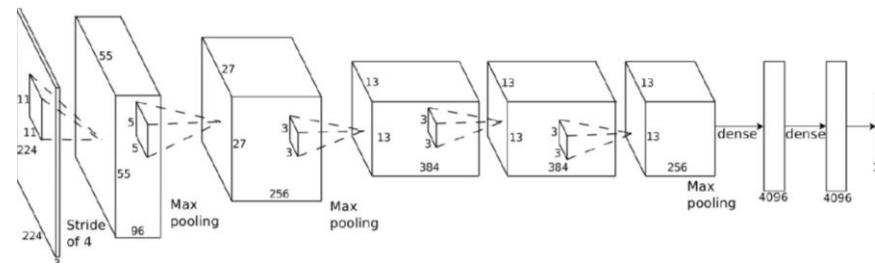
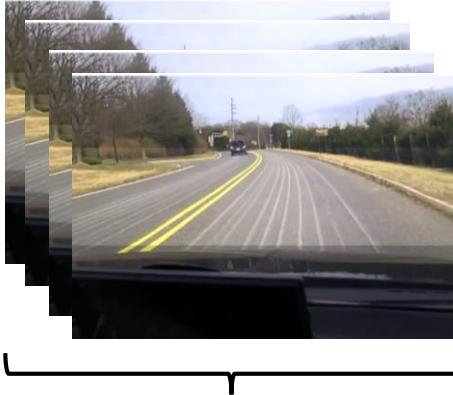
$$\pi_{\theta}(\mathbf{a}_t | \mathbf{o}_1, \dots, \mathbf{o}_t)$$


behavior depends on  
all past observations

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

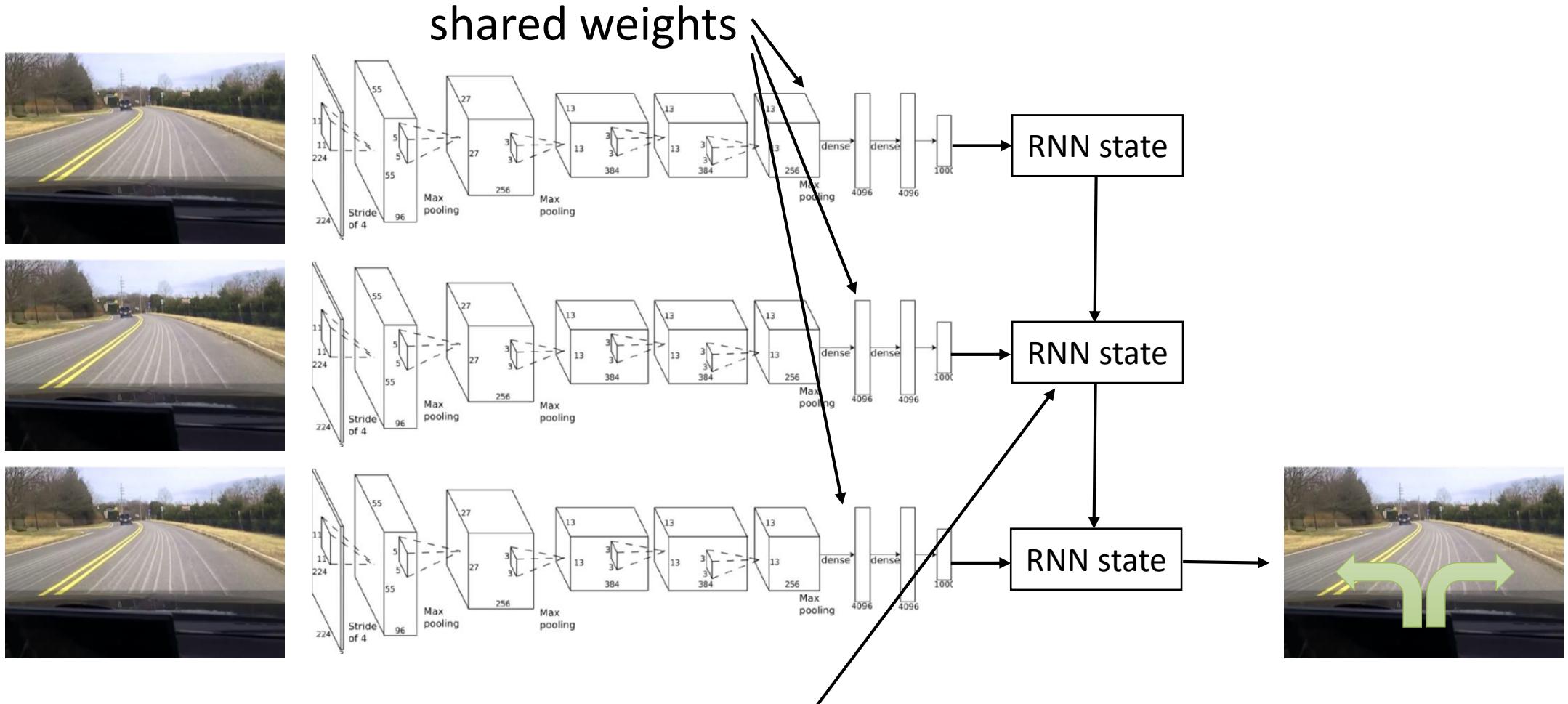
Often very unnatural for  
human demonstrators

# How can we use the whole history?



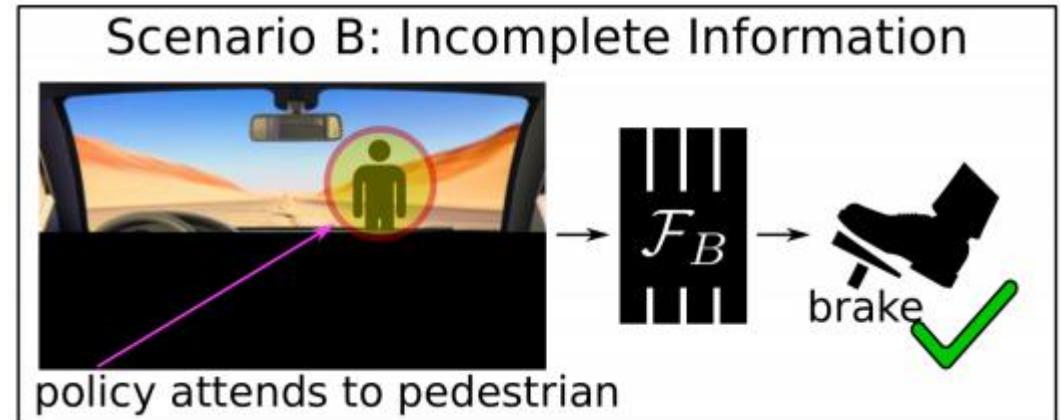
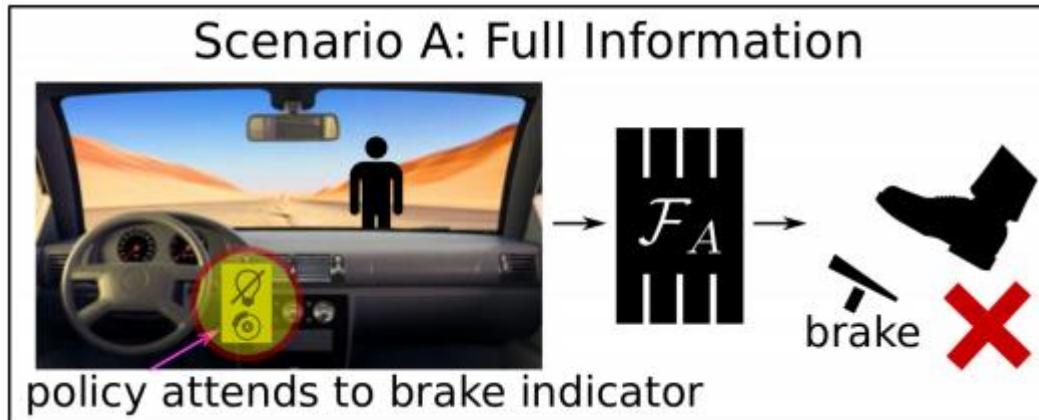
variable number of frames,  
too many weights

# How can we use the whole history?



Typically, LSTM cells work better here

# Aside: why might this work poorly?



“causal confusion”

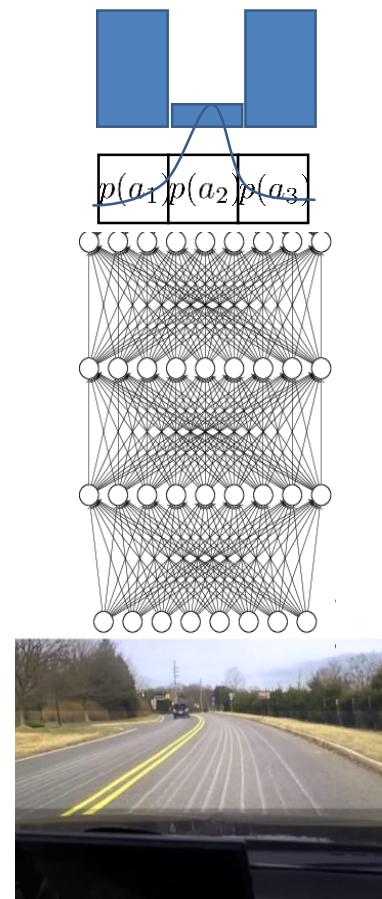
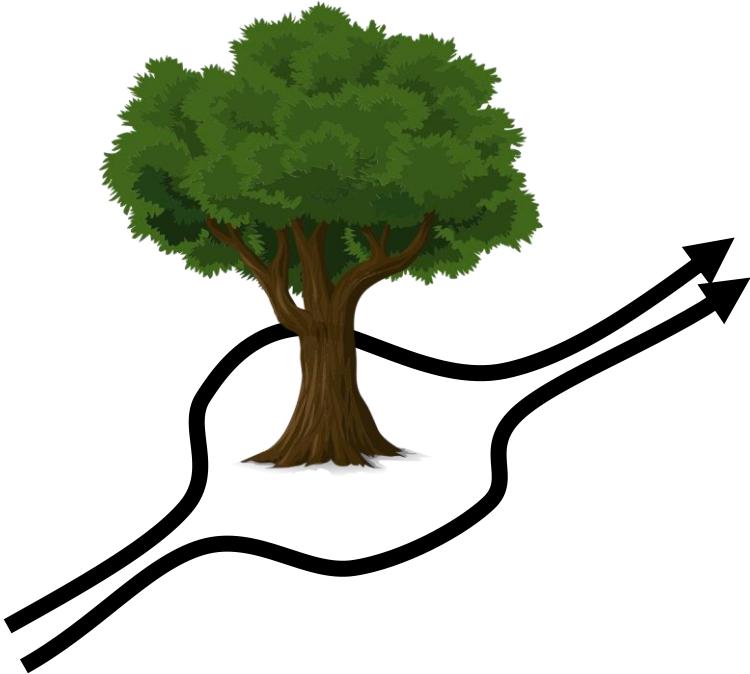
see: de Haan et al., “Causal Confusion in Imitation Learning”

**Question 1:** Does including history exacerbate causal confusion?

**Question 2:** Can DAgger mitigate causal confusion?

# Why might we fail to fit the expert?

1. Non-Markovian behavior
2. Multimodal behavior



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

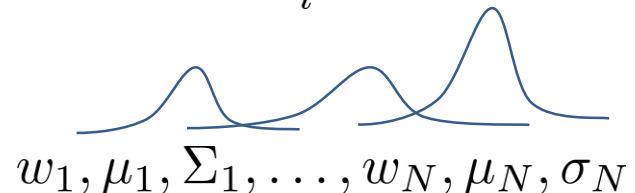


# Why might we fail to fit the expert?

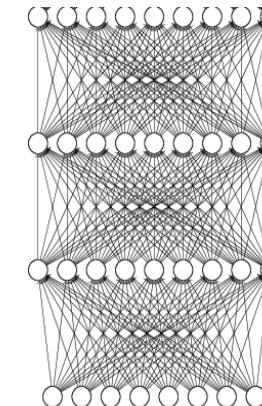


1. Output mixture of Gaussians
2. Latent variable models
3. Autoregressive discretization

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



$w_1, \mu_1, \Sigma_1, \dots, w_N, \mu_N, \sigma_N$

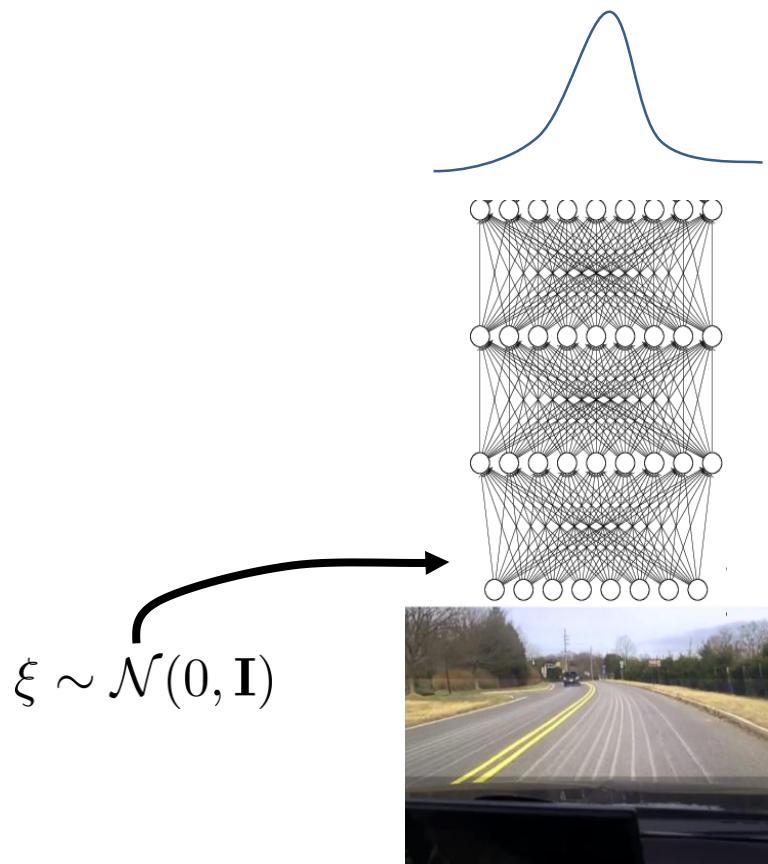


# Why might we fail to fit the expert?

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

Look up some of these:

- Conditional variational autoencoder
- Normalizing flow/realNVP
- Stein variational gradient descent

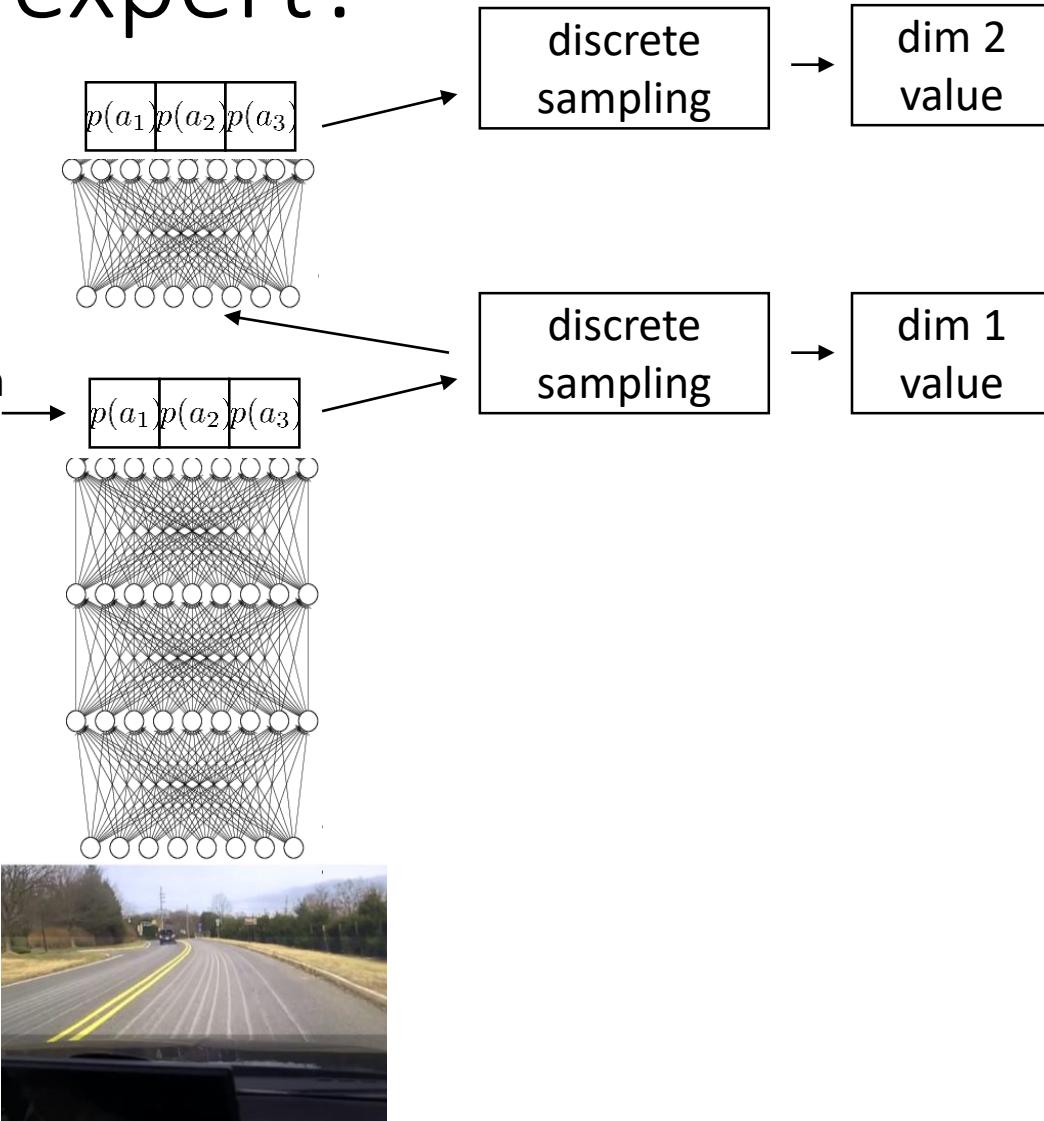


$$\xi \sim \mathcal{N}(0, \mathbf{I})$$

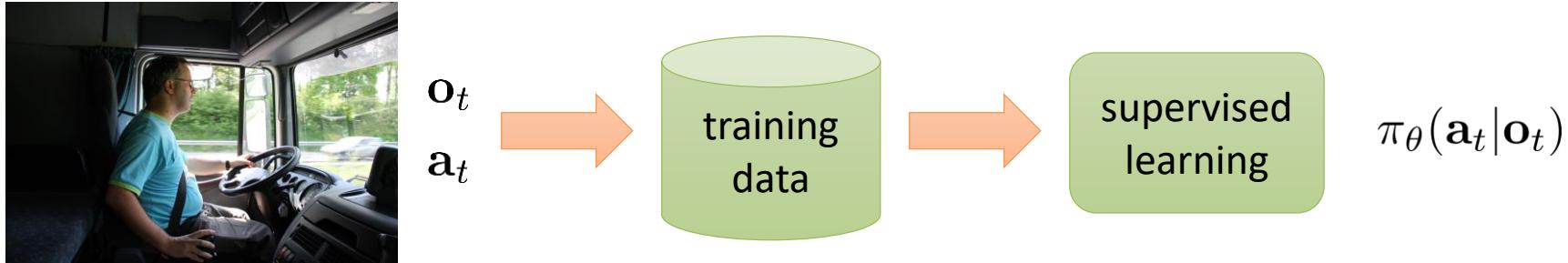
# Why might we fail to fit the expert?

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

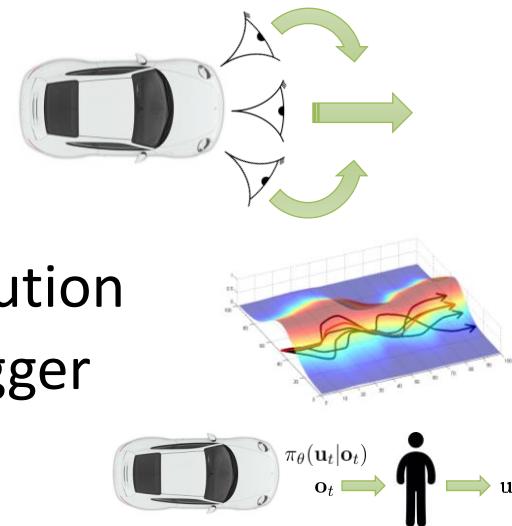
(discretized) distribution  
over dimension 1 **only**



# Imitation learning: recap



- Often (but not always) insufficient by itself
  - Distribution mismatch problem
- Sometimes works well
  - Hacks (e.g. left/right images)
  - Samples from a stable trajectory distribution
  - Add more **on-policy** data, e.g. using Dagger
  - Better models that fit more accurately



Break

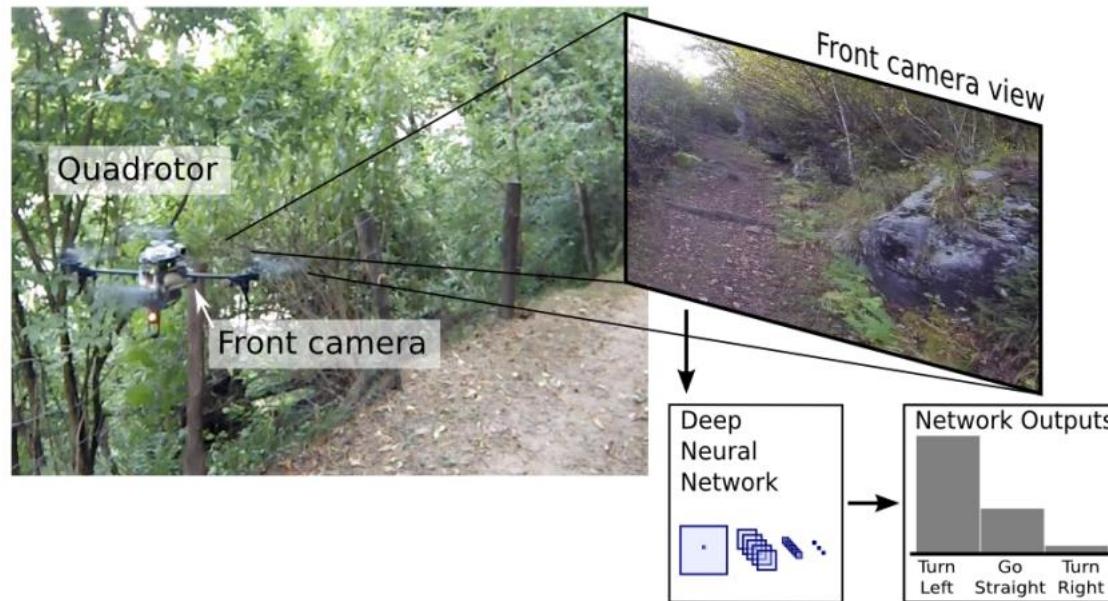
# Case study 1: trail following as classification

## A Machine Learning Approach to Visual Perception of Forest Trails for Mobile Robots

Alessandro Giusti<sup>1</sup>, Jérôme Guzzi<sup>1</sup>, Dan C. Cireşan<sup>1</sup>, Fang-Lin He<sup>1</sup>, Juan P. Rodríguez<sup>1</sup>

Flavio Fontana<sup>2</sup>, Matthias Faessler<sup>2</sup>, Christian Forster<sup>2</sup>

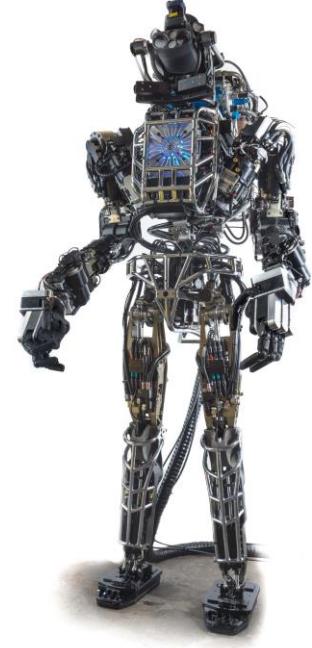
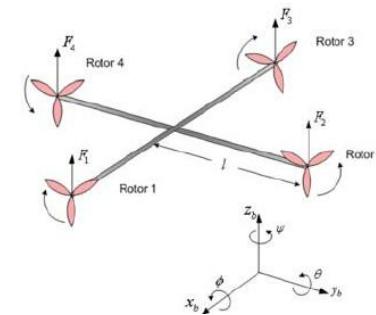
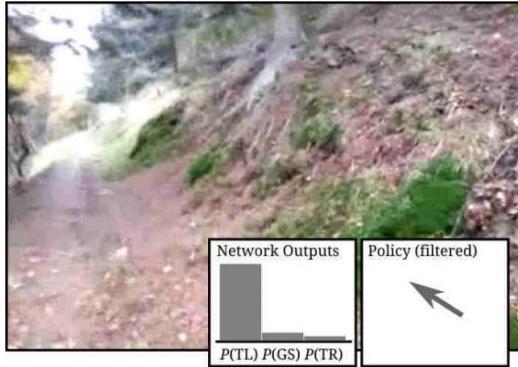
Jürgen Schmidhuber<sup>1</sup>, Gianni Di Caro<sup>1</sup>, Davide Scaramuzza<sup>2</sup>, Luca M. Gambardella<sup>1</sup>





# Imitation learning: what's the problem?

- Humans need to provide data, which is typically finite
  - Deep learning works best when data is plentiful
- Humans are not good at providing some kinds of actions



- Humans can learn autonomously; can our machines do the same?
  - Unlimited data from own experience
  - Continuous self-improvement

# Terminology & notation



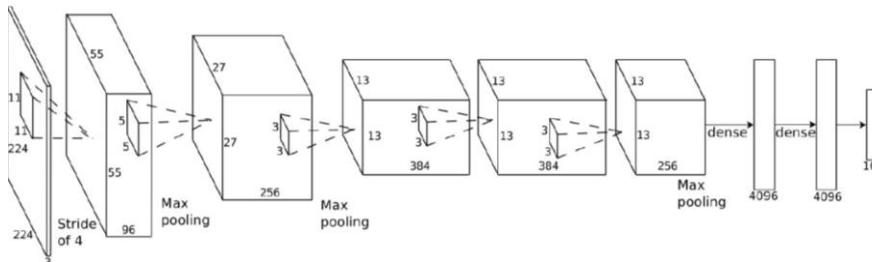
$\mathbf{o}_t$



$\mathbf{s}_t$  – state

$\mathbf{o}_t$  – observation

$\mathbf{a}_t$  – action



$\pi_{\theta}(\mathbf{a}_t | \mathbf{o}_t)$



$\mathbf{a}_t$



$c(\mathbf{s}_t, \mathbf{a}_t)$  – cost function

$r(\mathbf{s}_t, \mathbf{a}_t)$  – reward function

$$\min_{\theta} E_{\mathbf{a}_t \sim \pi_{\theta}(\mathbf{a}_t | \mathbf{o}_t), \mathbf{s}_t \sim p(\mathbf{s}_t | \mathbf{s}_{t-1}, \mathbf{a}_{t-1})} \left[ \sum_t \delta(s_t, a_t) r(s_t, a_t) + \beta \text{eyeiger}(s_t, a_t) \right]$$

[  
   $\sum_t$  [  
     $\delta(s_t, a_t) r(s_t, a_t)$ ]  
   $+ \beta \text{eyeiger}(s_t, a_t)$ ]  
]  
]

# Aside: notation

$s_t$  – state

$a_t$  – action

$r(s, a)$  – reward function



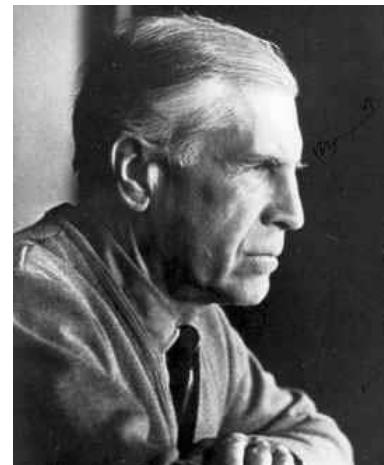
Richard Bellman

$x_t$  – state

$u_t$  – action

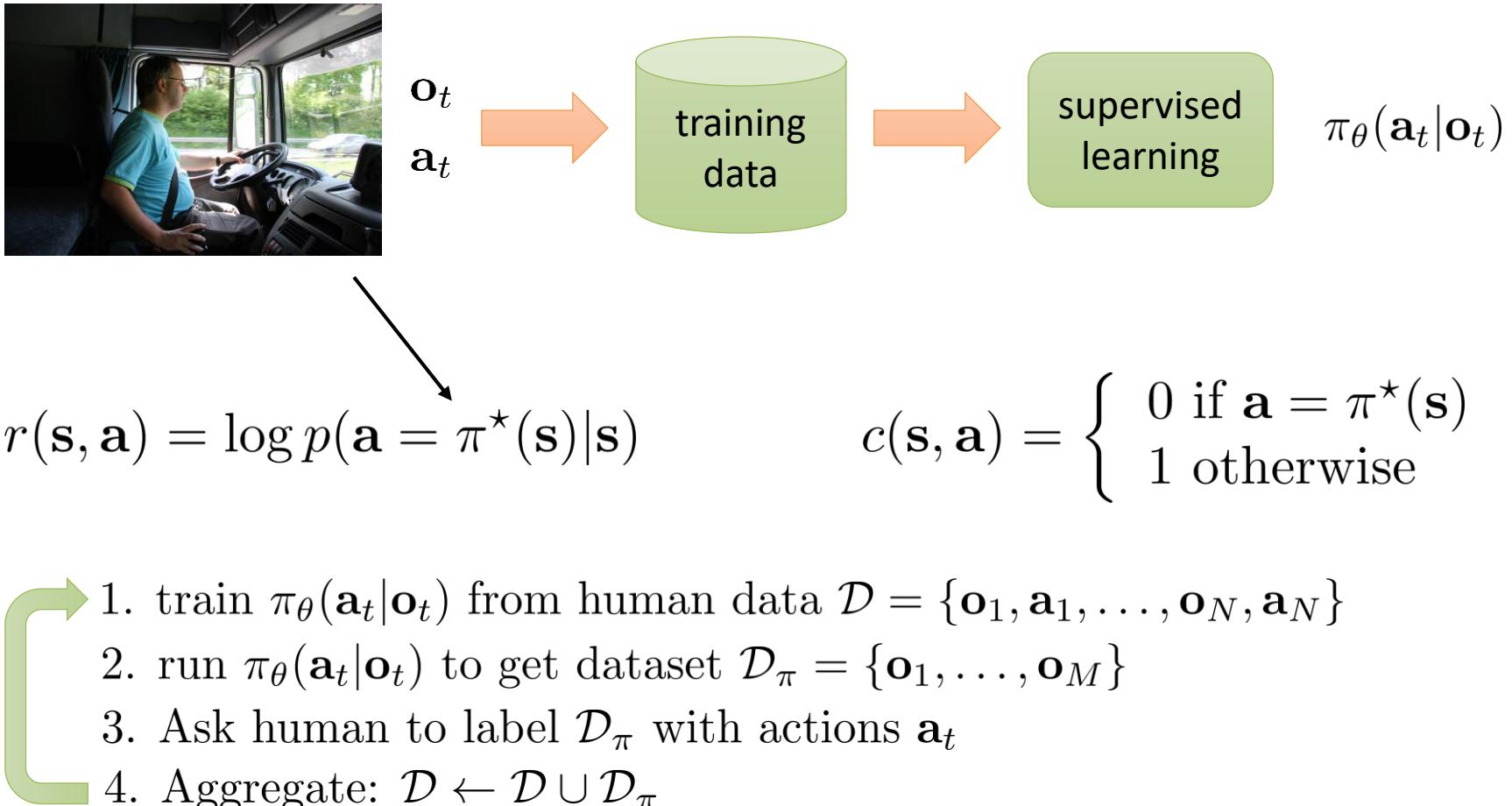
$c(x, u)$  – cost function

$$r(s, a) = -c(x, u)$$

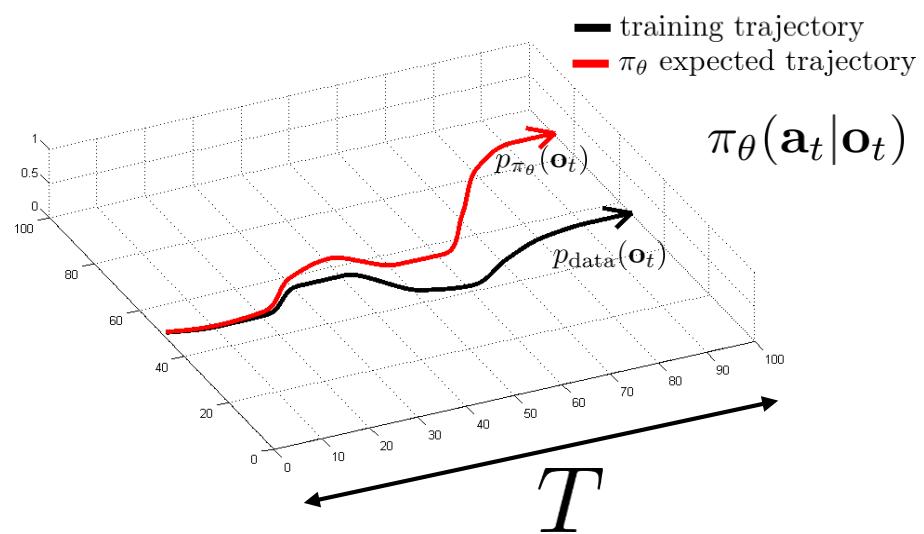


Lev Pontryagin

# A cost function for imitation?



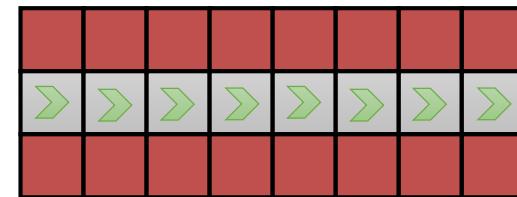
# Some analysis



## How bad is it?

$$c(\mathbf{s}, \mathbf{a}) = \begin{cases} 0 & \text{if } \mathbf{a} = \pi^*(\mathbf{s}) \\ 1 & \text{otherwise} \end{cases}$$

assume:  $\pi_\theta(\mathbf{a} \neq \pi^*(\mathbf{s}) | \mathbf{s}) \leq \epsilon$   
for all  $\mathbf{s} \in \mathcal{D}_{\text{train}}$



$$E \left[ \sum_t c(\mathbf{s}_t, \mathbf{a}_t) \right] \leq \underbrace{\epsilon T +}_{T \text{ terms, each } O(\epsilon T)} O(\epsilon T^2)$$

# More general analysis

$$c(\mathbf{s}, \mathbf{a}) = \begin{cases} 0 & \text{if } \mathbf{a} = \pi^*(\mathbf{s}) \\ 1 & \text{otherwise} \end{cases}$$

assume:  $\pi_\theta(\mathbf{a} \neq \pi^*(\mathbf{s}) | \mathbf{s}) \leq \epsilon$

~~for all  $\mathbf{s} \in \mathcal{D}_{\text{train}}$  for  $\mathbf{s} \sim p_{\text{train}}(\mathbf{s})$~~

if  $p_{\text{train}}(\mathbf{s}) \neq p_\theta(\mathbf{s})$ :

$$p_\theta(\mathbf{s}_t) = \underbrace{(1 - \epsilon)^t p_{\text{train}}(\mathbf{s}_t)}_{\text{probability we made no mistakes}} + \underbrace{(1 - (1 - \epsilon)^t)) p_{\text{mistake}}(\mathbf{s}_t)}_{\text{some } \textit{other} \text{ distribution}}$$

probability we made no mistakes

some *other* distribution

$$|p_\theta(\mathbf{s}_t) - p_{\text{train}}(\mathbf{s}_t)| = (1 - (1 - \epsilon)^t) |p_{\text{mistake}}(\mathbf{s}_t) - p_{\text{train}}(\mathbf{s}_t)| \leq 2(1 - (1 - \epsilon)^t)$$

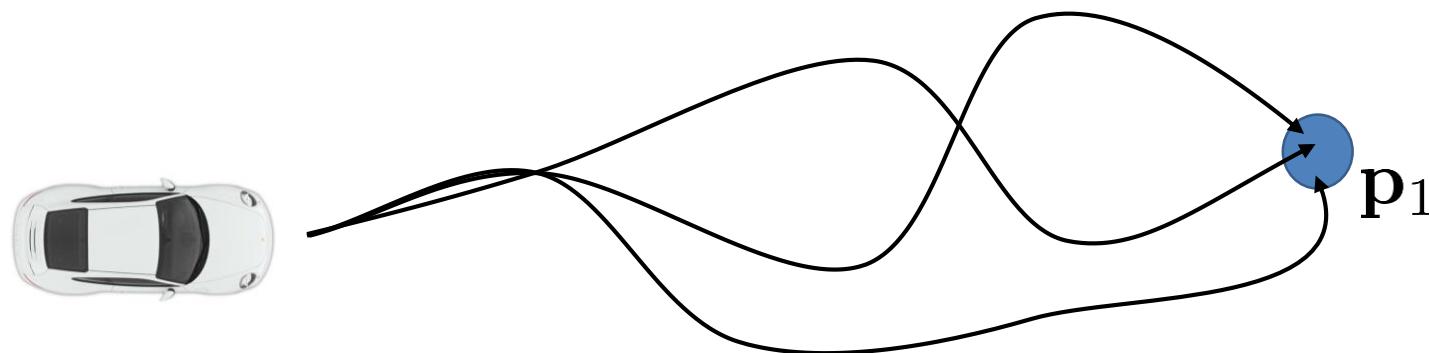
$$\text{useful identity: } (1 - \epsilon)^t \geq 1 - \epsilon t \text{ for } \epsilon \in [0, 1] \quad \leq 2\epsilon t$$

$$\begin{aligned} \sum_t E_{p_\theta(\mathbf{s}_t)}[c_t] &= \sum_t \sum_{\mathbf{s}_t} p_\theta(\mathbf{s}_t) c_t(\mathbf{s}_t) \leq \sum_t \sum_{\mathbf{s}_t} p_{\text{train}}(\mathbf{s}_t) c_t(\mathbf{s}_t) + |p_\theta(\mathbf{s}_t) - p_{\text{train}}(\mathbf{s}_t)| c_{\max} \\ &\leq \sum_t \epsilon + 2\epsilon t \quad O(\epsilon T^2) \end{aligned}$$

with DAgger,  $p_{\text{train}}(\mathbf{s}) \rightarrow p_\theta(\mathbf{s})$

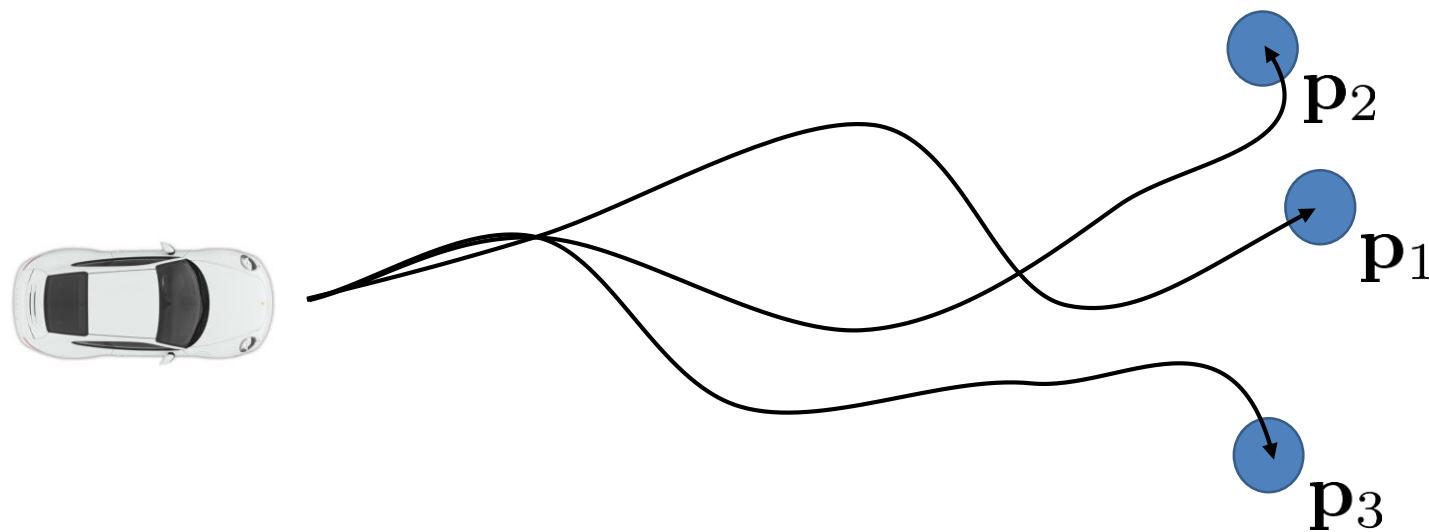
$$E \left[ \sum_t c(\mathbf{s}_t, \mathbf{a}_t) \right] \leq \epsilon T$$

# Another imitation idea



$$\pi_\theta(a|s)$$

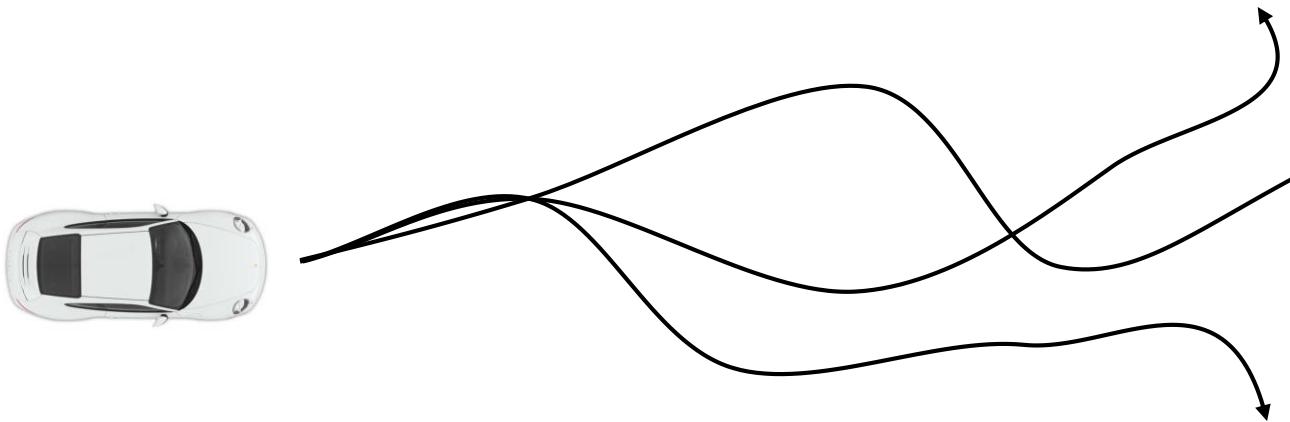
policy for reaching  $p_1$



$$\pi_\theta(a|s, p)$$

policy for reaching *any*  $p$

# Goal-conditioned behavioral cloning



training time:

demo 1:  $\{s_1, a_t, \dots, s_{T-1}, a_{T-1}, s_T\}$  ← successful demo for reaching  $s_T$

demo 2:  $\{s_1, a_t, \dots, s_{T-1}, a_{T-1}, s_T\}$

learn  $\pi_\theta(a|s, g)$

demo 3:  $\{s_1, a_t, \dots, s_{T-1}, a_{T-1}, s_T\}$

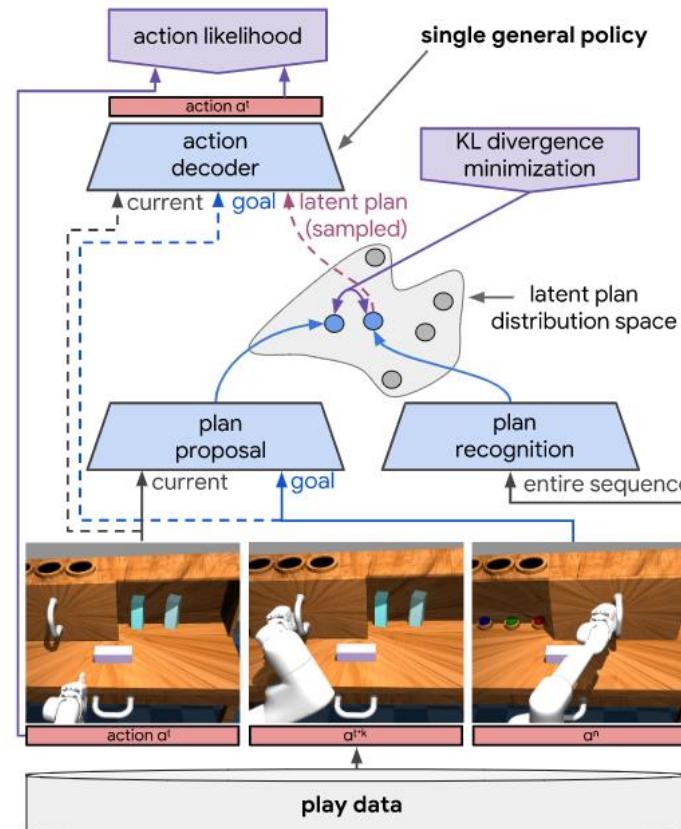
goal state

for each demo  $\{s_1^i, a_1^i, \dots, s_{T-1}^i, a_{T-1}^i, s_T^i\}$

maximize  $\log \pi_\theta(a_t^i | s_t^i, g = s_T^i)$

# Learning Latent Plans from Play

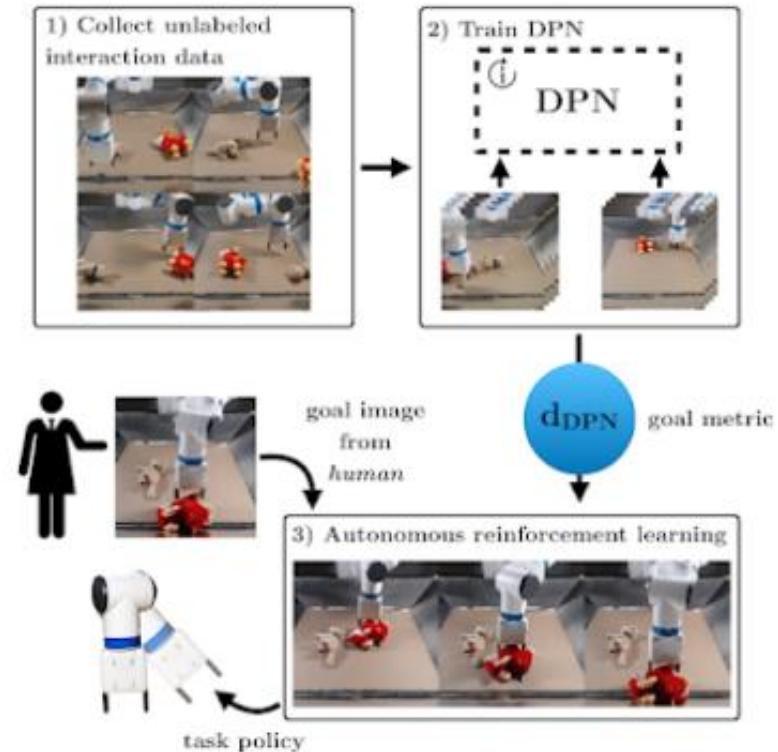
COREY LYNCH MOHI KHANSARI TED XIAO VIKASH KUMAR JONATHAN TOMPSON SERGEY LEVINE PIERRE SERMANET  
Google Brain Google X Google Brain Google Brain Google Brain Google Brain Google Brain



## Unsupervised Visuomotor Control through Distributional Planning Networks

Tianhe Yu, Gleb Shevchuk, Dorsa Sadigh, Chelsea Finn

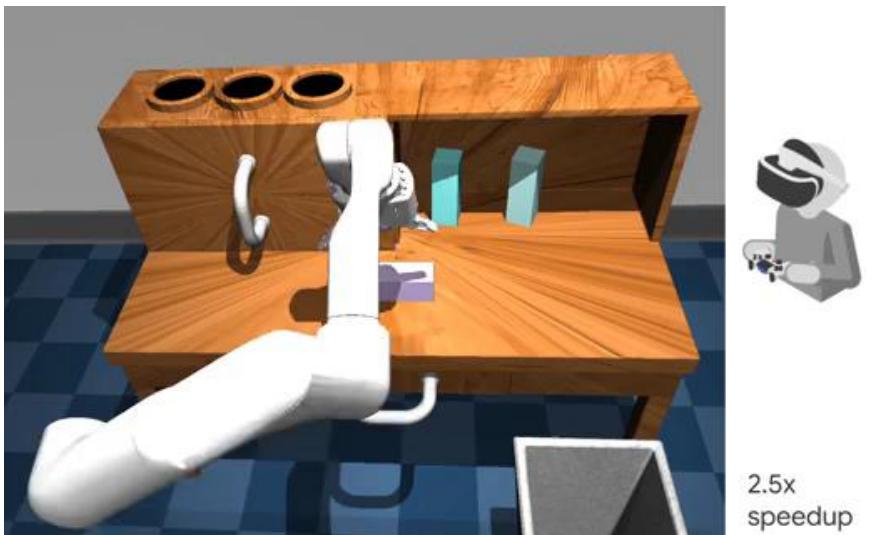
Stanford University



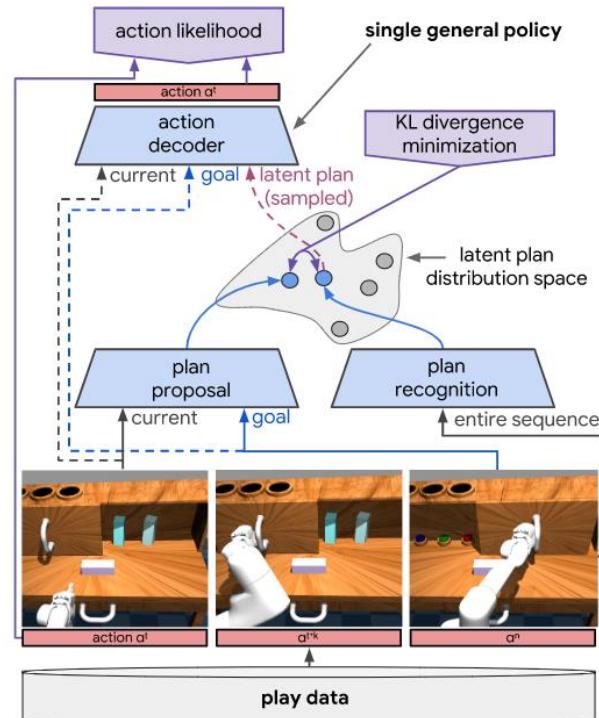
# Learning Latent Plans from Play

COREY LYNCH MOHI KHANSARI TED XIAO VIKASH KUMAR JONATHAN TOMPSON SERGEY LEVINE PIERRE SERMANET  
Google Brain Google X Google Brain Google Brain Google Brain Google Brain Google Brain

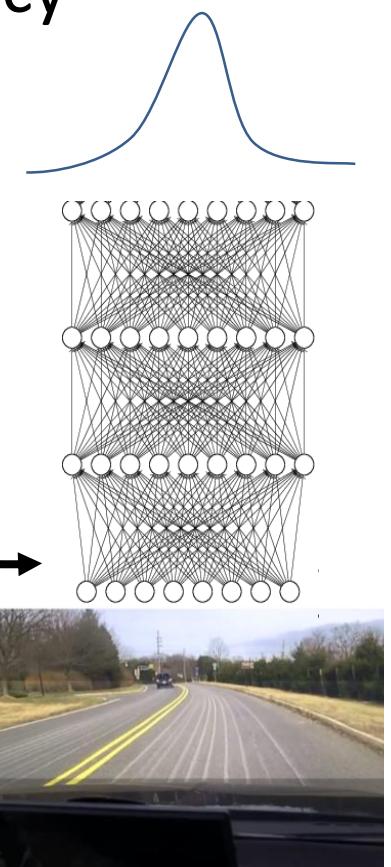
## 1. Collect data



## 2. Train goal conditioned policy



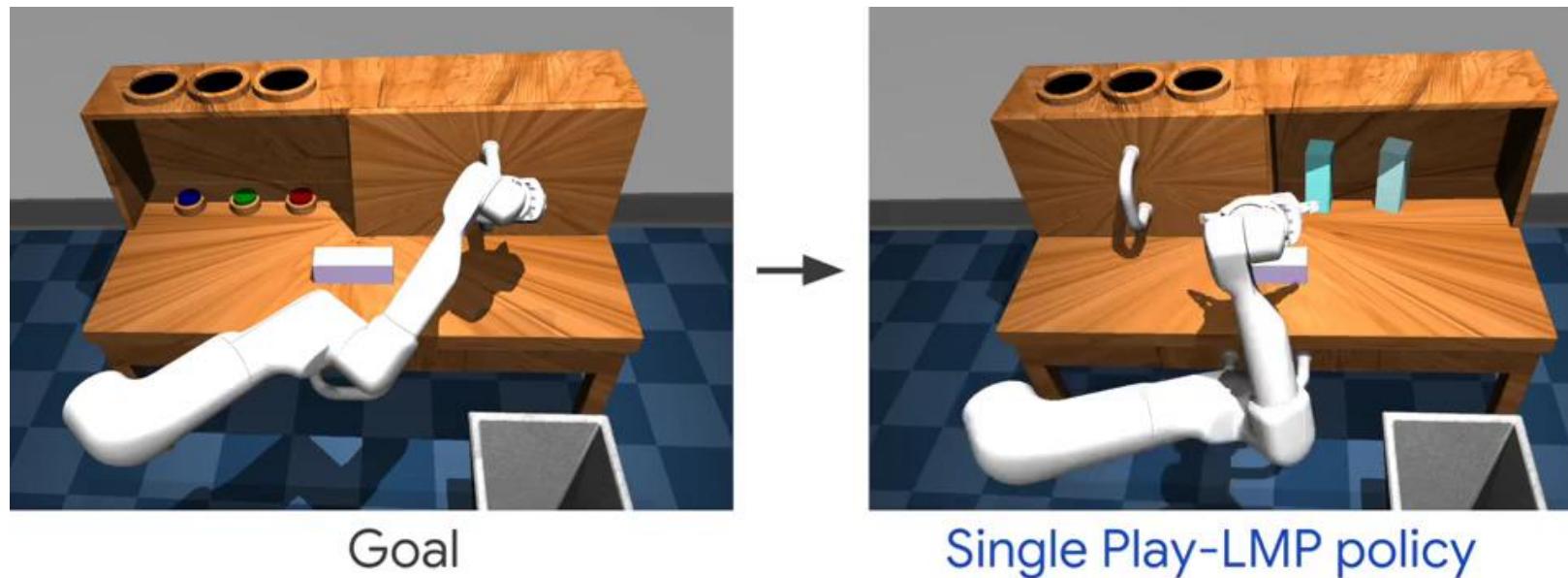
$$\xi \sim \mathcal{N}(0, \mathbf{I})$$



# Learning Latent Plans from Play

COREY LYNCH MOHI KHANSARI TED XIAO VIKASH KUMAR JONATHAN TOMPSON SERGEY LEVINE PIERRE SERMANET  
Google Brain Google X Google Brain Google Brain Google Brain Google Brain Google Brain

## 3. Reach goals



# Terminology & notation



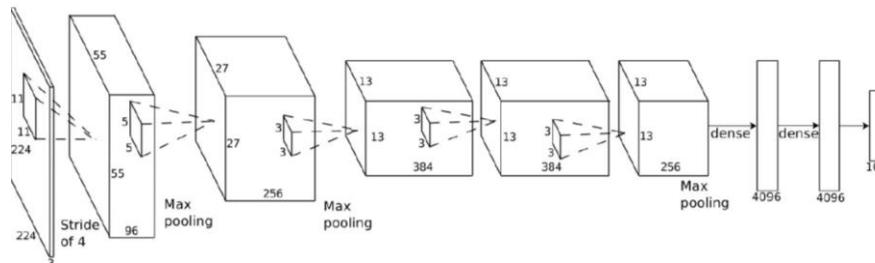
$\mathbf{o}_t$



$\mathbf{s}_t$  – state

$\mathbf{o}_t$  – observation

$\mathbf{a}_t$  – action



$\pi_\theta(\mathbf{a}_t | \mathbf{o}_t)$



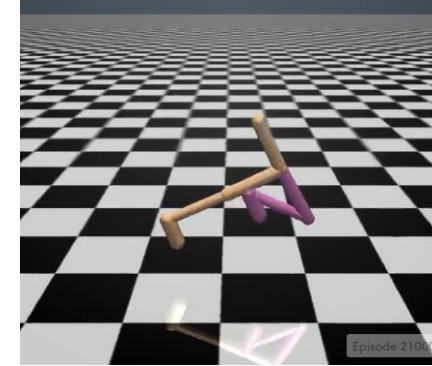
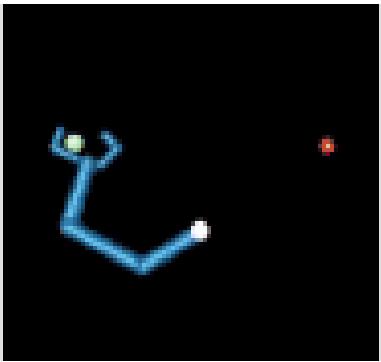
$\mathbf{a}_t$

$c(\mathbf{s}_t, \mathbf{a}_t)$  – cost function

$r(\mathbf{s}_t, \mathbf{a}_t)$  – reward function

$$\min_{\theta} E_{\mathbf{s}_{1:T}, \mathbf{a}_{1:T}} \left[ \sum_t c(\mathbf{s}_t, \mathbf{a}_t) \right]$$

# Cost/reward functions in theory and practice



$$r(\mathbf{s}, \mathbf{a}) = \begin{cases} 1 & \text{if object at target} \\ 0 & \text{otherwise} \end{cases}$$

$$r(\mathbf{s}, \mathbf{a}) = \begin{cases} 1 & \text{if walker is running} \\ 0 & \text{otherwise} \end{cases}$$

$$\begin{aligned} r(\mathbf{s}, \mathbf{a}) = & -w_1 \|p_{\text{gripper}}(\mathbf{s}) - p_{\text{object}}(\mathbf{s})\|^2 + \\ & -w_2 \|p_{\text{object}}(\mathbf{s}) - p_{\text{target}}(\mathbf{s})\|^2 + \\ & -w_3 \|\mathbf{a}\|^2 \end{aligned}$$

$$\begin{aligned} r(\mathbf{s}, \mathbf{a}) = & w_1 v(\mathbf{s}) + \\ & w_2 \delta(|\theta_{\text{torso}}(\mathbf{s})| < \epsilon) + \\ & w_3 \delta(h_{\text{torso}}(\mathbf{s}) \geq h) \end{aligned}$$