

# Supervised Learning of Behaviors

CS 294-112: Deep Reinforcement Learning

Sergey Levine

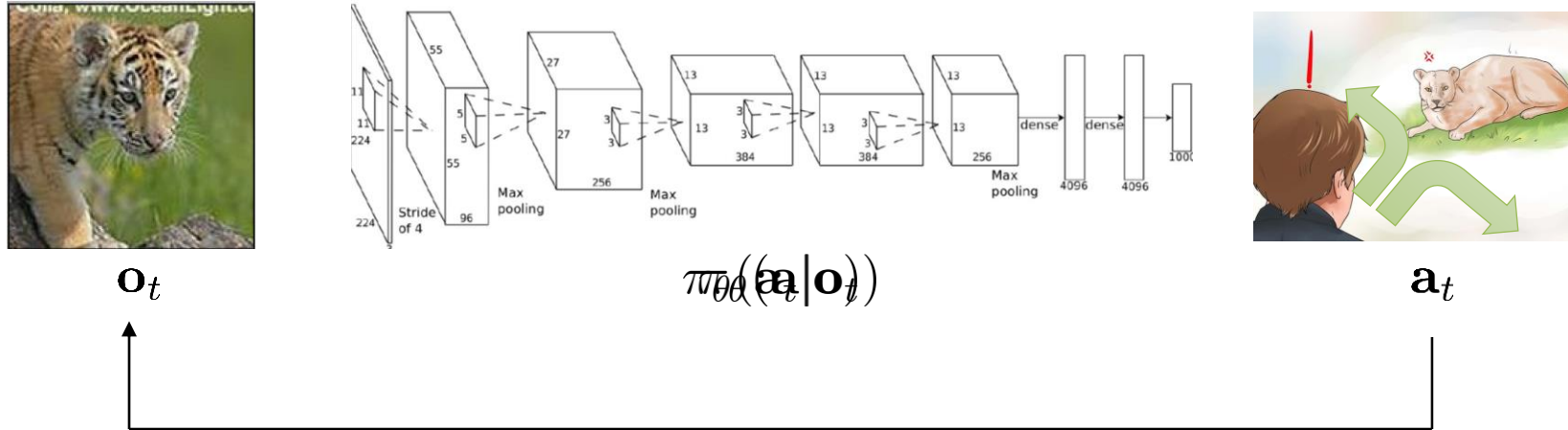
# Class Notes

1. Make sure you sign up for Piazza!
2. Homework 1 is now out
3. Remember to start forming final project groups

# 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. Case studies of recent work in (deep) imitation learning
  4. What is missing from imitation learning?
- Goals:
    - Understand definitions & notation
    - Understand basic imitation learning algorithms
    - Understand their strengths & weaknesses

# Terminology & notation



$\mathbf{s}_t$  – 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)

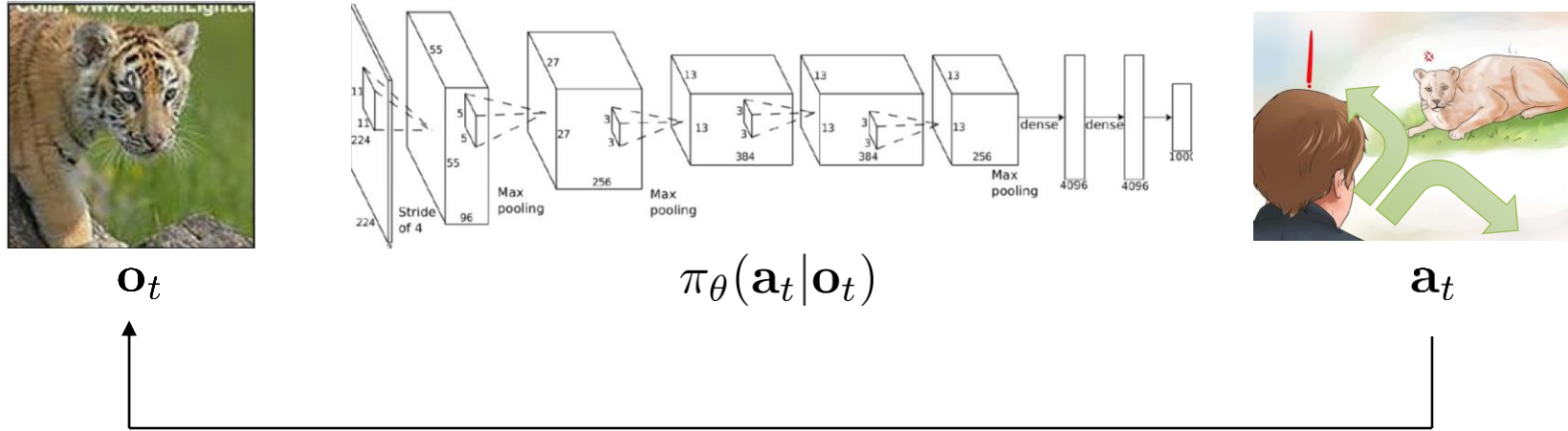


$\mathbf{o}_t$  – observation



$\mathbf{s}_t$  – state

# Terminology & notation



$\mathbf{o}_t$

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

$\mathbf{a}_t$

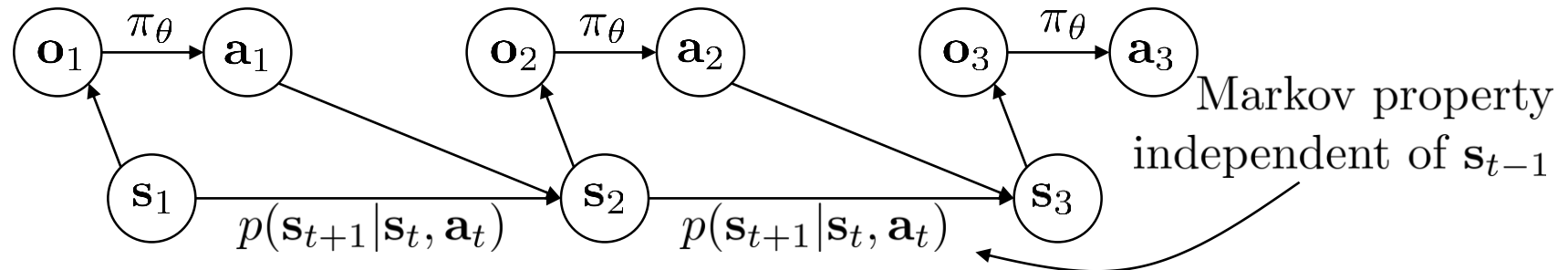
$\mathbf{s}_t$  – state

$\mathbf{o}_t$  – 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)



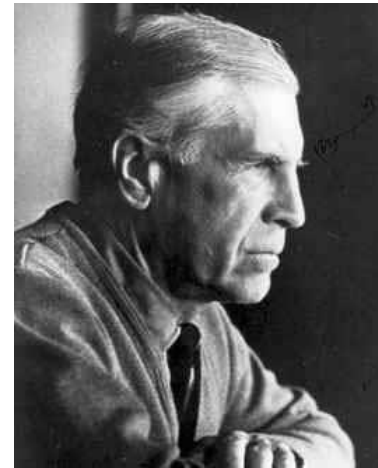
# Aside: notation

$\mathbf{s}_t$  – state  
 $\mathbf{a}_t$  – action

$\mathbf{x}_t$  – state  
 $\mathbf{u}_t$  – action    управление

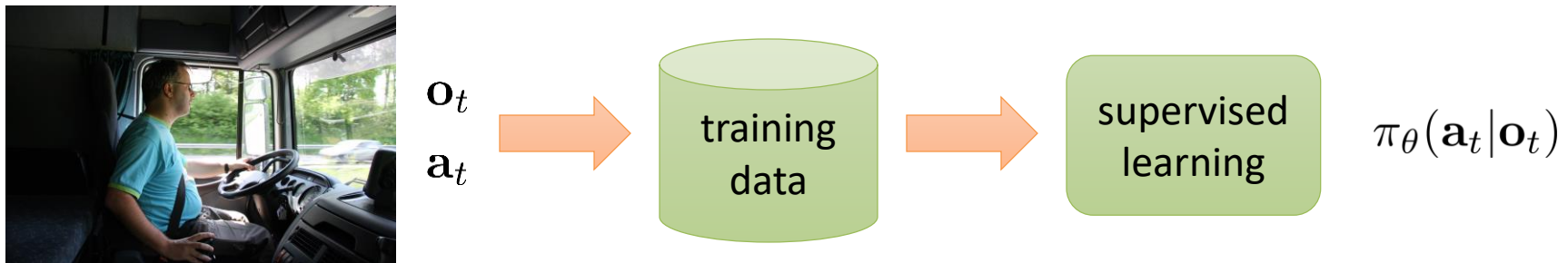
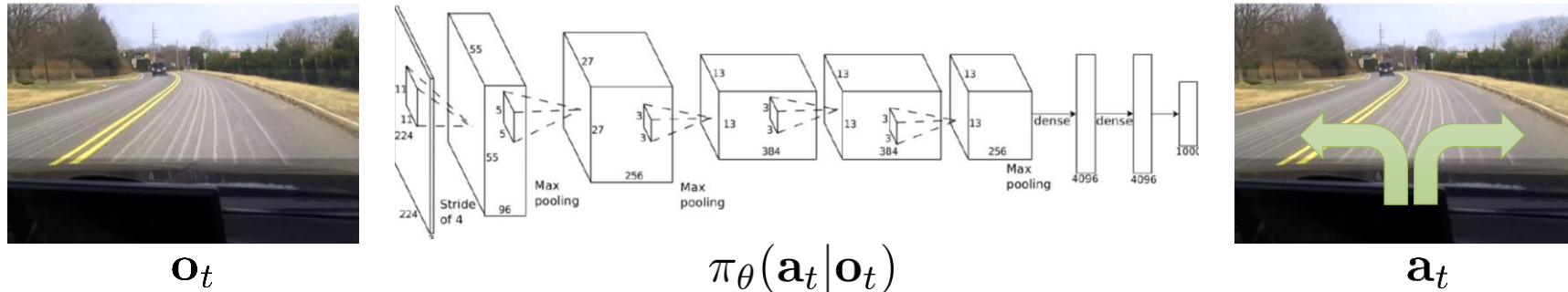


Richard Bellman



Lev Pontryagin

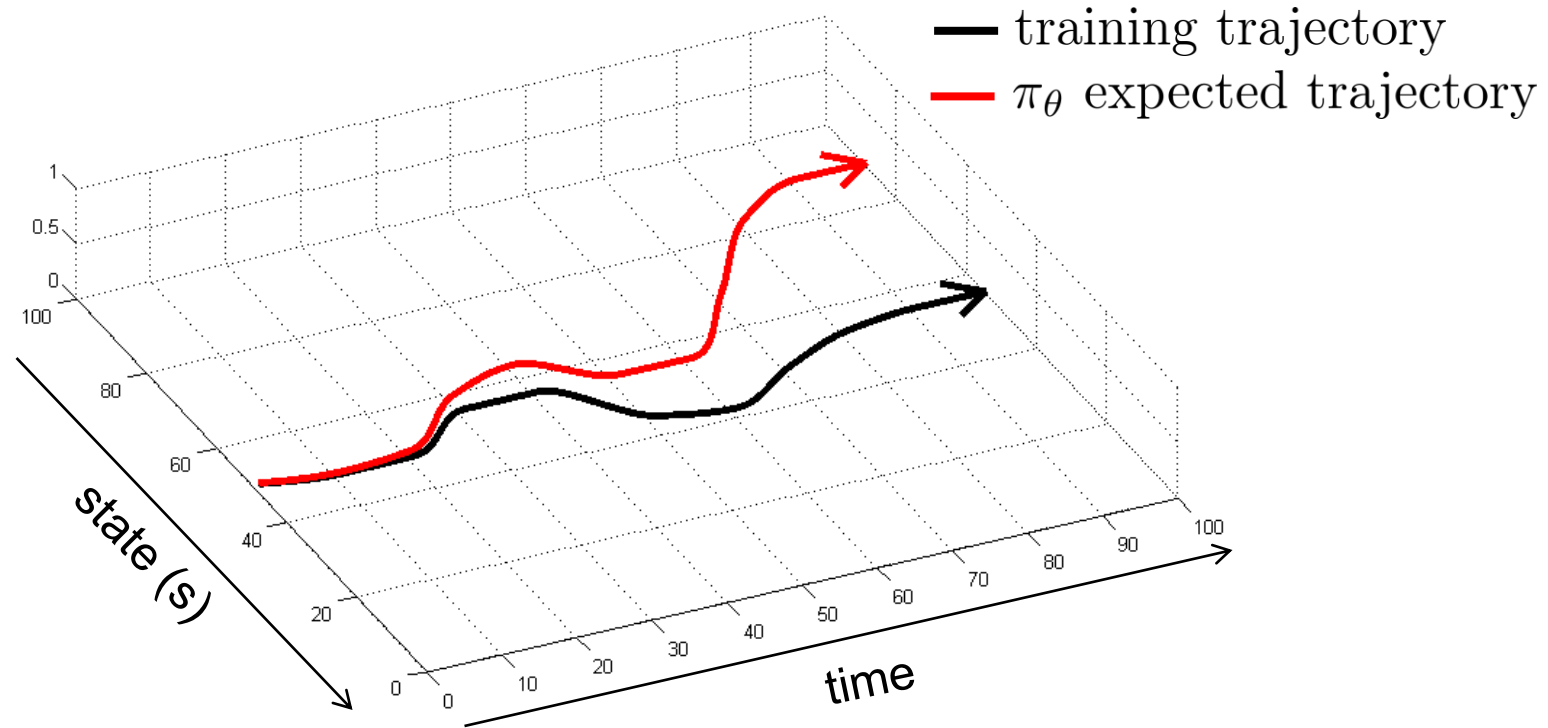
# Imitation Learning



behavior cloning

Does it work?

No!



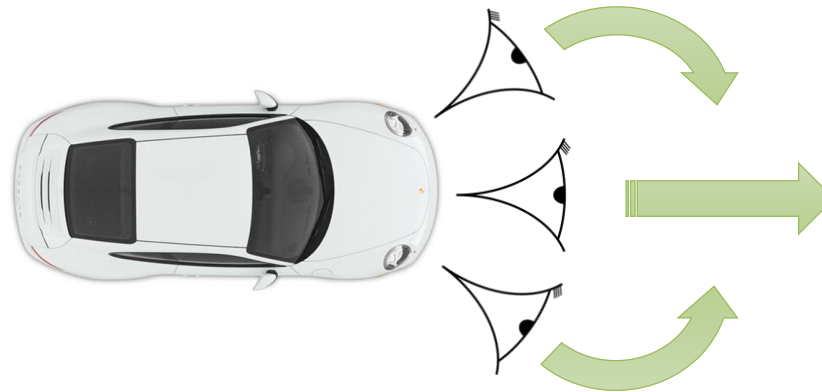
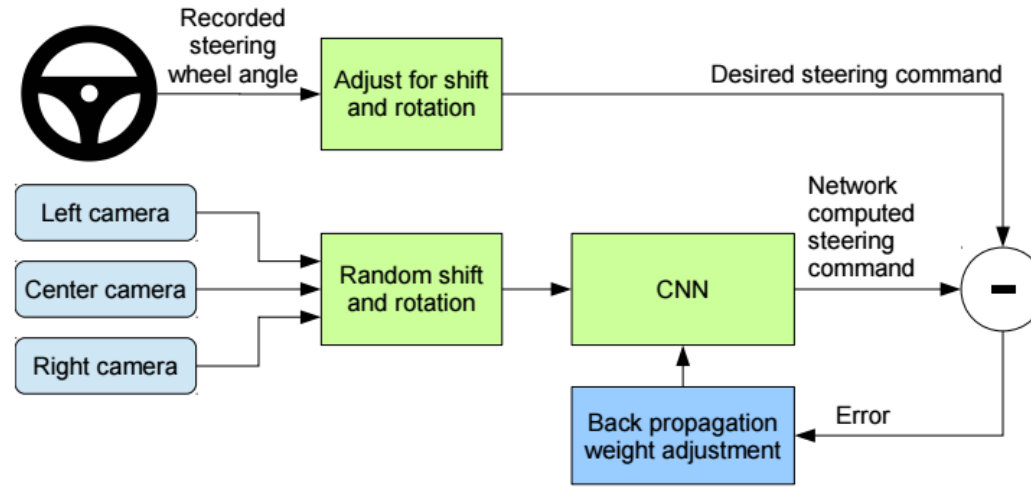


Does it work?

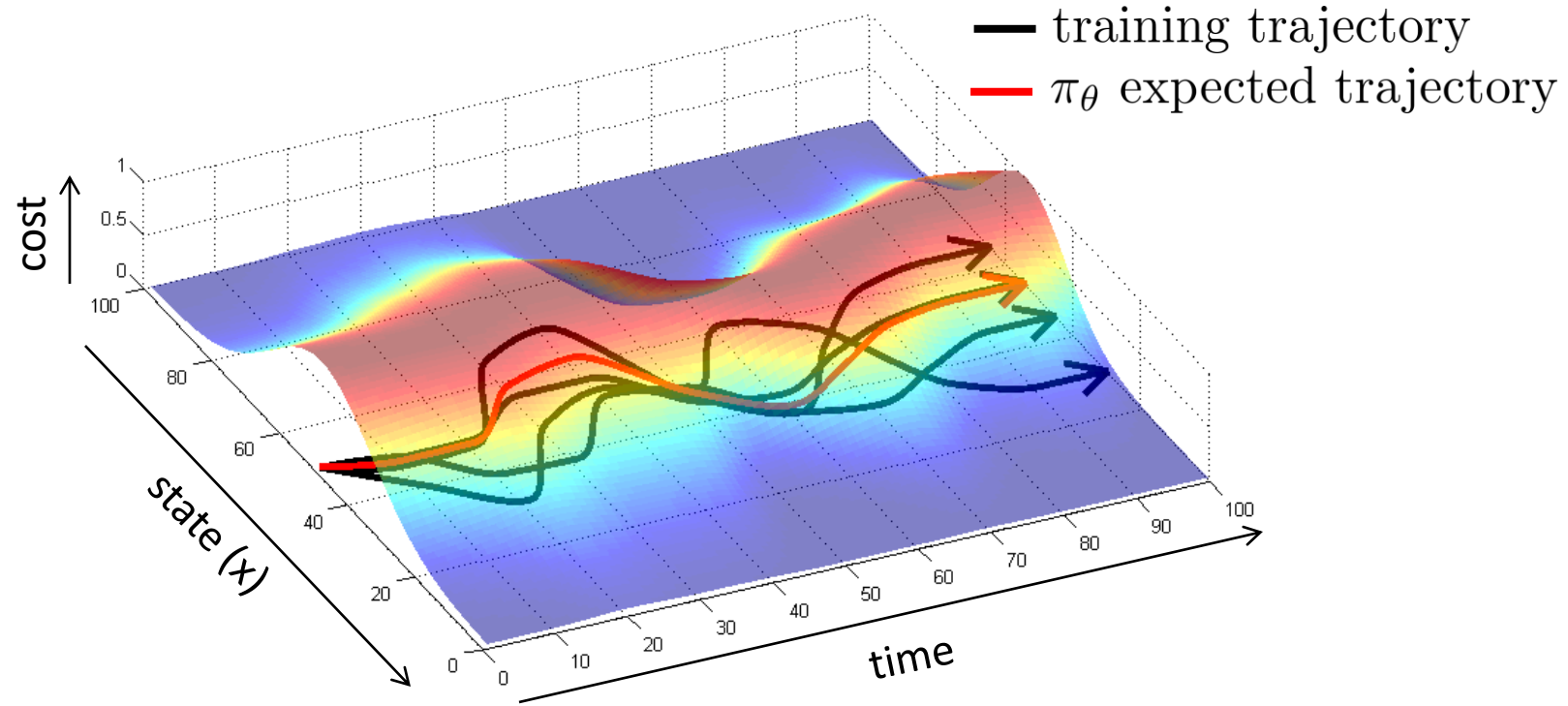
Yes!



# Why did that work?



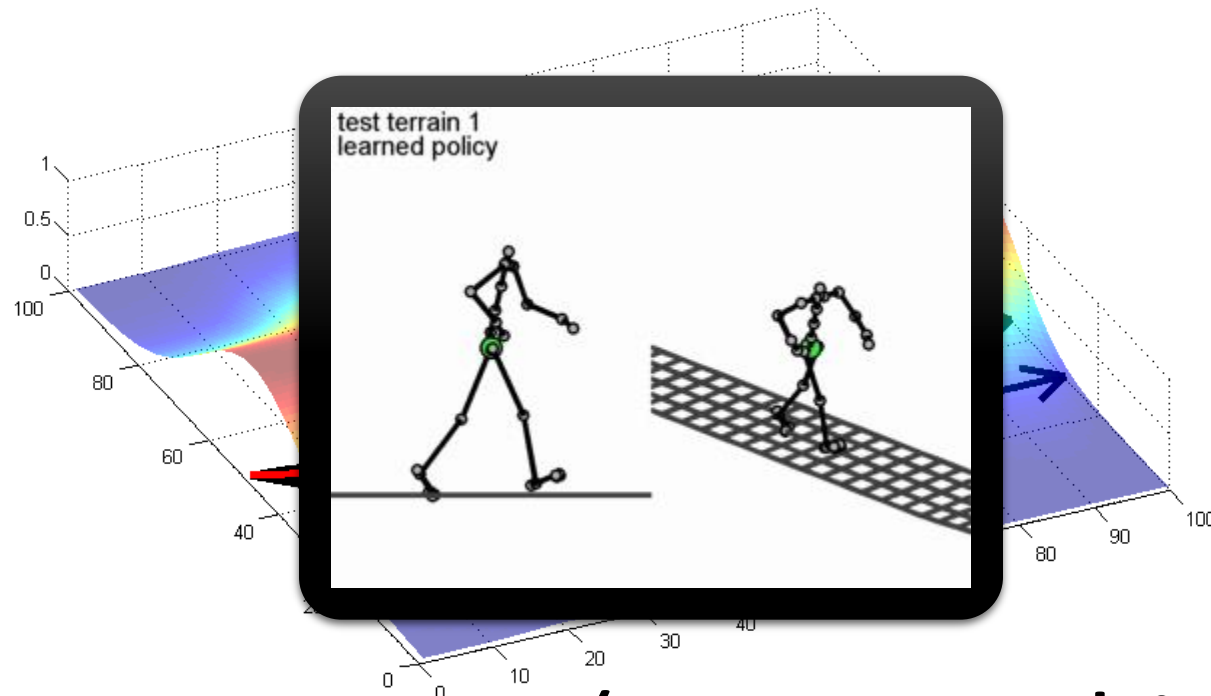
Can we make it work more often?



stability

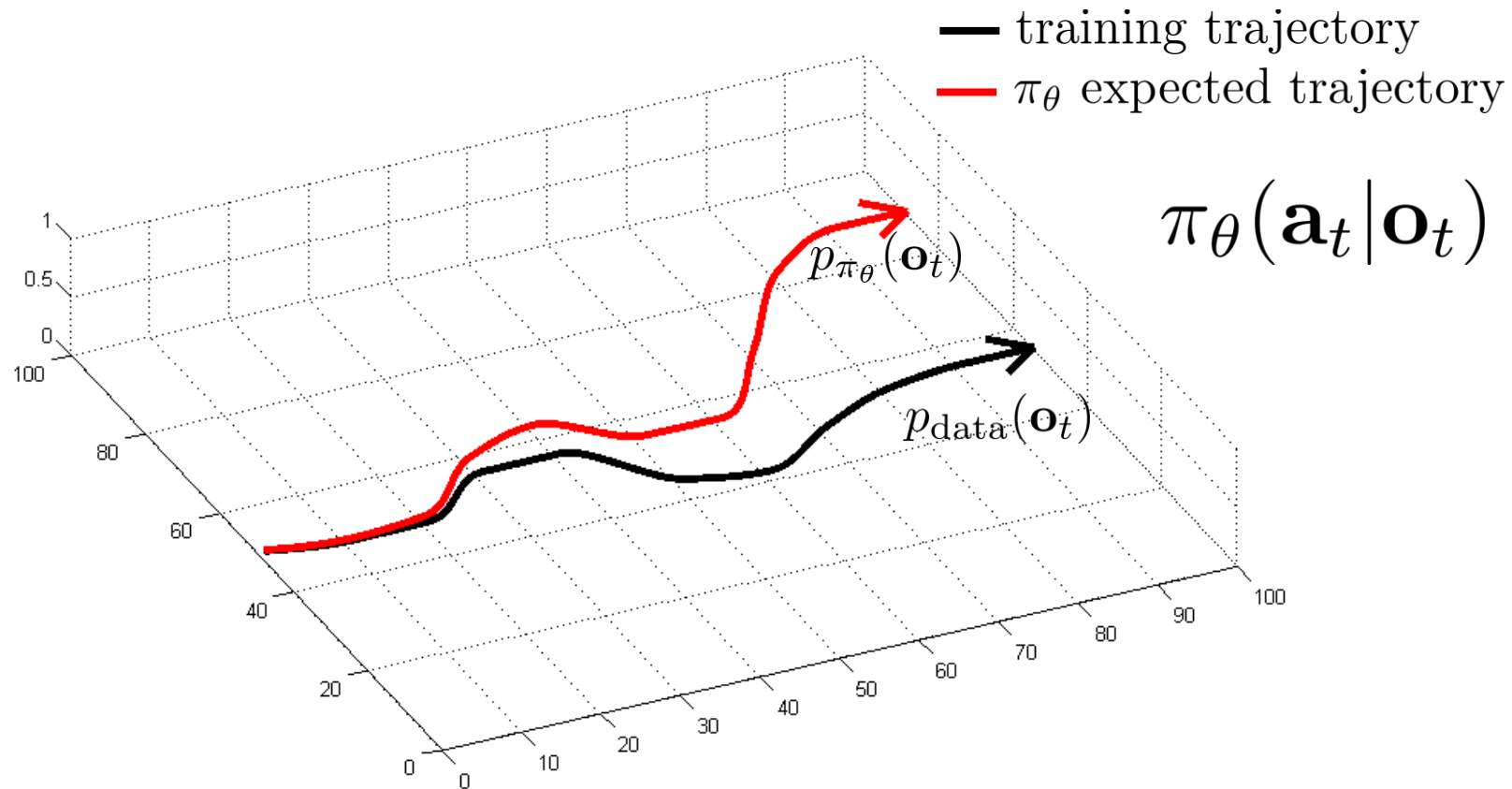
# Learning from a stabilizing controller

$p(\mathbf{s}_1)$ , a Gaussian distribution obtained using variant of iterative LQR



(more on this later)

# Can we make it work more often?



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

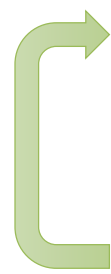
idea: instead of being clever about  $p_{\pi_\theta}(\mathbf{o}_t)$ , be clever about  $p_{\text{data}}(\mathbf{o}_t)$ !

## **D**Agger: 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  $\mathbf{a}_t$ !

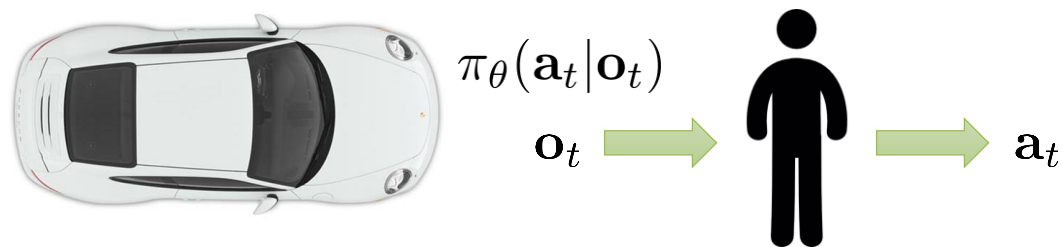
- 
1. train  $\pi_\theta(\mathbf{a}_t|\mathbf{o}_t)$  from human data  $\mathcal{D} = \{\mathbf{o}_1, \mathbf{a}_1, \dots, \mathbf{o}_N, \mathbf{a}_N\}$
  2. run  $\pi_\theta(\mathbf{a}_t|\mathbf{o}_t)$  to get dataset  $\mathcal{D}_\pi = \{\mathbf{o}_1, \dots, \mathbf{o}_M\}$
  3. Ask human to label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

# Dagger Example



# What's the problem?

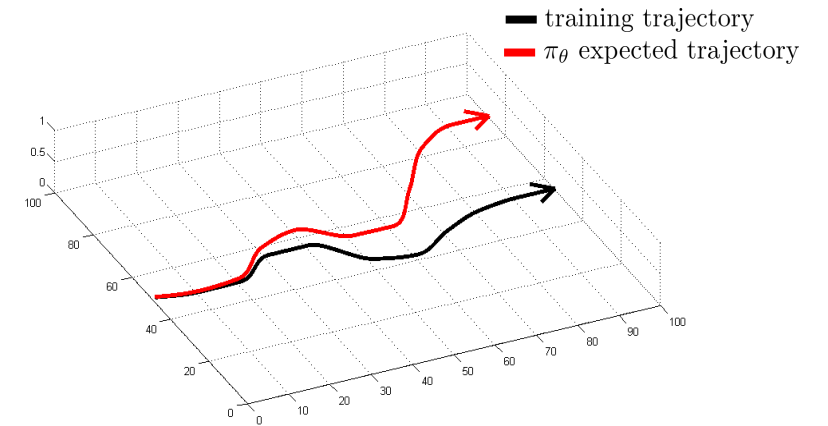
1. train  $\pi_\theta(\mathbf{a}_t|\mathbf{o}_t)$  from human data  $\mathcal{D} = \{\mathbf{o}_1, \mathbf{a}_1, \dots, \mathbf{o}_N, \mathbf{a}_N\}$
2. run  $\pi_\theta(\mathbf{a}_t|\mathbf{o}_t)$  to get dataset  $\mathcal{D}_\pi = \{\mathbf{o}_1, \dots, \mathbf{o}_M\}$
3. Ask human to label  $\mathcal{D}_\pi$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



# 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

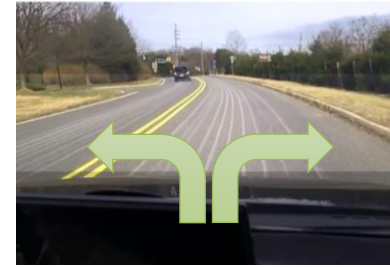
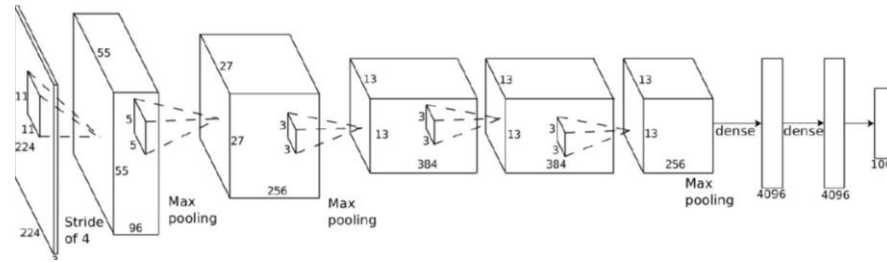
If we see the same thing  
twice, we do the same thing  
twice, regardless of what  
happened before

Often very unnatural for  
human demonstrators

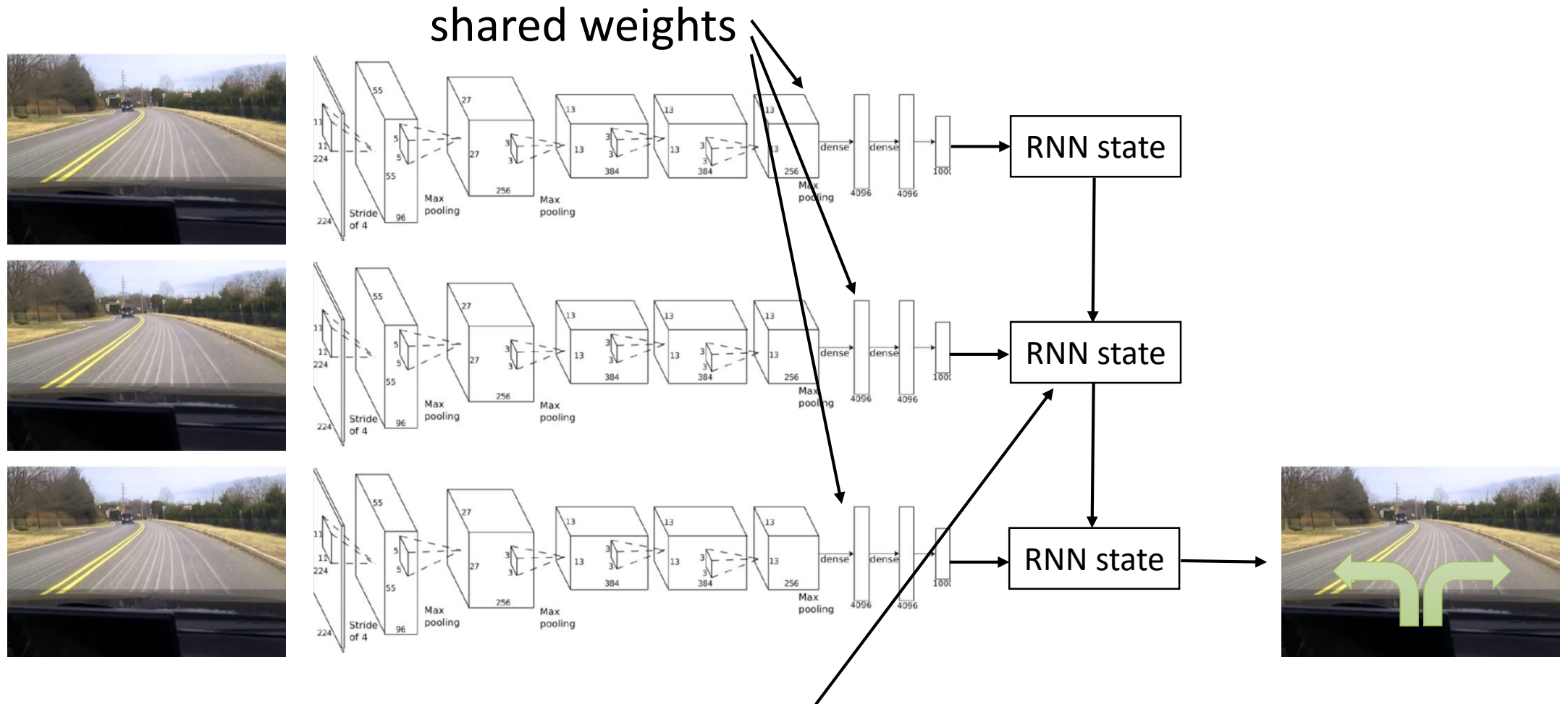
# How can we use the whole history?



variable number of frames,  
too many weights



# How can we use the whole history?

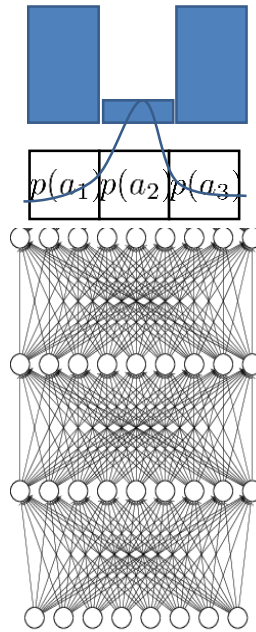
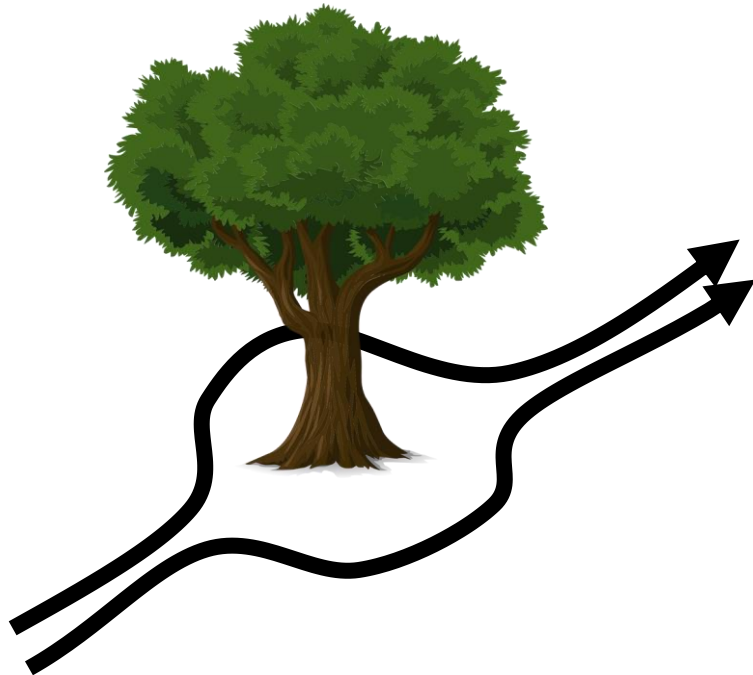


Typically, LSTM cells work better here

# Why might we fail to fit the expert?

1. Non-Markovian behavior

➔ 2. Multimodal behavior



1. Output mixture of Gaussians

2. Latent variable models

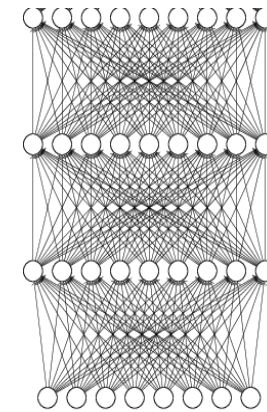
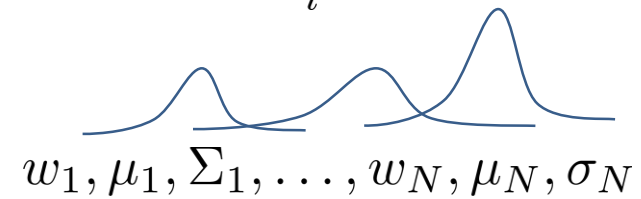
3. Autoregressive discretization



# Why might we fail to fit the expert?

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

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



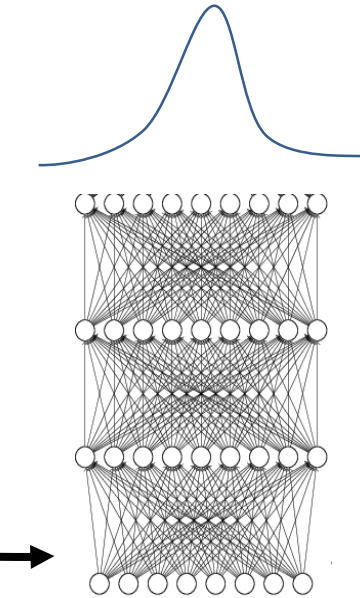
# Why might we fail to fit the expert?

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

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



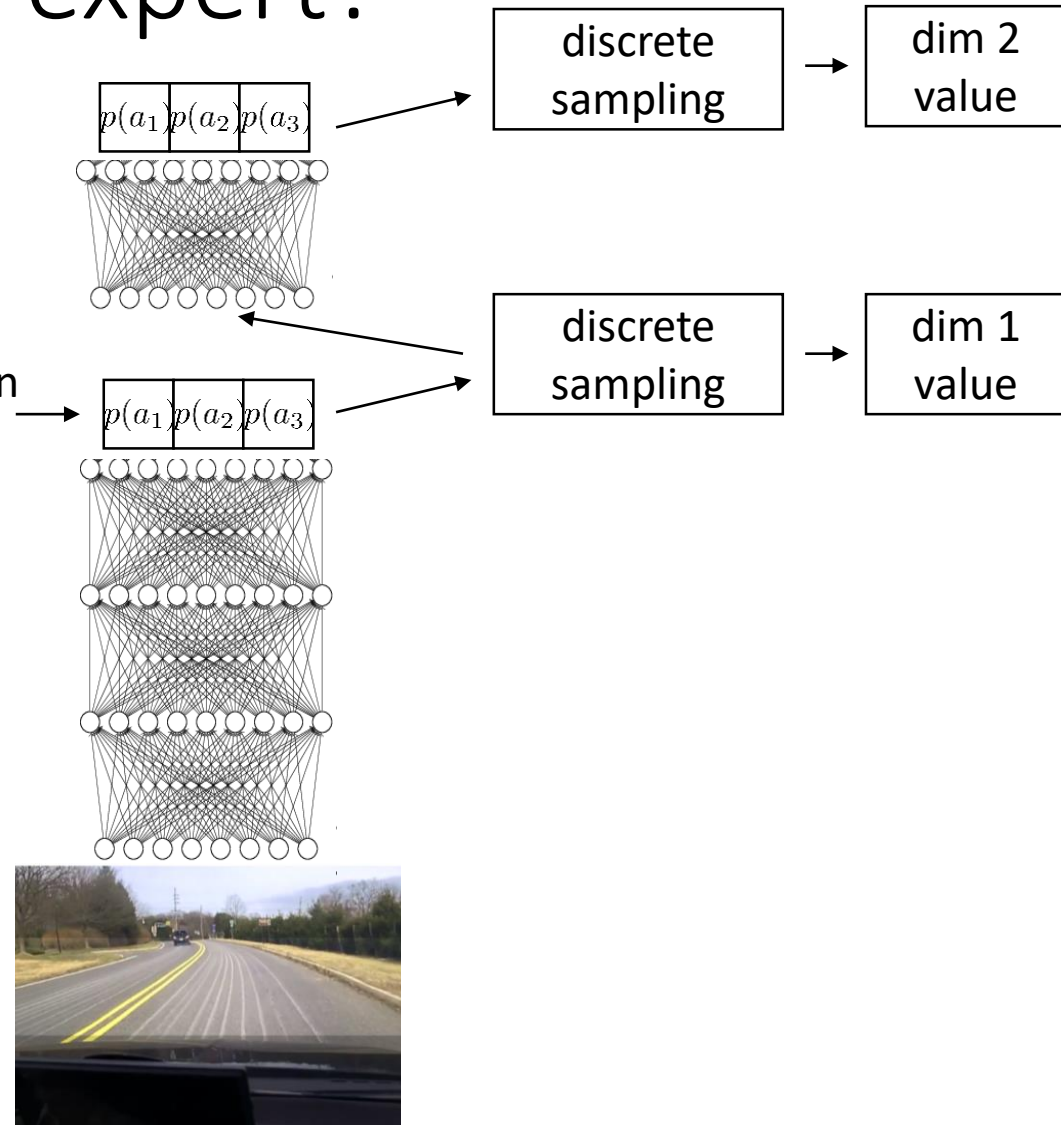
# Why might we fail to fit the expert?

1. Output mixture of Gaussians

2. Latent variable models

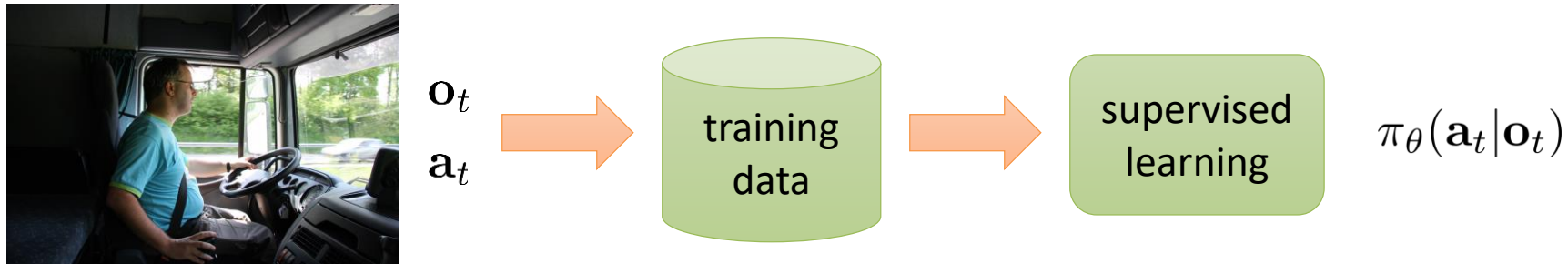
➔ 3. Autoregressive discretization

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

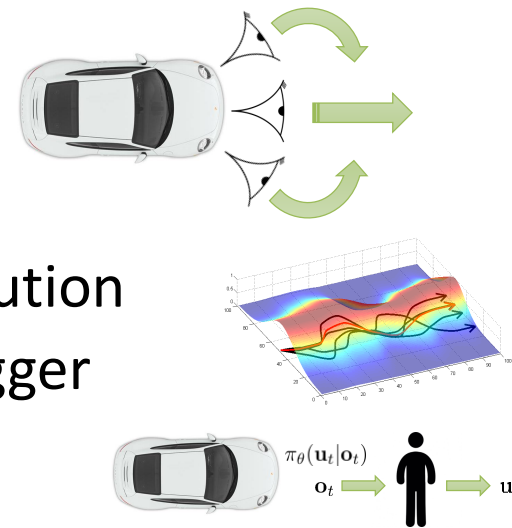




# Imitation learning: recap



- 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

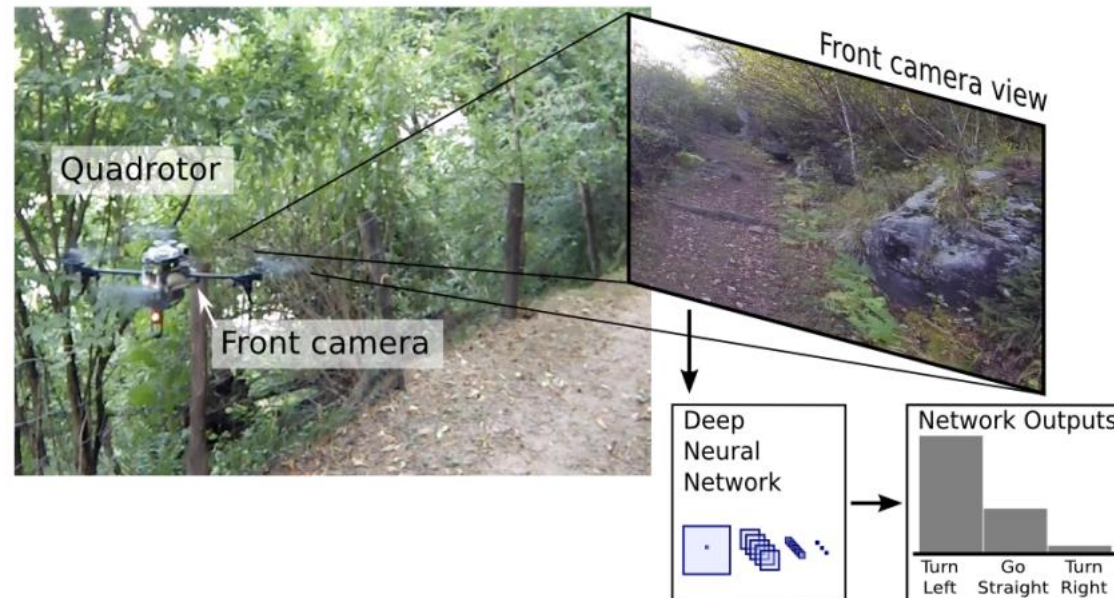


Break

# 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>

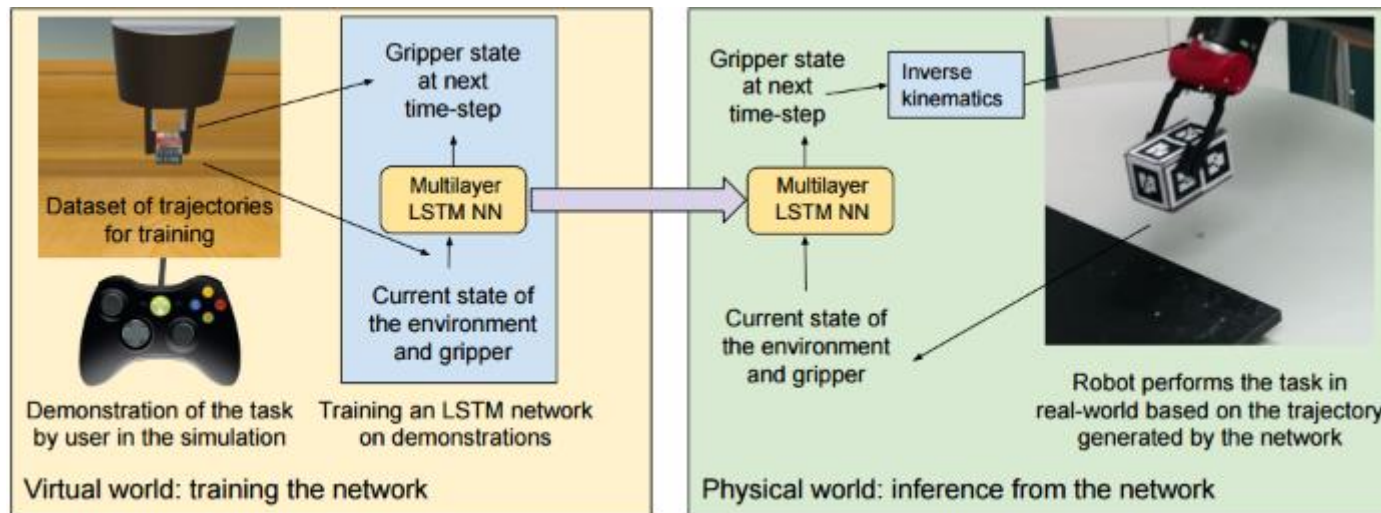




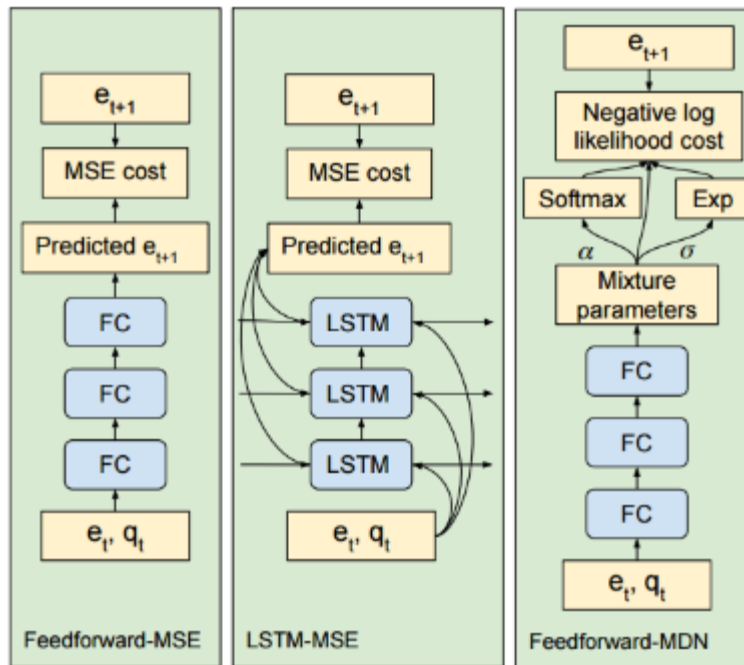
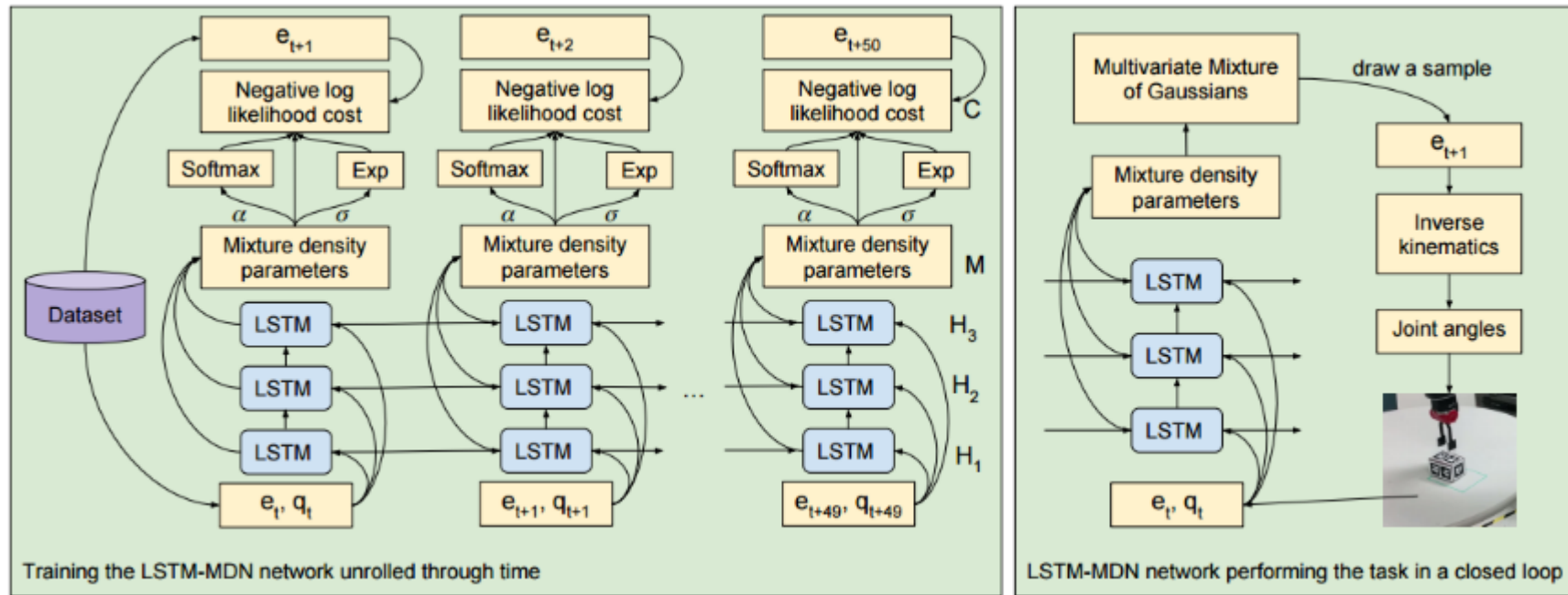
# Case study 2: Imitation with LSTMs

Learning real manipulation tasks from virtual demonstrations using LSTM

Rouhollah Rahmatizadeh<sup>1</sup>, Pooya Abolghasemi<sup>1</sup>, Aman Behal<sup>2</sup> and Ladislau Bölöni<sup>1</sup>



# Learning Manipulation Trajectories Using Recurrent Neural Networks

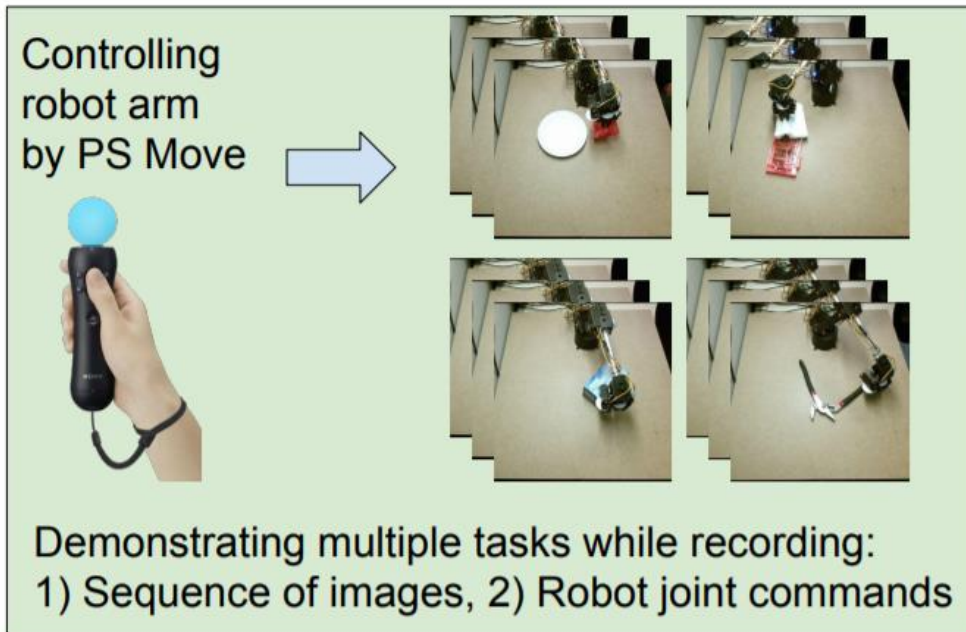


| Controller      | Pick and place | Push to pose |
|-----------------|----------------|--------------|
| Feedforward-MSE | 0%             | 0%           |
| LSTM-MSE        | 85%            | 0%           |
| Feedforward-MDN | 95%            | 15%          |
| LSTM-MDN        | <b>100%</b>    | <b>95%</b>   |

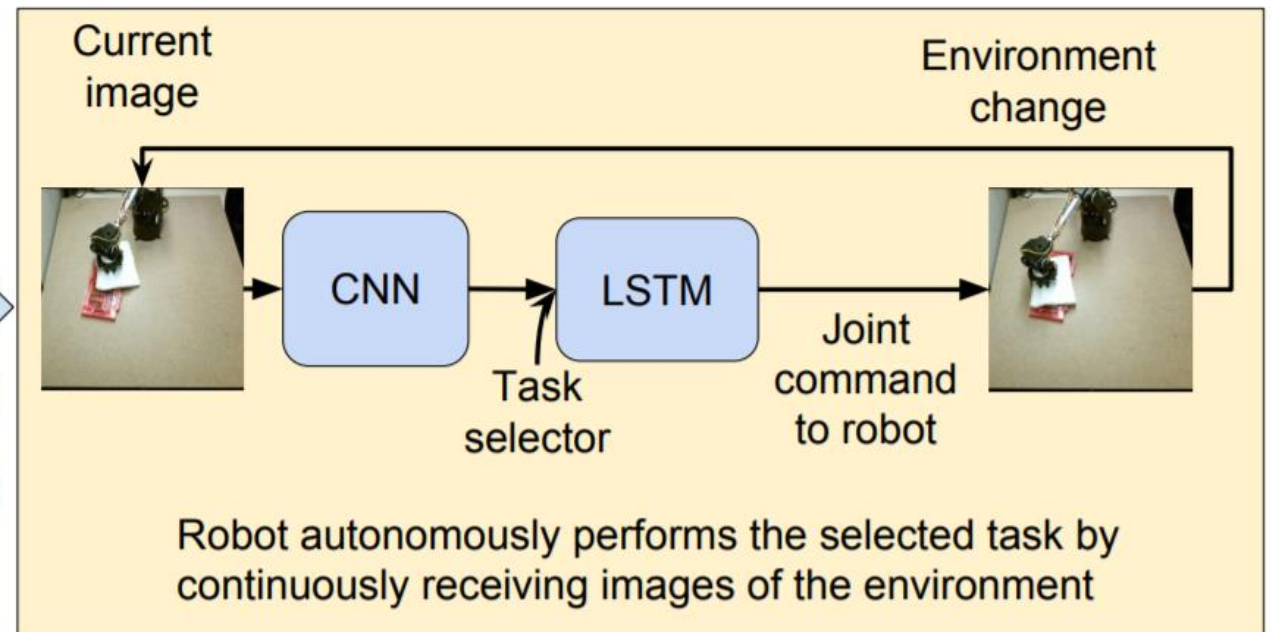
| Environment    | Pick and place | Push to pose |
|----------------|----------------|--------------|
| Virtual world  | 100%           | 95%          |
| Physical world | 80%            | 60%          |

# Follow-up: adding vision

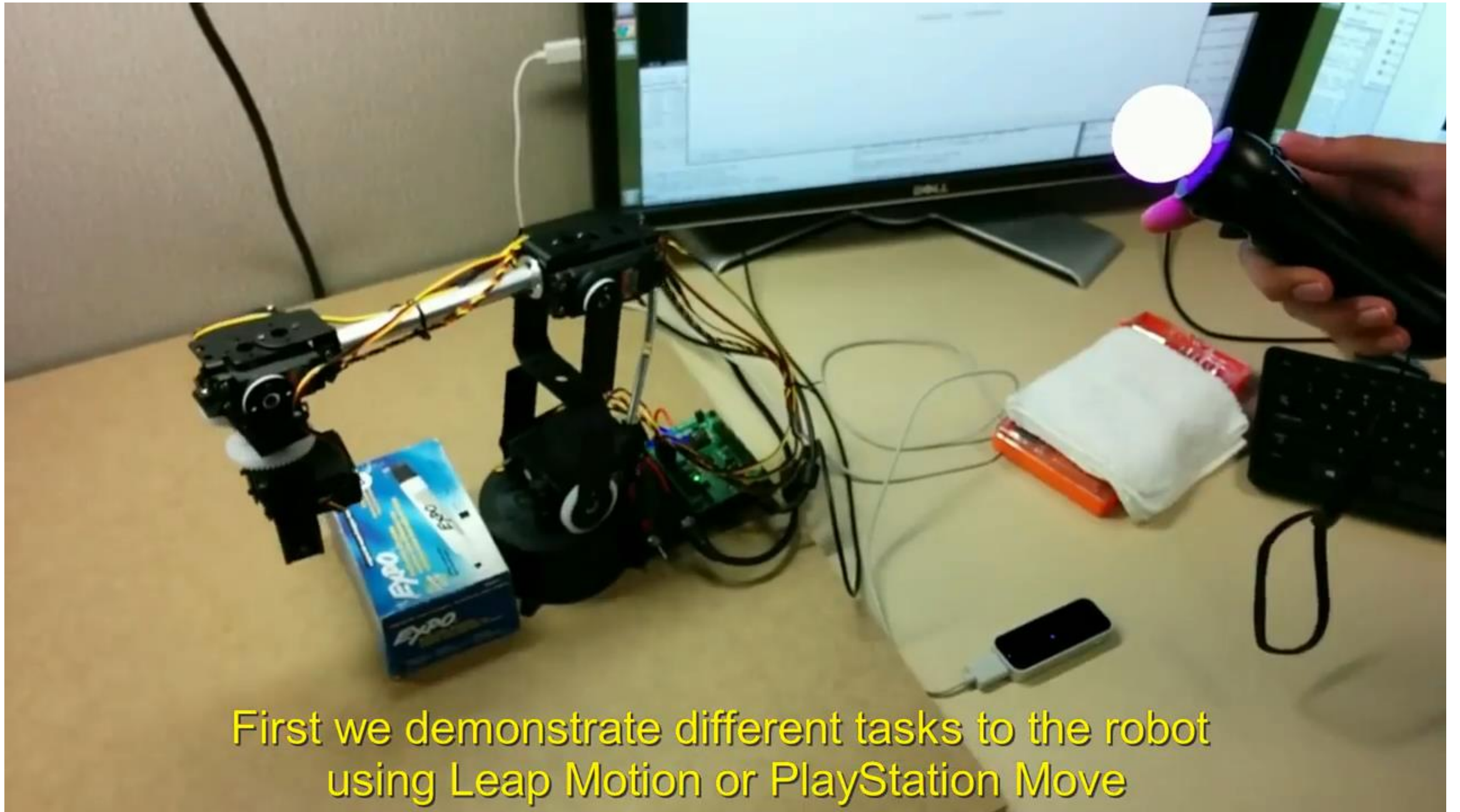
## Vision-Based Multi-Task Manipulation for Inexpensive Robots Using End-To-End Learning from Demonstration



Training neural network







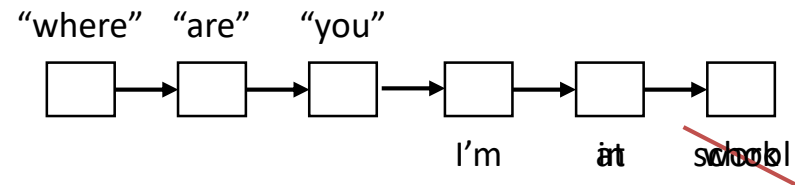
First we demonstrate different tasks to the robot using Leap Motion or PlayStation Move

# Other topics in imitation learning

- Structured prediction

x: where are you

y: I'm at work

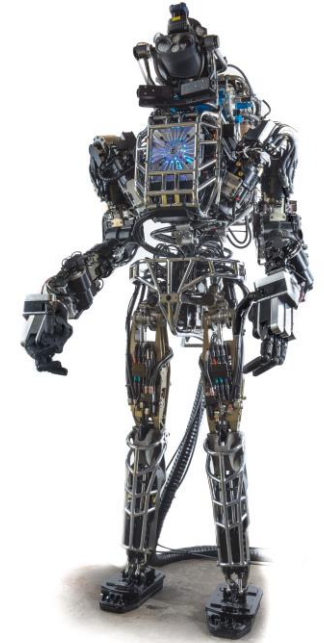
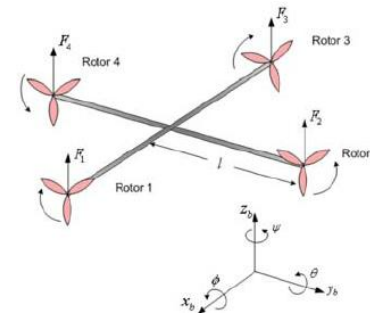
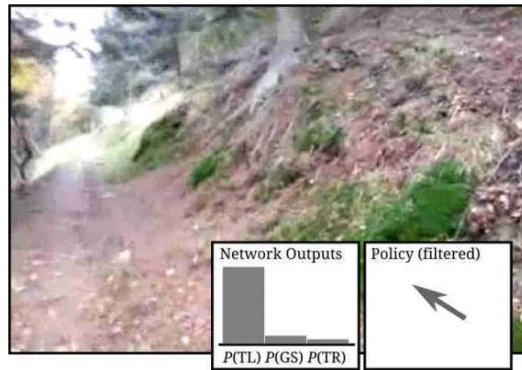


- Inverse reinforcement learning

- Instead of copying the demonstration, figure out the *goal*
- Will be covered later in this course

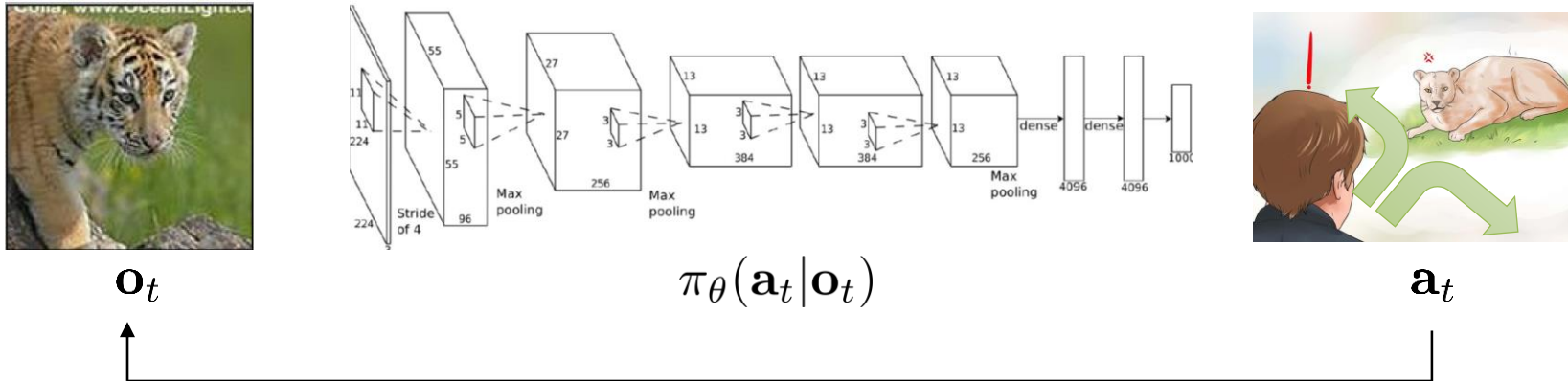
# Imitation learning: what's the problem?

- 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



$\mathbf{o}_t$

$\mathbf{a}_t$

$\mathbf{s}_t$  – state

$\mathbf{o}_t$  – observation

$\mathbf{a}_t$  – action

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

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

$$\min_{\mathbf{a}_1, \dots, \mathbf{a}_T} \sum_{t=1}^T \log p(\mathbf{s}_t | \mathbf{s}_{t-1}, \mathbf{a}_{t-1}) + \log p(\mathbf{o}_t | \mathbf{s}_t, \mathbf{a}_t) + c(\mathbf{s}_t, \mathbf{a}_t) - r(\mathbf{s}_t, \mathbf{a}_t)$$

# Aside: notation

$\mathbf{s}_t$  – state

$\mathbf{a}_t$  – action

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



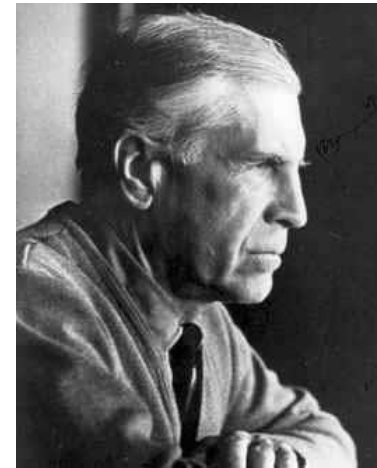
Richard Bellman

$\mathbf{x}_t$  – state

$\mathbf{u}_t$  – action управление

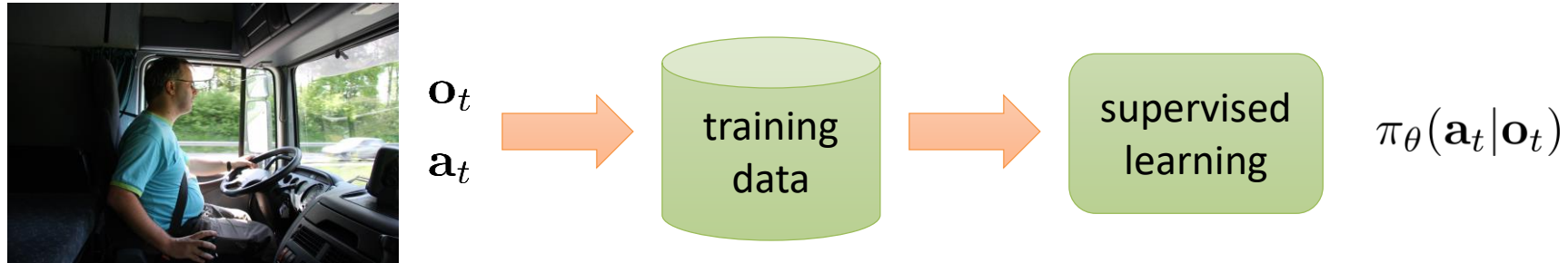
$c(\mathbf{x}, \mathbf{u})$  – cost function

$$r(\mathbf{s}, \mathbf{a}) = -c(\mathbf{x}, \mathbf{u})$$



Lev Pontryagin

# A cost function for imitation?

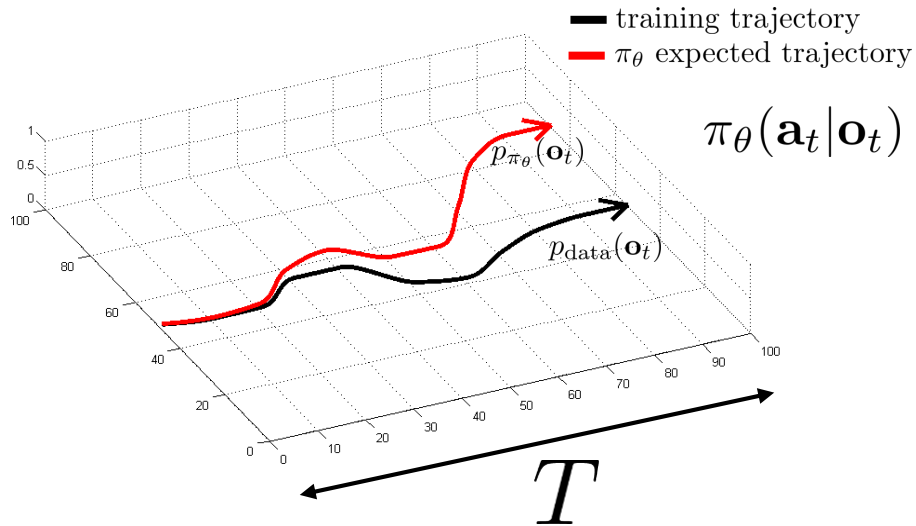


$$r(\mathbf{s}, \mathbf{a}) = \log p(\mathbf{a} = \pi^*(\mathbf{s}) | \mathbf{s})$$

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

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

# Some analysis

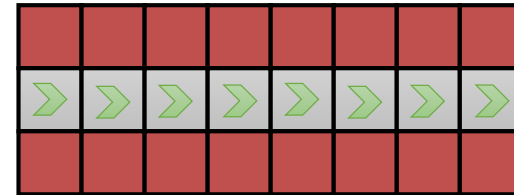


## How bad is it?

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

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

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

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

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

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

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

probability we made no mistakes

some *other* distribution

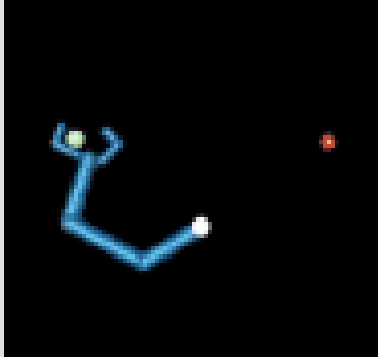
$$|p_\theta(\mathbf{s}_t) - p_{\text{train}}(\mathbf{s}_t)| = (1 - (1 - \epsilon)^t) |p_{\text{mistake}}(\mathbf{s}_t) - p_{\text{train}}(\mathbf{s}_t)| \leq 2(1 - (1 - \epsilon)^t)$$

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

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

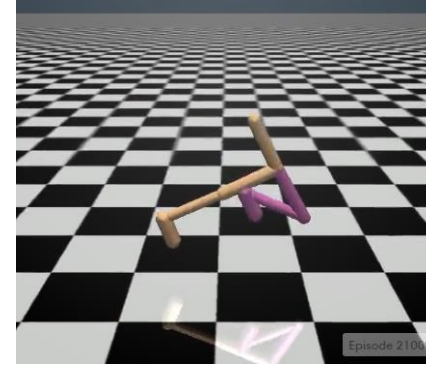


# Cost/reward functions in theory and practice



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

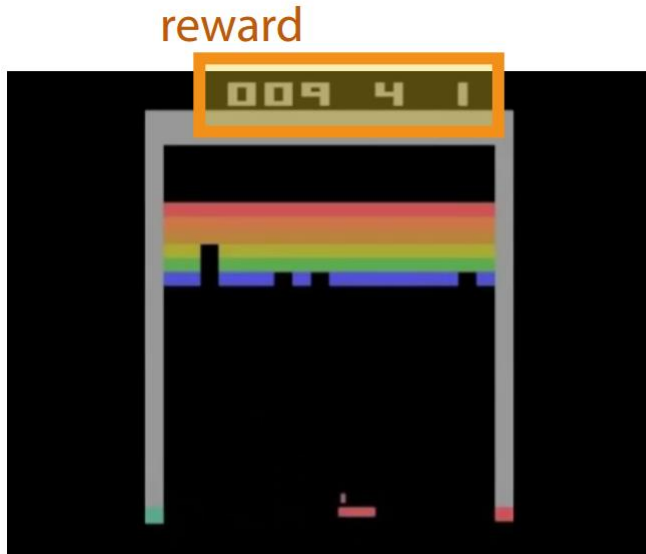
$$r(\mathbf{s}, \mathbf{a}) = -w_1 \|p_{\text{gripper}}(\mathbf{s}) - p_{\text{object}}(\mathbf{s})\|^2 + \\ -w_2 \|p_{\text{object}}(\mathbf{s}) - p_{\text{target}}(\mathbf{s})\|^2 + \\ -w_3 \|\mathbf{a}\|^2$$



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

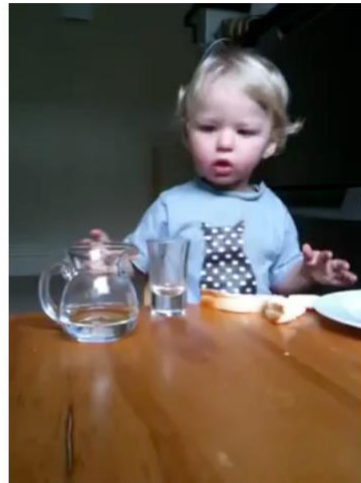
$$r(\mathbf{s}, \mathbf{a}) = w_1 v(\mathbf{s}) + \\ w_2 \delta(|\theta_{\text{torso}}(\mathbf{s})| < \epsilon) + \\ w_3 \delta(h_{\text{torso}}(\mathbf{s}) \geq h)$$

# The trouble with cost & reward functions



Mnih et al. '15

reinforcement learning agent



what is the **reward**?

---

## Sim-to-Real Robot Learning from Pixels with Progressive Nets

---

Andrei A. Rusu, Matej Vecerik, Thomas Rothörl, Nicolas Heess,  
Razvan Pascanu, Raia Hadsell

Google DeepMind  
London, UK

{andreirusu, matejvecerik, tcr, heess, razp, raia}@google.com



Rewards are given automatically by tracking the colored target

More on this later...

# A note about terminology...

## the “R” word

a bit of history...

reinforcement learning  
(the **problem** statement)

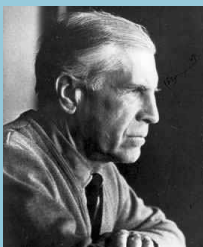
$$\max \sum_{t=1}^T E[r(\mathbf{s}_t, \mathbf{a}_t)]$$

$$\mathbf{s}_{t+1} \sim p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)$$

reinforcement learning  
(the **method**)

without using the **model**

$$\mathbf{s}_{t+1} \sim p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)$$



Lev Pontryagin



Richard Bellman



Andrew Barto



Richard Sutton