

Supervised Learning of Behaviors: Deep Learning, Dynamical Systems, and Behavior Cloning

CS 294-112: Deep Reinforcement Learning

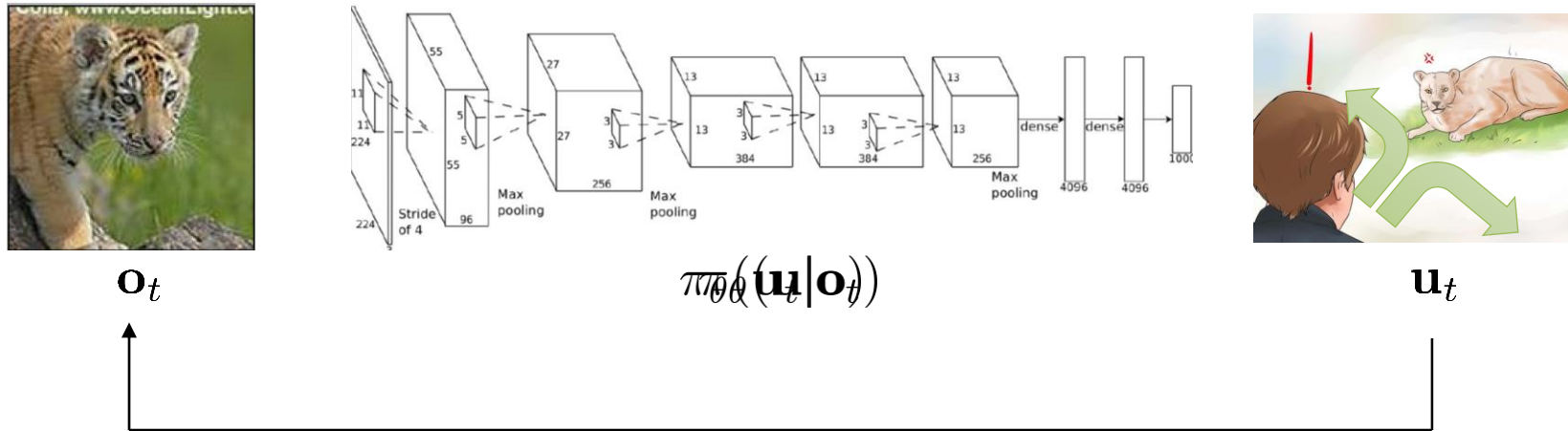
Week 2, Lecture 1

Sergey Levine

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{x}_t – state

\mathbf{o}_t – observation

\mathbf{u}_t – action

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

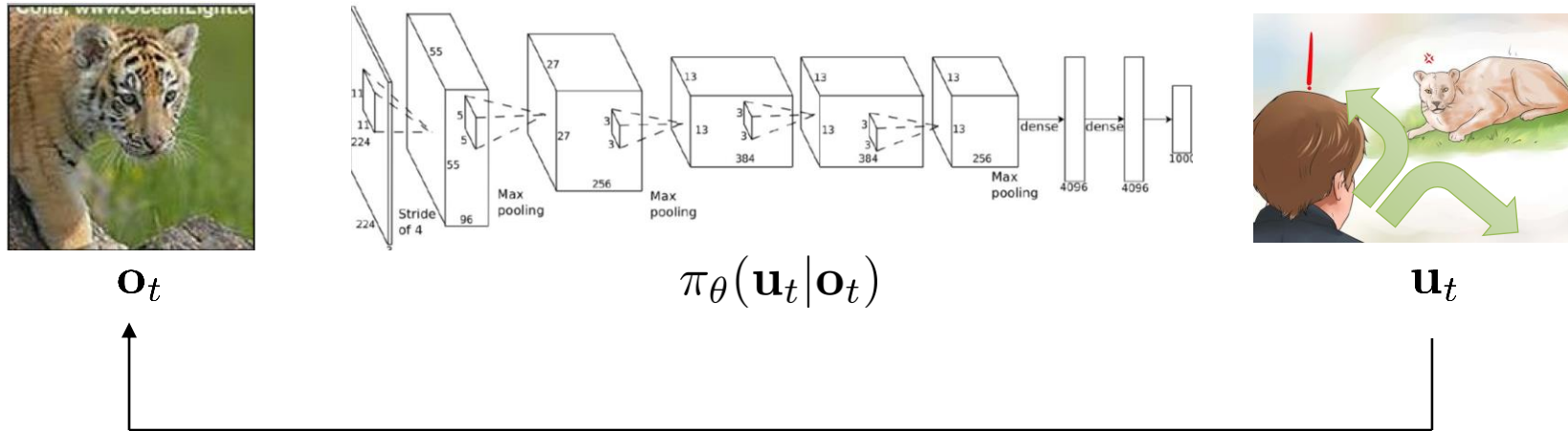


\mathbf{o}_t – observation



\mathbf{x}_t – state

Terminology & notation

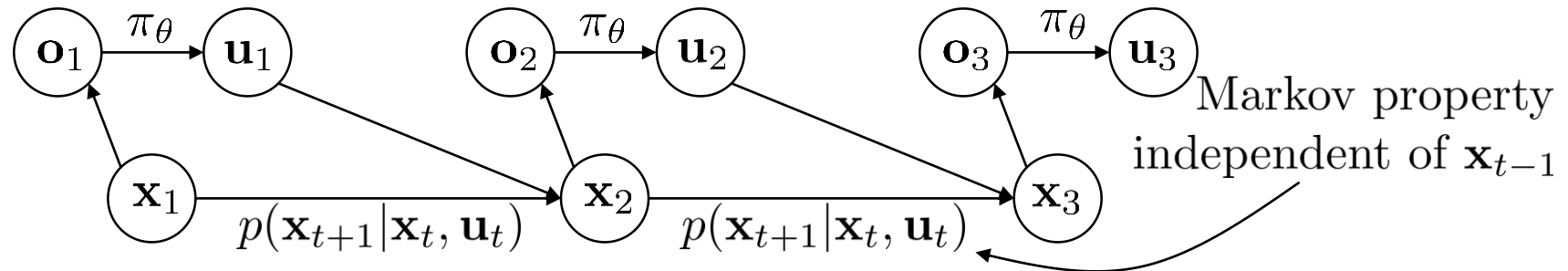


\mathbf{x}_t – state

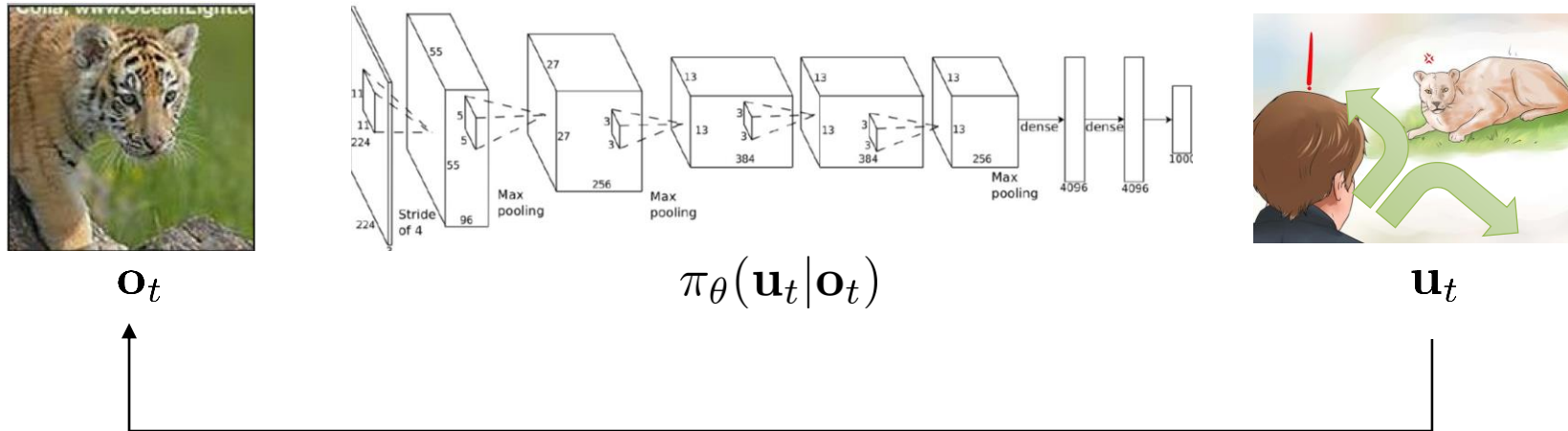
\mathbf{o}_t – observation

\mathbf{u}_t – action

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



Terminology & notation



\mathbf{x}_t – state

\mathbf{o}_t – observation

\mathbf{u}_t – action

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

a bit of history...

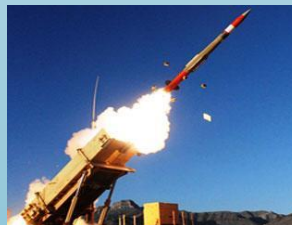
\mathbf{x}_t – state

\mathbf{u}_t – action

управление



Lev Pontryagin



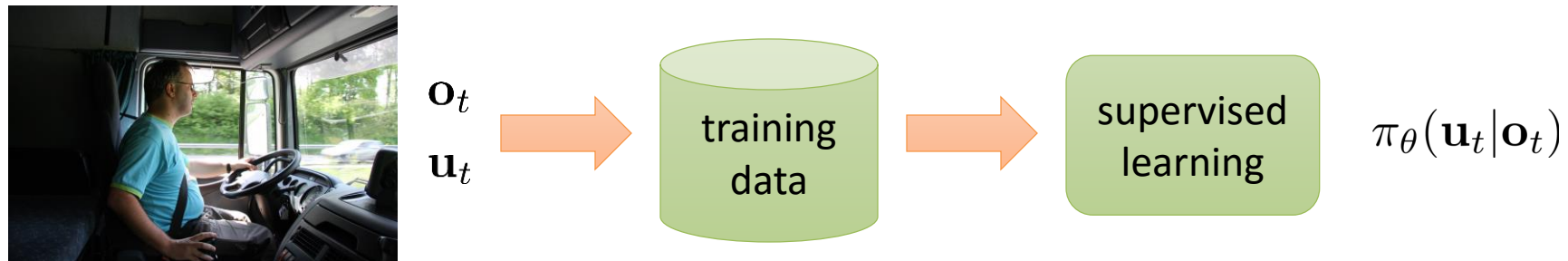
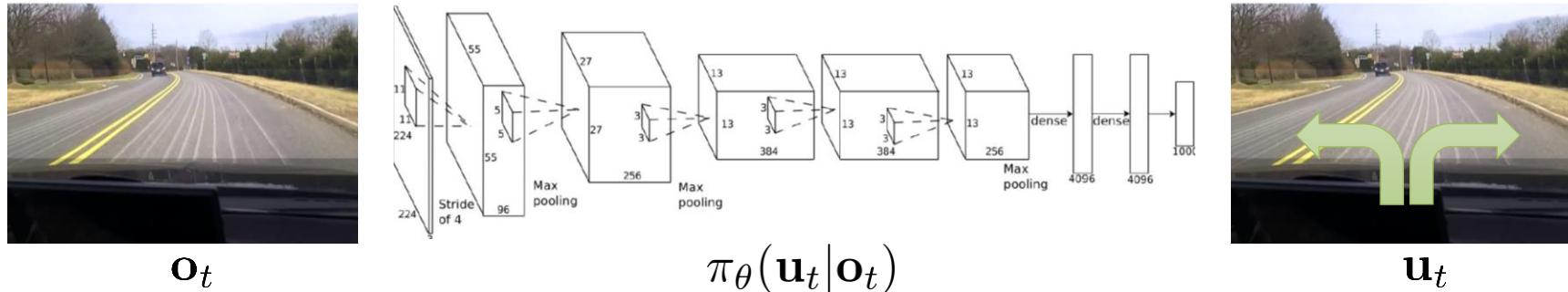
\mathbf{s}_t – state

\mathbf{a}_t – action



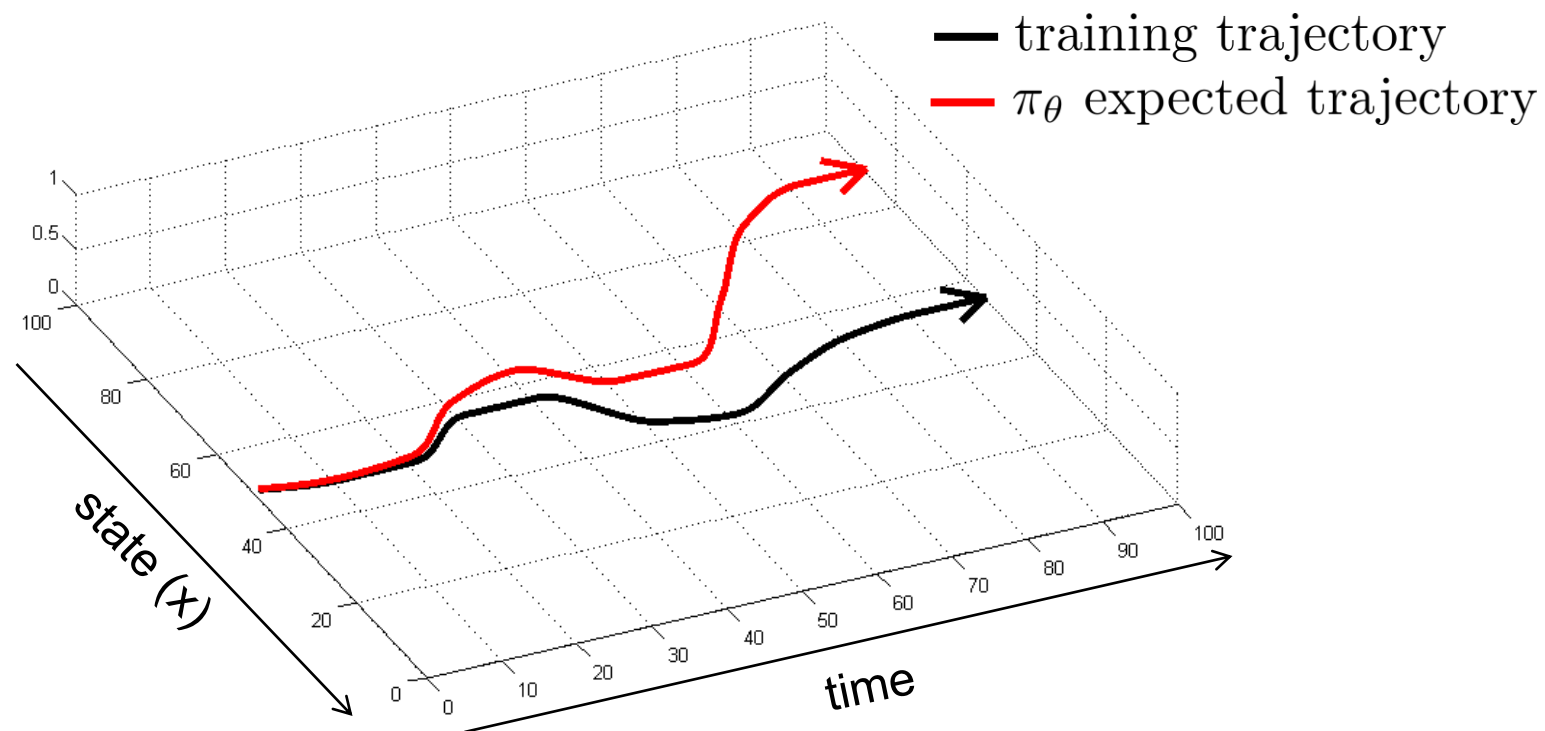
Richard Bellman

Imitation Learning



Does it work?

No!



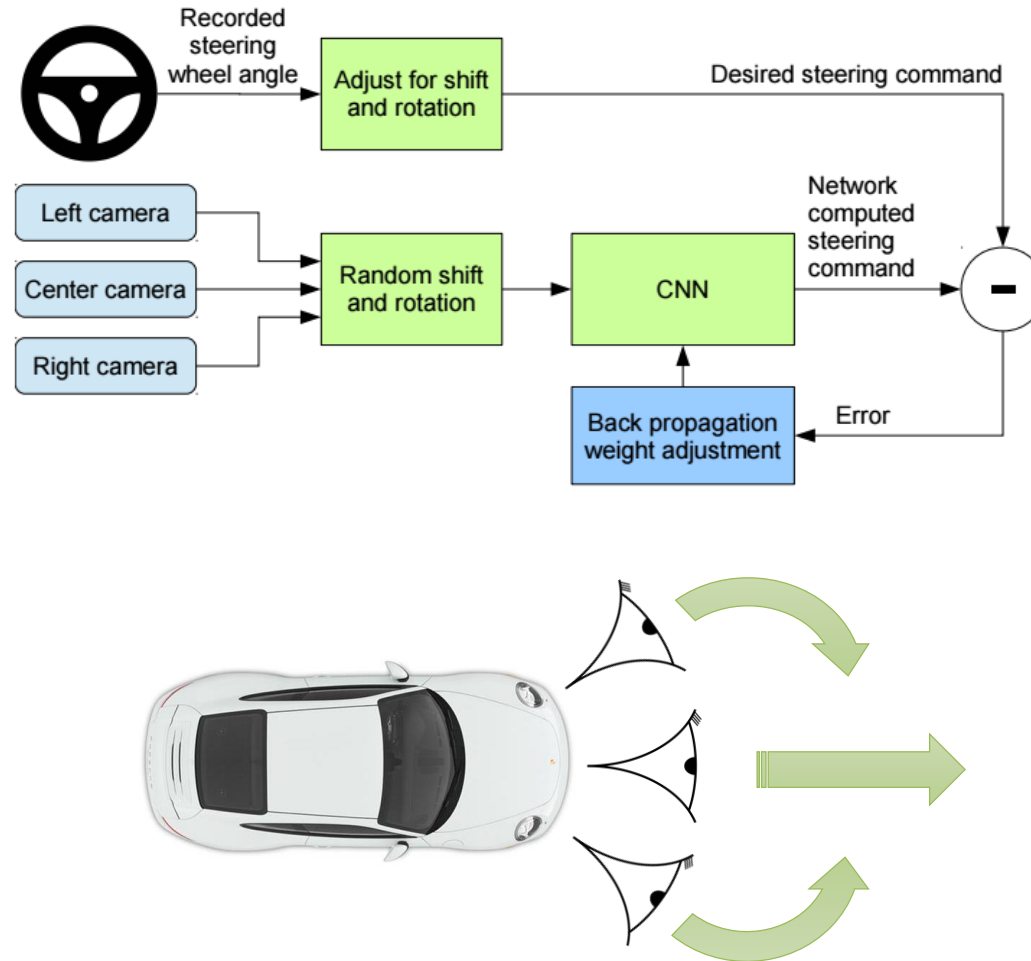
Does it work?

Yes!

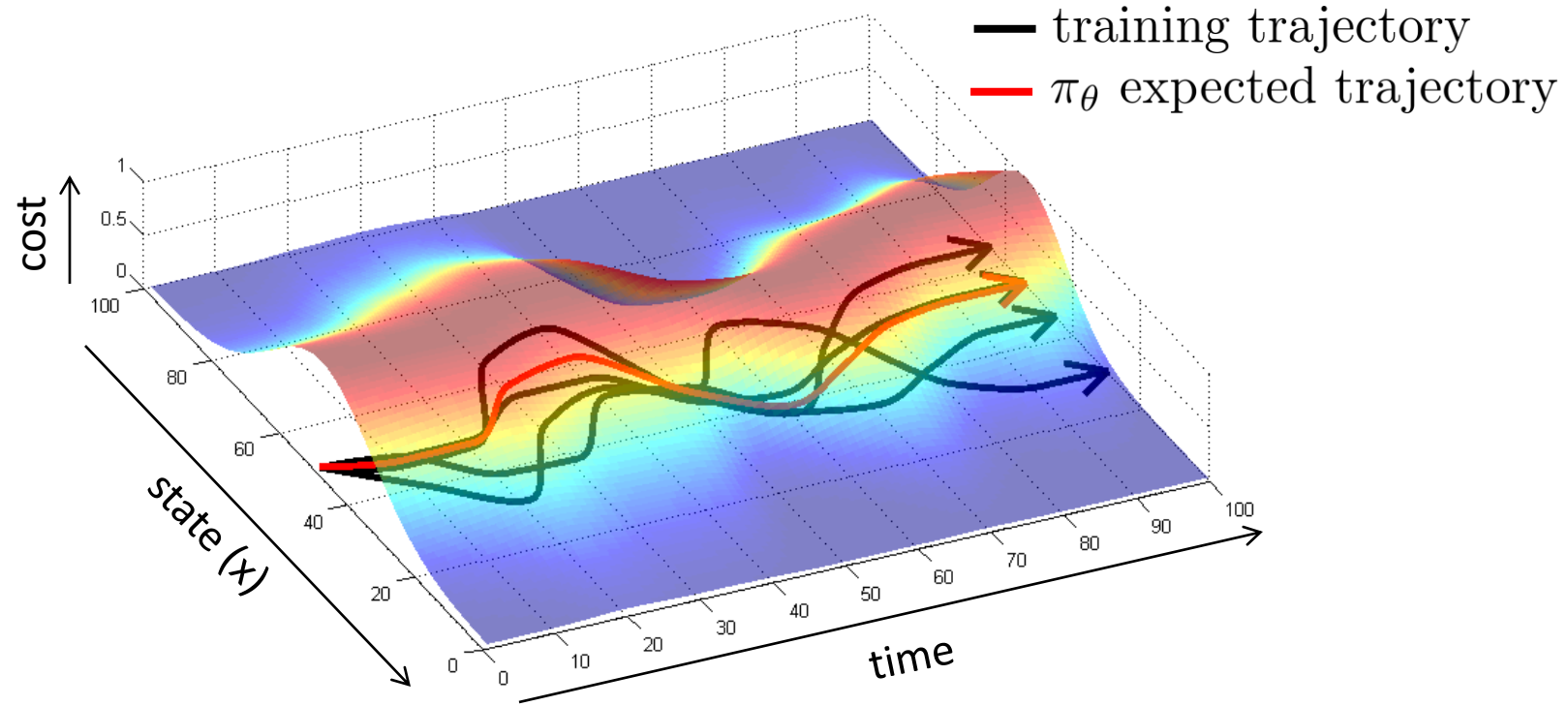


Video: Bojarski et al. '16, NVIDIA

Why did that work?



Can we make it work more often?

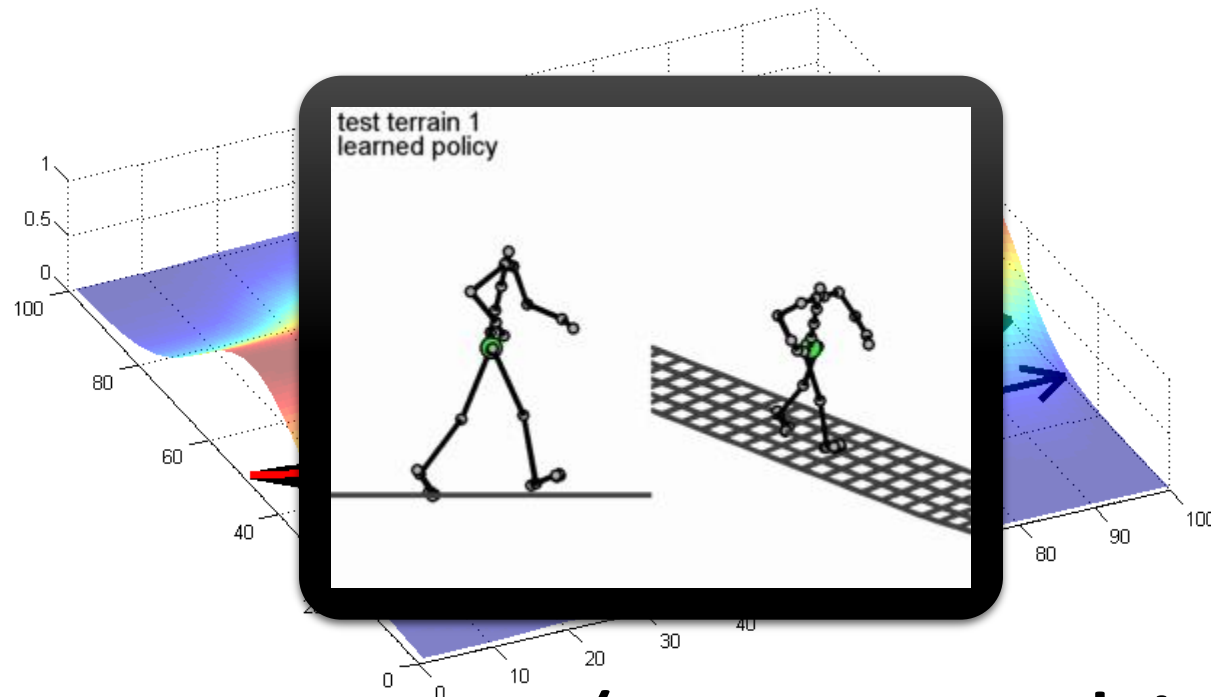


stability

Learning from a stabilizing controller

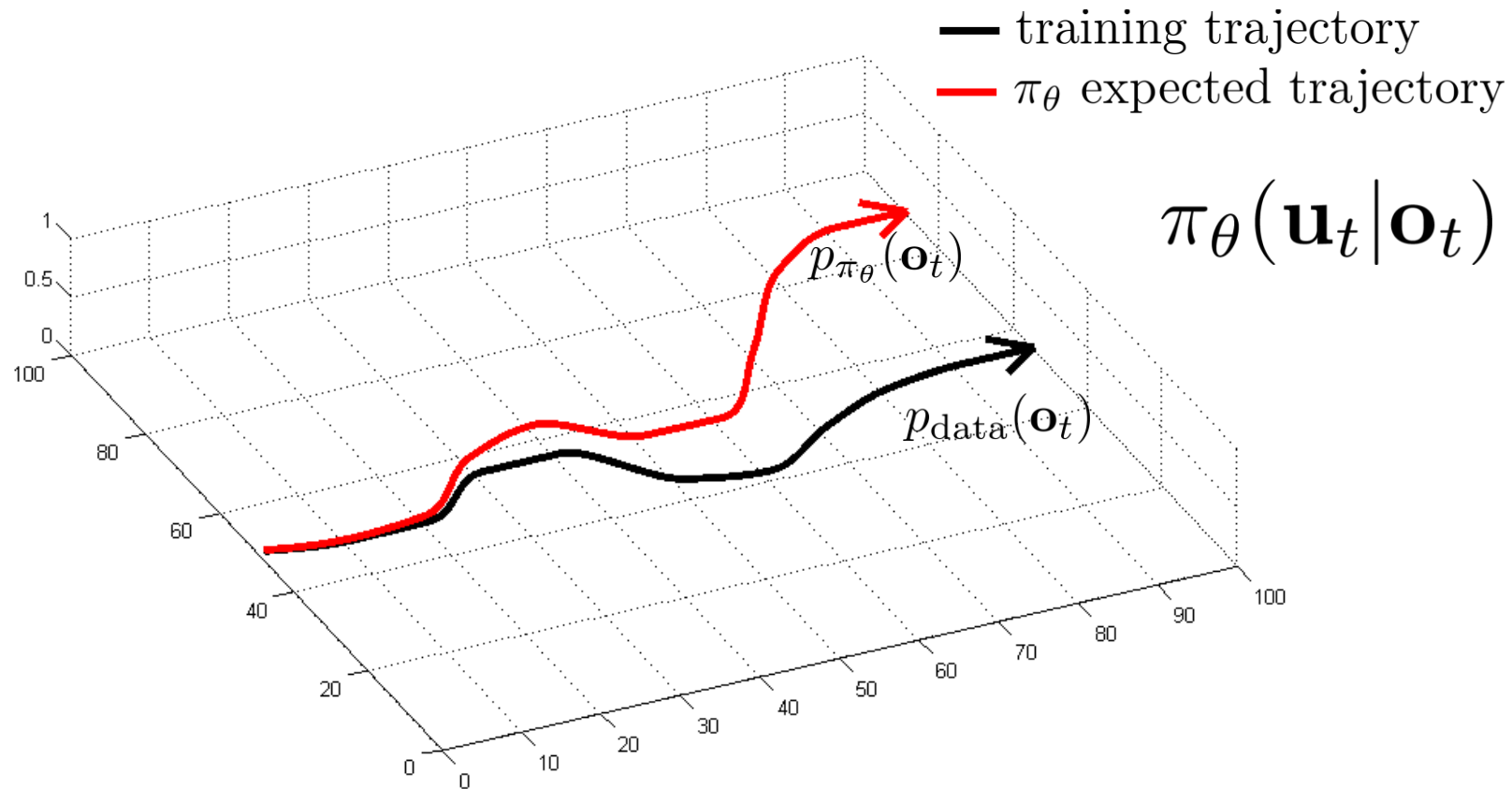
$p(\mathbf{x}_1), u_{\text{Gaussian distribution}}$ obtained using variant of iterative LQR

τ



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

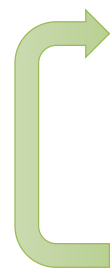
idea: instead of being clever about $p_{\pi_\theta}(\mathbf{o}_t)$, be clever about $p_{\text{data}}(\mathbf{o}_t)$!

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

but need labels \mathbf{u}_t !

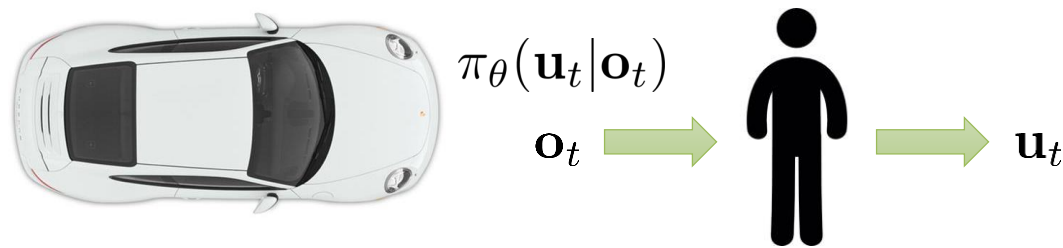
- 
1. train $\pi_\theta(\mathbf{u}_t|\mathbf{o}_t)$ from human data $\mathcal{D} = \{\mathbf{o}_1, \mathbf{u}_1, \dots, \mathbf{o}_N, \mathbf{u}_N\}$
 2. run $\pi_\theta(\mathbf{u}_t|\mathbf{o}_t)$ to get dataset $\mathcal{D}_\pi = \{\mathbf{o}_1, \dots, \mathbf{o}_M\}$
 3. Ask human to label \mathcal{D}_π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

Dagger Example

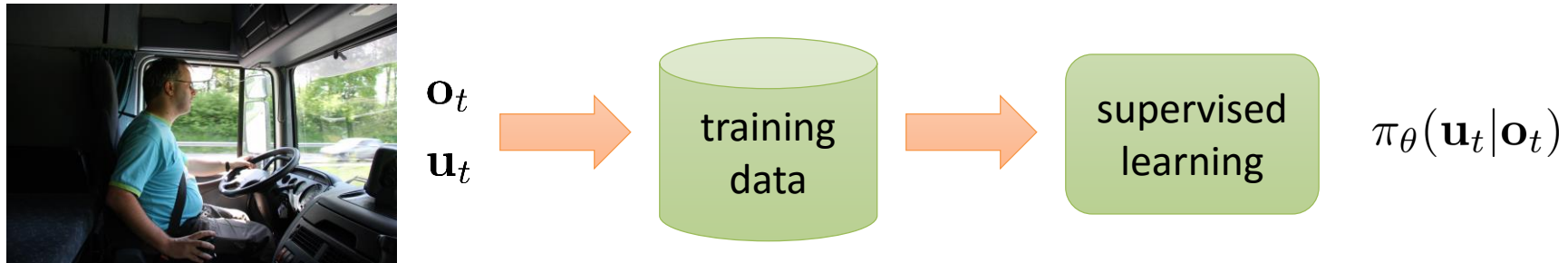


What's the problem?

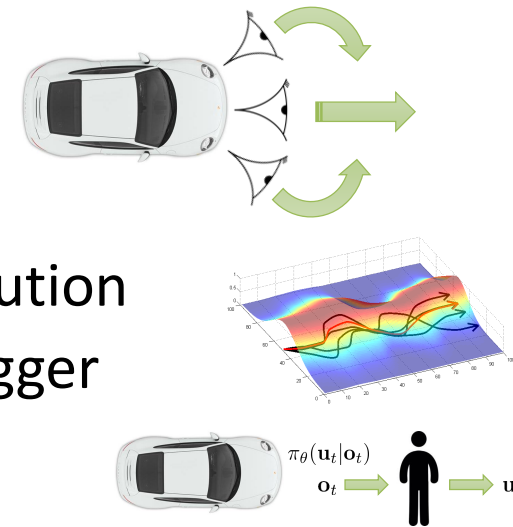
1. train $\pi_{\theta}(\mathbf{u}_t|\mathbf{o}_t)$ from human data $\mathcal{D} = \{\mathbf{o}_1, \mathbf{u}_1, \dots, \mathbf{o}_N, \mathbf{u}_N\}$
2. run $\pi_{\theta}(\mathbf{u}_t|\mathbf{o}_t)$ to get dataset $\mathcal{D}_{\pi} = \{\mathbf{o}_1, \dots, \mathbf{o}_M\}$
3. Ask human to label \mathcal{D}_{π}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



Imitation learning: recap



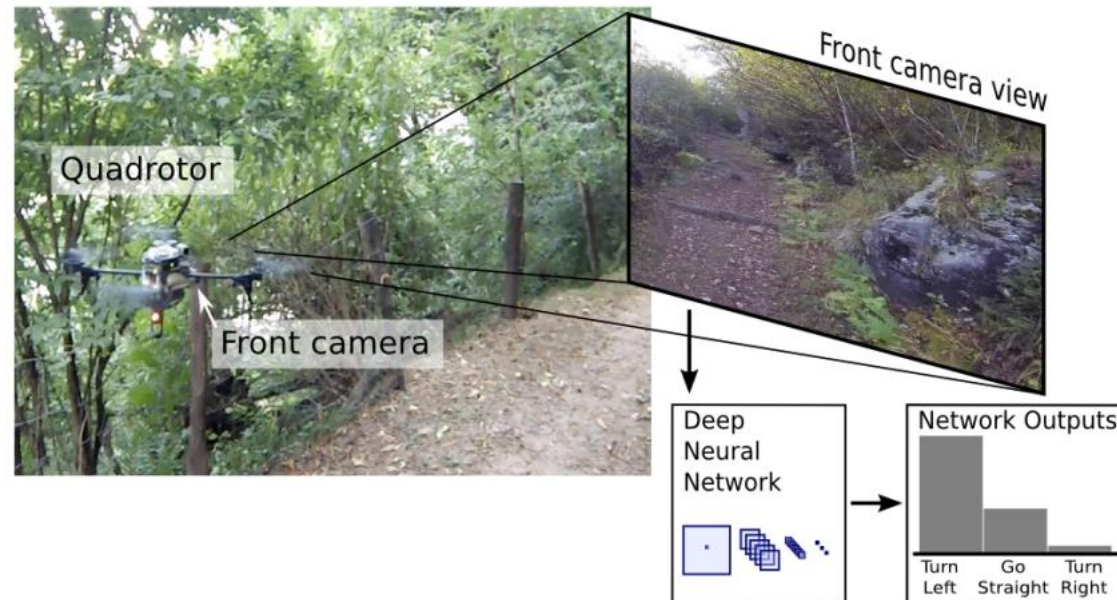
- Often (but not always) insufficient by itself
 - Distribution mismatch problem
- Sometimes works well
 - Hacks (e.g. left/right images)
 - Samples from a stable trajectory distribution
 - Add more **on-policy** data, e.g. using DAgger



Case study 1: trail following as classification

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

Alessandro Giusti¹, Jérôme Guzzi¹, Dan C. Cireşan¹, Fang-Lin He¹, Juan P. Rodríguez¹
Flavio Fontana², Matthias Faessler², Christian Forster²
Jürgen Schmidhuber¹, Gianni Di Caro¹, Davide Scaramuzza², Luca M. Gambardella¹

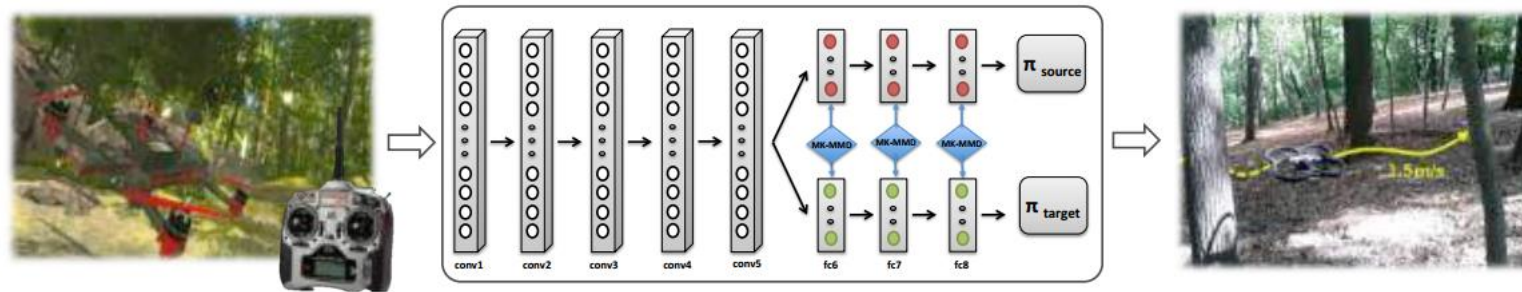


Case study 2: DAgger & domain adaptation

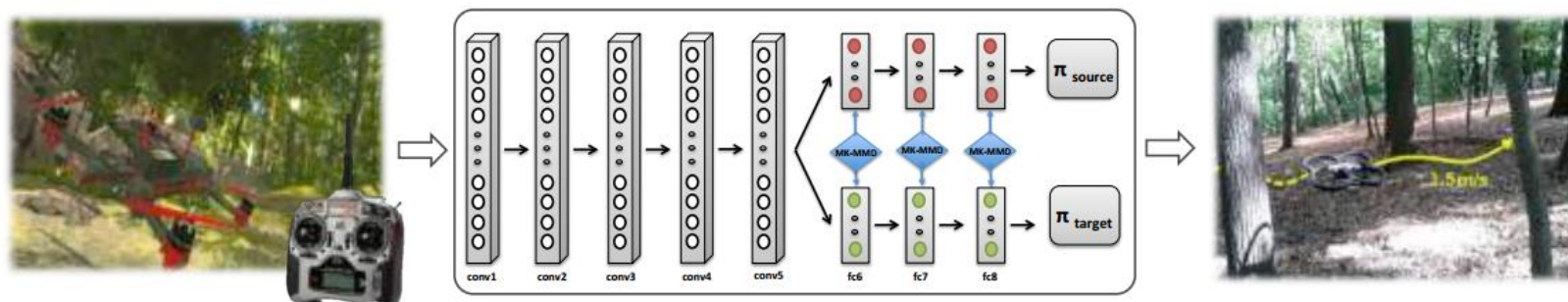
Learning Transferable Policies for Monocular Reactive MAV Control

Shreyansh Daftry, J. Andrew Bagnell, and Martial Hebert

Robotics Institute, Carnegie Mellon University, Pittsburgh, USA
{daftry, dbagnell, hebert}@ri.cmu.edu



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



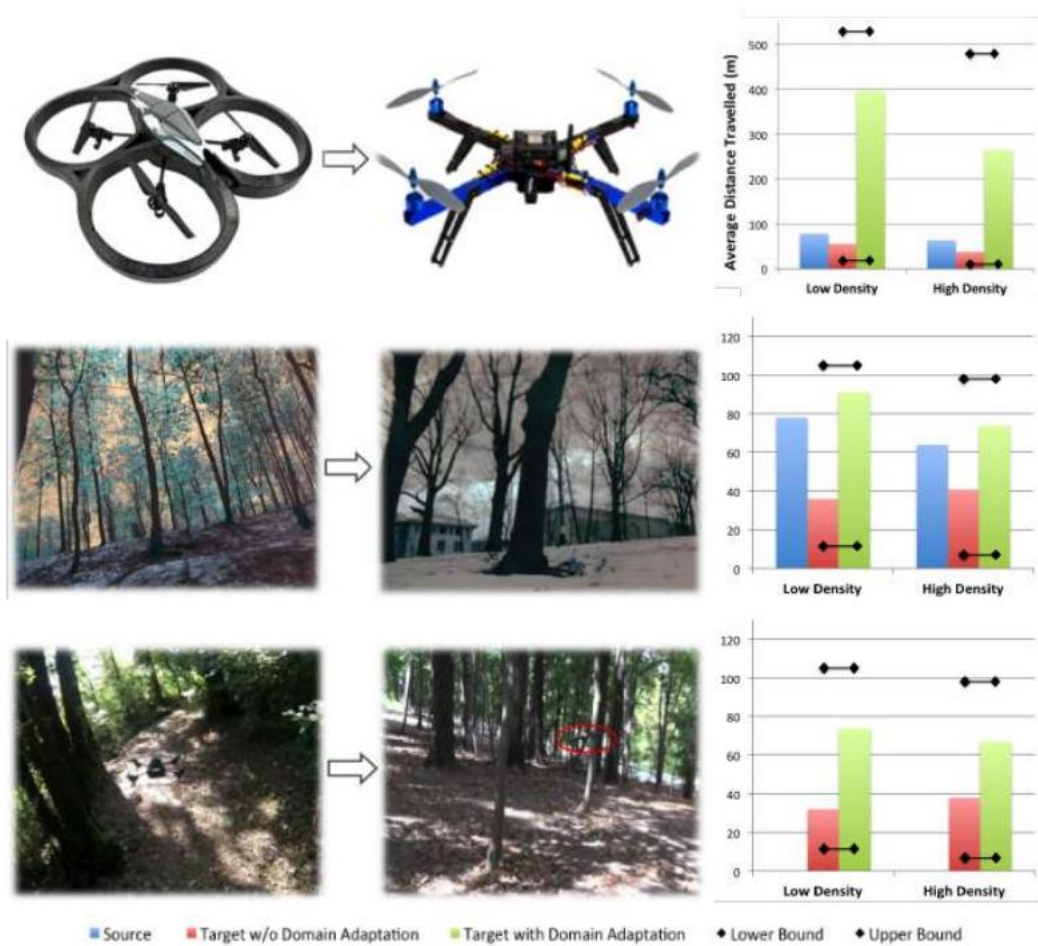
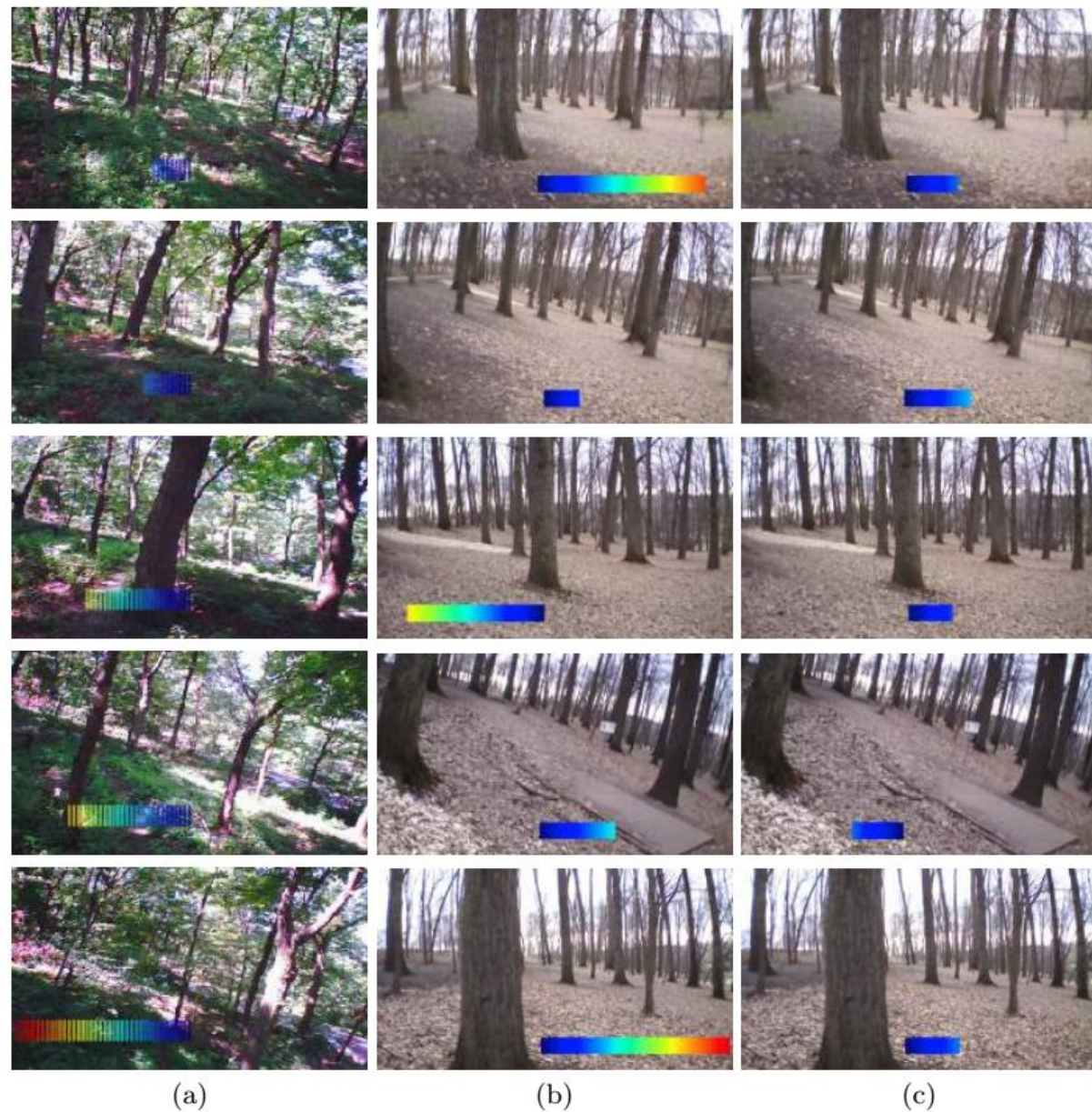


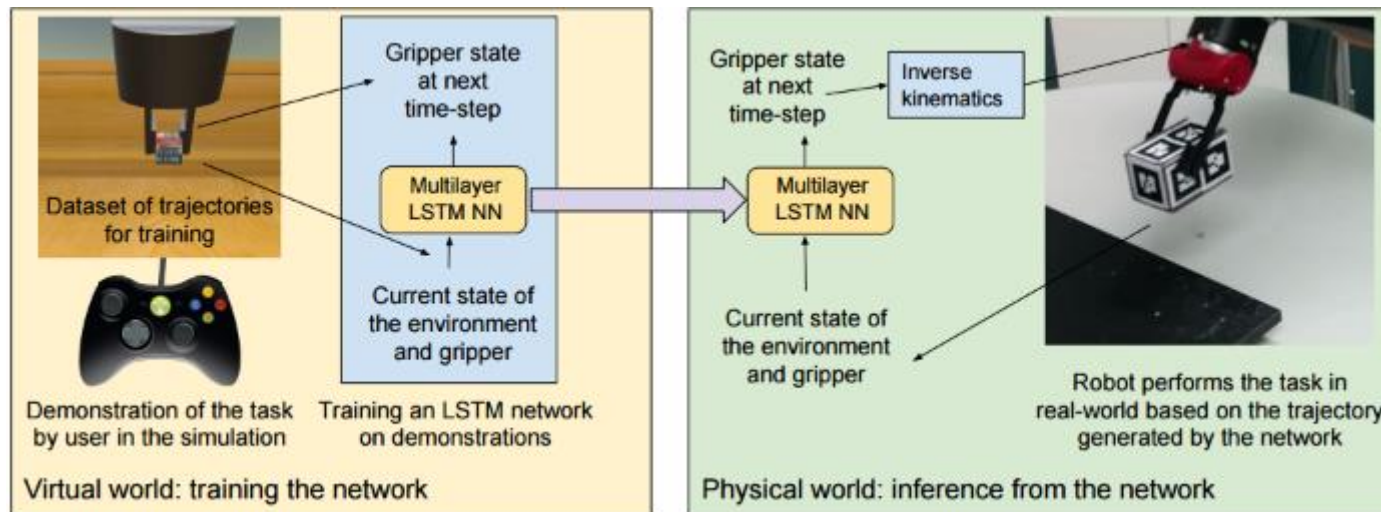
Fig. 2. Experiments and Results for (Row-1) Transfer across physical systems from ARDrone to ArduCopter, (Row-2) Transfer across weather conditions from summer to winter and (Row-3) Transfer across environments from Univ. of Zurich to CMU.



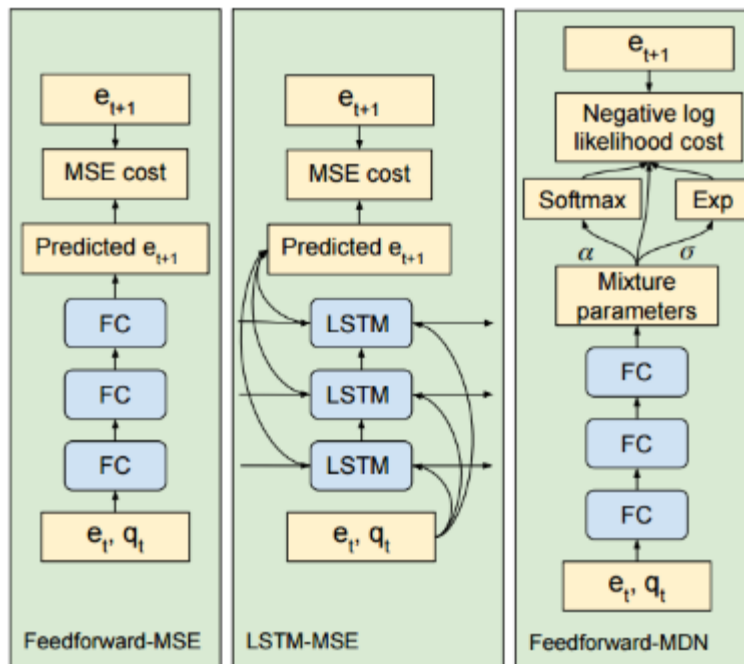
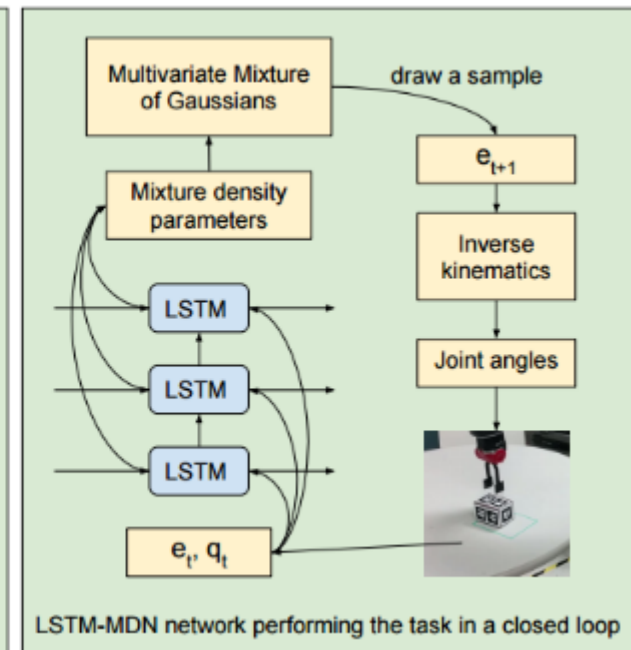
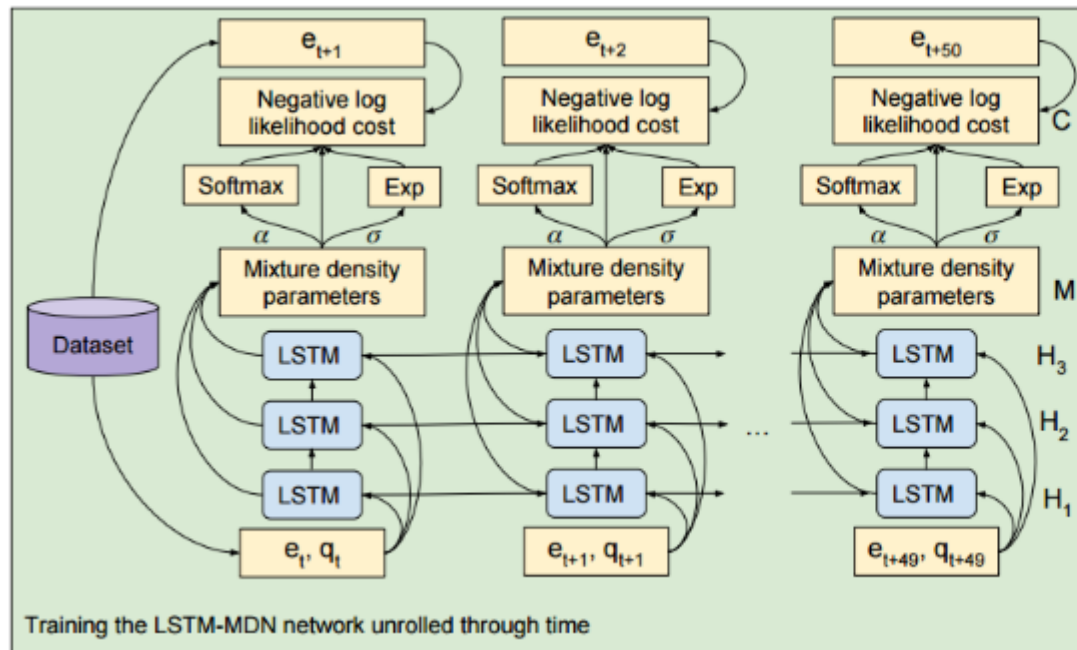
Case study 3: Imitation with LSTMs

Learning real manipulation tasks from virtual demonstrations using LSTM

Rouhollah Rahmatizadeh¹, Pooya Abolghasemi¹, Aman Behal² and Ladislau Bölöni¹



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 | 100% | 95% |

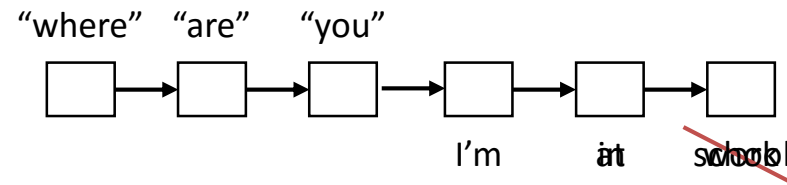
| Environment | Pick and place | Push to pose |
|----------------|----------------|--------------|
| Virtual world | 100% | 95% |
| Physical world | 80% | 60% |

Other topics in imitation learning

- Structured prediction

x: where are you

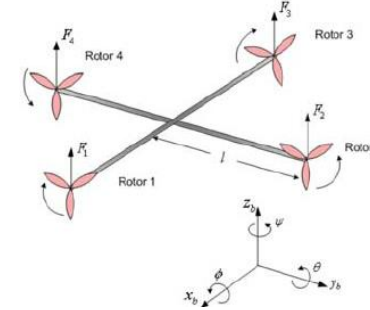
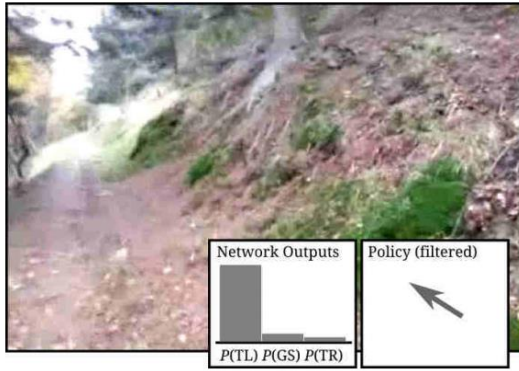
y: I'm at work



- See Mohammad Norouzi's lecture in April!
- Interaction & active learning
- Inverse reinforcement learning
 - Instead of copying the demonstration, figure out the *goal*
 - Will be covered later in this course

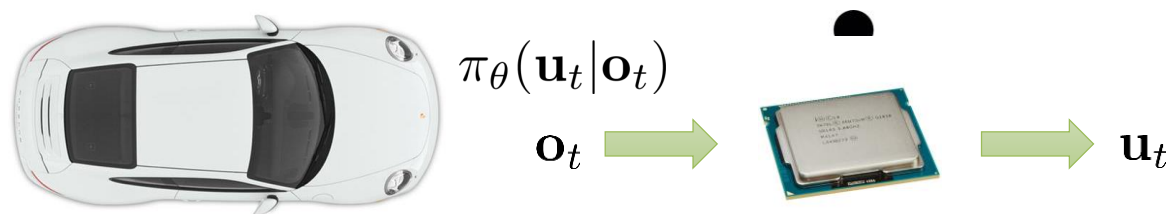
Imitation learning: what's the problem?

- 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

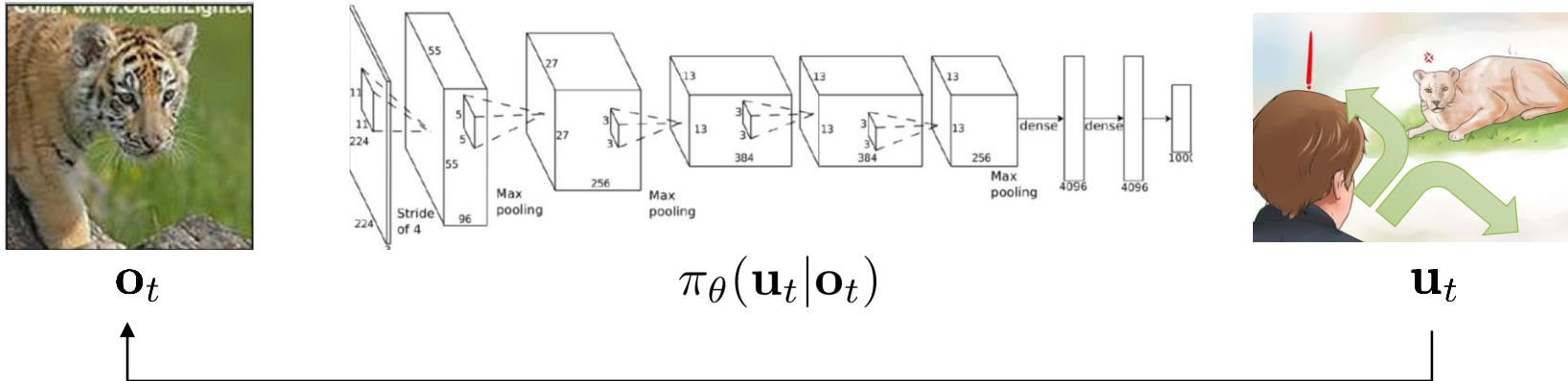


- Humans can learn autonomously; can our machines do the same?
 - Unlimited data from own experience
 - Continuous self-improvement

Next time: learning without humans



Terminology & notation



\mathbf{x}_t – state

\mathbf{o}_t – observation

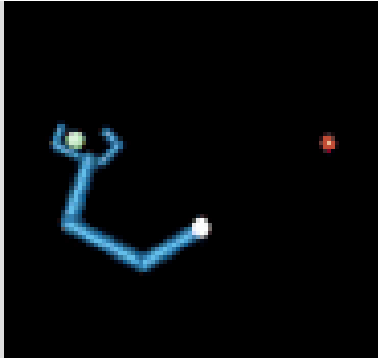
\mathbf{u}_t – action

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

$r(\mathbf{x}_t, \mathbf{u}_t)$ – reward function

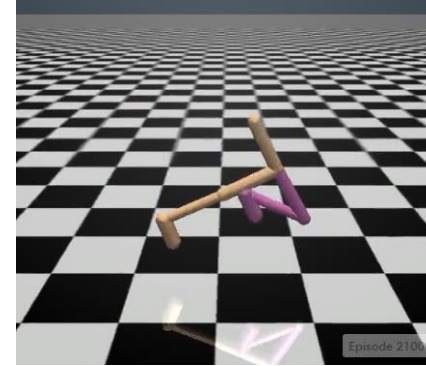
$$\min_{\mathbf{u}_1, \dots, \mathbf{u}_T} \sum_{t=1}^T \log p(\mathbf{o}_t | \mathbf{x}_t, \mathbf{u}_t) + \sum_{t=1}^T f(\mathbf{x}_t, \mathbf{u}_t)$$

Cost/reward functions in theory and practice



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

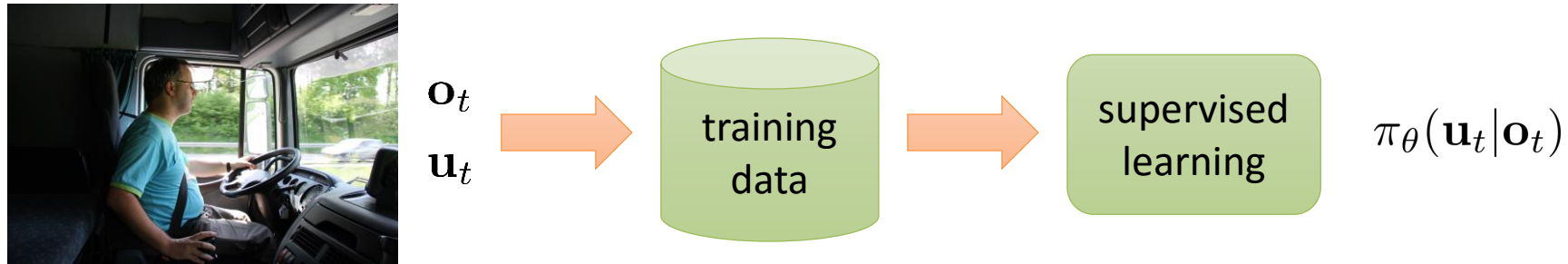
$$\begin{aligned} r(\mathbf{x}, \mathbf{u}) = & -w_1 \|p_{\text{gripper}}(\mathbf{x}) - p_{\text{object}}(\mathbf{x})\|^2 + \\ & -w_2 \|p_{\text{object}}(\mathbf{x}) - p_{\text{target}}(\mathbf{x})\|^2 + \\ & -w_3 \|\mathbf{u}\|^2 \end{aligned}$$



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

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

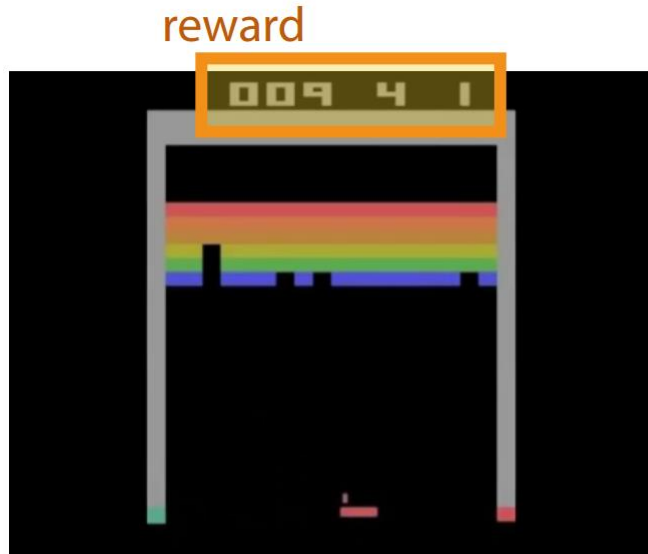
A cost function for imitation?



$$c(\mathbf{x}, \mathbf{u}) = -\log p(\mathbf{u} = \pi^*(\mathbf{x})|\mathbf{x})$$

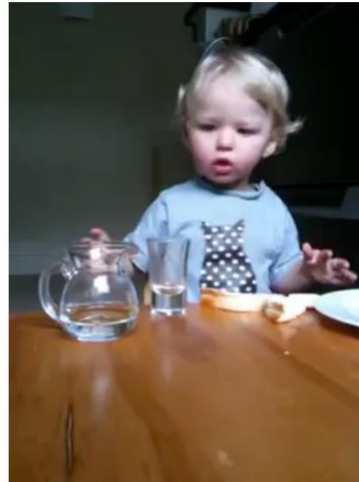
1. train $\pi_\theta(\mathbf{u}_t|\mathbf{o}_t)$ from human data $\mathcal{D} = \{\mathbf{o}_1, \mathbf{u}_1, \dots, \mathbf{o}_N, \mathbf{u}_N\}$
2. run $\pi_\theta(\mathbf{u}_t|\mathbf{o}_t)$ to get dataset $\mathcal{D}_\pi = \{\mathbf{o}_1, \dots, \mathbf{o}_M\}$
3. Ask human to label \mathcal{D}_π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

The trouble with cost & reward functions



Mnih et al. '15

reinforcement learning agent



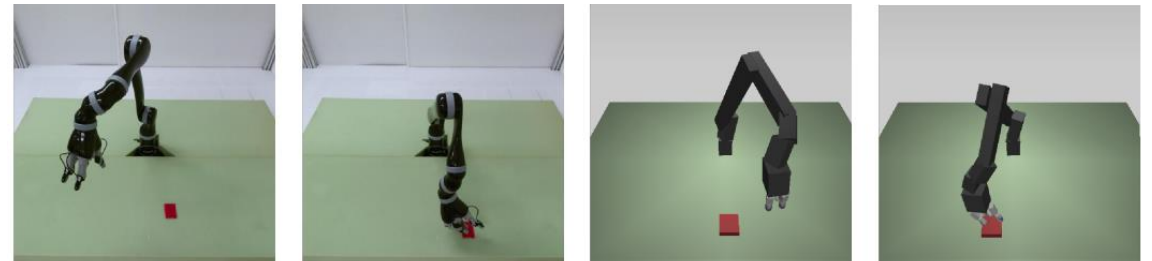
what is the **reward**?

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...

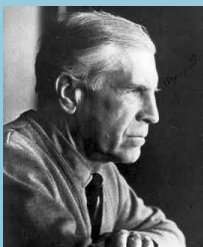
reinforcement learning
(the **problem** statement)

$$\min \sum_{t=1}^T E[c(\mathbf{x}_t, \mathbf{u}_t)] \quad \mathbf{x}_{t+1} \sim p(\mathbf{x}_{t+1} | \mathbf{x}_t, \mathbf{u}_t)$$

reinforcement learning
(the **method**)

without using the **model**

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



Lev Pontryagin



Richard Bellman



Andrew Barto



Richard Sutton