

# CS 277: Control and Reinforcement Learning

## Winter 2022

# Lecture 2: Imitation Learning

Roy Fox

Department of Computer Science

Bren School of Information and Computer Sciences

University of California, Irvine



# Logistics

---

logistics

- Follow announcements and discussions on [ed](#)
- See [website](#) for schedule, recordings, resources, etc.

assignments

- Quiz 1 due [this Friday](#)
- Assignment 1 to be published soon

# Imitation Learning (IL)

- How can we **teach** an agent to perform a task?
- Often there is an **expert** that already knows how to perform the task
  - ▶ A **human** operator who controls a robot
  - ▶ A **black-box** artificial agent that we can observe but not copy
  - ▶ An agent with different **representation** or **embodiment**
- The expert can **demonstrate** the task to create a training dataset  $\mathcal{D} = \{\xi^{(i)}\}_i$ 
  - ▶ Each demonstration is a trajectory  $\xi = s_0, a_0, s_1, a_1, \dots$
  - ▶ Then the learner **imitates** these demonstrations



# IL = Learning from Demonstrations (LfD)

- Teacher provides **demonstration** trajectories  $\mathcal{D} = \{\xi^{(1)}, \dots, \xi^{(m)}\}$
- Learner trains a policy  $\pi_\theta$  to **minimize a loss**  $\mathcal{L}(\theta)$
- For example, **negative log-likelihood (NLL)**:

$$\begin{aligned} \arg \min_{\theta} \mathcal{L}_\theta(\xi) &= \arg \min_{\theta} (-\log p_\theta(\xi)) \\ &= \arg \max_{\theta} \left( \log p(s_0) + \sum_{t=0}^{T-1} \log \pi_\theta(a_t | s_t) + \log p(s_{t+1} | s_t, a_t) \right) \\ &= \arg \max_{\theta} \sum_{t=0}^{T-1} \log \pi_\theta(a_t | s_t) \end{aligned}$$

**model-free**  
= no need to know the environment dynamics  $p$

# Today's lecture

---

**Behavior Cloning**

**Better behavior modeling**

**Alleviating train–test mismatch**

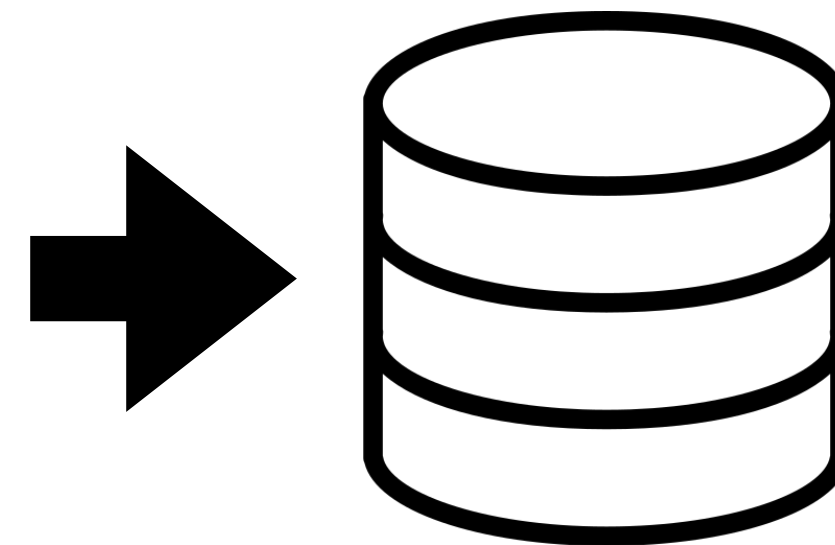
# Behavior Cloning (BC)

- Behavior Cloning:

- Break down trajectories  $\{\xi^{(1)}, \dots, \xi^{(m)}\}$  into steps  $\{(s_0^{(1)}, a_0^{(1)}), \dots, (s_{T_m-1}^{(m)}, a_{T_m-1}^{(m)})\}$
- Train  $\pi_\theta : s \mapsto a$  using supervised learning



observations  
+  
actions



training  
data

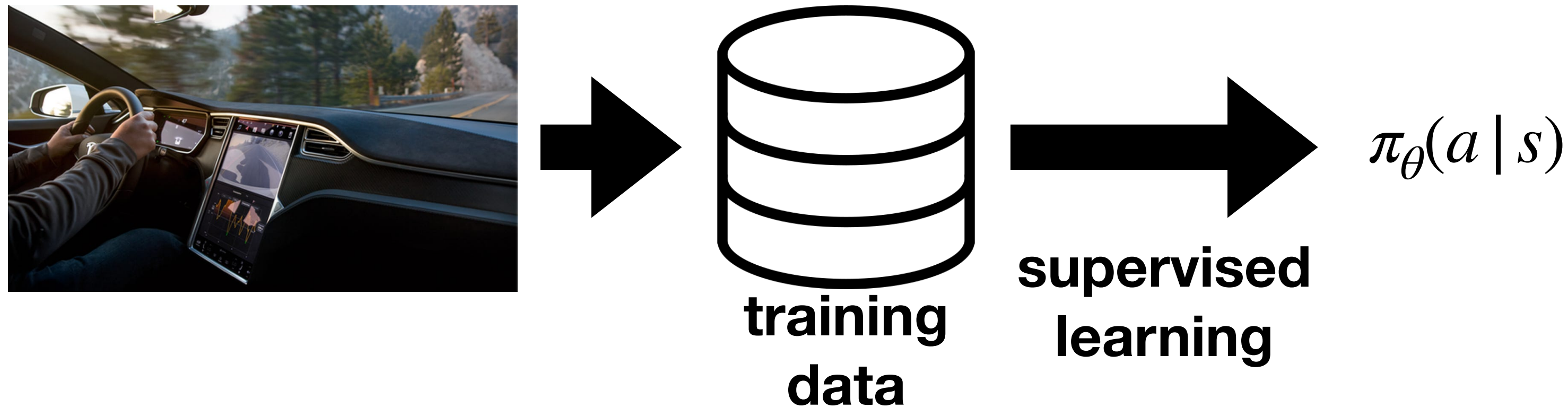
$$\mathcal{D} = \{(s_t^{(i)}, a_t^{(i)})\}_{i,t}$$



$\pi_\theta(a | s)$

$$\max_{\theta} \frac{1}{|\mathcal{D}|} \sum_{(s,a) \in \mathcal{D}} \log \pi_\theta(a | s)$$

# Behavior Cloning (BC)



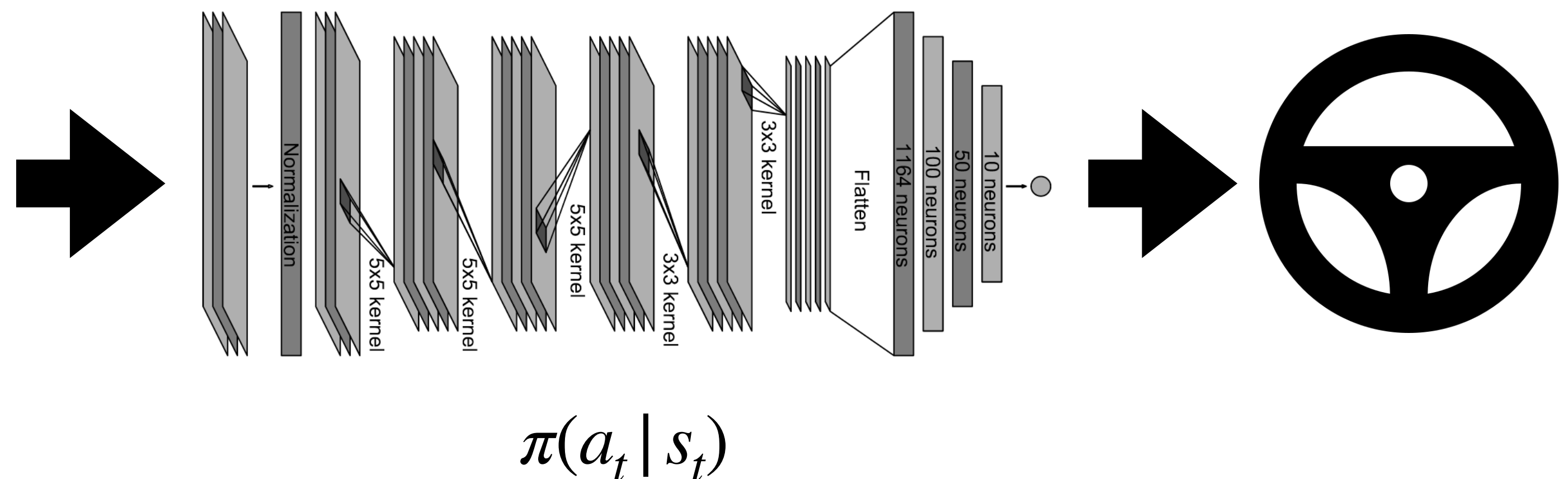
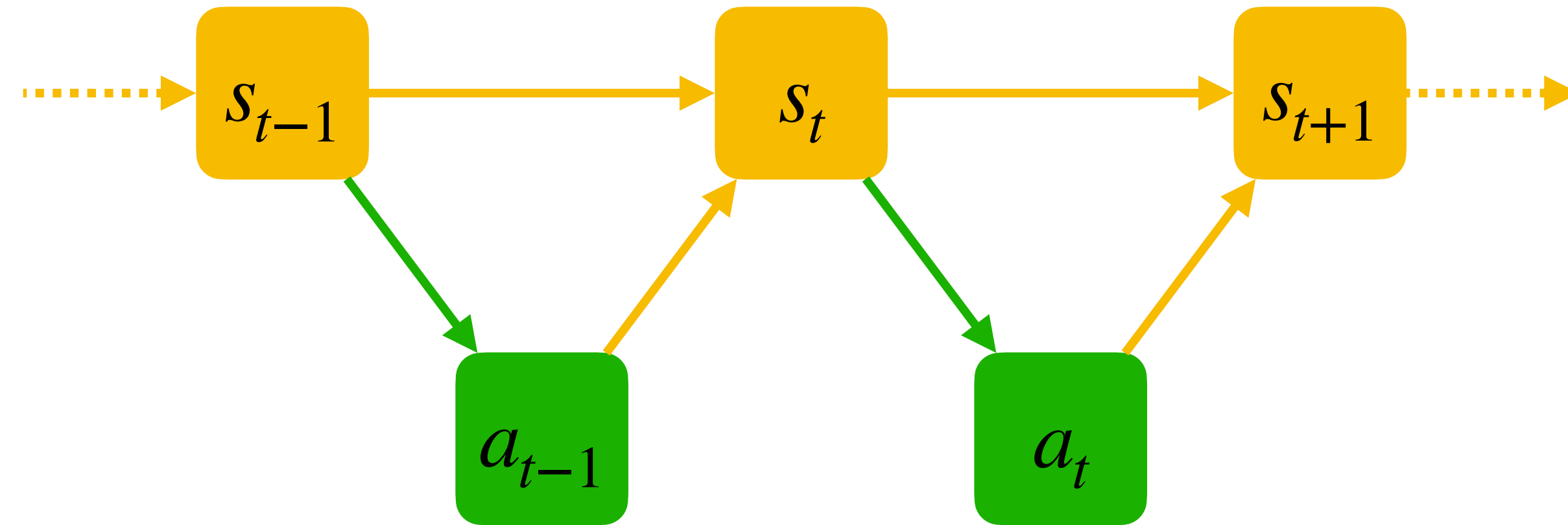
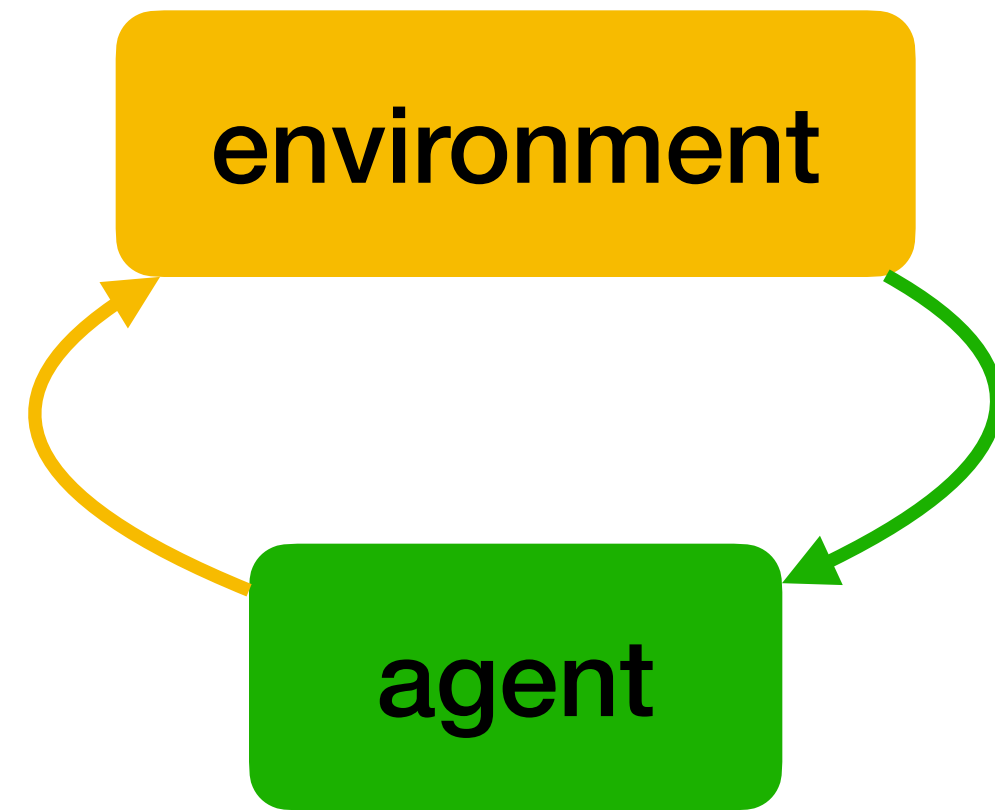
- **Benefits:**

- Simple, flexible — can use any learning algorithm
- **Model-free** — never need to know the environment

- **Drawbacks:**

- Only as good as the demonstrator
- Only good in **demonstrated states** — how do we evaluate?
  - Validation loss (on held out data)? Task success rate in rollouts?

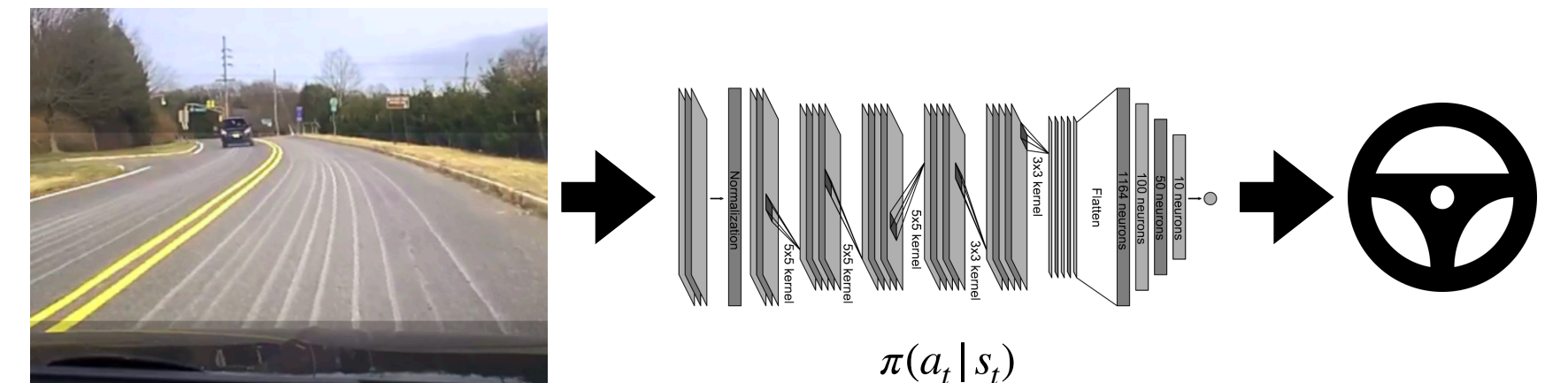
# A policy is a (stochastic) function



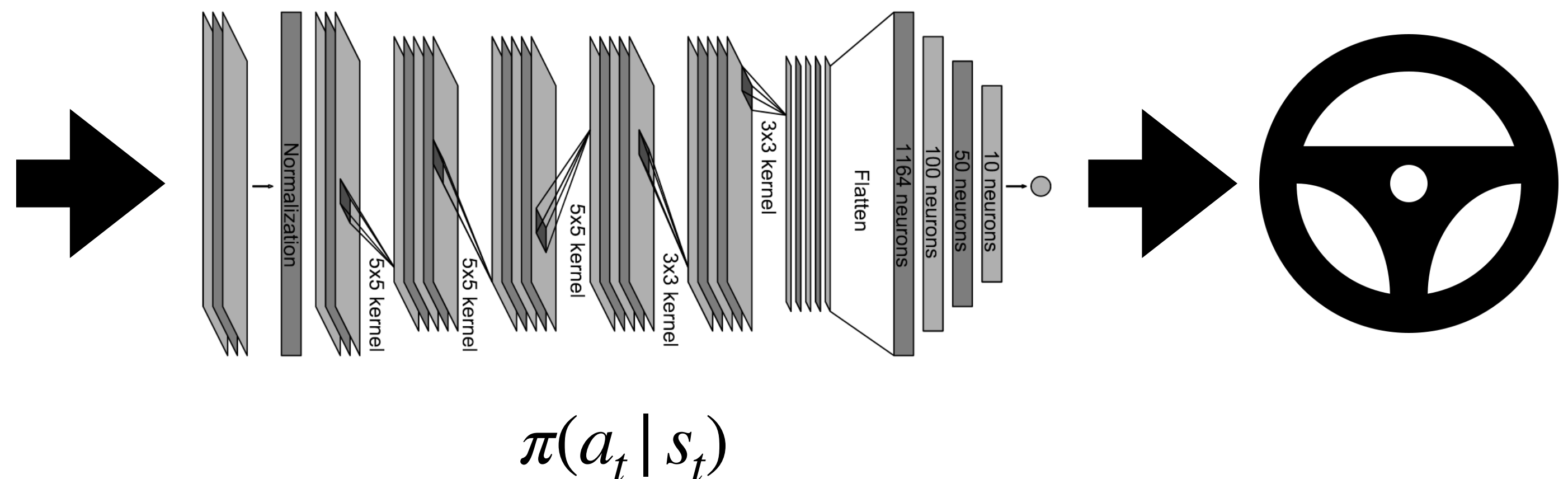
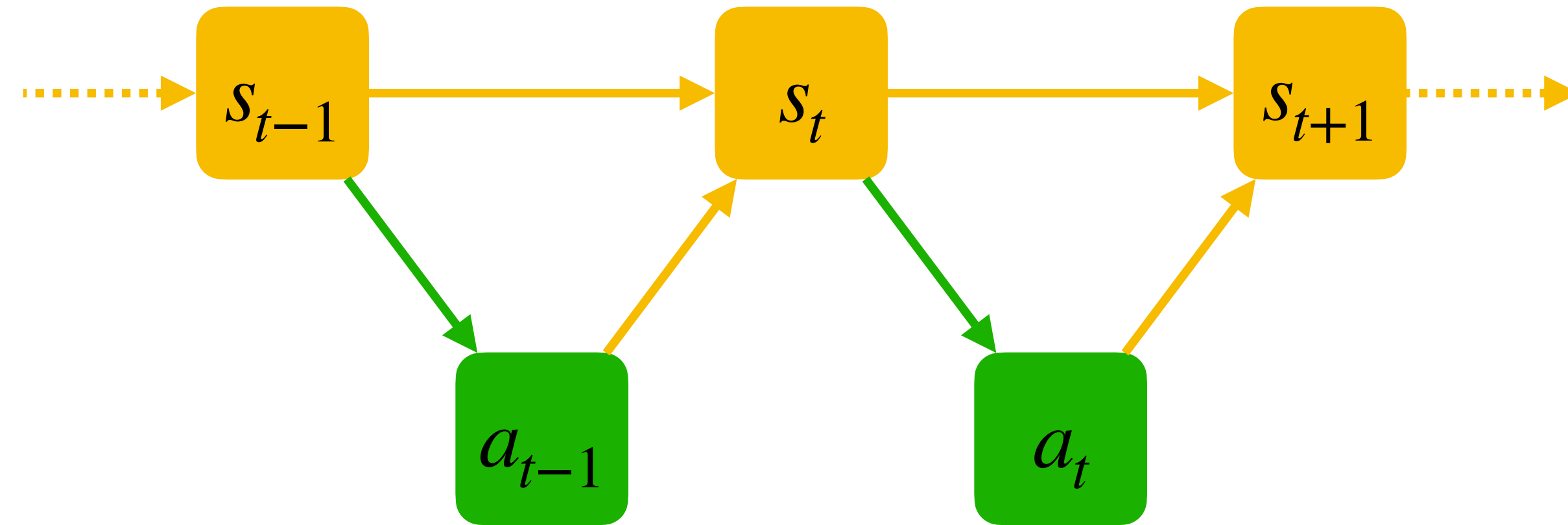
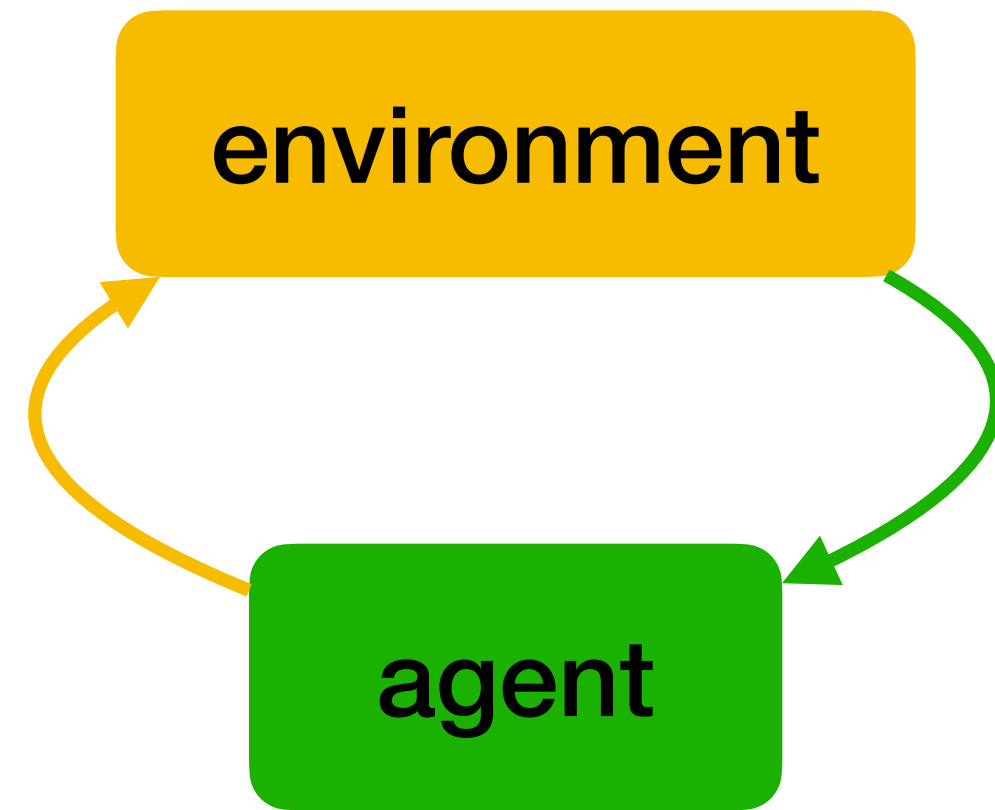


# Stochastic policies

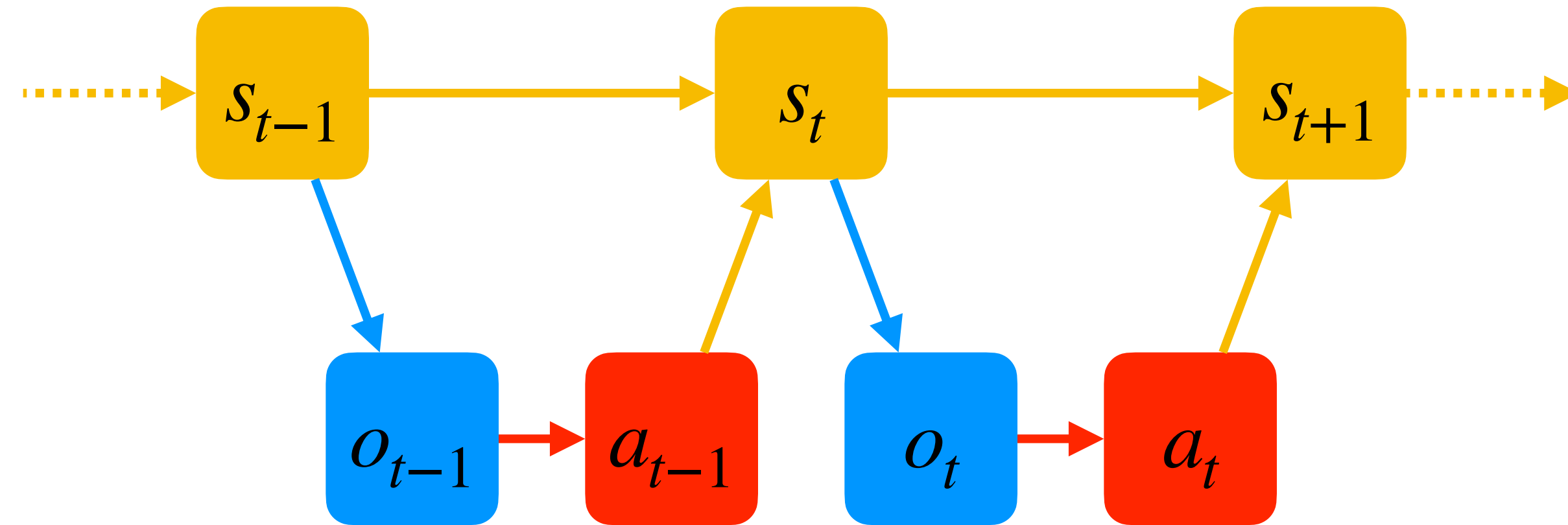
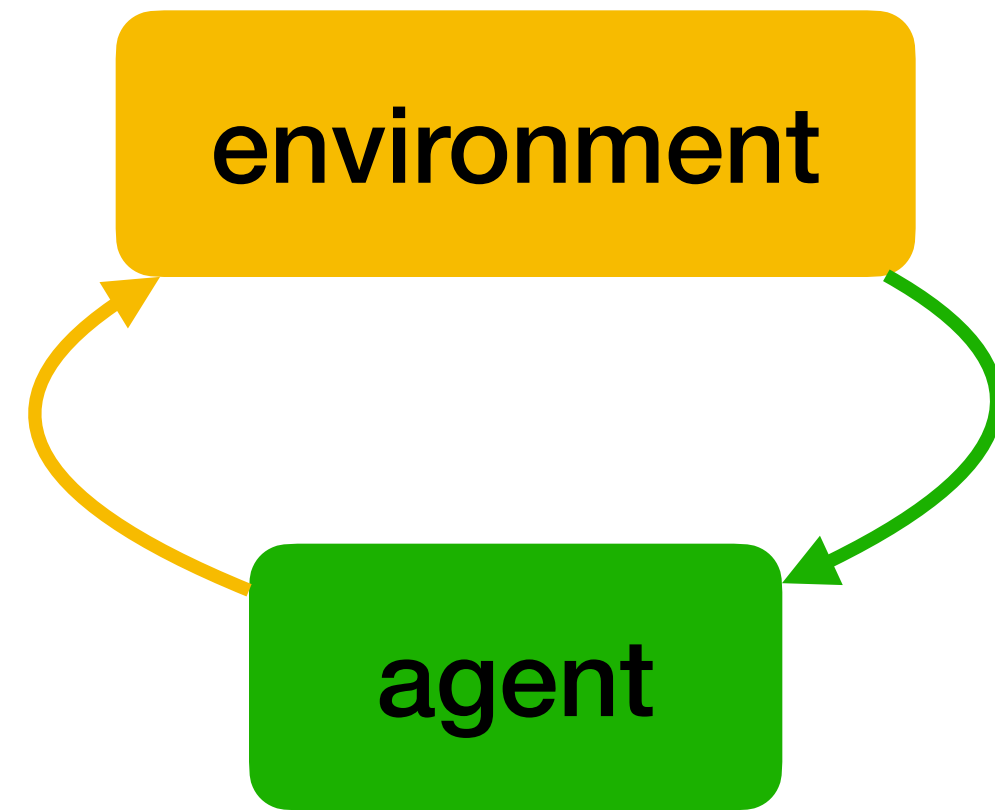
- Learned models are often **deterministic** functions  $f_\theta : x \mapsto y$
- To implement a stochastic policy: output **distribution parameters**
- Examples:
  - ▶ Discrete action space: **categorical** distribution
    - $\pi_\theta : s \mapsto \{\lambda_a\}_a; \pi_\theta(a | s) = \text{softmax}_a \lambda_a \propto \exp \lambda_a$
  - ▶ Continuous action space: **Gaussian** distribution
    - $\pi_\theta : s \mapsto (\mu, \Sigma); \pi_\theta(a | s) = \mathcal{N}(\mu, \Sigma)$



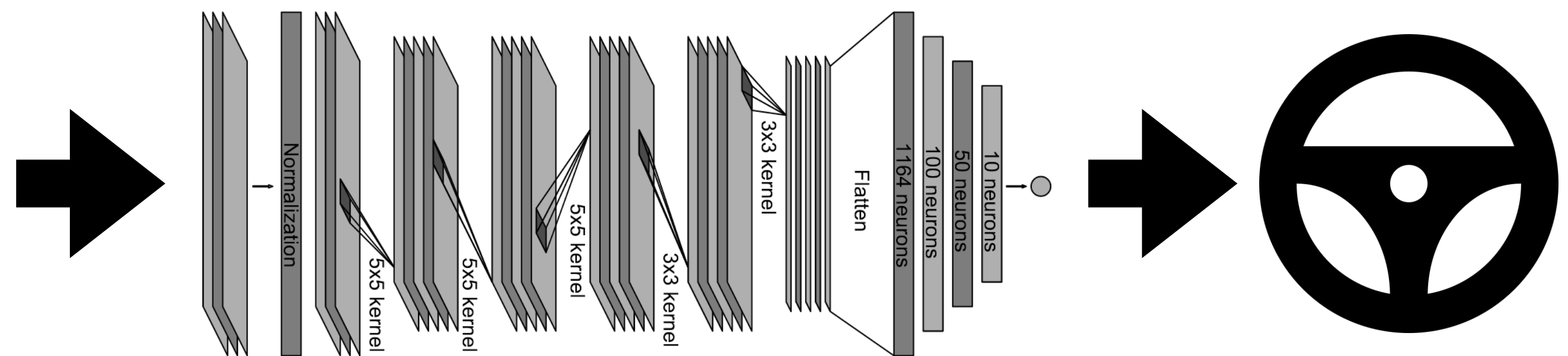
# A policy is a (stochastic) function



# A policy is a (stochastic) function



observation

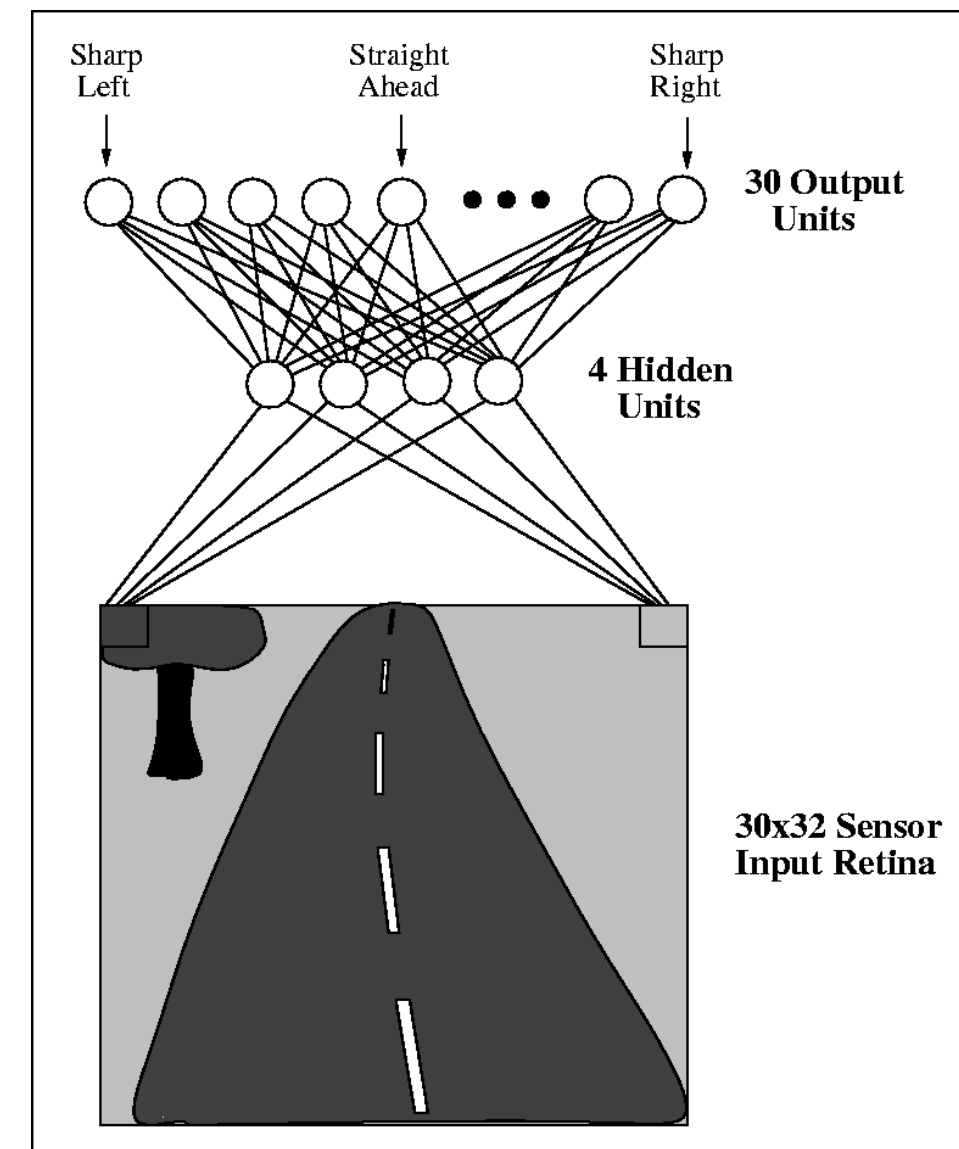


$$\pi(a_t | o_t)$$

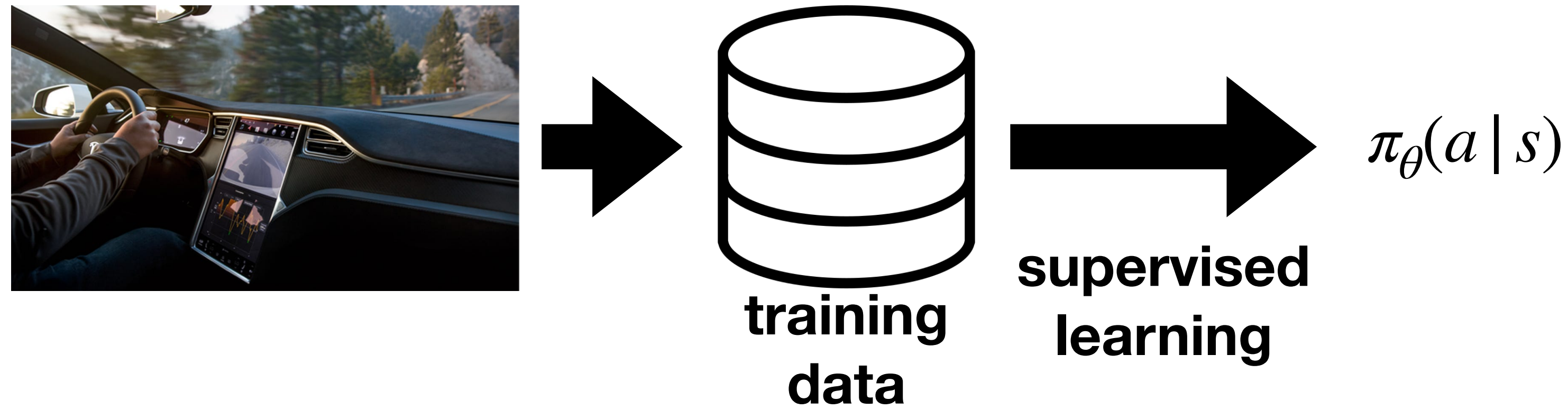
action

# ALVINN

- Autonomous Land Vehicle in a Neural Network (ALVINN, 1989)



# Inaccuracy in BC



- We could evaluate on held out teacher data, but really interested in **using**  $\pi_{\theta}$
- If the policy approximates the teacher  $\pi_{\theta}(a_t | s_t) \approx \pi^*(a_t | s_t)$ 
  - The trajectory distribution will also **approximate teacher behavior**  $p_{\theta}(\xi) \approx p^*(\xi)$
- But **errors accumulate** over time
  - May reach states **not seen** in the training dataset

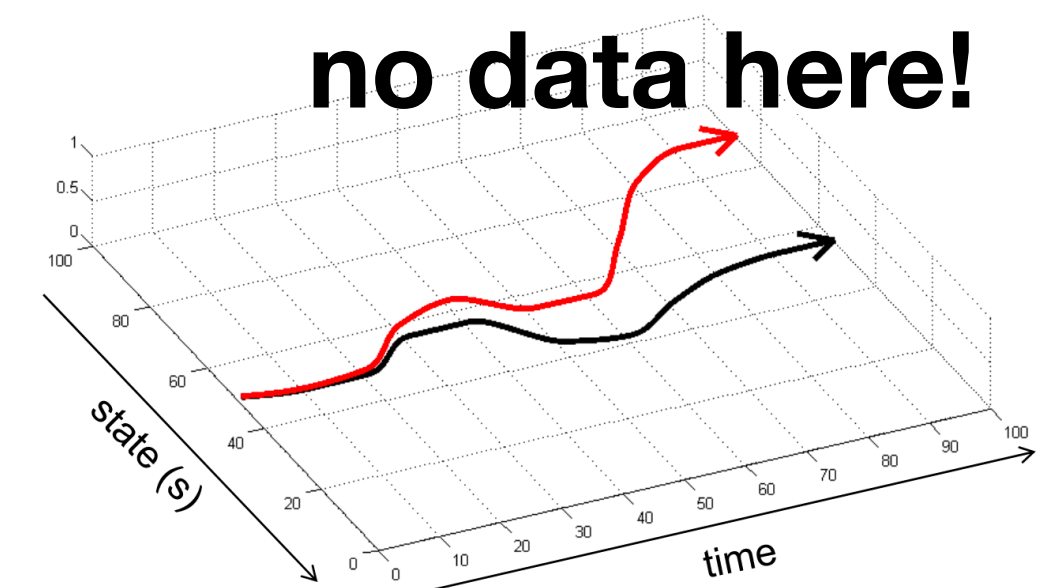
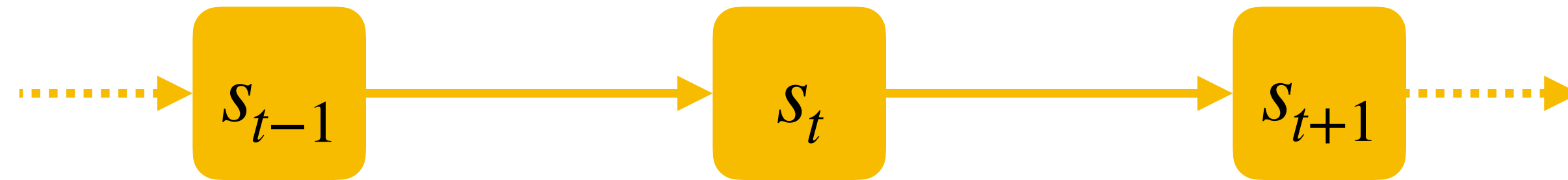
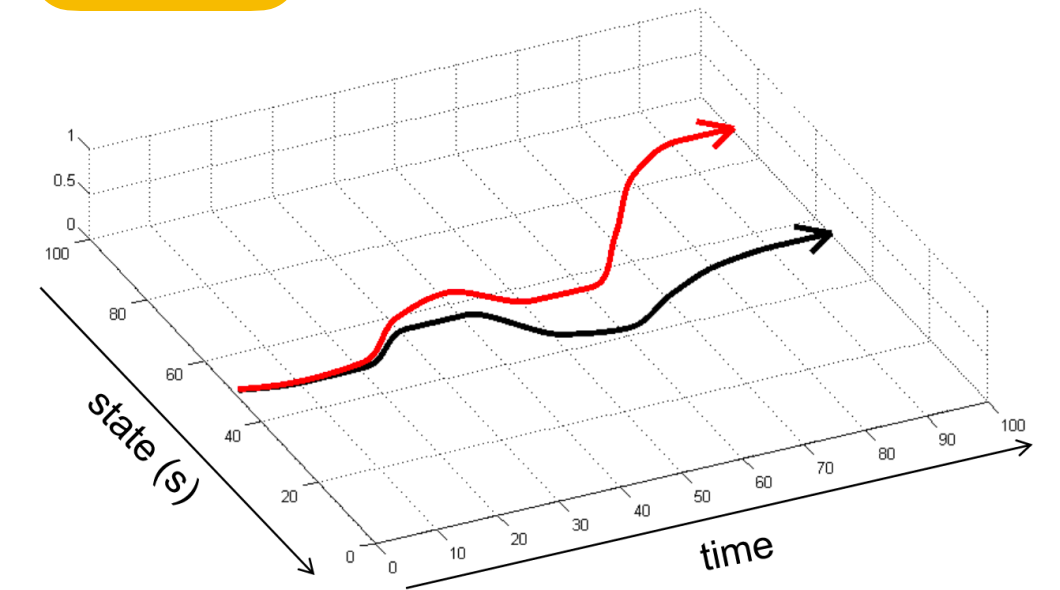


Image: Sergey Levine

# The impact of inaccurate dynamics



- **Errors** in learning are unavoidable
- What impact do they have on **sequential** behavior?



- Bounded **one-step error** in a dynamical model  $\sum_{s'} \left| p_{\theta}(s' | s) - p^*(s' | s) \right| \leq \epsilon$

▶ Can lead to growing error **over time**  $\sum_{s_t} \left| p_{\theta}(s_t) - p^*(s_t) \right| \leq \epsilon t$

▶ Not too bad by itself, but can **drift** outside training distribution  $\mathcal{D}$

# Today's lecture

---

Behavior Cloning

**Better behavior modeling**

Alleviating train–test mismatch

# Modeling other agents is hard

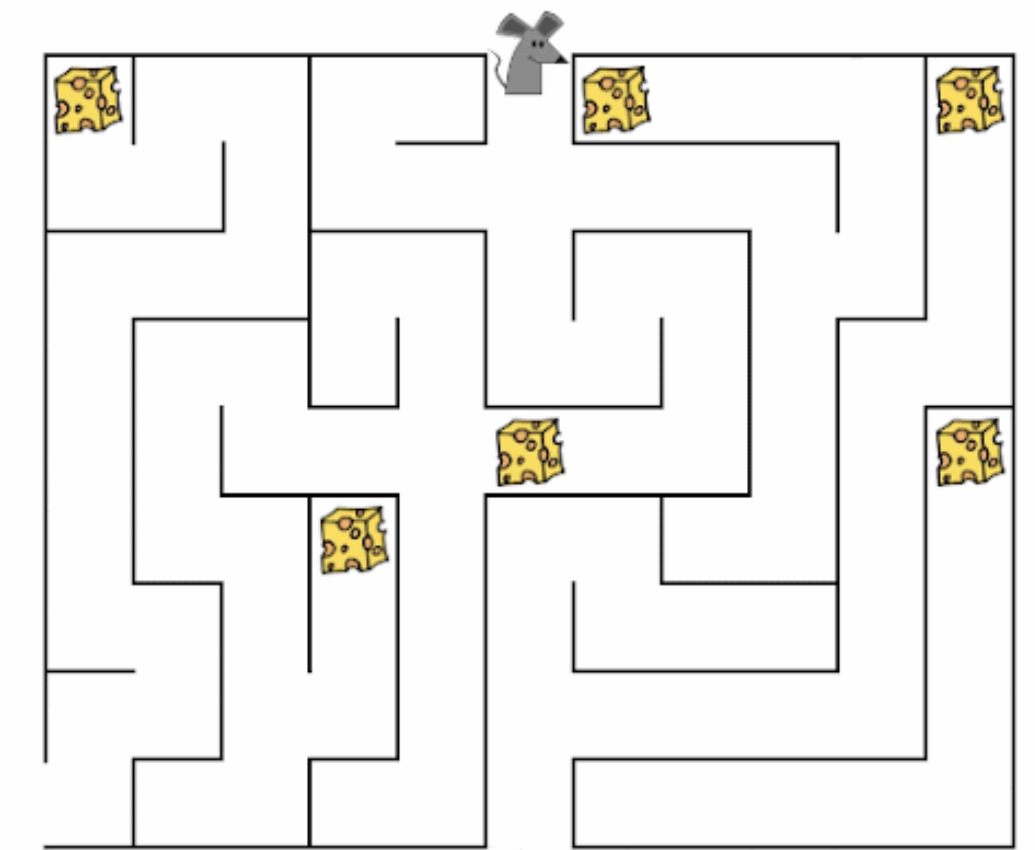
---

- Is there sufficient **data**? Demonstrating puts a burden on the teacher
- Are demonstrations **correct**? Humans are fallible, some supervision is hard
- Are demonstrations **consistent**? Some tasks can be done in multiple ways
- Is the state **partially observable**?  $o_t \stackrel{?}{=} s_t$
- Are the learner and teacher **observations** the same?  $o_t \stackrel{?}{=} o_t^*$



# Inconsistent demonstrations: multiple goals

- What if the task is to reach one of **multiple goals**?
  - Different episodes can successfully reach different goals
  - We need to train one policy to reach multiple goals
- If we know the goal, **condition** on it
  - **Goal-conditioned** policy:  $\pi_{\theta}(a_t | s_t, g)$
- More generally: **task-conditioned** policy  $\pi_{\theta}(a_t | s_t, \tau)$ 
  - **Goal** = desired final state; but how to represent other kinds of tasks?



# Goal-conditioned Behavior Cloning

- Can we train a goal-conditioned policy  $\pi_{\theta}(a_t | s_t, g)$  from demonstrations?
  - Assume **goal** = state that the agent should reach
- How can we know the goal in demonstrations  $\xi = s_0, a_0, s_1, a_1, \dots$ ?
  - Manual **labeling**?  $\mathcal{D} = \{(\xi^{(i)}, g^{(i)})\}_i$
- **Hindsight**: take each  $s_t$  as the goal of the trajectory leading to it

$$s_0, a_0, \dots, s_{t-1}, a_{t-1}, s_t = g$$

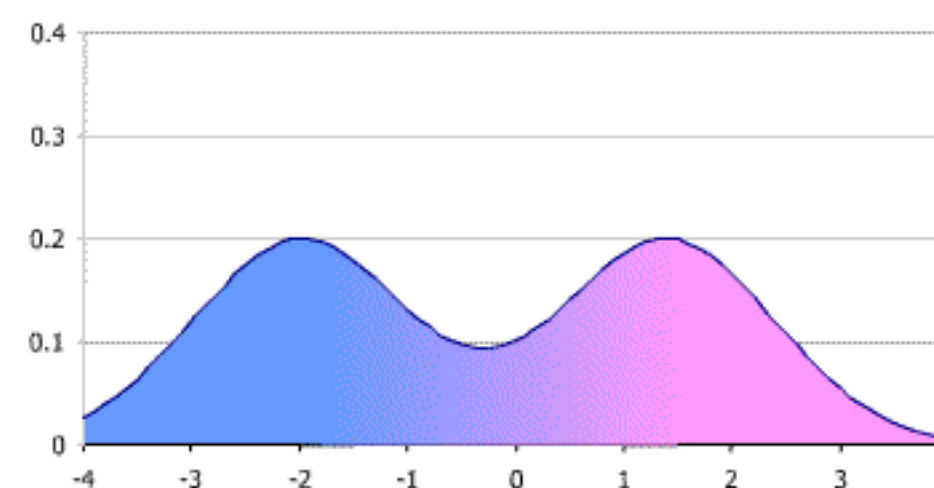
- Supervised learning of  $\pi(a | s, g)$  from data points  $((s_t, g = s_{t'}), a_t)$  for  $t' > t$

# Inconsistency due to multimodal behavior

- Goal-conditioning assumes **known goals**
  - More generally, known behavior modifiers
- Usually, the behavior **mode** is unknown



- Need **multimodal policy**  $\pi(a | s)$ 
  - Mixture models (e.g. GMM)
  - Latent-variable models (e.g. normalizing flows)

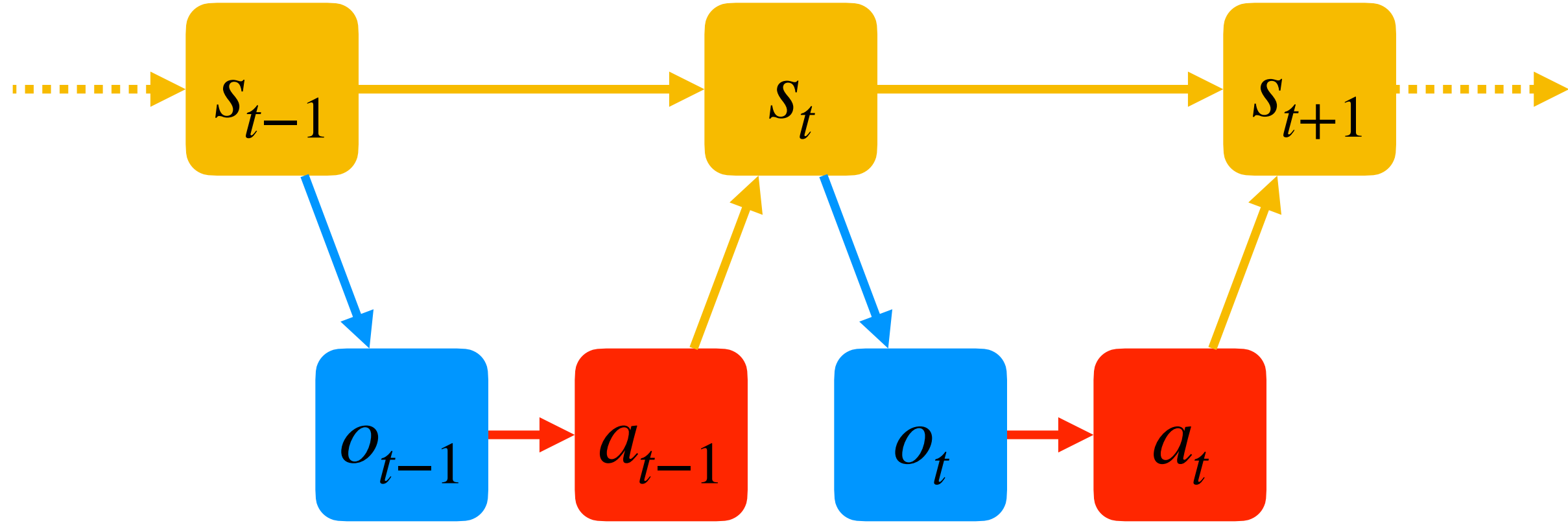
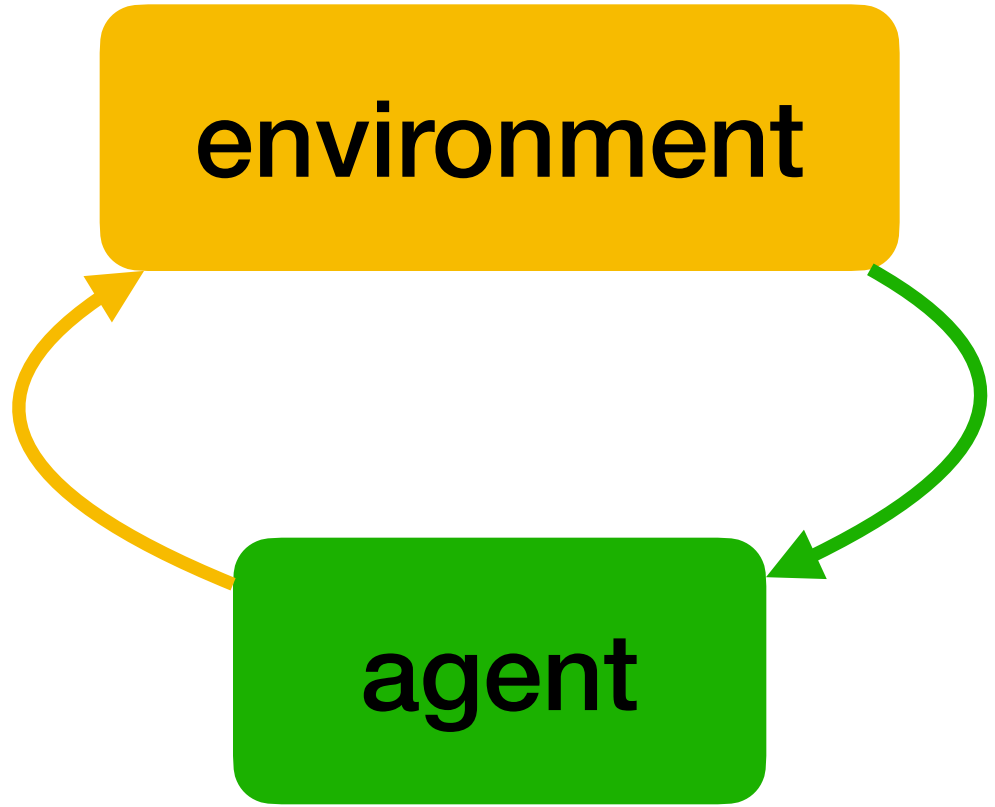


- Need to be **consistent** along a trajectory
  - Condition the policy on **memory of past actions**  $\pi(a_t | s_t, a_{\leq t})$

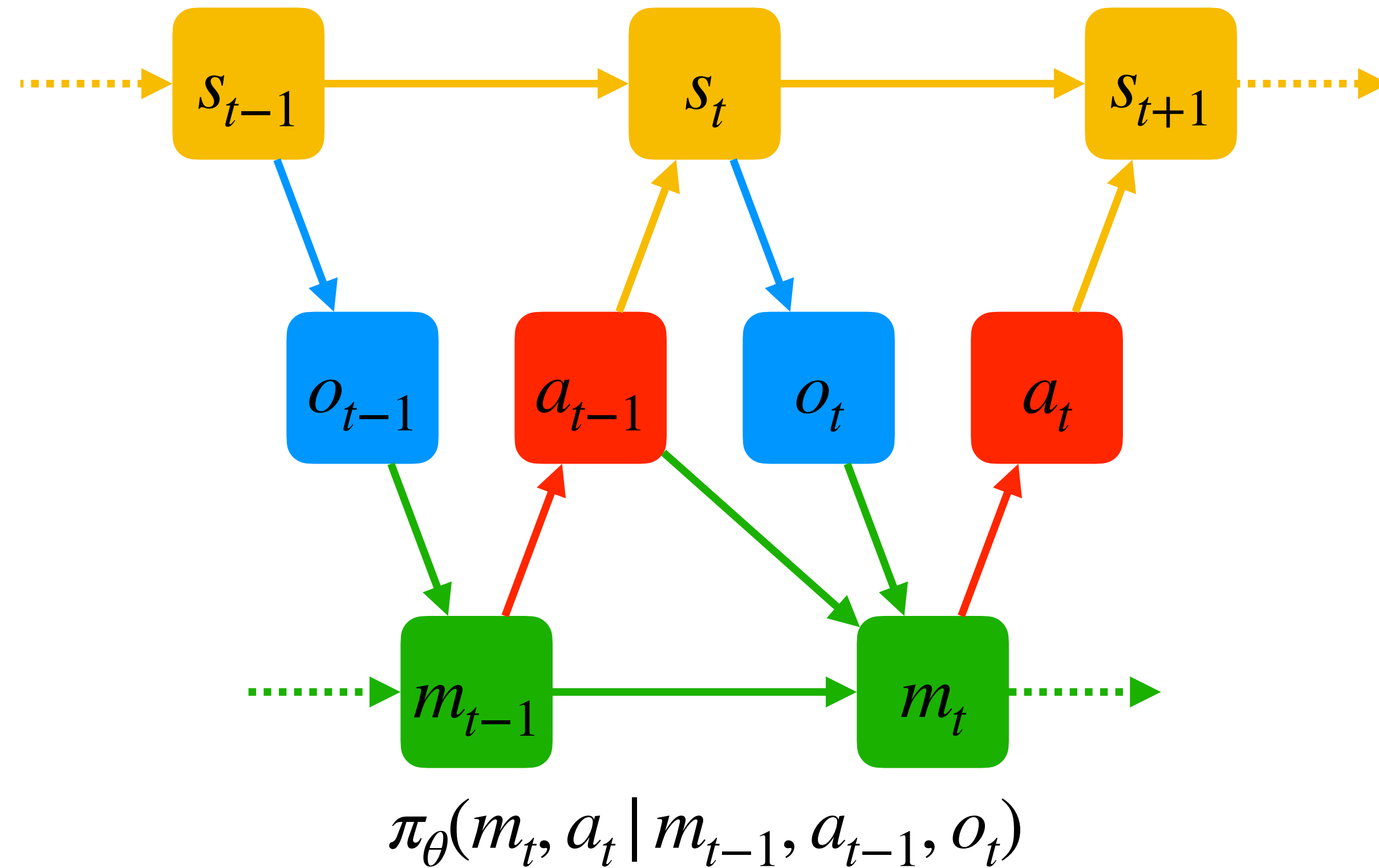
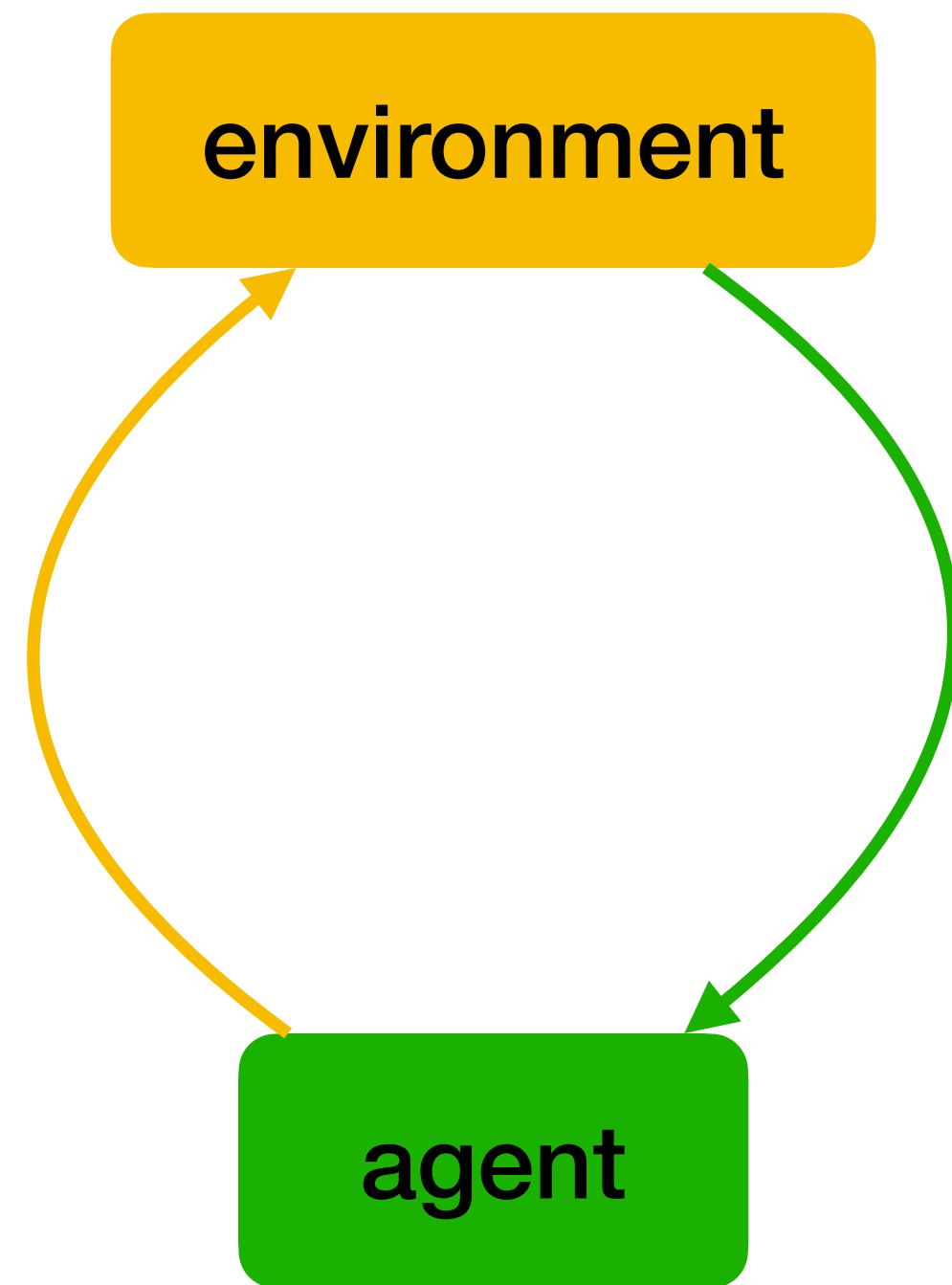
# Modeling partially observable behavior

- Partial observations are **not Markov**
  - Generally, this means  $p(o_{t+1} | o_t, a_t) \neq p(o_{t+1} | o_{\leq t}, a_{\leq t})$
  - **Reactive policy**  $\pi_{\theta}(a_t | o_t)$  may not be optimal
    - May need  $\pi_{\theta}(a_t | o_{\leq t})$ , or even  $\pi_{\theta}(a_t | o_{\leq t}, a_{<t})$ ; but how?
- Can use **RNNs**  $f_{\theta} : (h_{t-1}, a_{t-1}, o_t) \mapsto h_t$ , or other memory models
- But memory state is **latent** in demonstrations
  - Modeling memory is hard  $\rightarrow$  **prior structure** may help; more on this later

# Modeling memory



# Modeling memory



- A common architecture:

- A **recurrent** model  $m_t = f_{\theta}(m_{t-1}, a_{t-1}, o_t)$ ; and an **action** model  $\pi_{\theta}(a_t | m_t)$

# Today's lecture

---

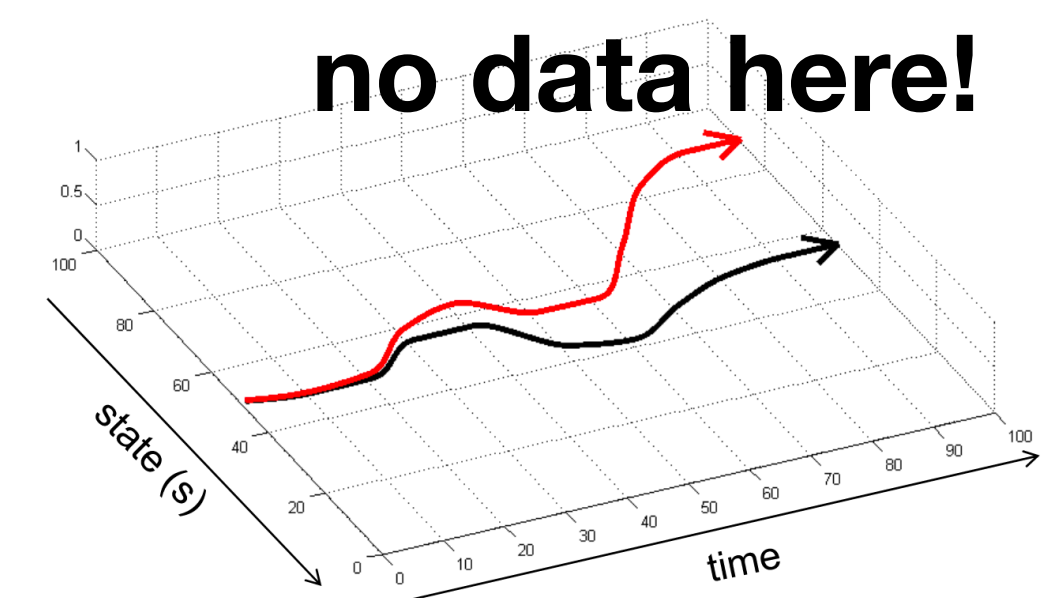
Behavior Cloning

Better behavior modeling

**Alleviating train–test mismatch**

# Alleviating train–test mismatch

- ML promises **generalization** when training distribution = test distribution
  - ▶ But this is challenging in IL: errors accumulate
  - ▶ We can quickly get to error states that we haven't seen fixed
  - ▶ Train–test distribution mismatch = **covariate shift**
- Ideas:
  - ▶ **Augment** the training dataset to expand the distribution
  - ▶ **Update** train distribution → test distribution
  - ▶ **Intervene** during demonstrations to expand the distribution

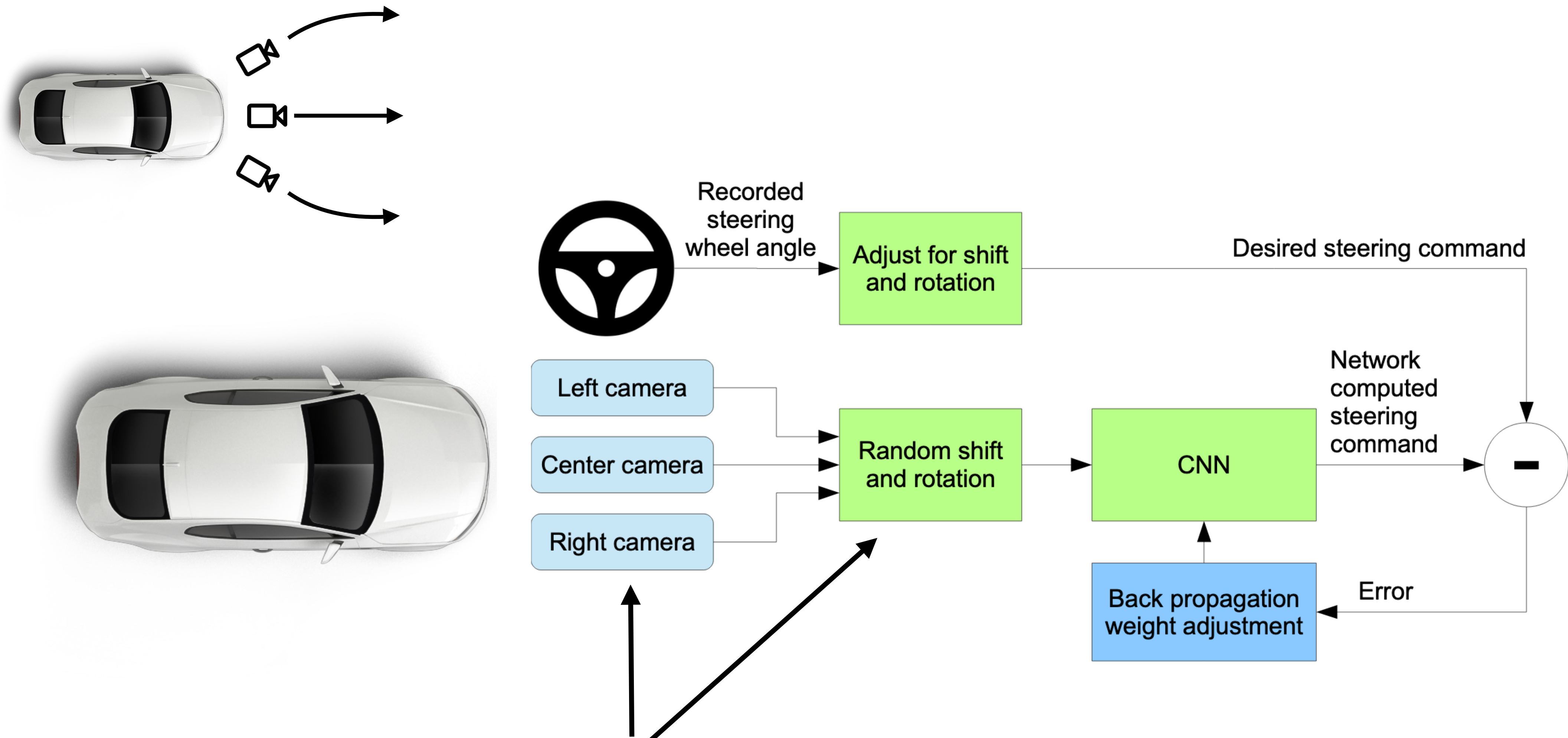




# Imitation Learning can work



# How did they do it?



**augmented data to better cover test distribution**

# DAgger: Dataset Aggregation

- Can we collect demonstration data from the **test distribution**?
  - We don't know  $p_\theta(\xi)$  until we're **done training**  $\theta$
  - But we **get closer** and closer during training

---

## Algorithm DAgger

---

Collect dataset  $\mathcal{D}$  of teacher demonstrations  $\xi \sim p^*$

**repeat**

Train  $\pi_\theta$  on  $\mathcal{D}$

Execute  $\pi_\theta$  to get  $\xi \sim p_\theta$

Ask teacher to label  $(a_t^* | s_t) \sim \pi^*$

Aggregate  $\{(s_t, a_t^*)\}_t$  into  $\mathcal{D}$

---

**but how? challenging...**

# DAgger demo



It turns automatically to avoid trees  
based on what its camera sees

# DART: Disturbances Augmenting Robot Training

- Off-policy vs. on-policy

- ▶ **On-policy** = data comes from the learner's current policy

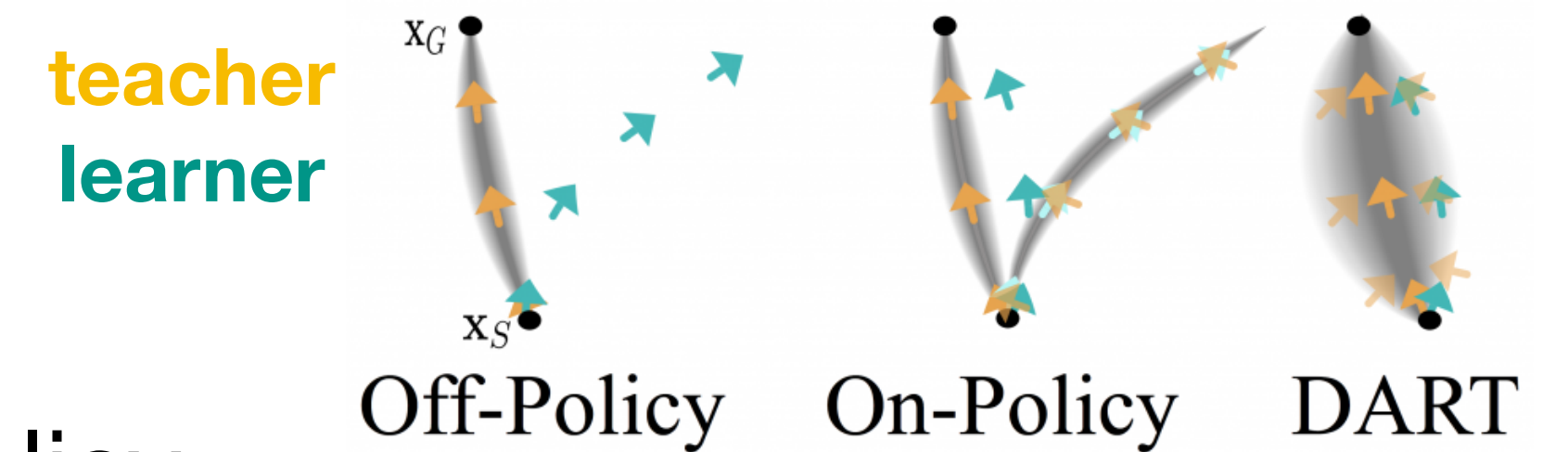
- ▶ **Off-policy** = data comes from another policy (another agent or past learner)

- In off-policy IL (e.g. BC) learner may go off the teacher's support

- In on-policy IL (e.g. DAgger) learner initially goes off, until corrected

- **DART**: increase the data support by injecting noise during demonstrations

- ▶ Force teacher into slight-error states, to see how they are fixed



# DART

- **Noise** = perturbation of actions  $q(\tilde{a} | a)$

▶ New effective dynamics:  $\tilde{p}(s' | s, a) = \sum_{\tilde{a}} q(\tilde{a} | a) p(s' | s, \tilde{a})$

▶ For example, in continuous actions:  $\tilde{a} = a + \epsilon; \quad \epsilon \sim \mathcal{N}(0, \Sigma)$

---

## Algorithm DART

---

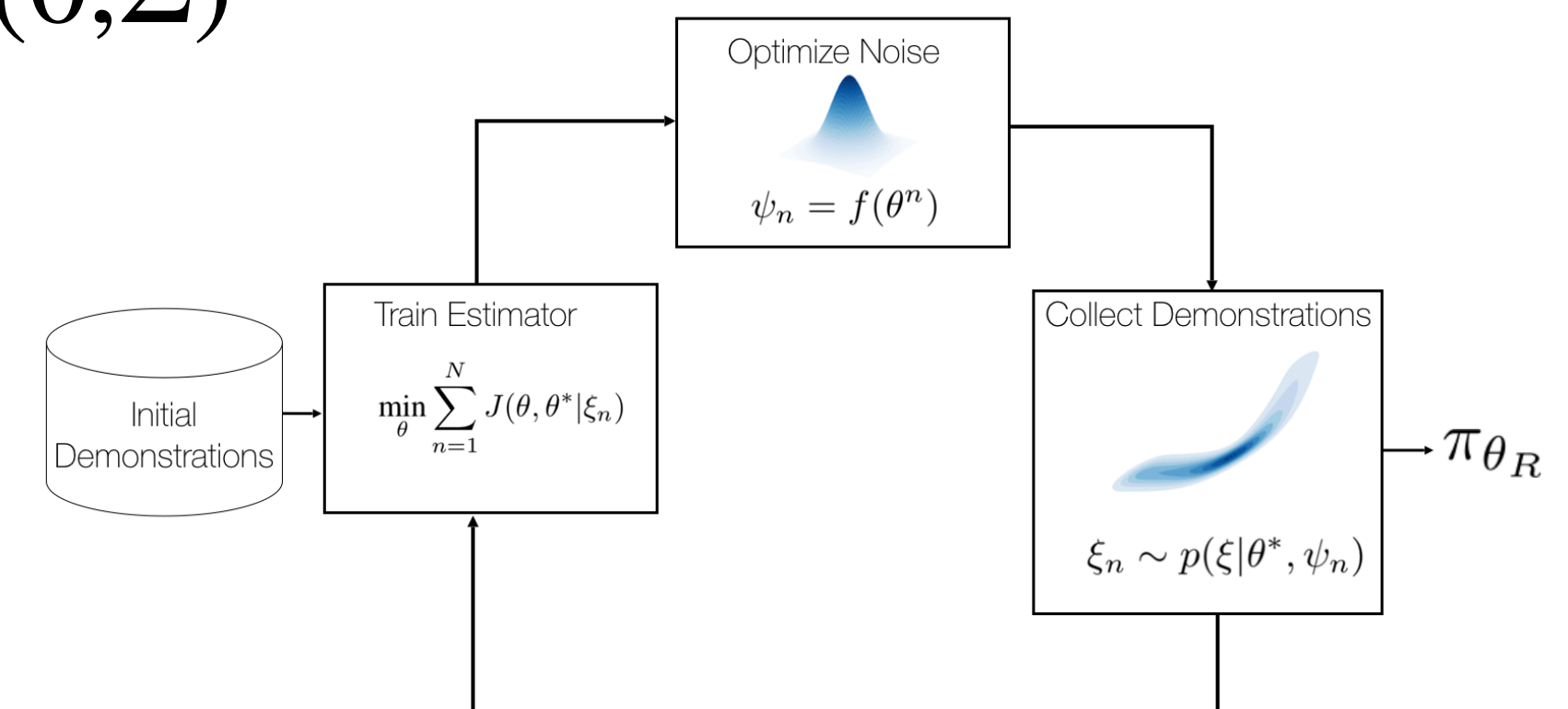
**repeat**

Collect dataset  $\mathcal{D}$  of teacher demonstrations  $\xi \sim \tilde{p}^*$

Train  $\pi_\theta$  on  $\mathcal{D}$

Update noise  $q$  such that  $p_\theta$  is better supported by  $\tilde{p}^*$

---

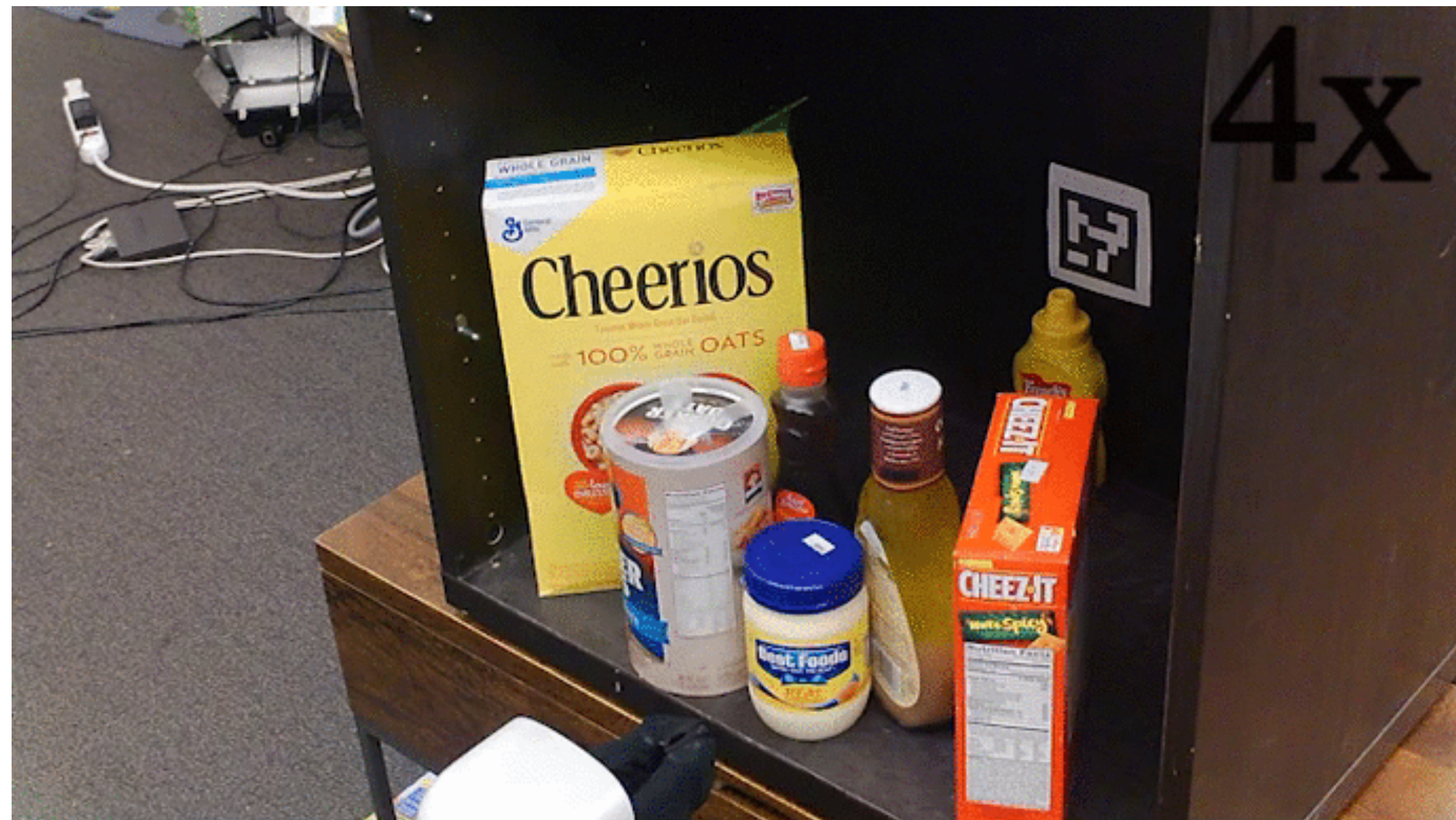


# Grasping task



Behavior Cloning

DART



[Laskey et al., 2017]

# Recap

- **Imitation Learning** = Learning from Demonstrations
  - Learn policy  $\pi(a | s)$  from teacher demonstrations
- **Behavior Cloning**: supervised learning
  - Minimize loss, e.g. NLL, on training set of trajectories
- Accurate imitation is crucial
  - Improve imitation through **goal-conditioning**, **multimodal** policies, **memory**, etc.
- Errors accumulate and cause **train–test distribution mismatch**
  - Can be alleviated through **augmentation**, **on-policy** data collection, **noise** injection

