## CSW182（2021）· 课程资料包 @ShowMeAI

**视频**
中英双语字幕

**课件**
一键打包下载

**笔记**
官方笔记翻译

**代码**
作业项目解析

**视频·B 站 [ 扫码或点击链接 ]**
https://www.bilibili.com/video/BV1Ff4y1n7ar

**课件 & 代码·博客 [ 扫码或点击链接 ]**
http://blog.showmeai.tech/berkeley-csw182

Berkeley
循环神经网络  可视化  梯度策略
Q-Learning  风格迁移  模仿学习  元学习
计算机视觉  机器学习基础  生成模型  卷积网络

Awesome AI Courses Notes Cheatsheets 是 **ShowMeAI** 资料库的分支系列，覆盖最具知名度的 **TOP50+** 门 AI 课程，旨在为读者和学习者提供一整套高品质中文学习笔记和速查表。

**点击**课程名称，跳转至课程**资料包**页面，**一键下载**课程全部资料！

| 机器学习 | 深度学习 | 自然语言处理 | 计算机视觉 |
|---|---|---|---|
| Stanford · CS229 | Stanford · CS230 | Stanford · CS224n | Stanford · CS231n |

**# Awesome AI Courses Notes Cheatsheets· 持续更新中**

| 知识图谱 | 图机器学习 | 深度强化学习 | 自动驾驶 |
|---|---|---|---|
| Stanford · CS520 | Stanford · CS224W | UCBerkeley · CS285 | MIT · 6.S094 |

**微信公众号**

资料下载方式 2：扫码点击底部菜单栏

称为 **AI 内容创作者？** 回复 [ 添砖加瓦 ]
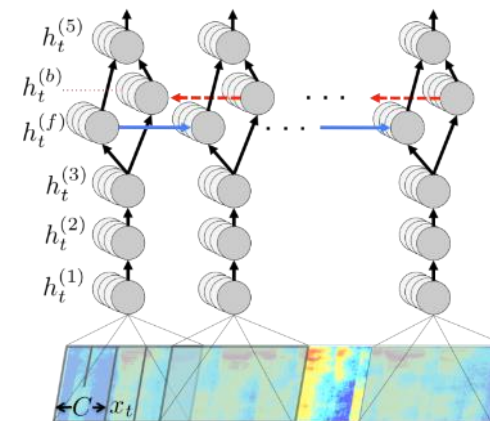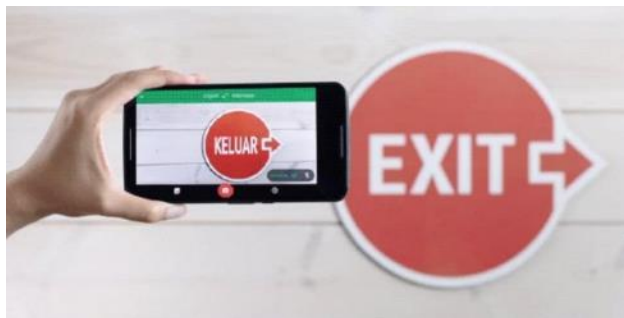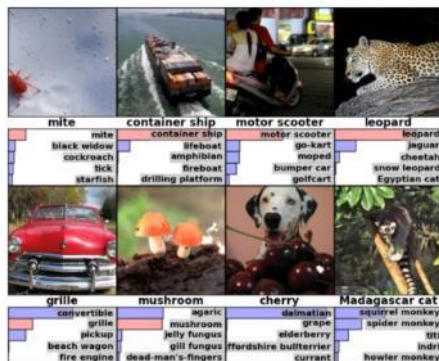
# Learning-Based Control & Imitation

Designing, Visualizing and Understanding Deep Neural Networks

# CS W182/282A

Instructor: Sergey Levine
UC Berkeley

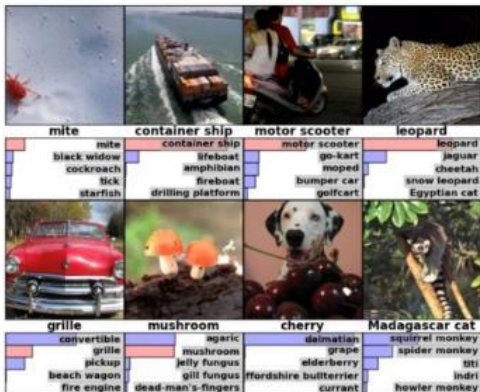# So far: learning to *predict*

What about learning to **control**?

# From *prediction* to *control*: challenges



i.i.d.: $p(\mathcal{D}) = \prod_i p(y_i|x_i)p(x_i)$

output $y_1$ does not change $x_2$

this is **very** important, because it allows us to just focus on getting the highest **average** accuracy over the whole dataset

making the wrong choice here is a disaster



making the wrong choice here is perhaps OK

# From *prediction* to *control*: challenges



**Ground truth labels:**



"puppy"

**Abstract goals:**



"drive to the grocery store"

> what steering command is that?

# From *prediction* to *control*: challenges



- i.i.d. distributed data (each datapoint is independent)
- ground truth supervision
- objective is to predict the right label

These are not **just** issues for control: in many cases, real-world deployment of ML has these same **feedback** issues
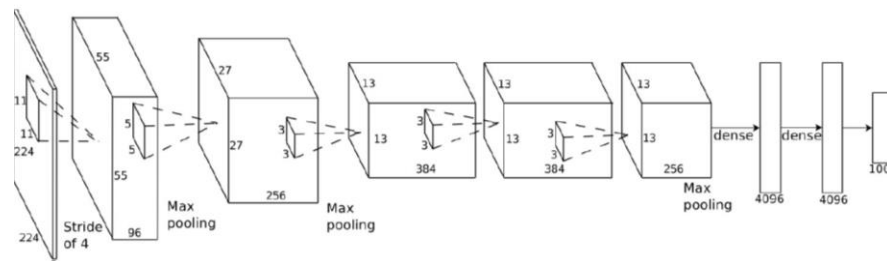
**Example:** decisions made by a traffic prediction system might affect the route that people take, which changes traffic



- each decision can change future inputs (not independent)
- supervision may be high-level (e.g., a goal)
- objective is to accomplish the task

We will **build up** toward a **reinforcement learning** system that addresses all of these issues, but we'll do so one piece at a time...

# Terminology



$\mathbf{o}_t$

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

used to be $p_\theta(y|x)$

$\mathbf{a}_t$

used to be $x$

used to be $y$

$\mathbf{s}_t$ − 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)

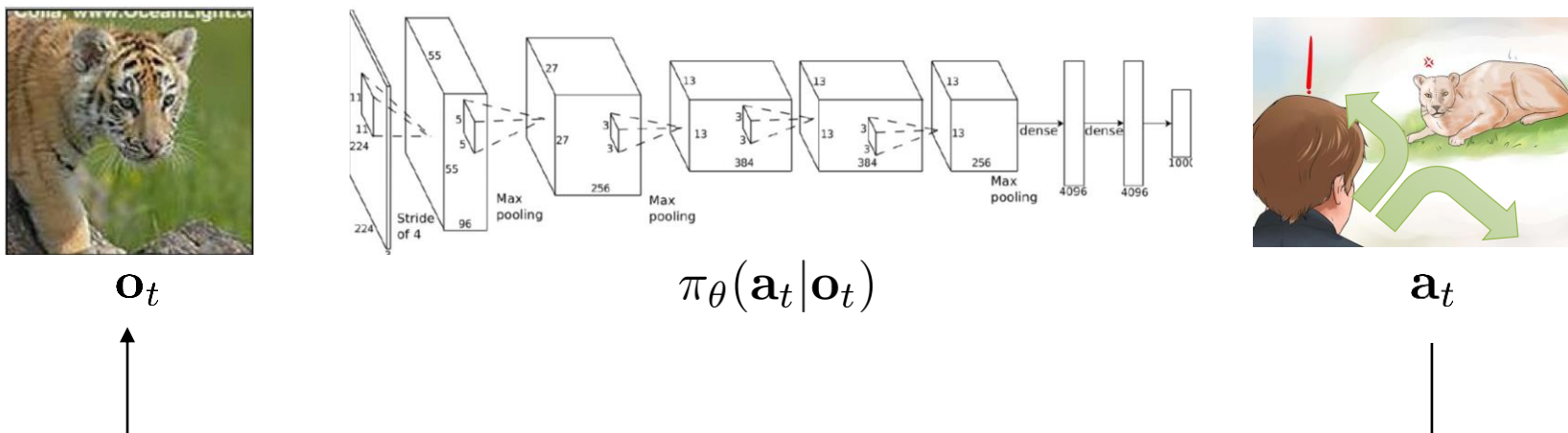This distinction will very important later, but is not so important today

$\mathbf{o}_t$ − observation

$\mathbf{s}_t$ − state

# Terminology



$$\mathbf{o}_t$$

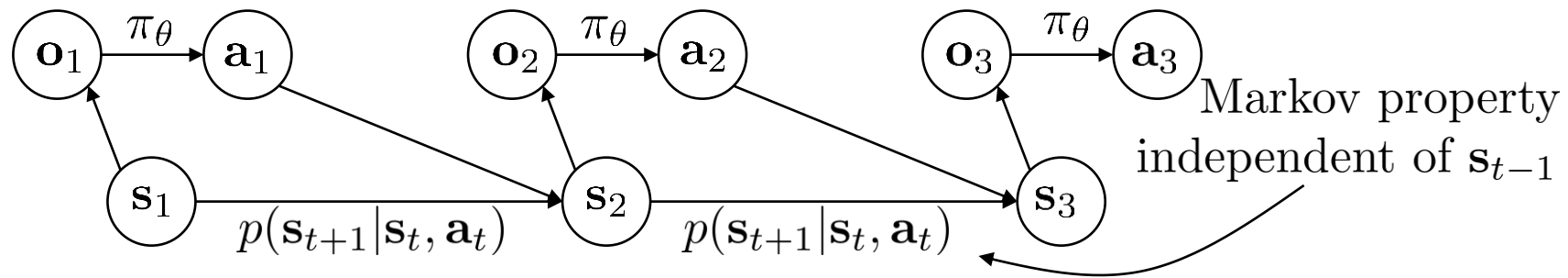$$\pi_\theta(\mathbf{a}_t|\mathbf{o}_t)$$

$$\mathbf{a}_t$$

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

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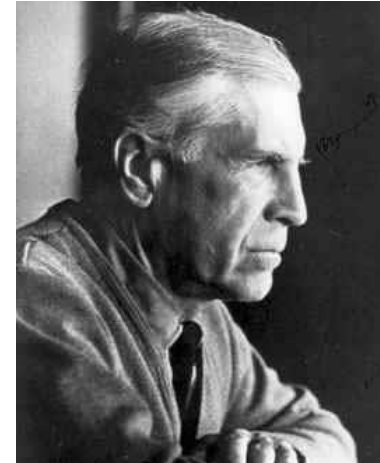$\pi_\theta(\mathbf{a}_t|\mathbf{s}_t)$ – policy (fully observed)



Markov property independent of $\mathbf{s}_{t-1}$

# 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



$\mathbf{o}_t$

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

$\mathbf{a}_t$

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

training data

supervised learning

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

behavioral cloning

# Does it work?

# No!



Where have we seen this before?

# Does it work?                    Yes!



Video: Bojarski et al. '16, NVIDIA
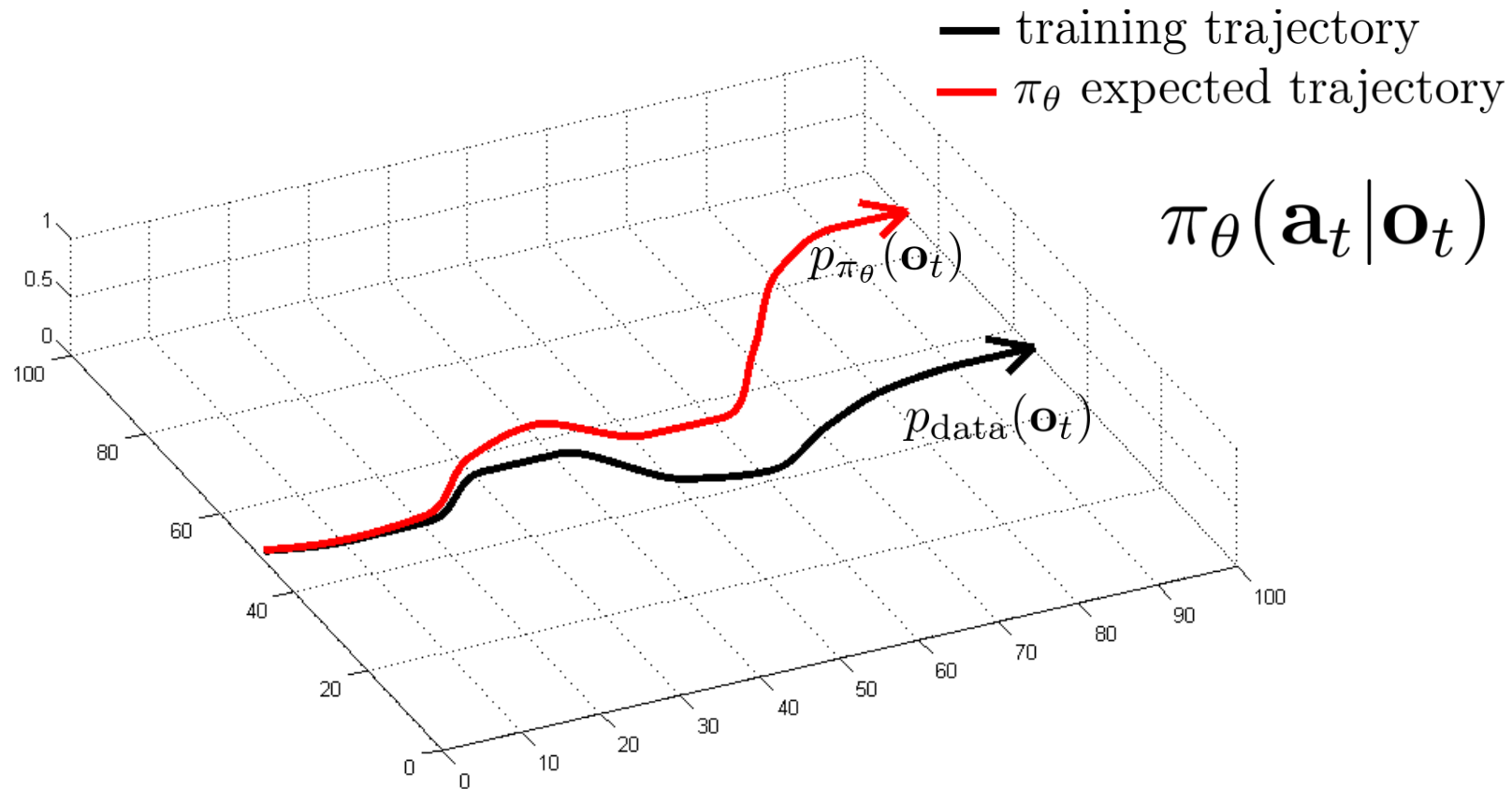
# Getting behavioral cloning to work

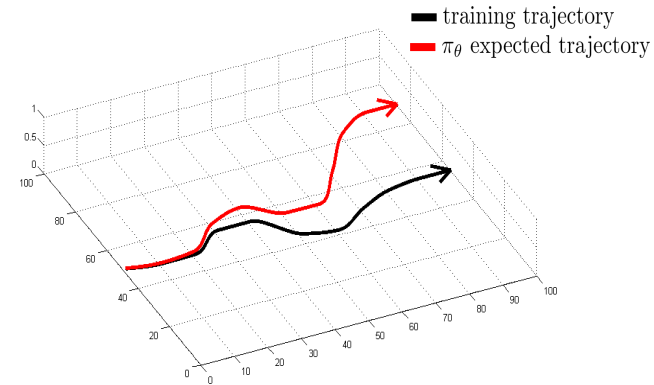# What is the problem?



the problem: $p_{\text{data}}(\mathbf{o}_t) \neq p_{\pi_\theta}(\mathbf{o}_t)$

# What is the problem?


training trajectory
$\pi_\theta$ expected trajectory

the problem: $p_{\text{data}}(\mathbf{o}_t) \neq p_{\pi_\theta}(\mathbf{o}_t)$

think: 0.6
like: 0.3
drive: 0.1

hippo: 0.6
paintbrush: 0.3
California: 0.1

complete nonsense,
because the network never
saw inputs remotely like this

This is the same problem!

unlikely but
possible
mistake

I          drive

we got unlucky, but now the
model is completely confused

it never saw "I drive" before

**The problem:** this is a training/test discrepancy:
the network always saw **true** sequences as
inputs, but at test-time it gets as input its own
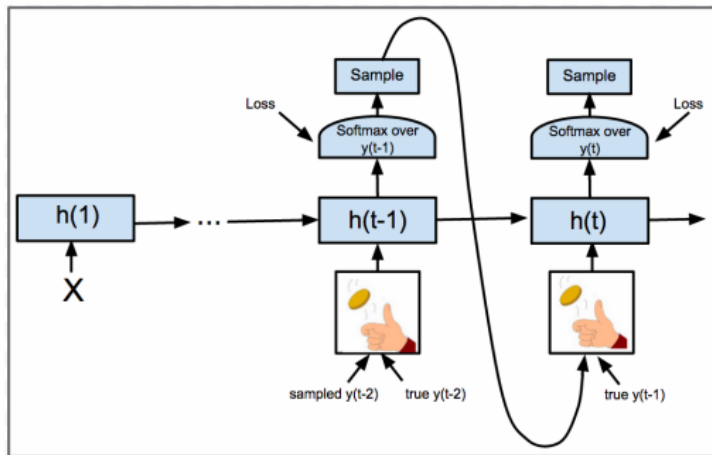(potentially incorrect) predictions

This is called **distributional shift**, because the
input distribution **shifts** from true strings (at
training) to synthetic strings (at test time)

# Why not use the same solution?

the problem: $p_{\text{data}}(\mathbf{o}_t) \neq p_{\pi_\theta}(\mathbf{o}_t)$
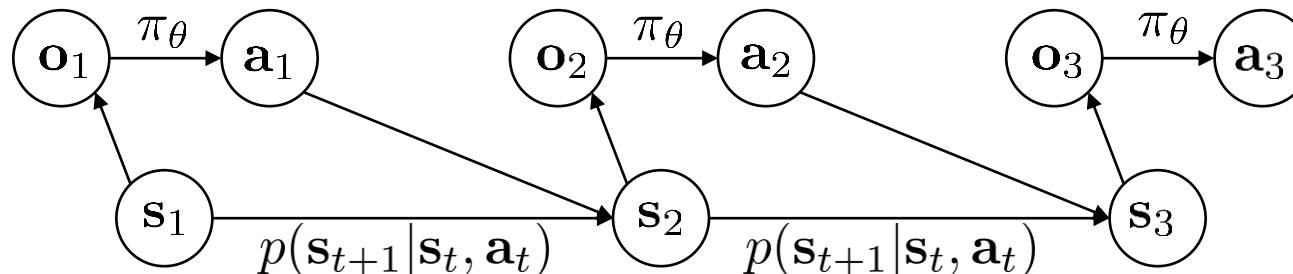
**Before:** scheduled sampling



**Now:** control

we *could* take the predicted action $\mathbf{a}_t \sim \pi_\theta(\mathbf{a}_t|\mathbf{o}_t)$ and observe the resulting $\mathbf{o}_{t+1}$

but this requires interacting with the world! why?

we don't know $p(\mathbf{s}_{t+1}|\mathbf{s}_t, \mathbf{a}_t)$!

# Can we **mitigate** the problem?

the problem: $p_{\text{data}}(\mathbf{o}_t) \neq p_{\pi_\theta}(\mathbf{o}_t)$

if $\pi_\theta(\mathbf{a}_t|\mathbf{o}_t)$ is *very* accurate
maybe $p_{\text{data}}(\mathbf{o}_t) \approx p_\theta(\mathbf{o}_t)$

Why **might** we fail to fit the expert?

➡ 1. Non-Markovian behavior

2. Multimodal behavior

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

behavior depends only on current observation
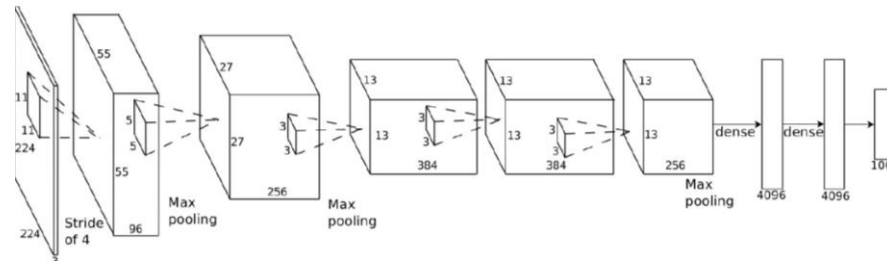
behavior depends on all past observations

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

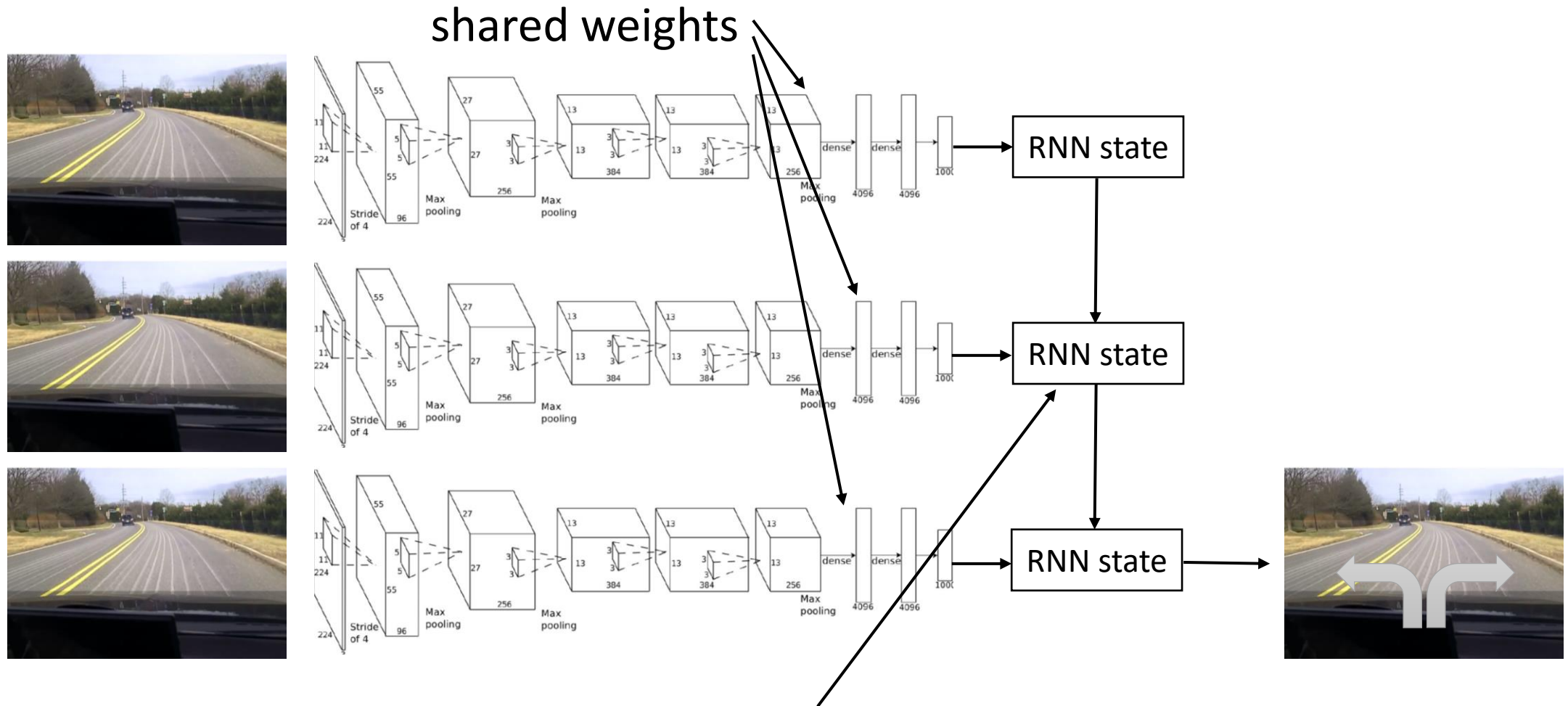Often very unnatural for human demonstrators

# How can we use the whole history?



variable number of frames,
too many weights

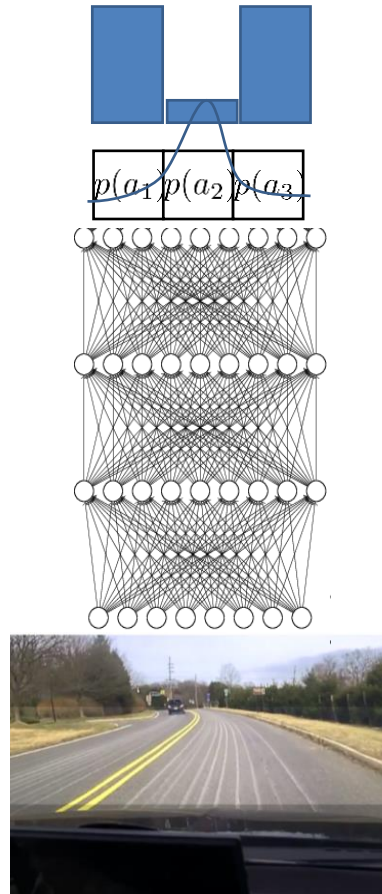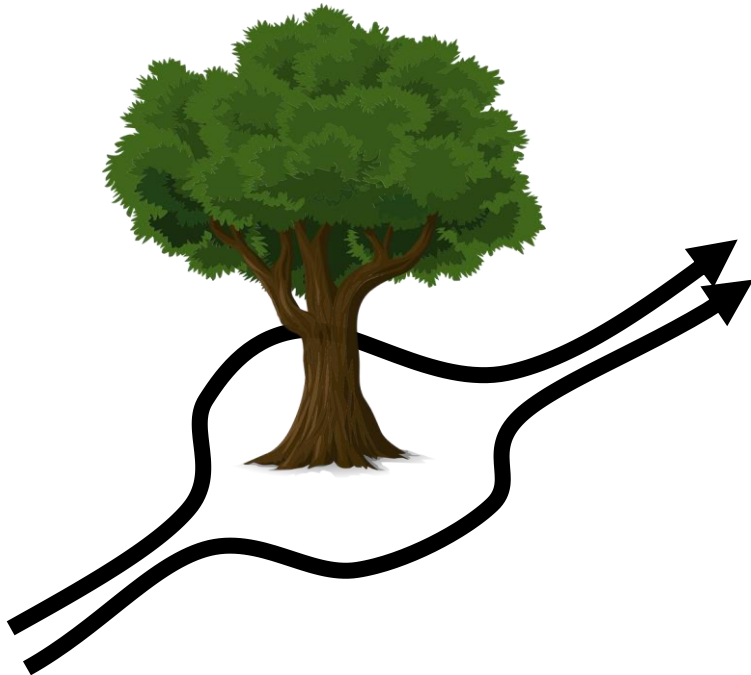# How can we use the whole history?



Typically, LSTM cells work better here

# Why might we fail to fit the expert?

1. Non-Markovian behavior
2. Multimodal behavior
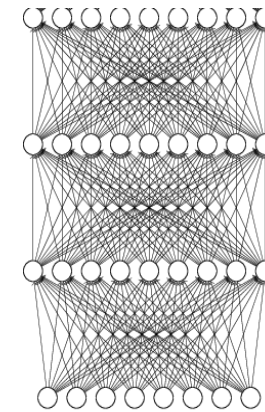
1. Output mixture of Gaussians
2. Latent variable models
3. Autoregressive discretization

# Why might we fail to fit the expert?

1. Output mixture of Gaussians

2. Latent variable models

3. Autoregressive discretization

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

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

# 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

$$\xi \sim \mathcal{N}(0, \mathbf{I})$$
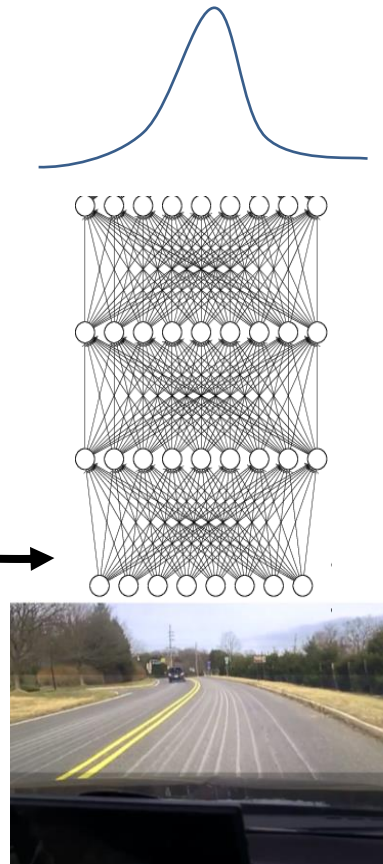
# Why might we fail to fit the expert?

1. Output mixture of Gaussians

2. Latent variable models

3. Autoregressive discretization

We'll learn more about better ways to model multi-modal distributions when we cover generative models later



$p(a_1)\ p(a_2)\ p(a_3)$

discrete sampling → dim 2 value

discrete sampling → dim 1 value

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

$p(a_1)\ p(a_2)\ p(a_3)$

# Does it work?          Yes!

# Why did that work?

# Summary



$\mathbf{o}_t$ → training data → supervised learning → $\pi_\theta(\mathbf{a}_t | \mathbf{o}_t)$

$\mathbf{a}_t$

- In principle it should not work
  - Distribution mismatch problem

- Sometimes works well
  - Hacks (e.g. left/right images)
  - Models with memory (i.e., RNNs)
  - Better distribution modeling
  - Generally taking care to get high accuracy

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

$\mathbf{o}_t$ → → $\mathbf{u}_t$

# A (perhaps) better approach

# Can we make it work more often?



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

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**: **D**ataset **A**ggregation

goal: collect training data from $p_{\pi_\theta}(\mathbf{o}_t)$ instead of $p_{\text{data}}(\mathbf{o}_t)$
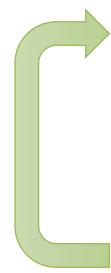
how? just run $\pi_\theta(\mathbf{a}_t|\mathbf{o}_t)$

but need labels $\mathbf{a}_t$!

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

Ross et al. '11

# DAgger Example



Ross et al. '11
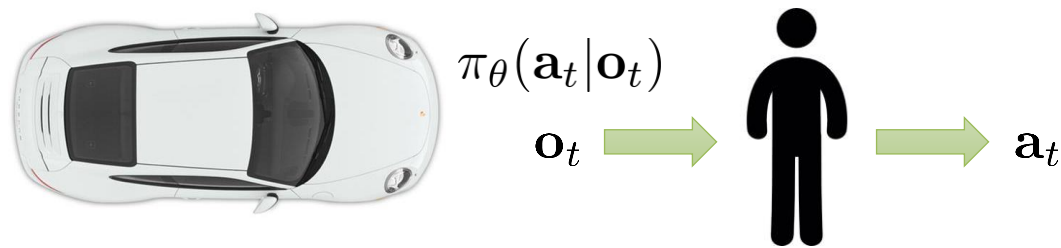
# What's the problem?

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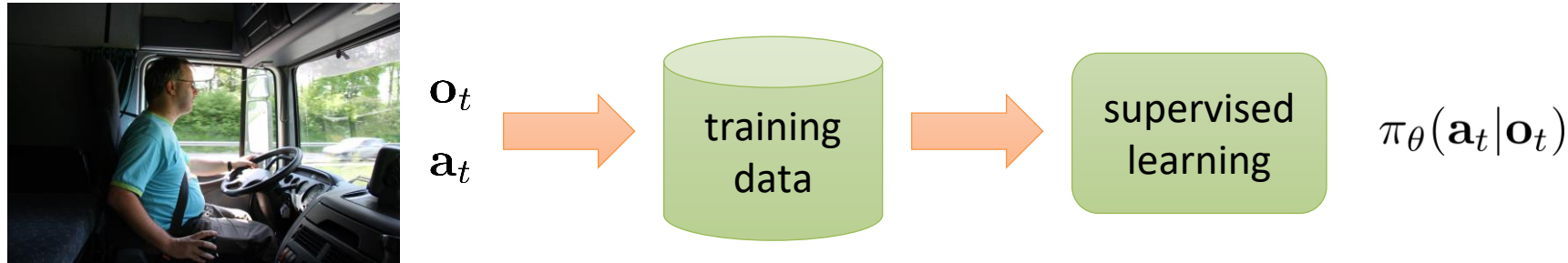4. Aggregate: $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_\pi$

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

$\mathbf{o}_t$ → → $\mathbf{a}_t$

Ross et al. '11

# Summary and takeaways



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

training data → supervised learning $\pi_\theta(\mathbf{a}_t | \mathbf{o}_t)$

- In principle it should not work
  - Distribution mismatch problem
  - **DAgger can address this, but requires costly data collection and labeling**
- Sometimes works well
  - Requires a bit of (heuristic) hacks, and very good (high-accuracy) models

**My recommendation:** try behavioral cloning first, but prepare to be disappointed

# Next time



- i.i.d. distributed data (each datapoint is independent)
- ground truth supervision
- objective is to predict the right label



- each decision can change future inputs (not independent)
- supervision may be high-level (e.g., a goal)
- objective is to accomplish the task

We'll tackle these issues with **reinforcement learning**

## CSW182 (2021)· 课程资料包 @ShowMeAI

**视频**
中英双语字幕

**课件**
一键打包下载

**笔记**
官方笔记翻译

**代码**
作业项目解析

**视频·B 站 [ 扫码或点击链接 ]**
https://www.bilibili.com/video/BV1Ff4y1n7ar

**课件 & 代码·博客 [ 扫码或点击链接 ]**
http://blog.showmeai.tech/berkeley-csw182

**Berkeley**
循环神经网络    可视化    梯度策略
Q-Learning    风格迁移    模仿学习    元学习
计算机视觉    机器学习基础    生成模型    卷积网络

Awesome AI Courses Notes Cheatsheets 是 **ShowMeAI** 资料库的分支系列，覆盖最具知名度的 **TOP50+** 门 AI 课程，旨在为读者和学习者提供一整套高品质中文学习笔记和速查表。

**点击**课程名称，跳转至课程**资料包**页面，**一键下载**课程全部资料！

| 机器学习 | 深度学习 | 自然语言处理 | 计算机视觉 |
|---|---|---|---|
| Stanford · CS229 | Stanford · CS230 | Stanford · CS224n | Stanford · CS231n |

**# Awesome AI Courses Notes Cheatsheets· 持续更新中**

| 知识图谱 | 图机器学习 | 深度强化学习 | 自动驾驶 |
|---|---|---|---|
| Stanford · CS520 | Stanford · CS224W | UCBerkeley · CS285 | MIT · 6.S094 |

**微信公众号**

资料下载方式 2：扫码点击底部菜单栏

称为 **AI 内容创作者？** 回复 [ 添砖加瓦 ]