# Imitation Learning

# Introduction

- Imitation Learning
  - Also known as learning by demonstration, apprenticeship learning
- An expert demonstrates how to solve the task
  - Machine can also interact with the environment, but cannot explicitly obtain reward.
  - It is hard to define reward in some tasks.
  - Hand-crafted rewards can lead to uncontrolled behavior
- Two approaches:
  - Behavior Cloning
  - Inverse Reinforcement Learning (inverse optimal control)

# Behavior Cloning

# Behavior Cloning

- Self-driving cars as example

observation



Expert (Human driver): 向前

Machine: 向前

機器學習expert的行為

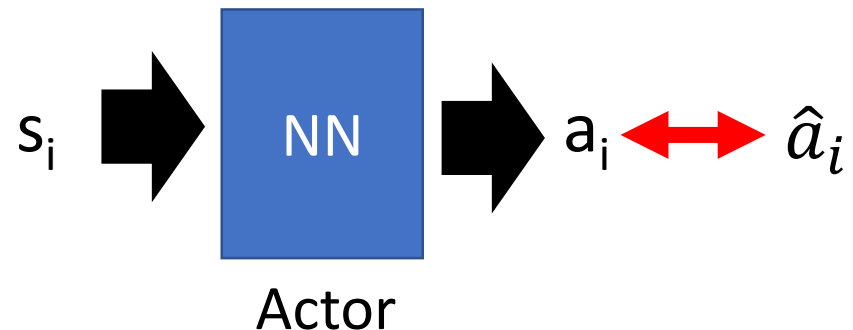Training data:

$(s_1, \hat{a}_1)$
$(s_2, \hat{a}_2)$
$(s_3, \hat{a}_3)$
......

$s_i$ → **NN** → $a_i$ ⟷ $\hat{a}_i$
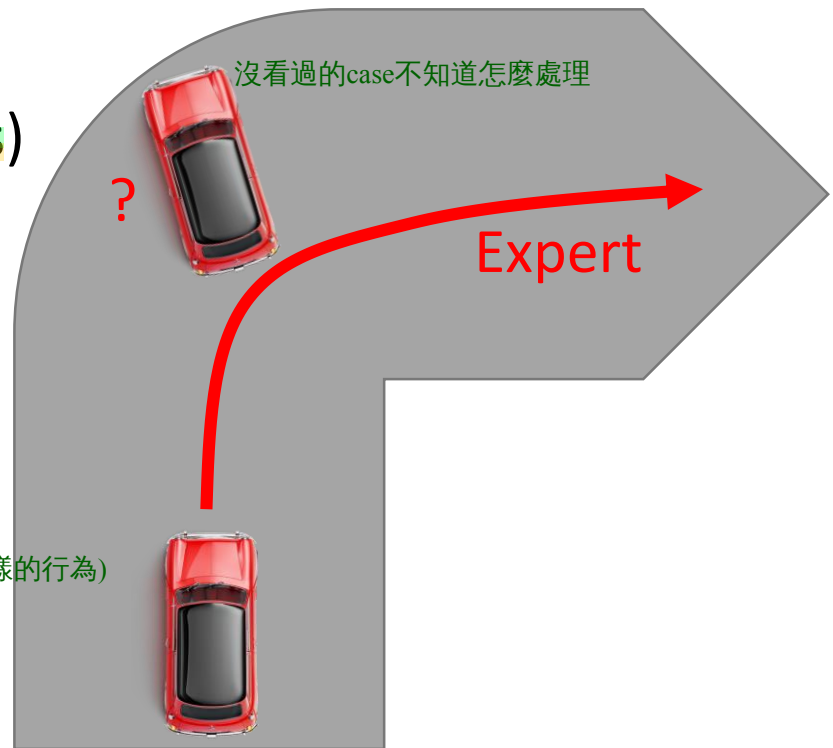
Actor

# Behavior Cloning

- Problem

Expert only samples limited observation (states)

看過的observation是有限的

Let the expert in the states seem by machine

需要蒐集更多樣性的data(在各種極端的情況下會做出怎樣的行為)

Dataset Aggregation

沒看過的case不知道怎麼處理
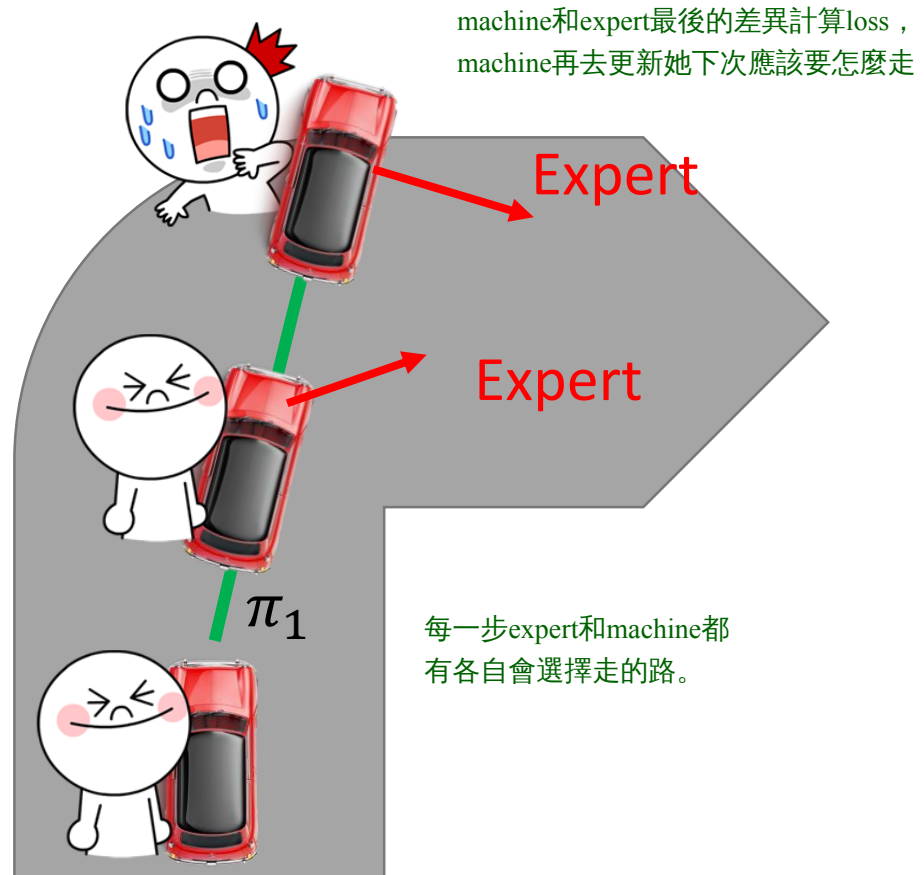
?

Expert

# Behavior Cloning

-

Get actor $\pi_1$ by behavior cloning

Using $\pi_1$ to interact with the environment

Ask the expert to label the observation of $\pi_1$

Using new data to train $\pi_2$

machine和expert最後的差異計算loss，machine再去更新她下次應該要怎麼走

Expert

Expert

$\pi_1$

每一步expert和machine都有各自會選擇走的路。

# Behavior Cloning

The agent will copy every behavior, even irrelevant actions.



https://www.youtube.com/watch?v=j2FSB3bseek

# Behavior Cloning

缺點: capacity有限: 因此甚麼該學甚麼不該學是該分重要程度的

- Major problem: if machine has limited capacity, it may choose the wrong behavior to copy.

$s_i$ ➡️ NN ➡️ $a_i$ ↗ speech ✖️ gesture

Actor 👍 machine自己去learn哪個部分比較重要~

$s_i$ ➡️ NN ➡️ $a_i$ ✖️ speech ↘ gesture

Actor 👎

- Some behavior must copy, but some can be ignored.
  - Supervised learning takes all errors equally

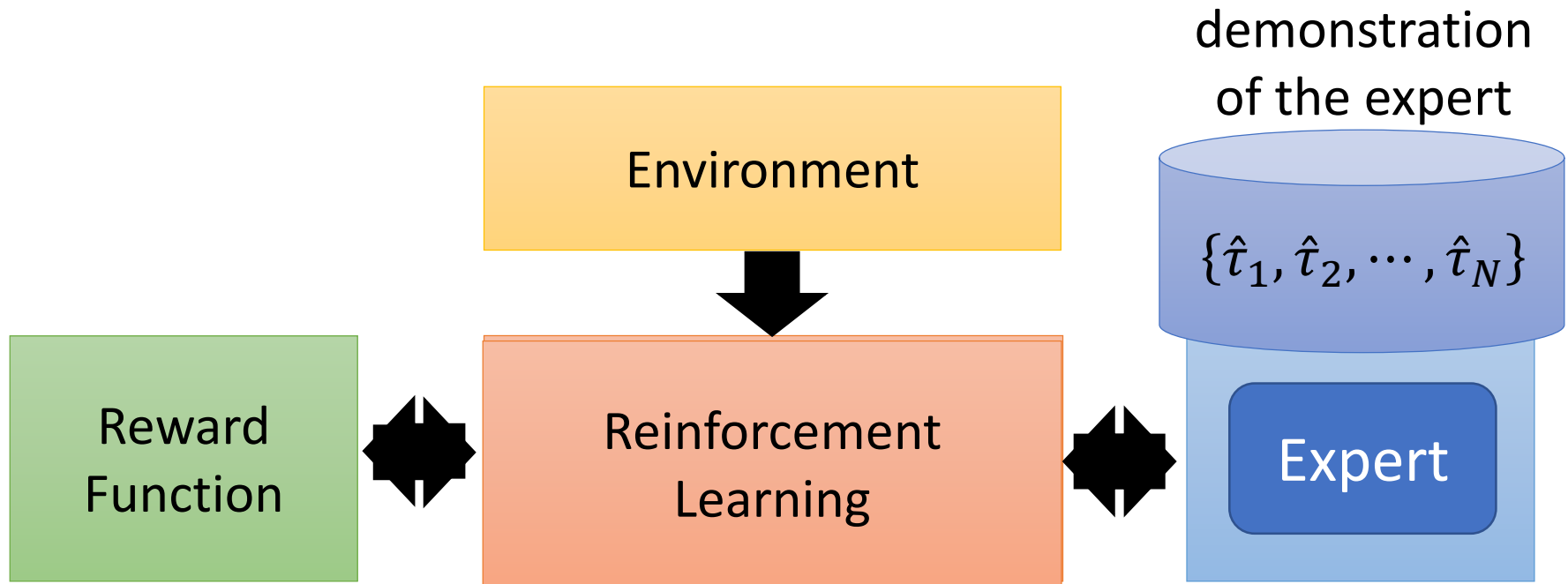沒有區分那些重要那些不重要

# Mismatch

$s_i$ ➡️ Actor ➡️ $a_i$

- In supervised learning, we expect training and testing data have the same distribution.

- In behavior cloning:
  - Training: $(s, a) \sim \hat{\pi}$ (expert)
    - ***Action a taken by actor influences the distribution of s***
  - Testing: $(s', a') \sim \pi^*$ (actor cloning expert)
    - If $\hat{\pi} = \pi^*$, $(s, a)$ and $(s', a')$ from the same distribution
    - If $\hat{\pi}$ and $\pi^*$ have difference, the distribution of $s$ and s$'$ can be very different. 失之毫釐差之千里

# Inverse Reinforcement Learning (IRL)
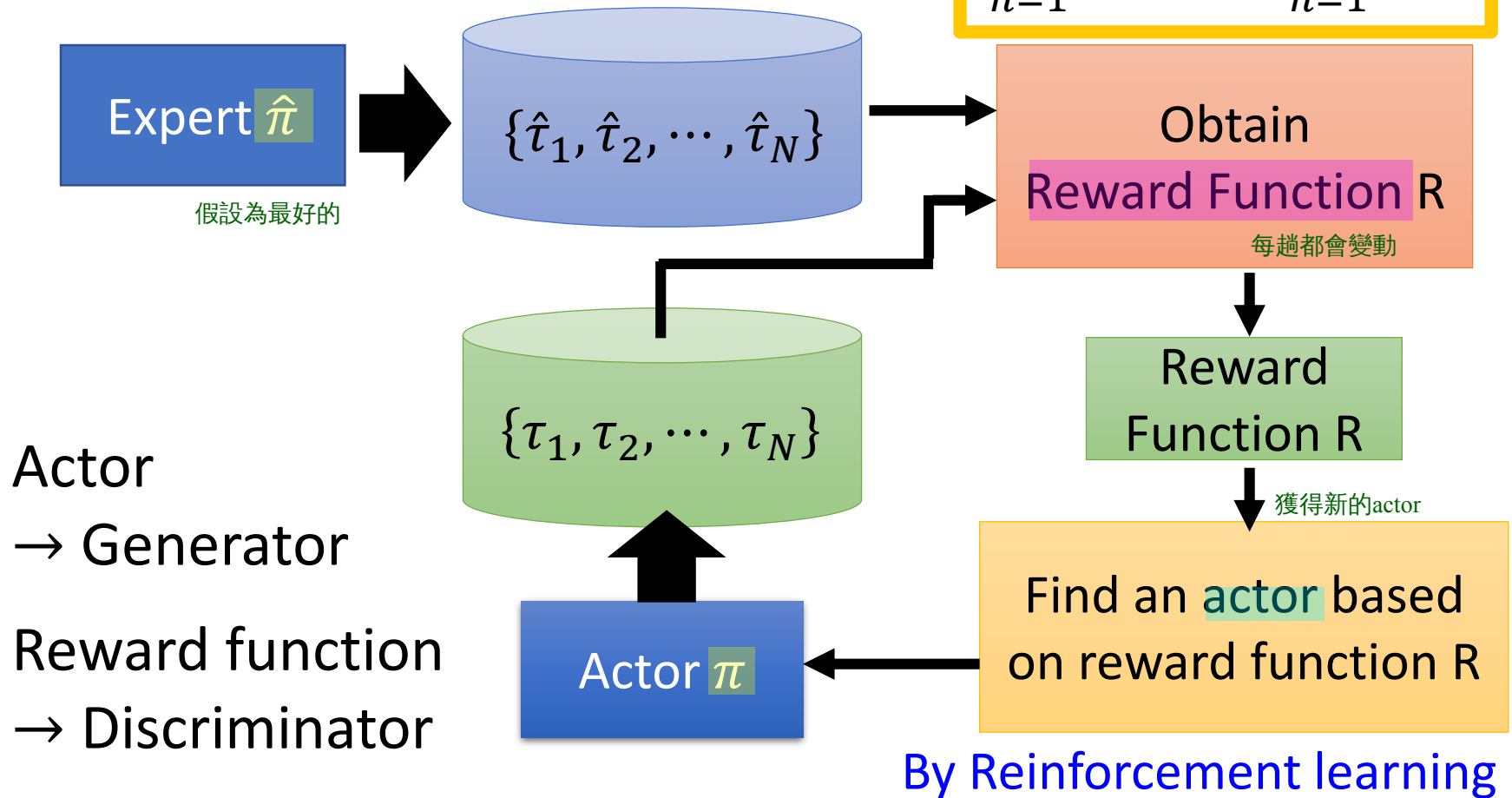
# Inverse Reinforcement Learning

Environment

demonstration
of the expert

$$\{\hat{\tau}_1, \hat{\tau}_2, \cdots, \hat{\tau}_N\}$$

Reward
Function

Reinforcement
Learning

Expert

也許reward function很簡單，卻可以呈現很複雜的行為

➢ Using the reward function to find the ***optimal actor***.

➢ Modeling reward can be easier. Simple reward function can lead to complex policy.
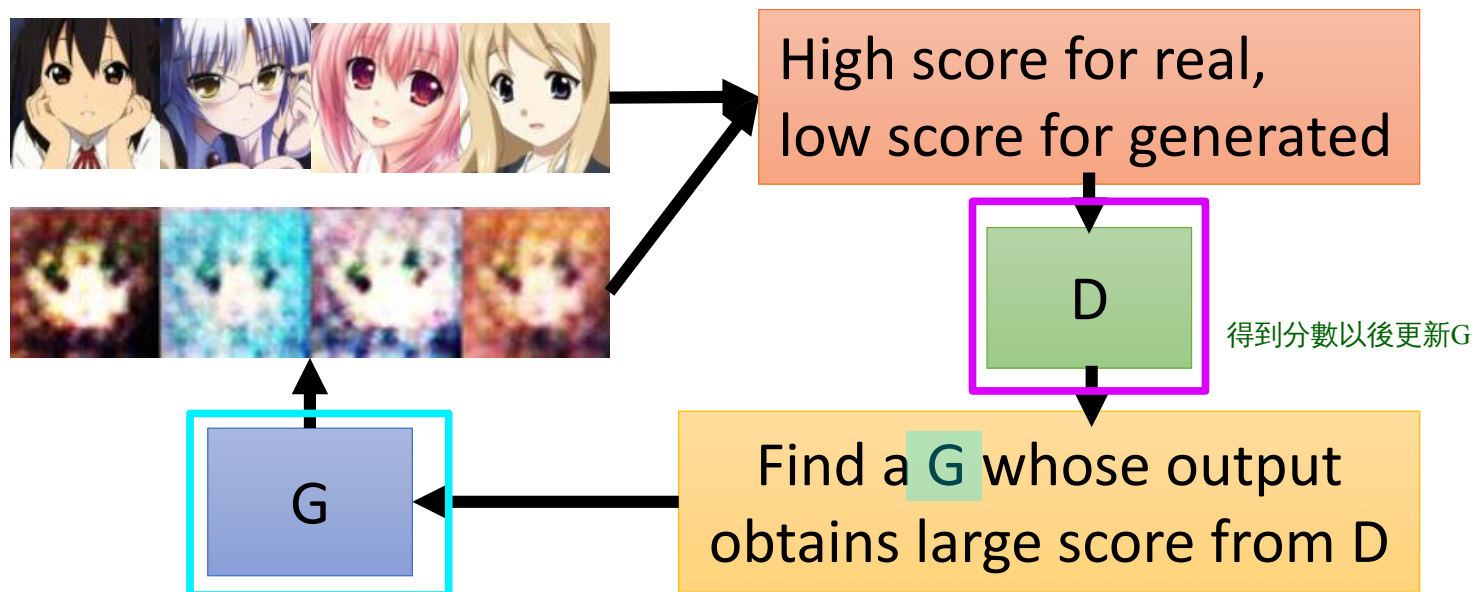
# Framework of IRL

$$\sum_{n=1}^{N} R(\hat{\tau}_n) > \sum_{n=1}^{N} R(\tau)$$

Expert $\hat{\pi}$

假設為最好的

$\{\hat{\tau}_1, \hat{\tau}_2, \cdots, \hat{\tau}_N\}$

Obtain
Reward Function R

每趟都會變動

Reward
Function R

獲得新的actor

$\{\tau_1, \tau_2, \cdots, \tau_N\}$

Actor $\pi$

Find an actor based
on reward function R

Actor
→ Generator

Reward function
→ Discriminator

By Reinforcement learning

# GAN

有！正確答案



High score for real,
low score for generated

D

得到分數以後更新G

Find a G whose output obtains large score from D

G

# IRL

Expert

可找出基於expert的最佳解
因為expert沒有看過所有state
因此沒有所謂的"完全正確答案"

$\{\hat{\tau}_1, \hat{\tau}_2, \cdots, \hat{\tau}_N\}$

$\{\tau_1, \tau_2, \cdots, \tau_N\}$

Larger reward for $\hat{\tau}_n$,
Lower reward for $\tau$

Reward
Function

Find a Actor obtains large reward

Actor

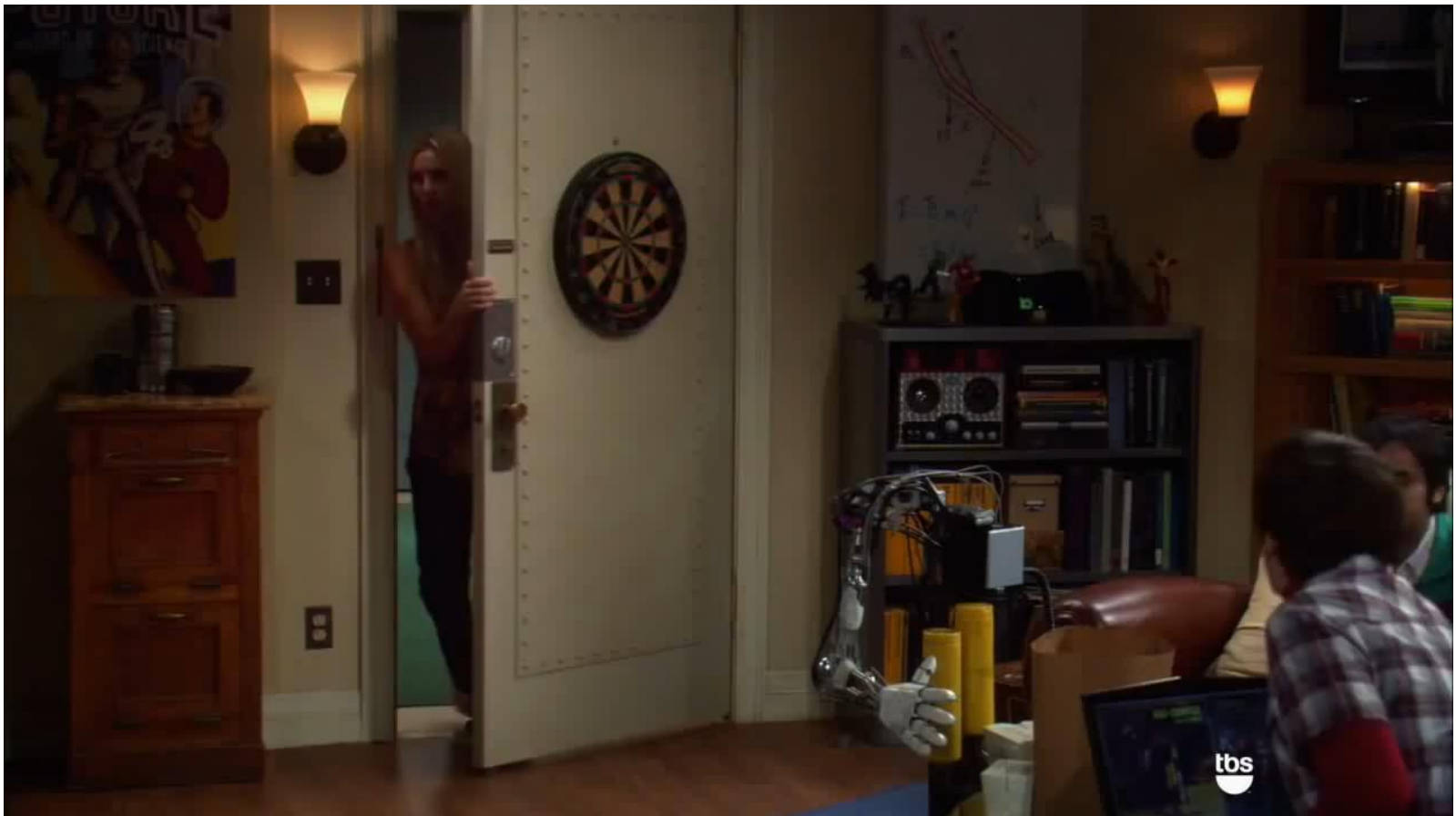# Parking Lot Navigation



- Reward function:
  - Forward vs. reverse driving
  - Amount of switching between forward and reverse
  - Lane keeping
  - On-road vs. off-road
  - Curvature of paths

# Robot

- How to teach robots?   https://www.youtube.com/watch?v=DEGbtjTOIB0

# Robot

Guided Cost Learning:
Deep Inverse Optimal Control via Policy Optimization

Chelsea Finn, Sergey Levine, Pieter Abbeel
UC Berkeley

# Third Person Imitation Learning

- Ref: Bradly C. Stadie, Pieter Abbeel, Ilya Sutskever, "Third-Person Imitation Learning", arXiv preprint, 2017

### First Person



http://lasa.epfl.ch/research_new/ML/index.php
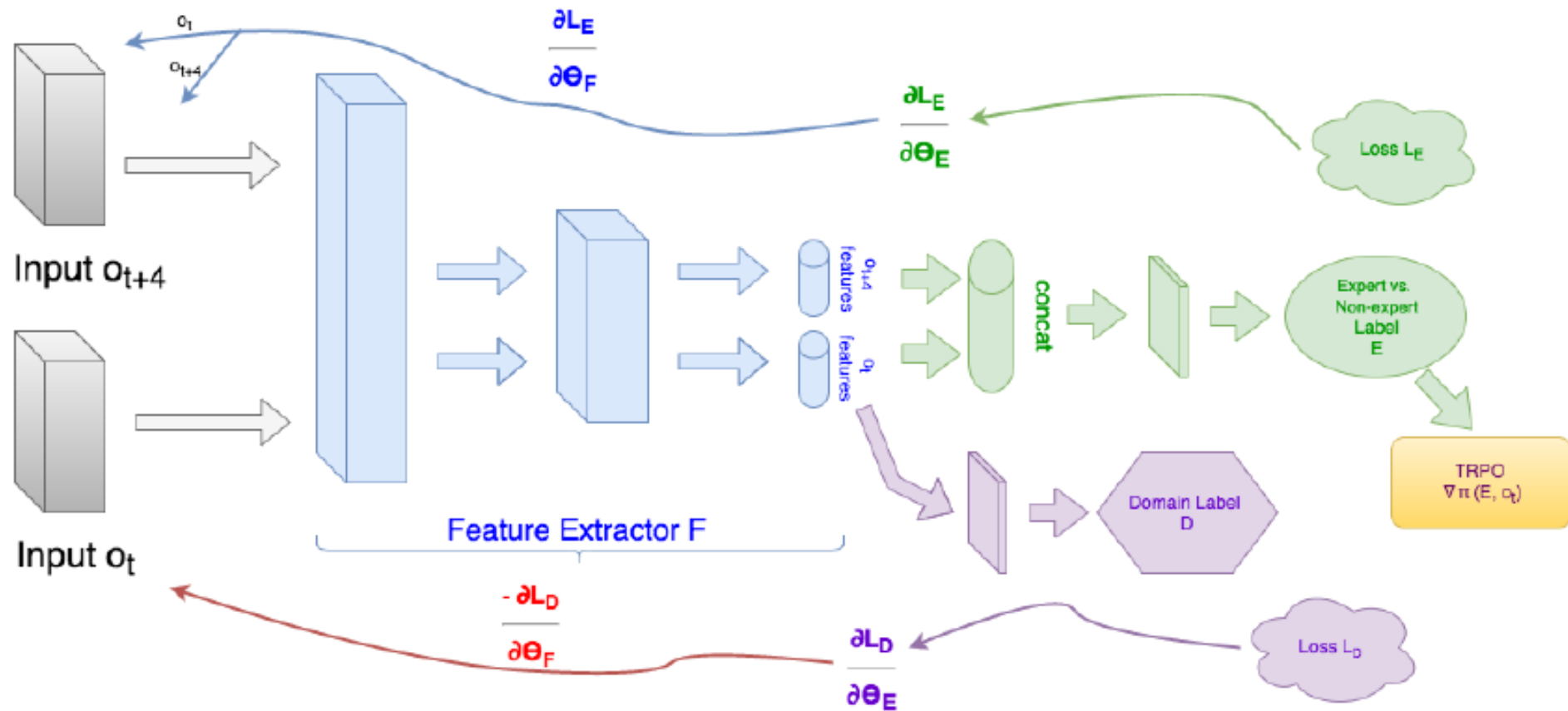
### Third Person



https://kknews.cc/sports/q5kbb8.html

http://sc.chinaz.com/Files/pic/icons/1913/%E6%9C%BA%E5%99%A8%E4%BA%BA%E5%9B%BE%E6%A0%87%E4%B8%8B%E8%BD%BD34.png

# Third Person Imitation Learning

# Recap: Sentence Generation & Chat-bot

### *Sentence Generation*

Expert trajectory:
床 前 明 月 光

$(s_1, a_1)$:  ("<BOS>",")床")

$(s_2, a_2)$:  ("床",")前")

$(s_3, a_3)$:  ("床前",")明")

⋮    ⋮

### *Chat-bot*

Expert trajectory:
input: how are you
Output: I am fine

$(s_1, a_1)$:  ("input, <BOS>",")I")

$(s_2, a_2)$:  ("input, I", "am")

$(s_3, a_3)$:  ("input, I am", "fine")

⋮    ⋮

Maximum likelihood is behavior cloning. Now we have better approach like SeqGAN.