

# imitation learning

1. Learning from reinforcement alone is hard

1. Exploration is hard

2. Credit assignment is hard

3. Designing rewards is hard

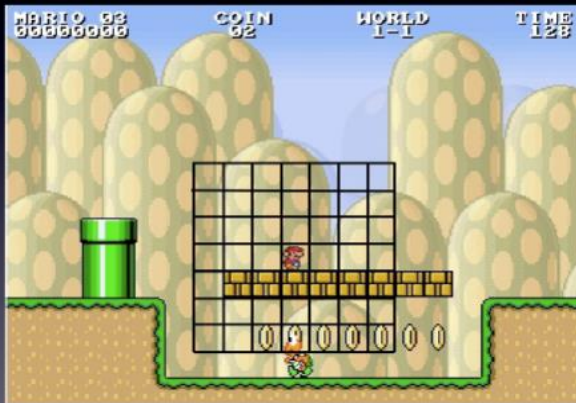
2. Yet people are pretty good at many tasks

3. Perhaps we can use them to help

# From Mario AI competition 2009

Input:

Interface



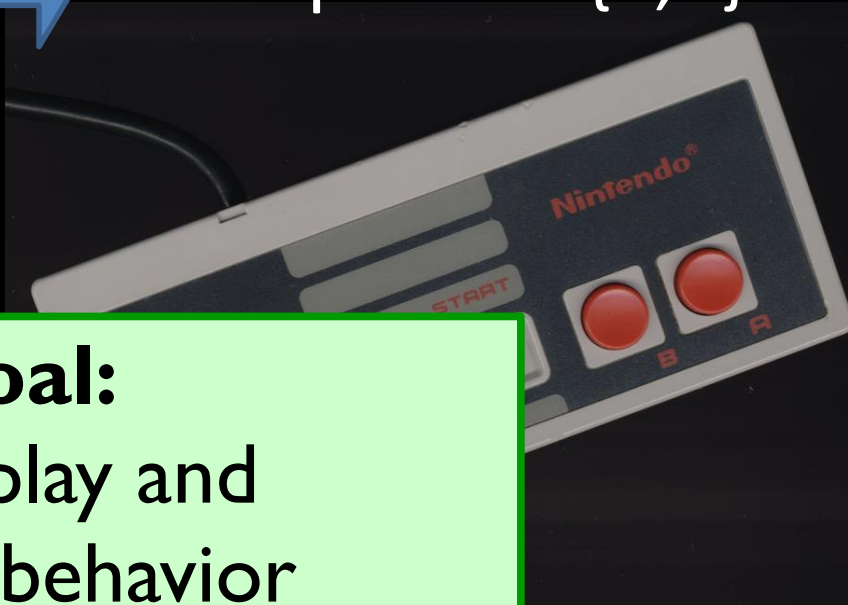
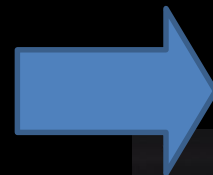
Output:

Jump in  $\{0,1\}$

Right in  $\{0,1\}$

Left in  $\{0,1\}$

Speed in  $\{0,1\}$



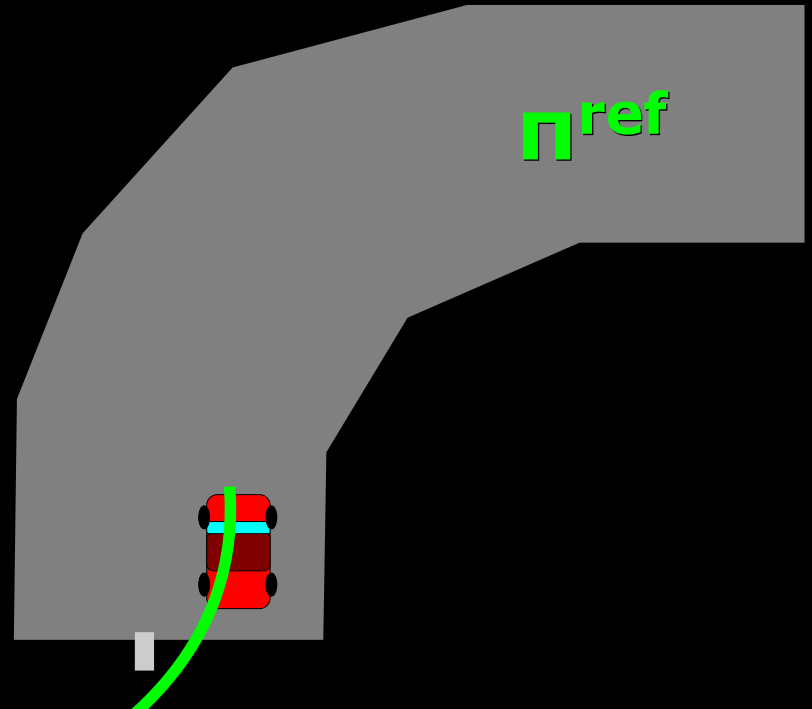
**High level goal:**

Watch an expert play and learn to mimic her behavior

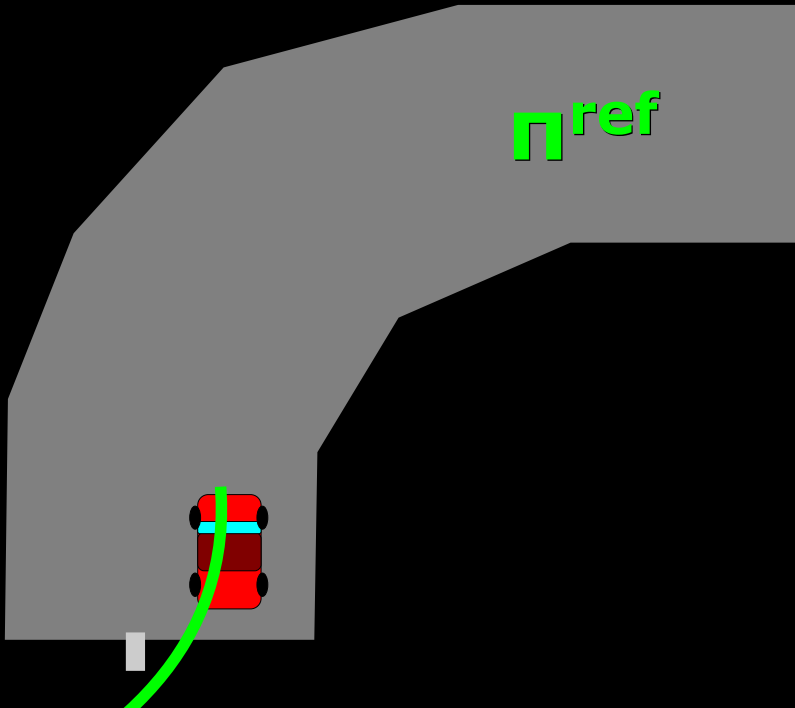
Video credit: Stéphane Ross, Geoff Gordon and Drew Bagnell



1. Collect trajectories from expert  $\pi^{\text{ref}}$
  2. Store dataset  $\mathbf{D} = \{ (o, \pi^{\text{ref}}(o)) \mid o \sim \pi^{\text{ref}} \}$
  3. Train classifier  $\pi$  on  $\mathbf{D}$
- Let  $\pi$  play the game!

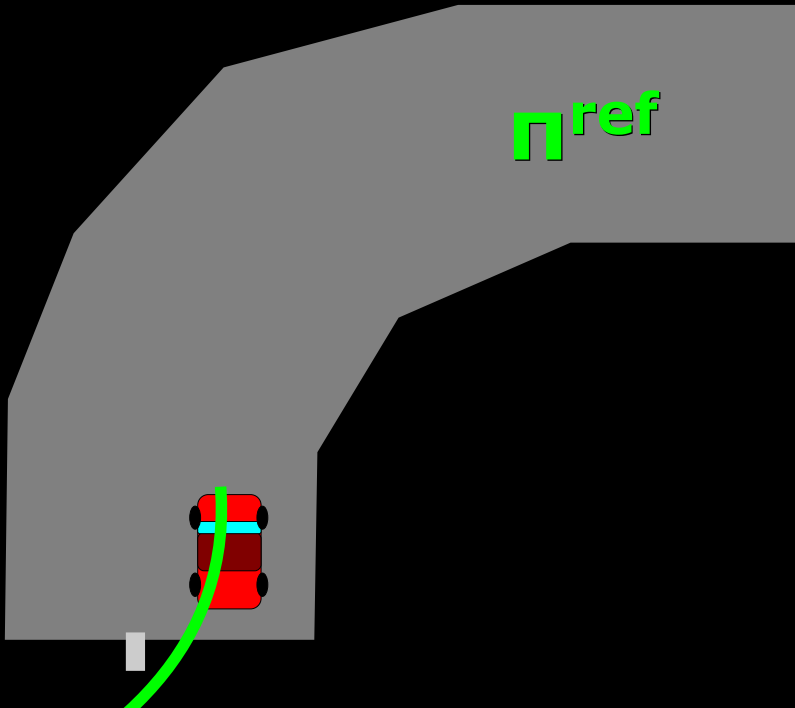


Video credit: Stéphane Ross, Geoff Gordon and Drew Bagnell







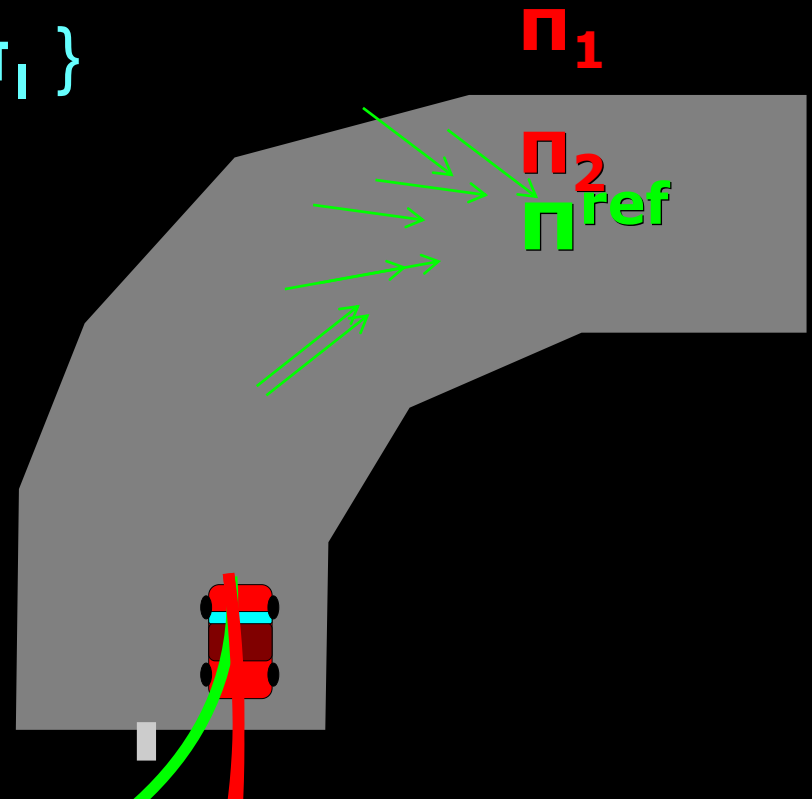


1. Collect trajectories from expert  $\pi^{\text{ref}}$
2. Dataset  $\mathbf{D}_0 = \{ (o, \pi^{\text{ref}}(o, y)) \mid o \sim \pi^{\text{ref}} \}$
3. Train  $\pi_1$  on  $\mathbf{D}_0$
4. Collect new trajectories from  $\pi_1$ 
  - › But let the *expert* steer!
5. Dataset  $\mathbf{D}_1 = \{ (o, \pi^{\text{ref}}(o, y)) \mid o \sim \pi_1 \}$
6. Train  $\pi_2$  on  $\mathbf{D}_0 \cup \mathbf{D}_1$

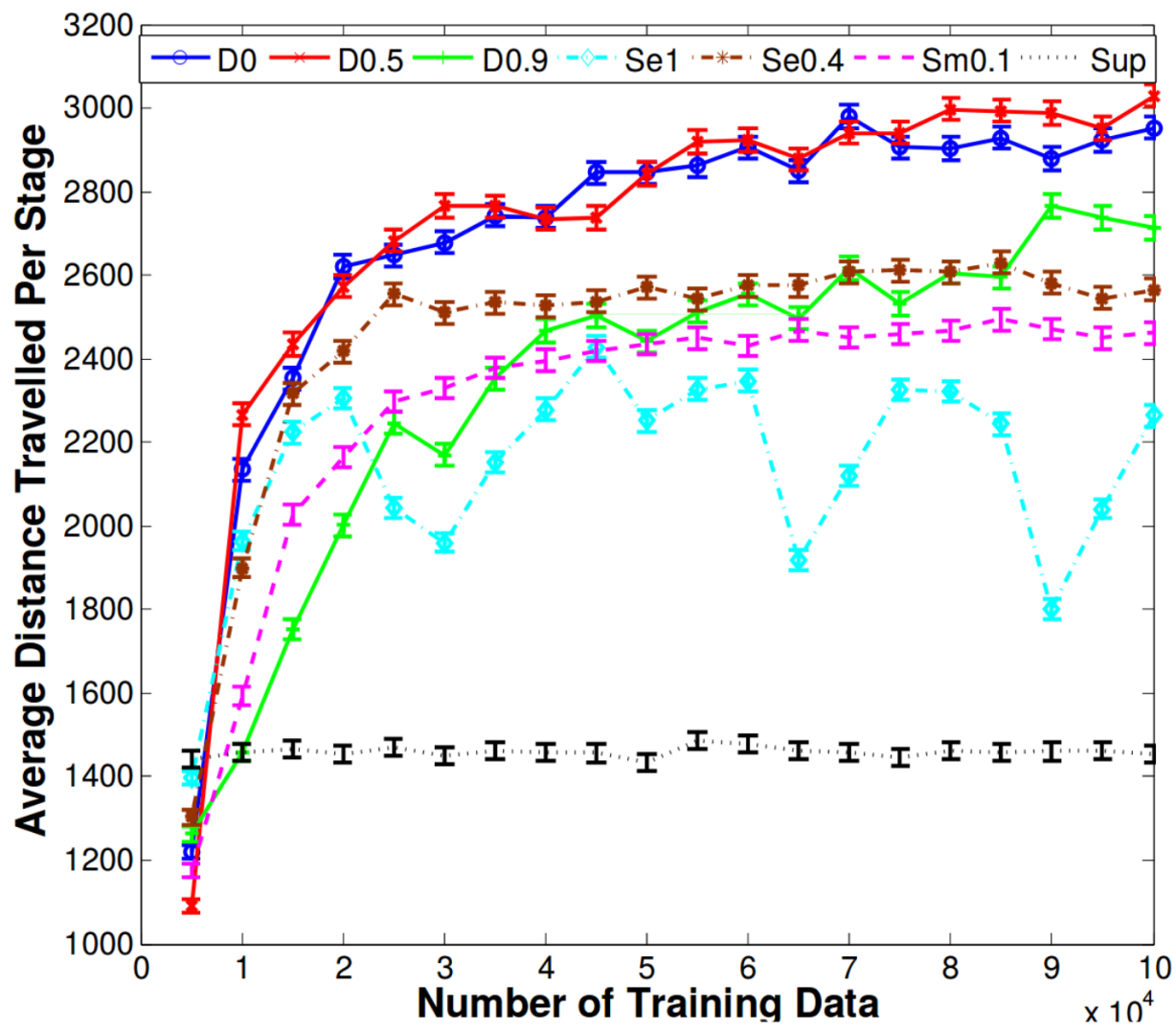
• In general:

- $\mathbf{D}_n = \{ (o, \pi^{\text{ref}}(o, y)) \mid o \sim \pi_n \}$
- Train  $\pi_{n+1}$  on  $\bigcup_{i \leq n} \mathbf{D}_i$

If  $N = T \log T$ ,  
 $\mathbf{L}(\pi_n) < T \epsilon_N + \mathbf{O}(1)$   
 for some  $n$

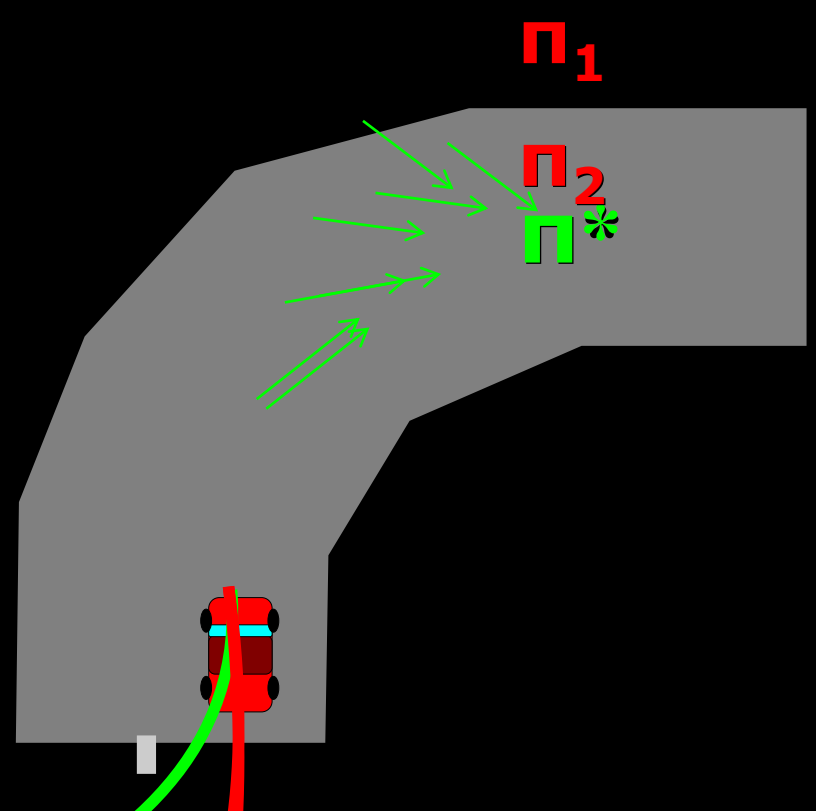


Video credit: Stéphane Ross, Geoff Gordon and Drew Bagnell



Ross+Bagnell, AISTats'10

- 
- 
- 



Classifier:  $h : x \rightarrow [K]$

- $(x, y) \in X \times [K]$
- $\min_h \Pr( h(x) \neq y )$

- $(x, c) \in X \times [0, \infty)^K$
- $\min_h E_{(x, c)} [ c_{h(x)} ]$

Classifier:  $h : x \rightarrow [K]$

- $(x, c) \in X \times [0, \infty)^K$
- $\min_h E_{(x, c)} [ c_{h(x)} ]$

Solution learn a K-dimensional regressor on costs;  
pick minimal cost



1. Let learned policy  $\pi$  drive for  $t$  timesteps to obs.  $o$

2. For each possible action  $a$ :

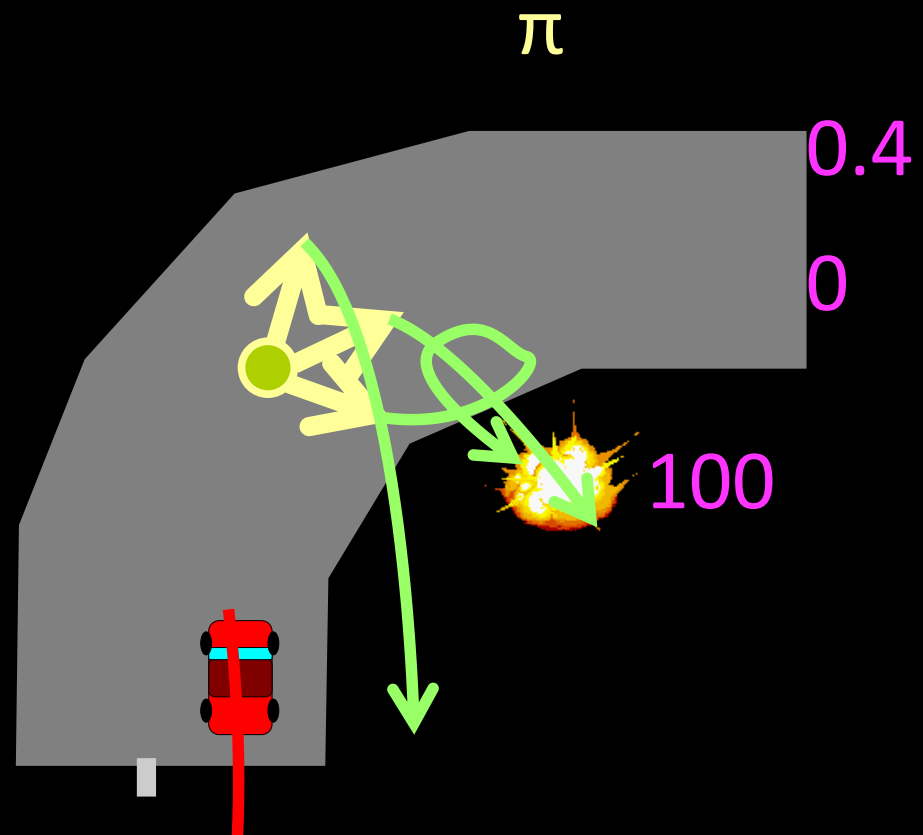
• Take action  $a$ , and let expert  $\pi^{\text{ref}}$  drive the rest

• Record the overall loss,  $c_a$

3. Update  $\pi$  based on example:

$$(o, \langle c_1, c_2, \dots, c_k \rangle)$$

4. Goto (1)



.From demonstrations  $\rightarrow$  expert decisions

.From expert decisions  $\rightarrow$  expert costs

observation

optimal(ish) action



$a_1, a_2, \dots, a_{t-1}$

*(expected) minimum achievable loss*

$$\min_{(a_t, a_{t+1}, \dots)} \mathbf{E} \text{ loss}(\vec{a})$$

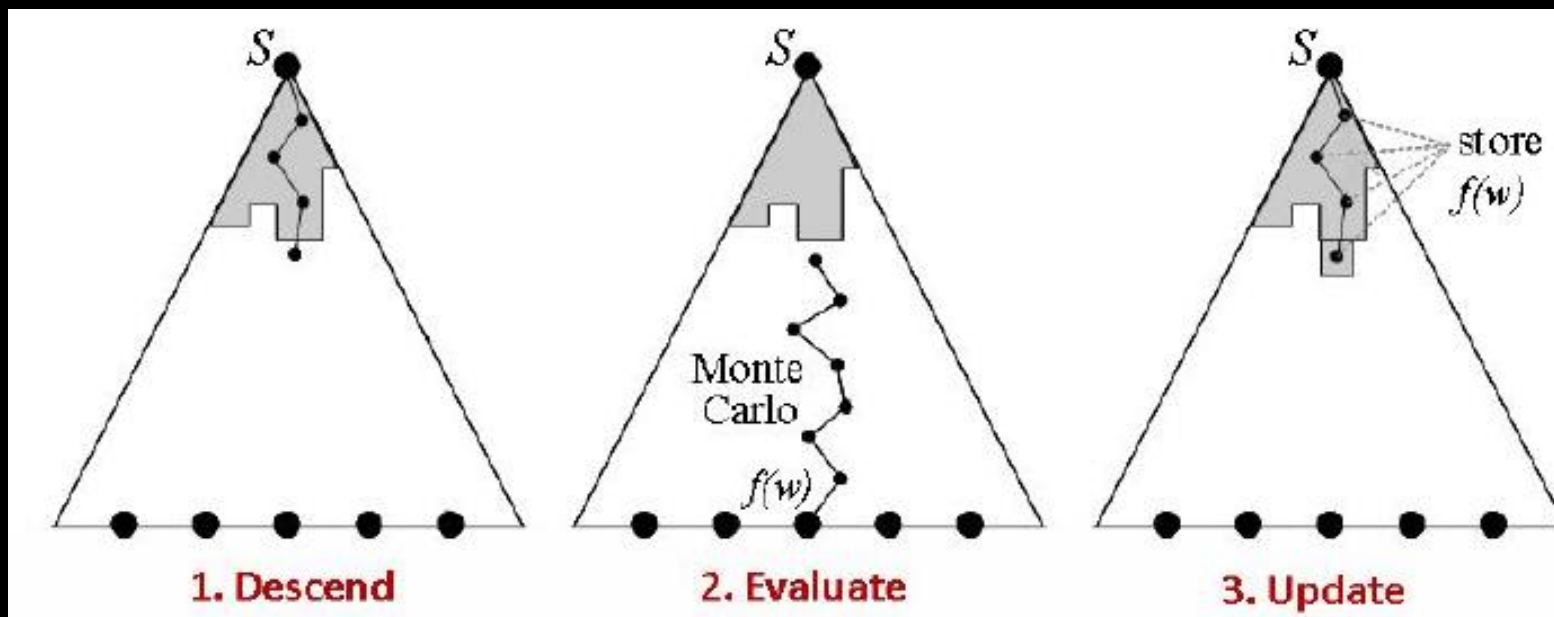
*optimal action*

$a_t$

*optimal Q values*

$a_t$





e.g., Monte Carlo  
Tree Search

Image credit: Michele Sebag  
and DeepMind



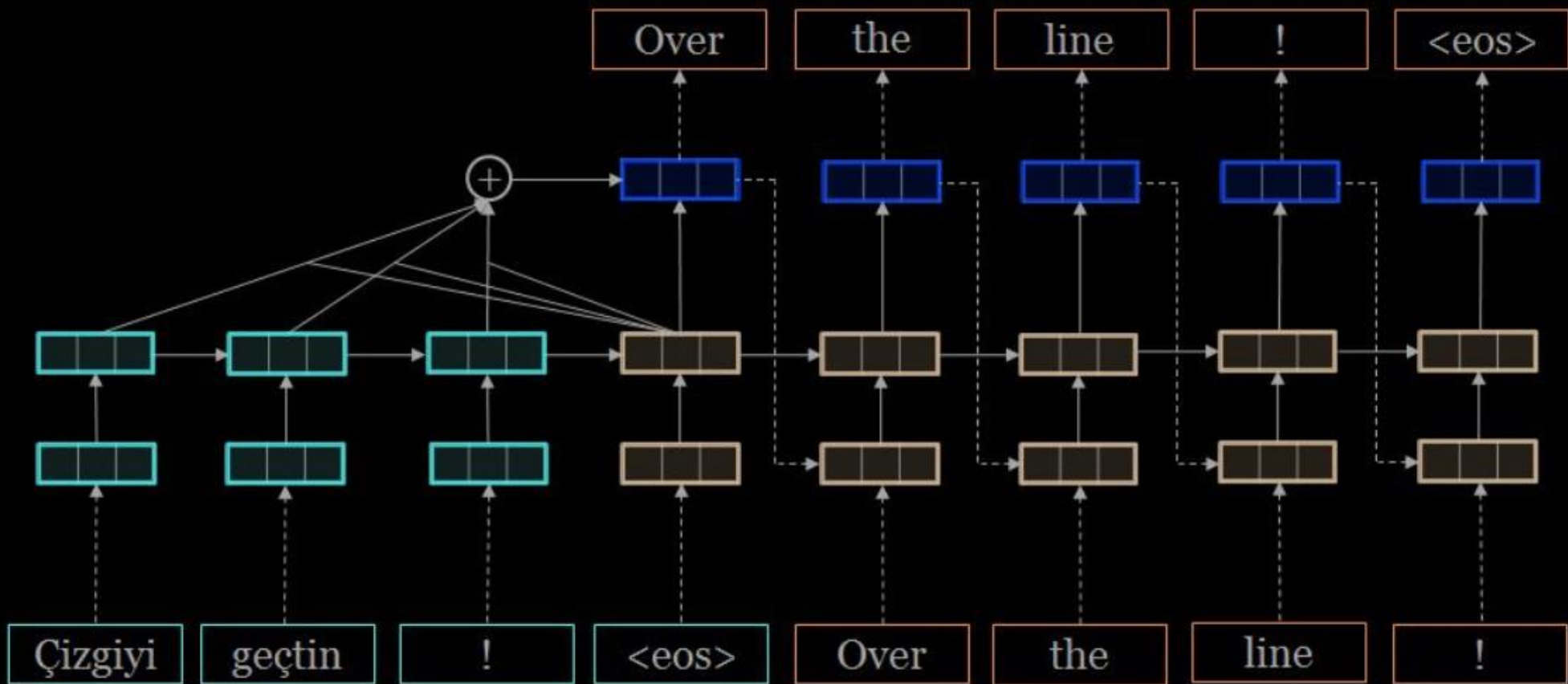
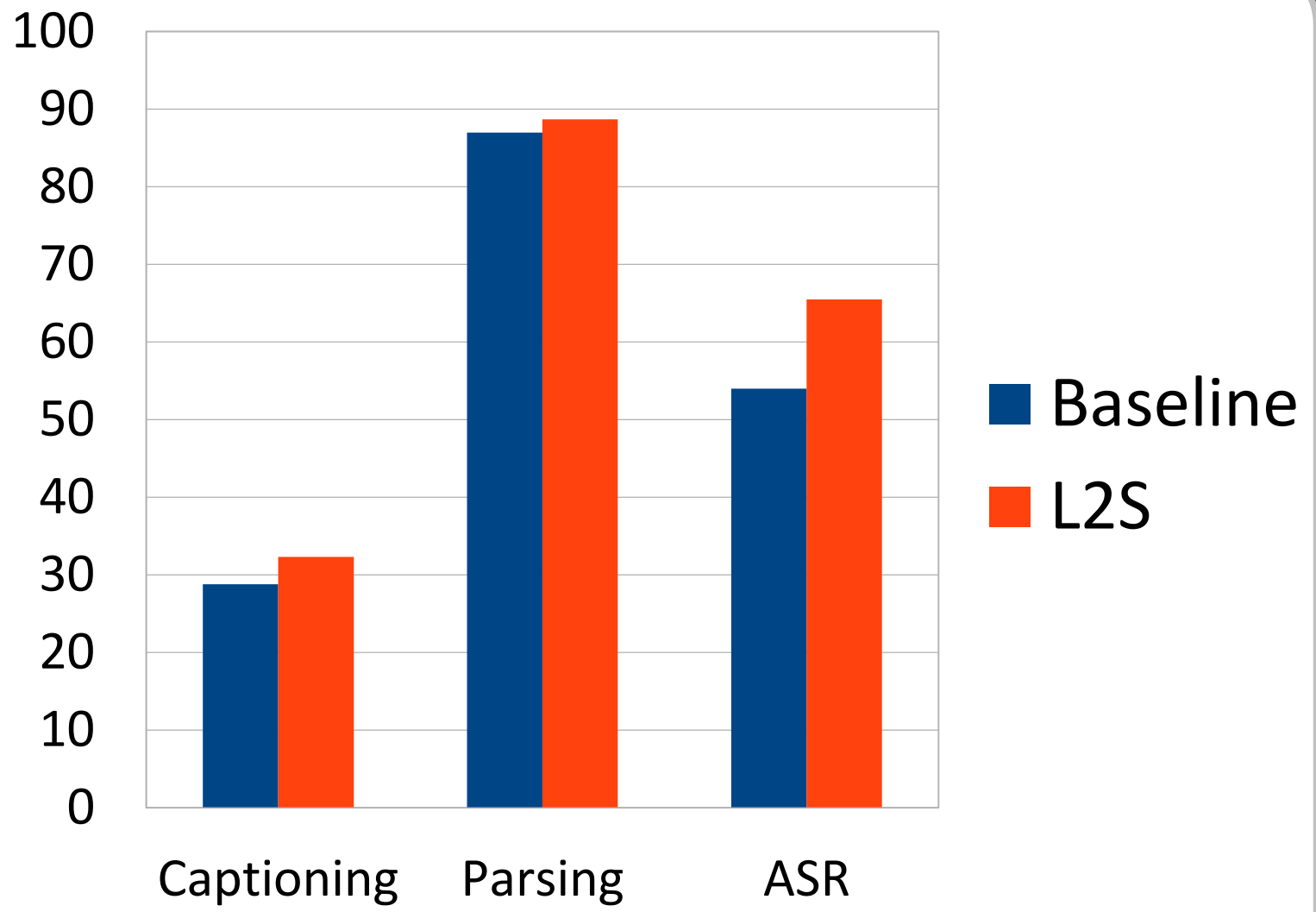
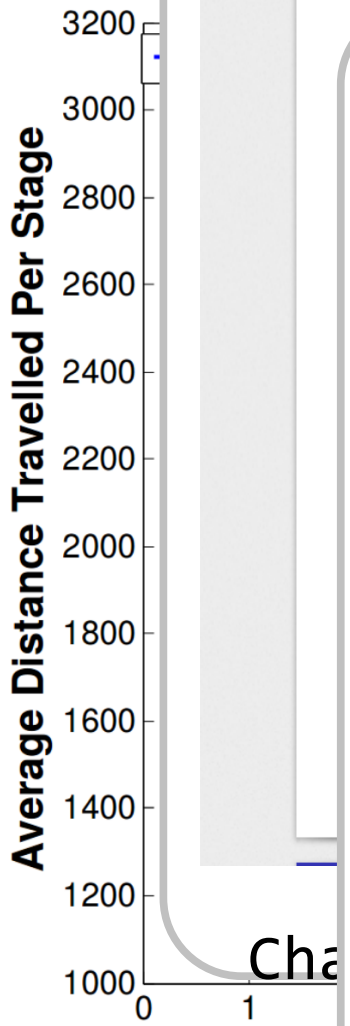
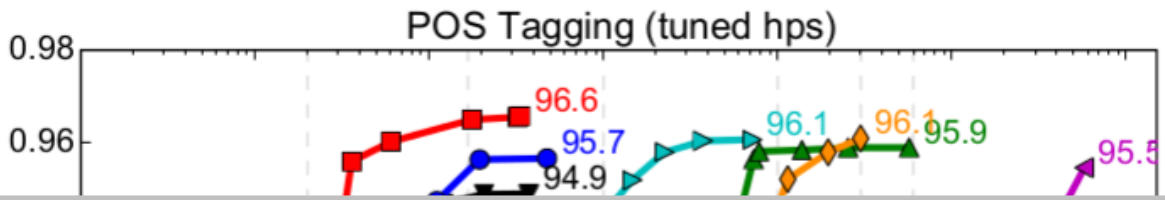


Image credit: Klein et al., 2017





Ross+Bagnoli

Bengio+Vinyals+Navdeep+Shazeer, NIPS'16

.From demonstrations  $\rightarrow$  expert decisions

.From expert decisions  $\rightarrow$  expert costs

.Whence the expert?



1. Let learned policy  $\pi^{\text{in}}$  drive for  $t$  timesteps to obs.  $o$

2. For each possible action  $a$ :

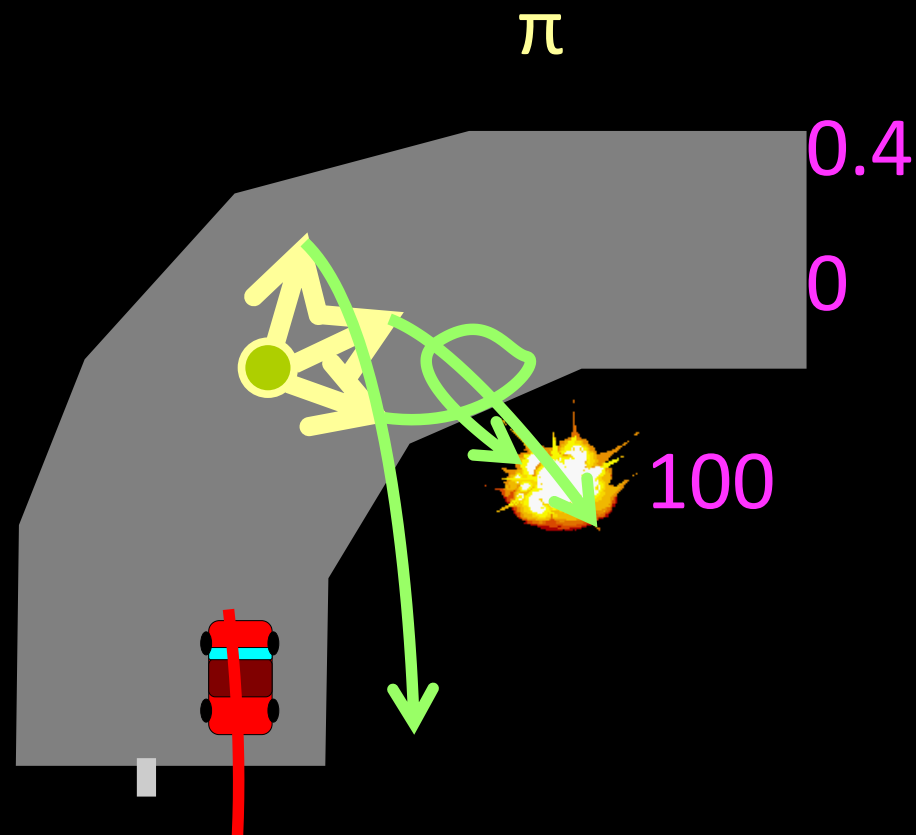
• Take action  $a$ , and let something  $\pi^{\text{out}}$  drive the rest

• Record the overall loss,  $c_a$

3. Update  $\pi$  based on example:

$(o, \langle c_1, c_2, \dots, c_k \rangle)$

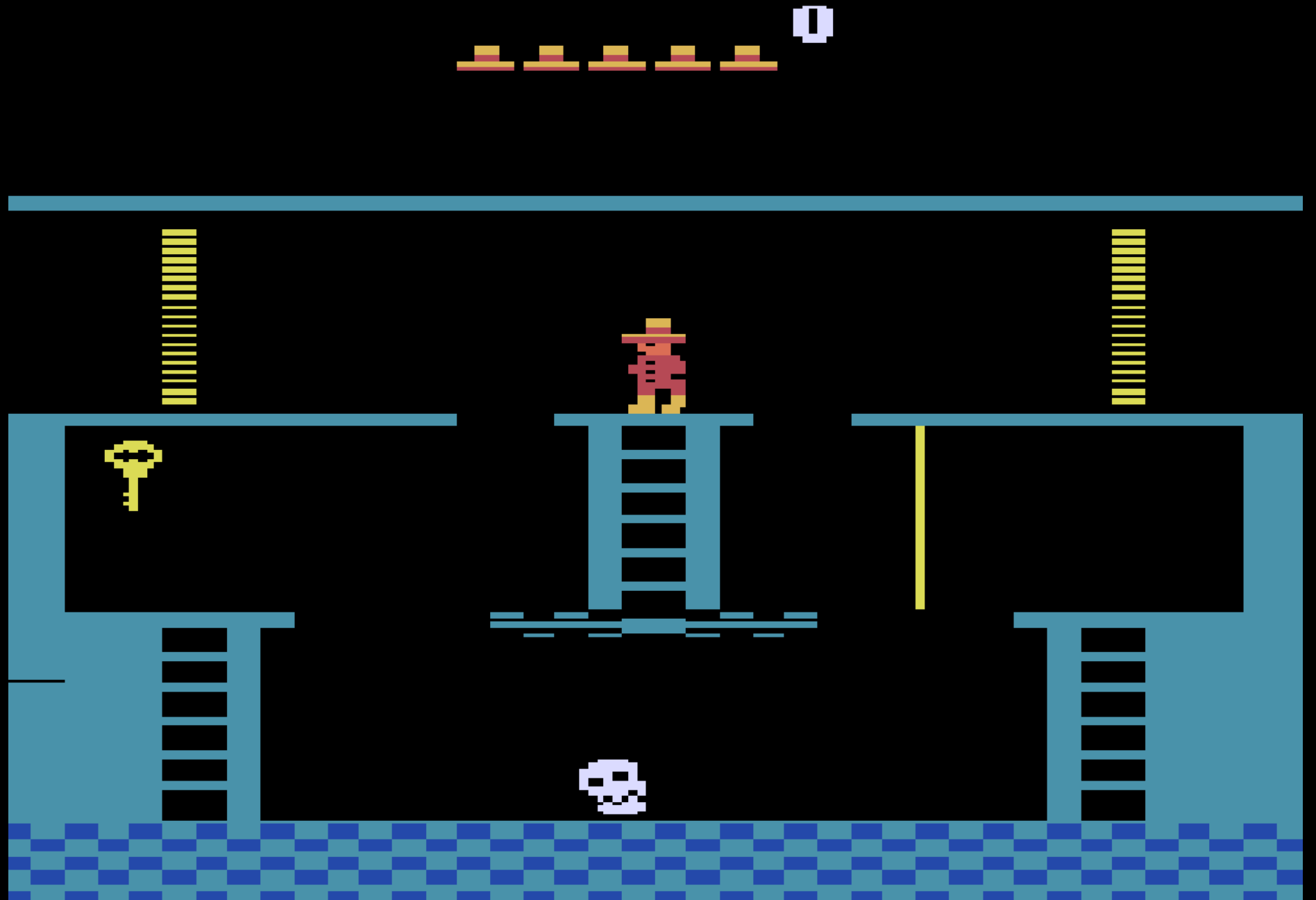
4. Goto (1)

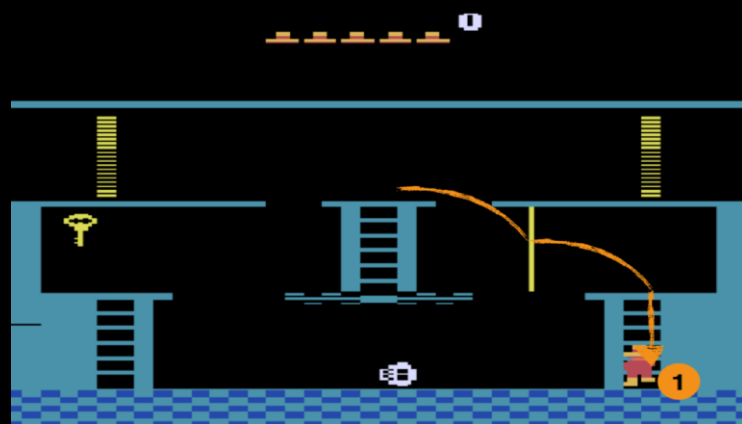


roll-out →	<b>Reference</b>	<b>Mixture</b>	<b>Learned</b>
↓ roll-in			
<b>Reference</b>	Inconsistent		
<b>Learned</b>	Not locally opt.	Good	RL

roll-out →	<b>Reference</b>	<b>Mixture</b>	<b>Learned</b>
↓ roll-in			
<b>Reference</b>	Inconsistent		
<b>Learned</b>	Not locally opt.	Good	RL

$$Regret = O \left( (KT)^{2/3} \sqrt[3]{\frac{\log(N|\Pi|)}{N}} + T\delta_{class} \right)$$

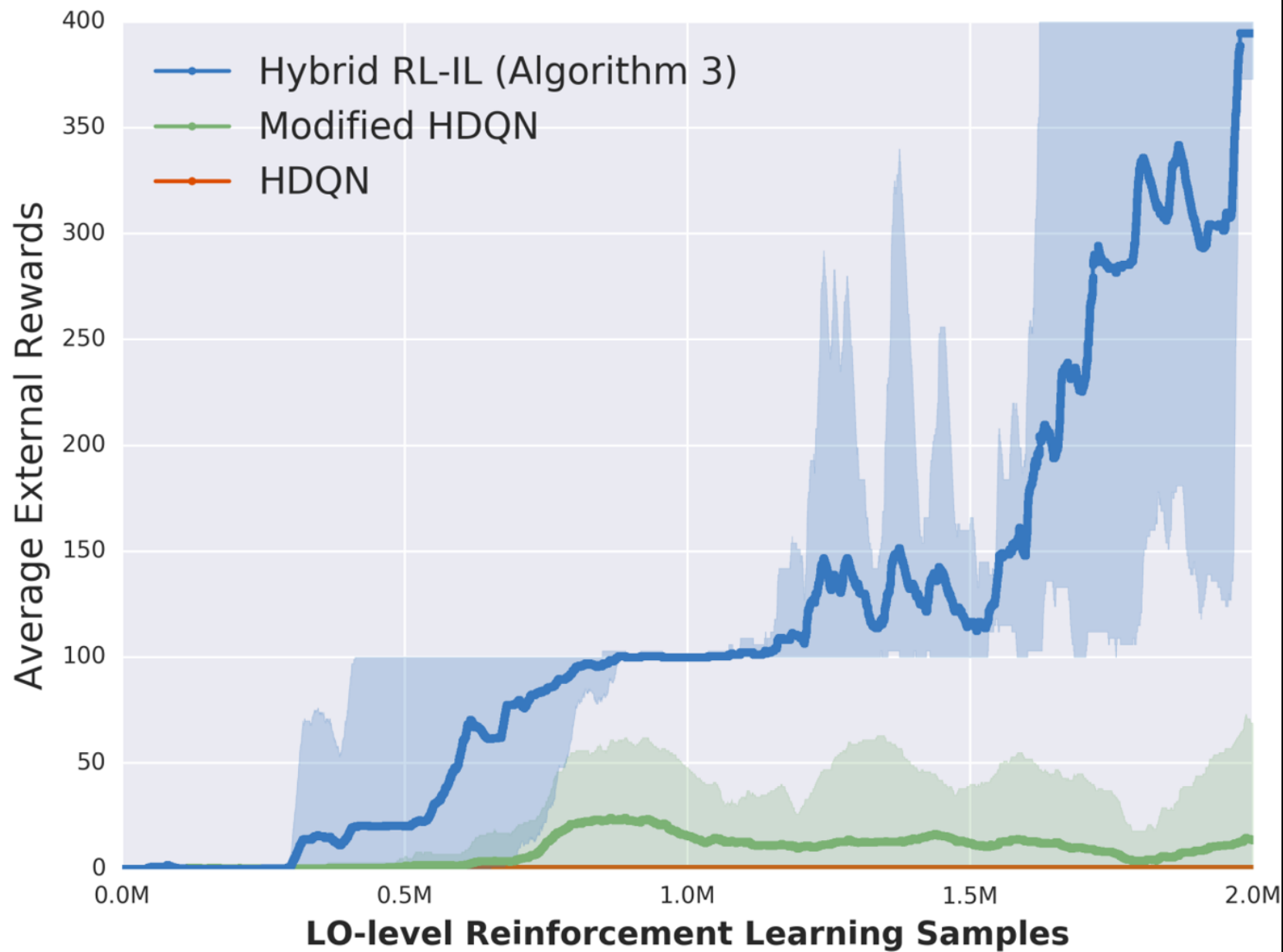




*Key insight:* **verifying if low-level trajectory is successful is cheaper than labeling low-level trajectory**

- **labeling effort** = high-level horizon + low-level horizon  
only a fraction of the full horizon (as low as sqrt of the full horizon)
- **subpolicies** are only learnt in the **relevant part** of the **state space**







# imitation learning summary

successes:

open problems:

Thank you! Queries!