

Reinforcement Learning:

A Primer, Multi-Task, Goal-Conditioned

CS 330

Introduction

Some background:

- Not a native English speaker so **please please** let me know if you don't understand something
- I like robots 😊
- Studied classical robotics first
- Got fascinated by deep RL in the middle of my PhD after a talk by Sergey Levine
- Research Scientist at Robotics @ Google

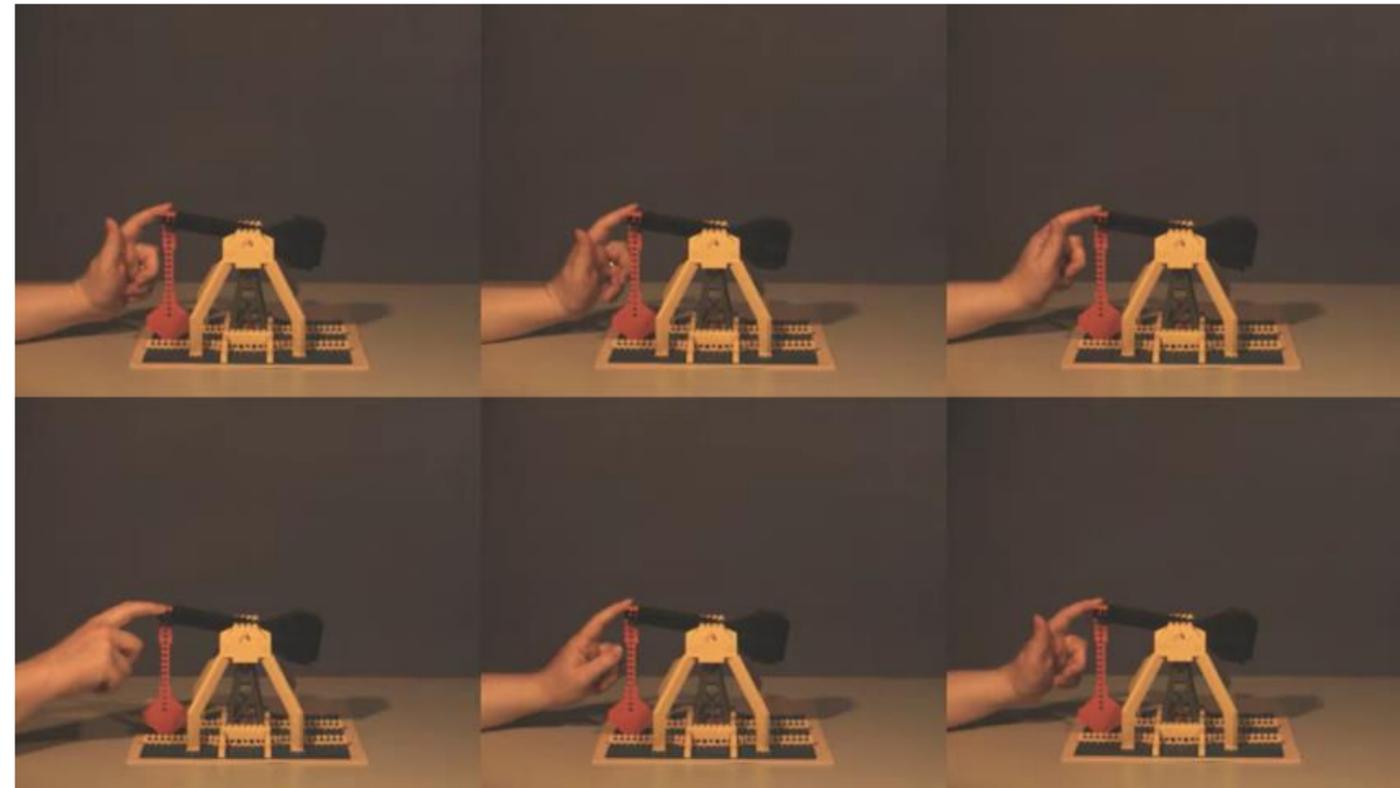


Karol
Hausman



Why Reinforcement Learning?

Isolated action that doesn't affect the future?



Why Reinforcement Learning?

Isolated action that doesn't affect the future?

Supervised learning?

Common applications



robotics



language & dialog



autonomous driving



business operations



finance

(most deployed ML systems)
+ a key aspect of intelligence

The Plan

Multi-task reinforcement learning problem

Policy gradients & their multi-task counterparts

Q-learning

<— should be review

Multi-task Q-learning

The Plan

Multi-task reinforcement learning problem

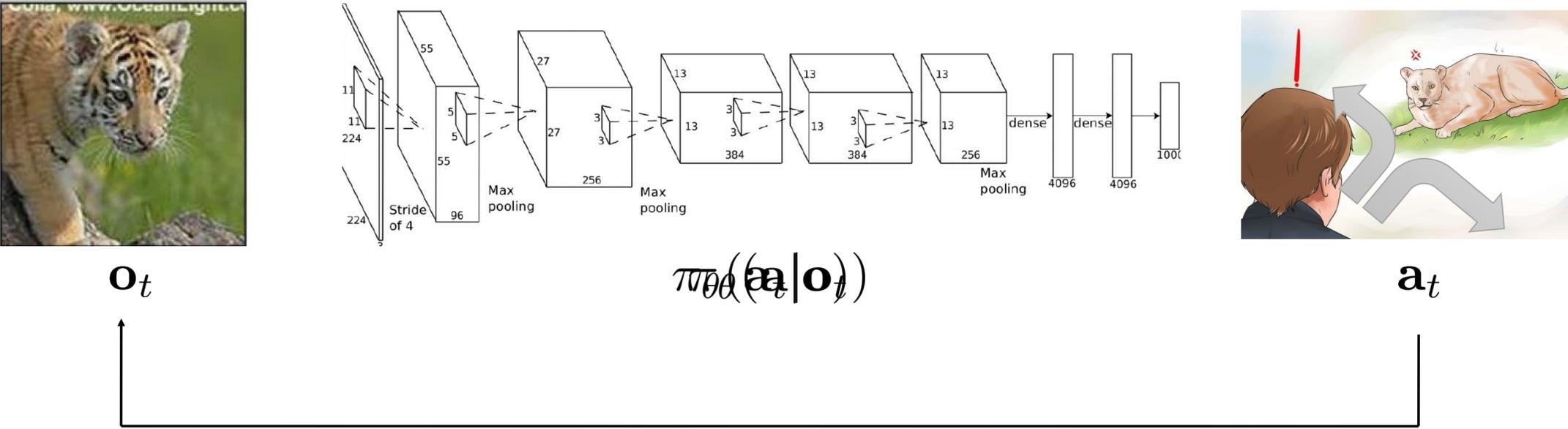
Policy gradients & their multi-task counterparts

Q-learning

<— should be review

Multi-task Q-learning

Terminology & notation



\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)

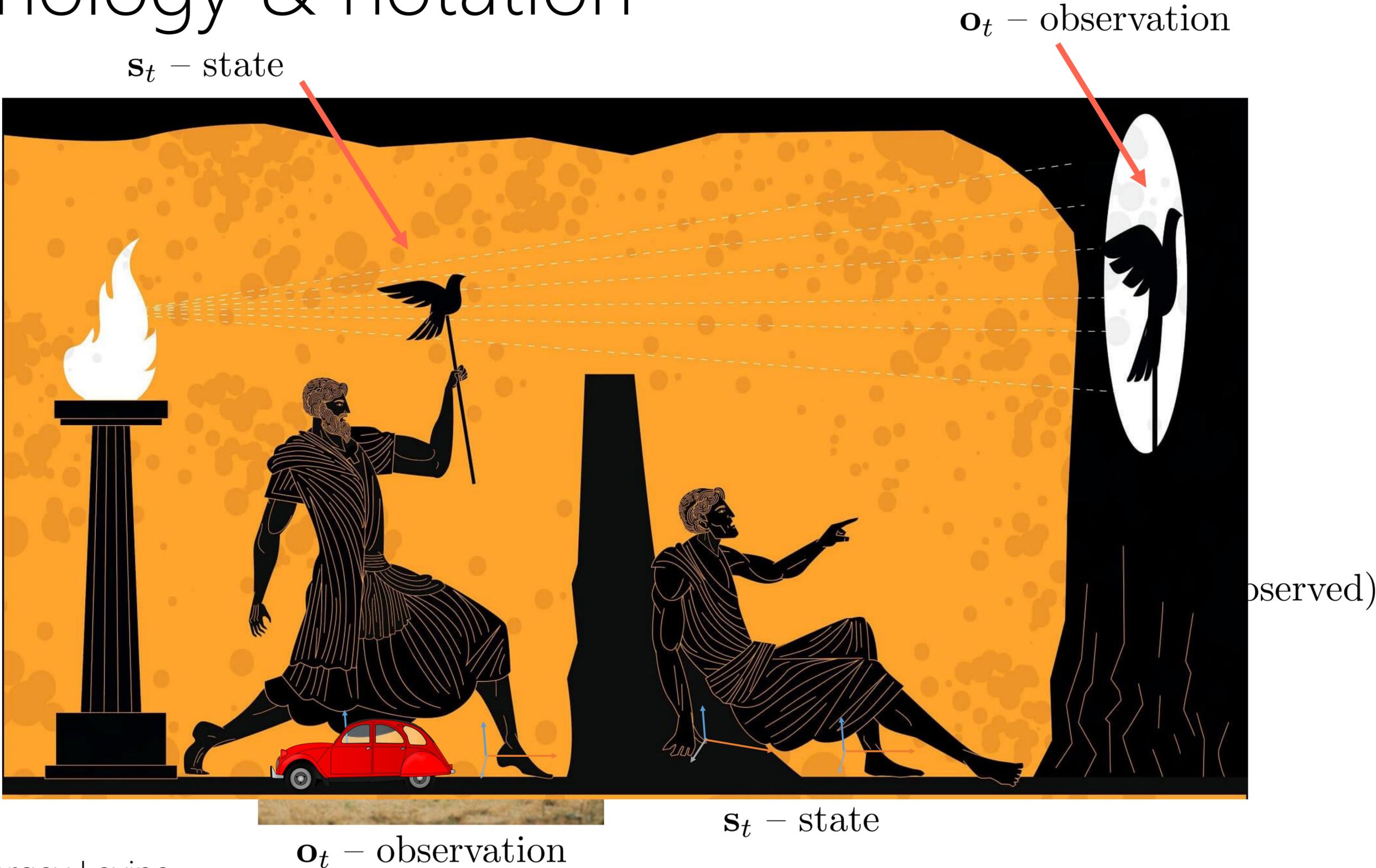


\mathbf{o}_t – observation

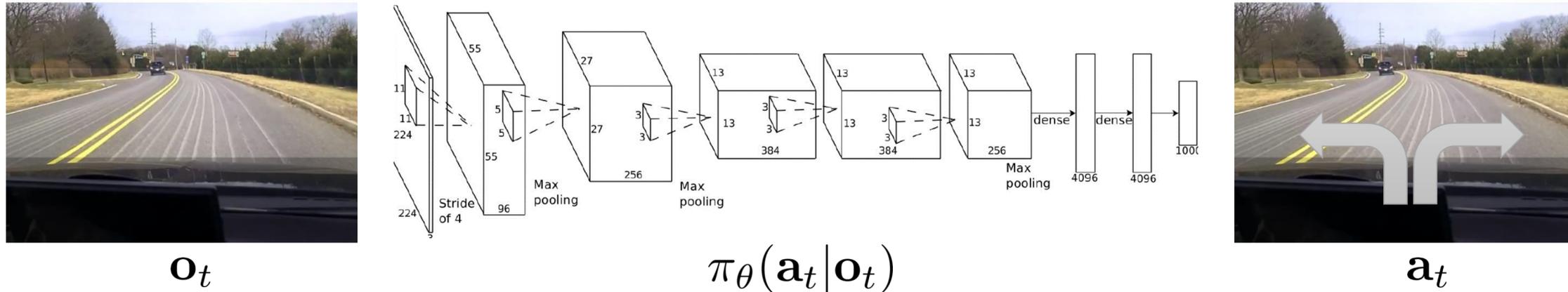


\mathbf{s}_t – state

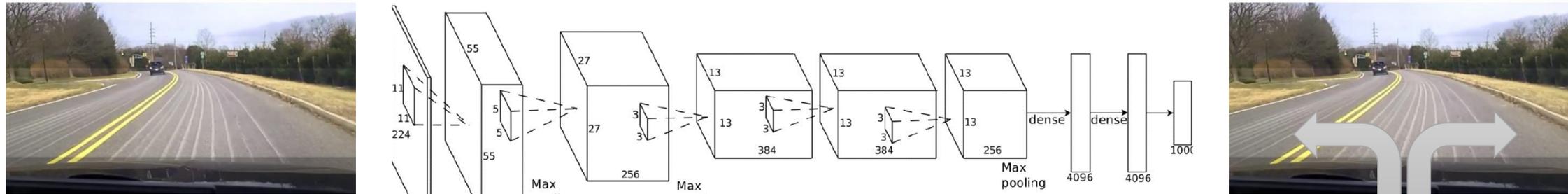
Terminology & notation



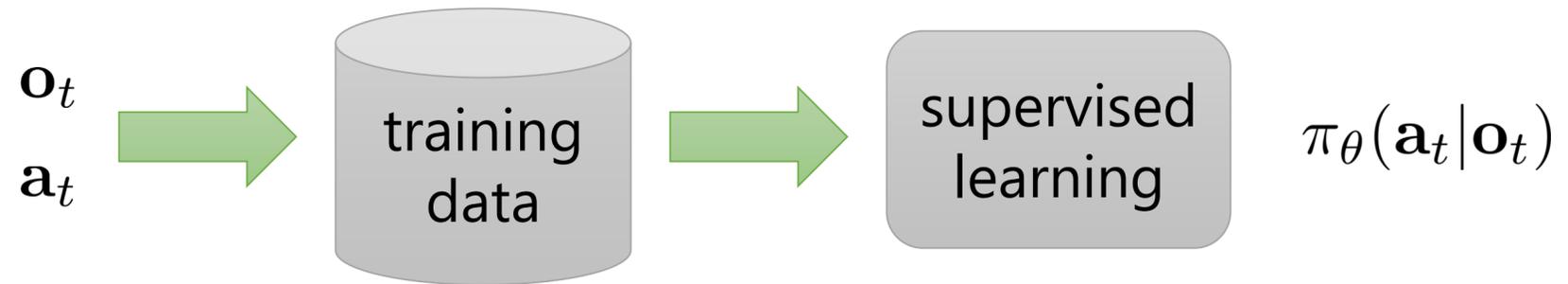
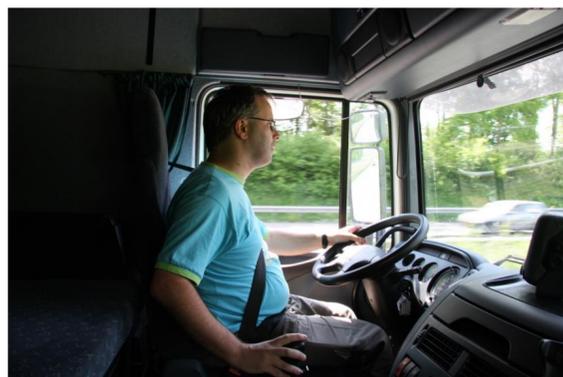
Imitation Learning



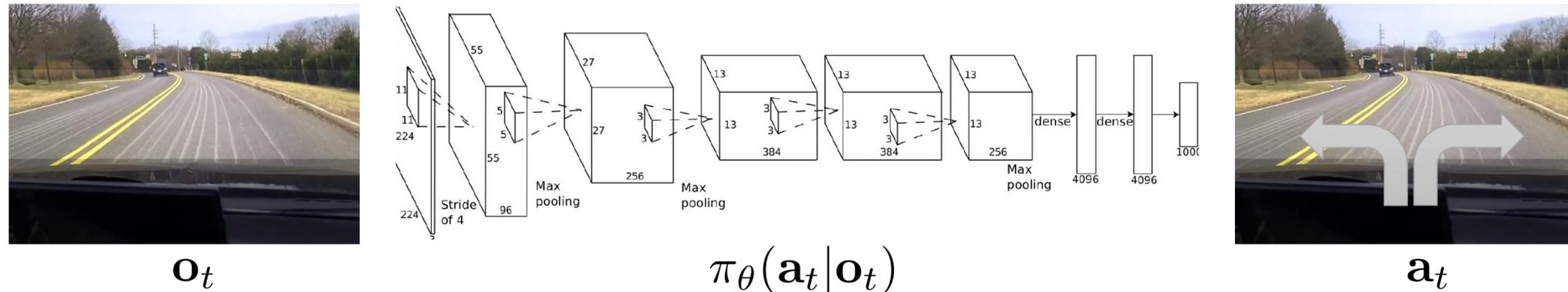
Imitation Learning



Imitation Learning vs Reinforcement Learning?



Reward functions



which action is better or worse?

$r(\mathbf{s}, \mathbf{a})$: reward function

tells us which states and actions are better

\mathbf{s} , \mathbf{a} , $r(\mathbf{s}, \mathbf{a})$, and $p(\mathbf{s}' | \mathbf{s}, \mathbf{a})$ define Markov decision process

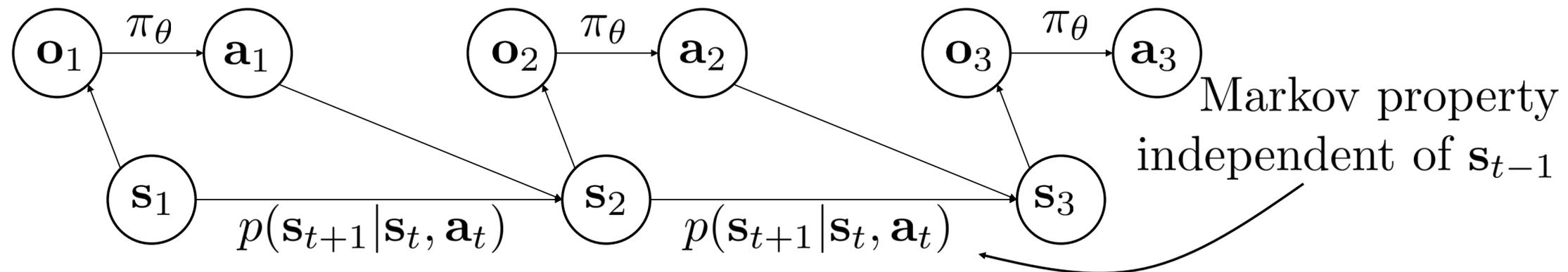
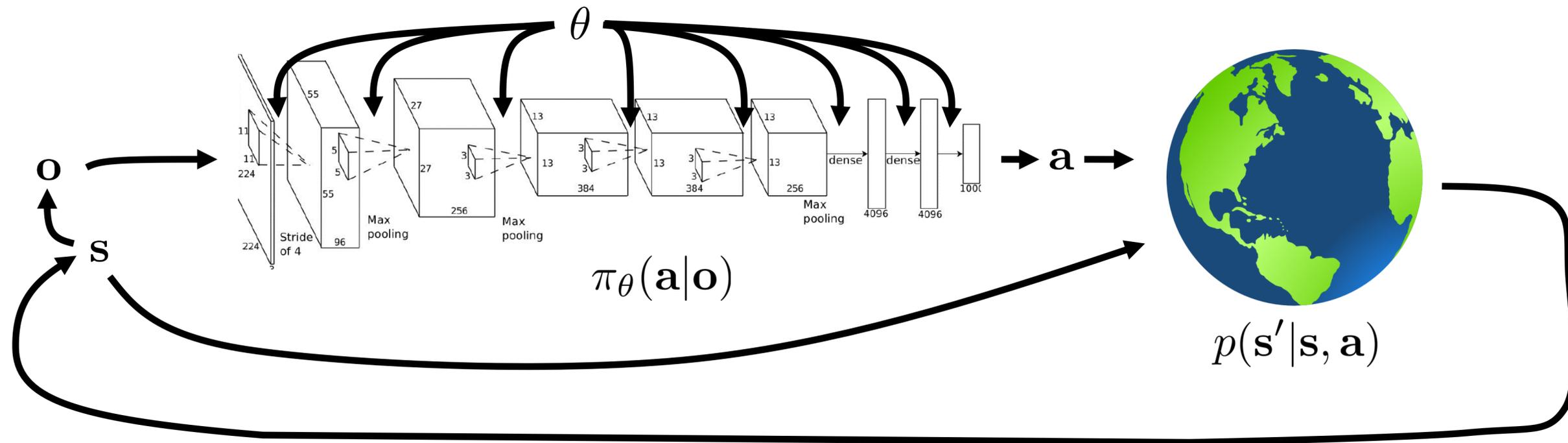


high reward



low reward

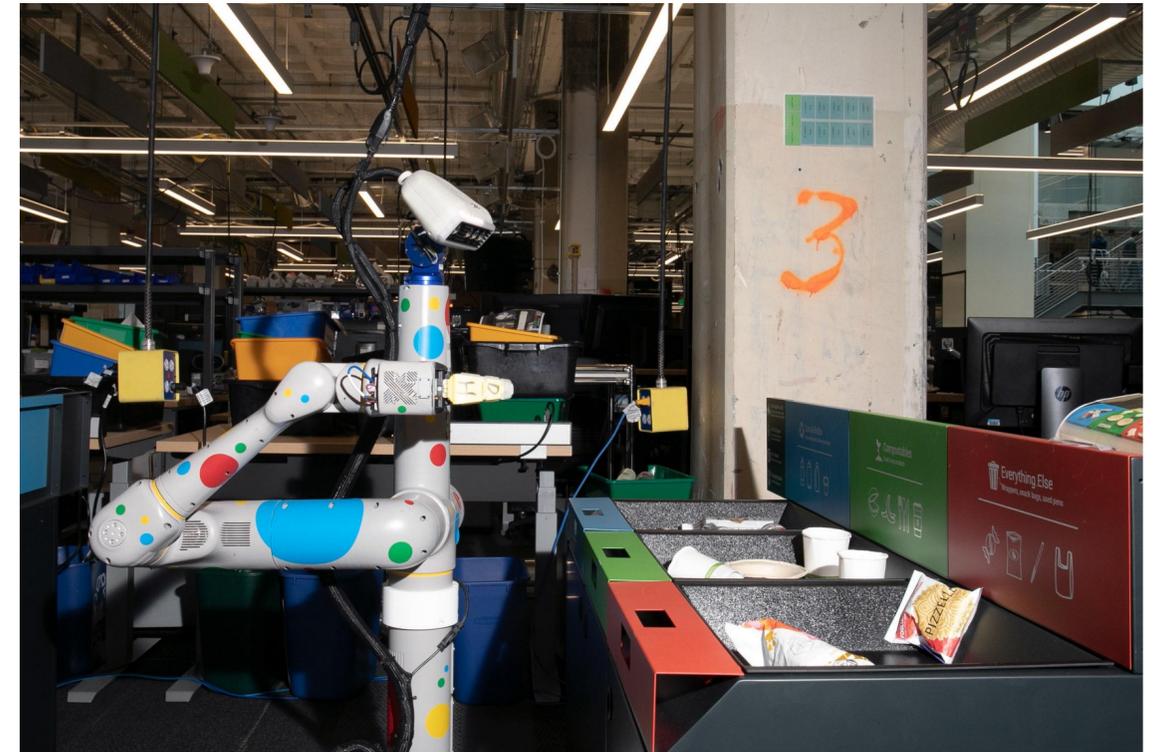
The goal of reinforcement learning



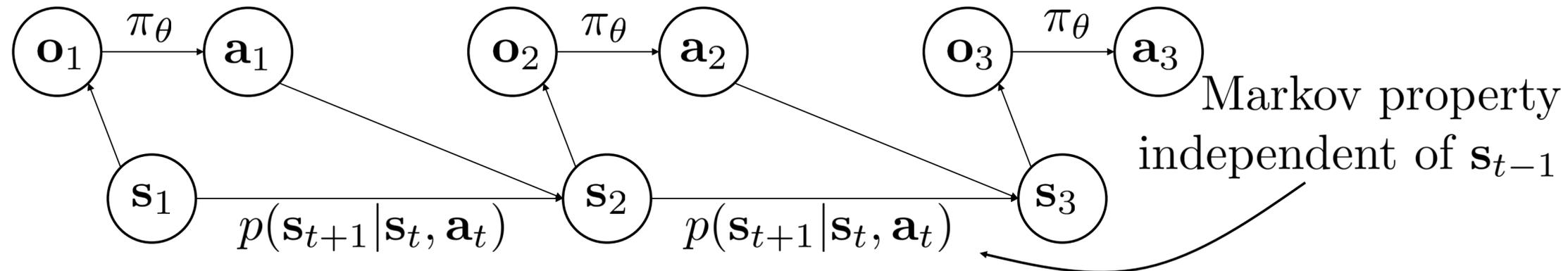
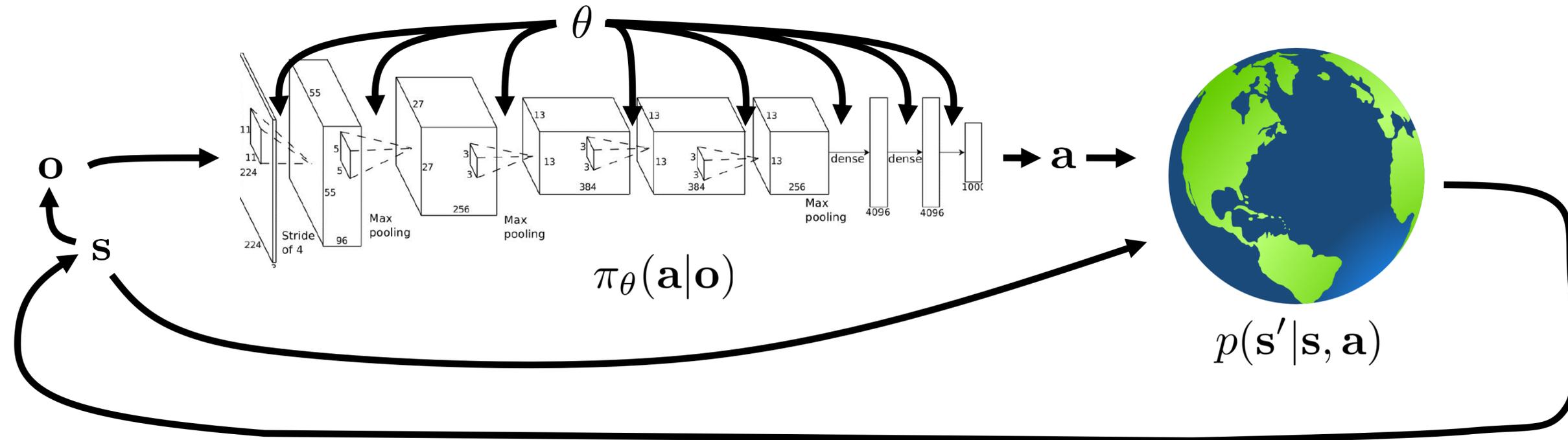
Partial observability

Fully observable?

- Simulated robot performing a reaching task given the goal position and positions and velocities of all of its joints
- Indiscriminate robotic grasping from a bin given an overhead image
- A robot sorting trash given a camera image



The goal of reinforcement learning



$$\theta^* = \arg \max_{\theta} E_{(s, a) \sim p_{\theta}(s, a)} [r(s, a)]$$

infinite horizon case

$$\theta^* = \arg \max_{\theta} \sum_{t=1}^T E_{(s_t, a_t) \sim p_{\theta}(s_t, a_t)} [r(s_t, a_t)]$$

finite horizon case

What is a reinforcement learning **task**?

Recall: supervised learning

data generating distributions, loss

A task: $\mathcal{T}_i \triangleq \{p_i(\mathbf{x}), p_i(\mathbf{y}|\mathbf{x}), \mathcal{L}_i\}$

Reinforcement learning

A task: $\mathcal{T}_i \triangleq \{\mathcal{S}_i, \mathcal{A}_i, p_i(\mathbf{s}_1), p_i(\mathbf{s}'|\mathbf{s}, \mathbf{a}), r_i(\mathbf{s}, \mathbf{a})\}$

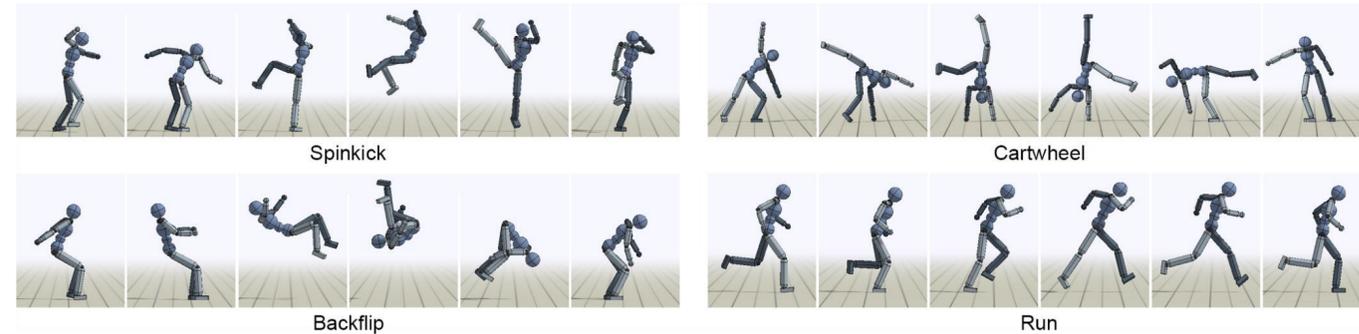
a Markov decision process

much more than the semantic meaning of task!

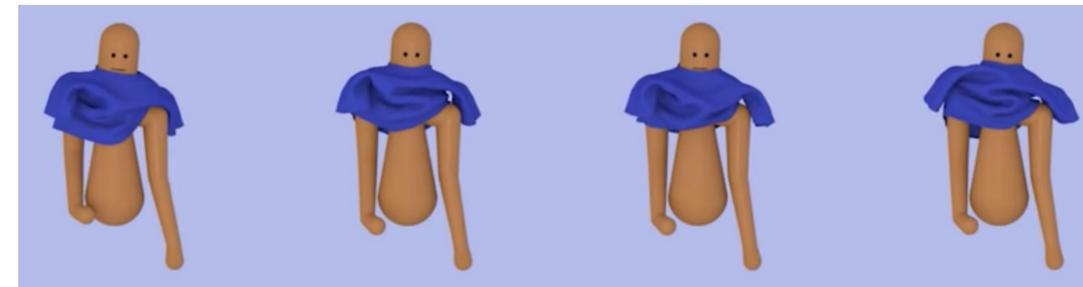
Examples Task Distributions

A task: $\mathcal{T}_i \triangleq \{\mathcal{S}_i, \mathcal{A}_i, p_i(\mathbf{s}_1), p_i(\mathbf{s}'|\mathbf{s}, \mathbf{a}), r_i(\mathbf{s}, \mathbf{a})\}$

Character animation: across maneuvers
 $r_i(\mathbf{s}, \mathbf{a})$ vary



across garments &
initial states
 $p_i(\mathbf{s}_1), p_i(\mathbf{s}'|\mathbf{s}, \mathbf{a})$ vary



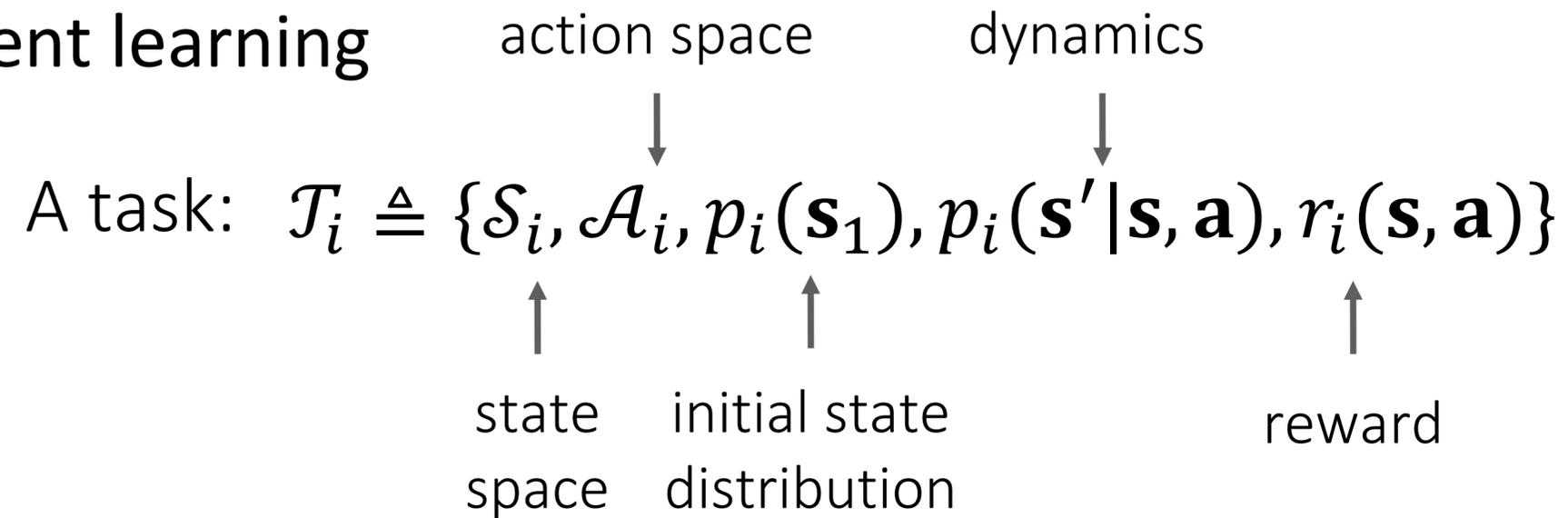
Multi-robot RL:



$\mathcal{S}_i, \mathcal{A}_i, p_i(\mathbf{s}_1), p_i(\mathbf{s}'|\mathbf{s}, \mathbf{a})$ vary

What is a reinforcement learning **task**?

Reinforcement learning



An alternative view:

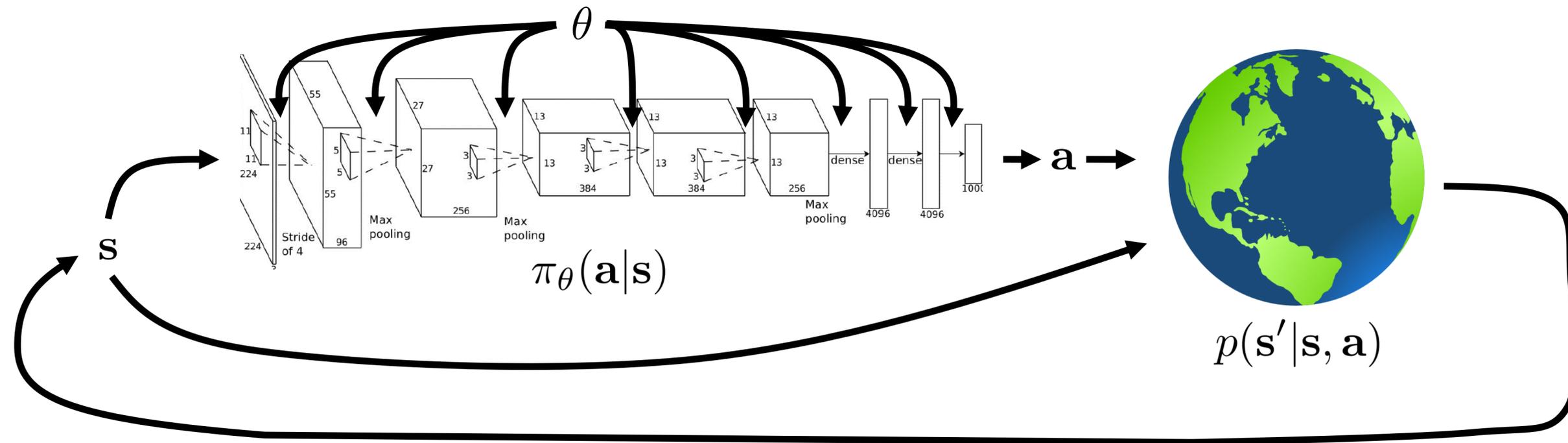
A task identifier is part of the state: $\mathbf{s} = (\bar{\mathbf{s}}, \mathbf{z}_i)$

original state

$$\mathcal{T}_i \triangleq \{\mathcal{S}_i, \mathcal{A}_i, p_i(\mathbf{s}_1), p(\mathbf{s}'|\mathbf{s}, \mathbf{a}), r(\mathbf{s}, \mathbf{a})\} \longrightarrow \{\mathcal{T}_i\} = \left\{ \cup \mathcal{S}_i, \cup \mathcal{A}_i, \frac{1}{N} \sum_i p_i(\mathbf{s}_1), p(\mathbf{s}'|\mathbf{s}, \mathbf{a}), r(\mathbf{s}, \mathbf{a}) \right\}$$

It can be cast as a standard **Markov decision process!**

The goal of **multi-task** reinforcement learning



Multi-task RL

The same as before, except:

a task identifier is part of the state: $\mathbf{s} = (\bar{\mathbf{s}}, \mathbf{z}_i)$

e.g. one-hot task ID

language description

desired goal state, $\mathbf{z}_i = \mathbf{s}_g \leftarrow$ “goal-conditioned RL”

If it's still a standard **Markov decision process**,

then, why not apply standard **RL algorithms**?

You can!

You can often do better.

What is the reward?

The same as before

Or, for **goal-conditioned** RL:

$$r(\mathbf{s}) = r(\bar{\mathbf{s}}, \mathbf{s}_g) = -d(\bar{\mathbf{s}}, \mathbf{s}_g)$$

Distance function d examples:

- Euclidean ℓ_2
- sparse 0/1

The Plan

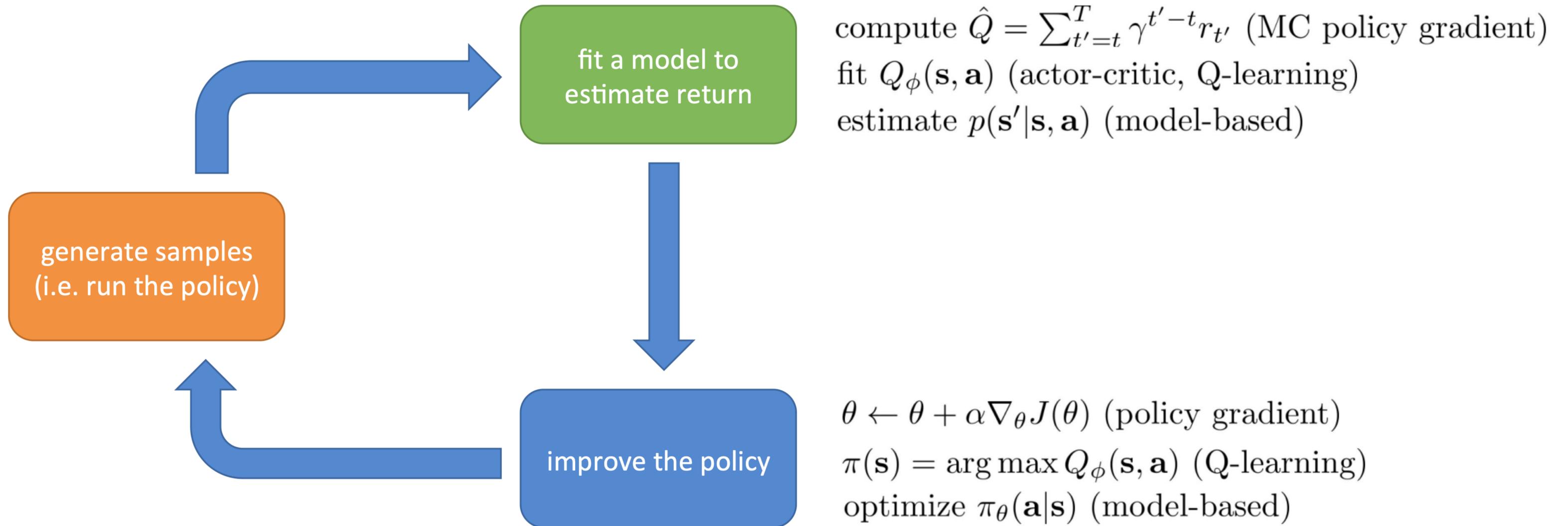
Multi-task reinforcement learning problem

Policy gradients & their multi-task counterparts

Q-learning

Multi-task Q-learning

The anatomy of a reinforcement learning algorithm



This lecture: focus on model-free RL methods (policy gradient, Q-learning)

10/19: focus on model-based RL methods

On-policy

vs

Off-policy

- Data comes from the current policy
- Compatible with all RL algorithms
- Can't reuse data from previous policies

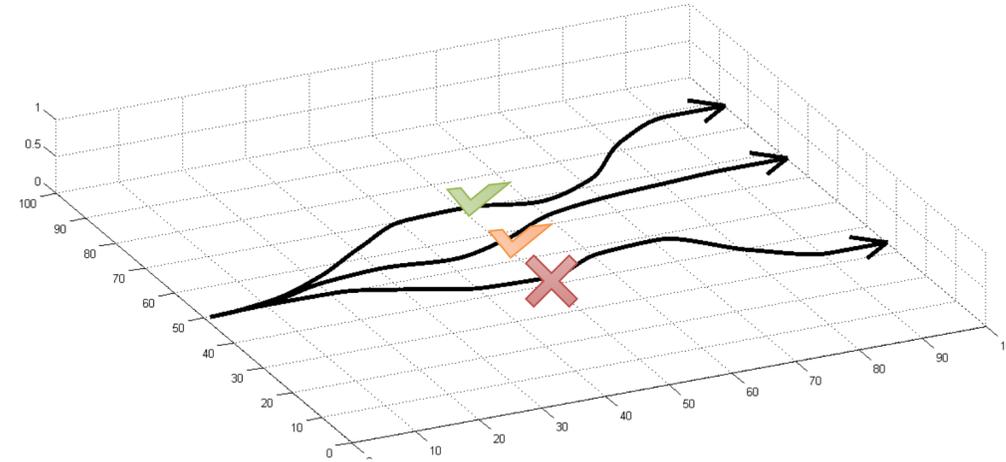
- Data comes from any policy
- Works with specific RL algorithms
- Much more sample efficient, can re-use old data

Evaluating the objective

$$\theta^* = \arg \max_{\theta} \underbrace{E_{\tau \sim p_{\theta}(\tau)} \left[\sum_t r(\mathbf{s}_t, \mathbf{a}_t) \right]}_{J(\theta)}$$

$$J(\theta) = E_{\tau \sim p_{\theta}(\tau)} \left[\sum_t r(\mathbf{s}_t, \mathbf{a}_t) \right] \approx \frac{1}{N} \sum_i \sum_t r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})$$

sum over samples from π_{θ}



Direct policy differentiation

$$\theta^* = \arg \max_{\theta} \underbrace{E_{\tau \sim p_{\theta}(\tau)} \left[\sum_t r(\mathbf{s}_t, \mathbf{a}_t) \right]}_{J(\theta)}$$

a convenient identity

$$\underline{\pi_{\theta}(\tau) \nabla_{\theta} \log \pi_{\theta}(\tau)} = \pi_{\theta}(\tau) \frac{\nabla_{\theta} \pi_{\theta}(\tau)}{\pi_{\theta}(\tau)} = \underline{\nabla_{\theta} \pi_{\theta}(\tau)}$$

$$J(\theta) = E_{\tau \sim \pi_{\theta}(\tau)} \left[\underbrace{r(\tau)}_{\sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t)} \right] = \int \pi_{\theta}(\tau) r(\tau) d\tau$$

$$\nabla_{\theta} J(\theta) = \int \underline{\nabla_{\theta} \pi_{\theta}(\tau)} r(\tau) d\tau = \int \underline{\pi_{\theta}(\tau) \nabla_{\theta} \log \pi_{\theta}(\tau)} r(\tau) d\tau = E_{\tau \sim \pi_{\theta}(\tau)} [\nabla_{\theta} \log \pi_{\theta}(\tau) r(\tau)]$$

Direct policy differentiation

$$\theta^* = \arg \max_{\theta} J(\theta)$$

$$J(\theta) = E_{\tau \sim \pi_{\theta}(\tau)} [r(\tau)]$$

$$\nabla_{\theta} J(\theta) = E_{\tau \sim \pi_{\theta}(\tau)} [\nabla_{\theta} \log \pi_{\theta}(\tau) r(\tau)]$$

$$\underbrace{\pi_{\theta}(\mathbf{s}_1, \mathbf{a}_1, \dots, \mathbf{s}_T, \mathbf{a}_T)}_{\pi_{\theta}(\tau)} = p(\mathbf{s}_1) \prod_{t=1}^T \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)$$

log of both sides

$$\log \pi_{\theta}(\tau) = \log p(\mathbf{s}_1) + \sum_{t=1}^T \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) + \log p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)$$

$$\nabla_{\theta} \left[\cancel{\log p(\mathbf{s}_1)} + \sum_{t=1}^T \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) + \cancel{\log p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)} \right]$$

$$\nabla_{\theta} J(\theta) = E_{\tau \sim \pi_{\theta}(\tau)} \left[\left(\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right) \left(\sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) \right) \right]$$

Evaluating the policy gradient

$$\text{recall: } J(\theta) = E_{\tau \sim p_{\theta}(\tau)} \left[\sum_t r(\mathbf{s}_t, \mathbf{a}_t) \right] \approx \frac{1}{N} \sum_i \sum_t r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})$$

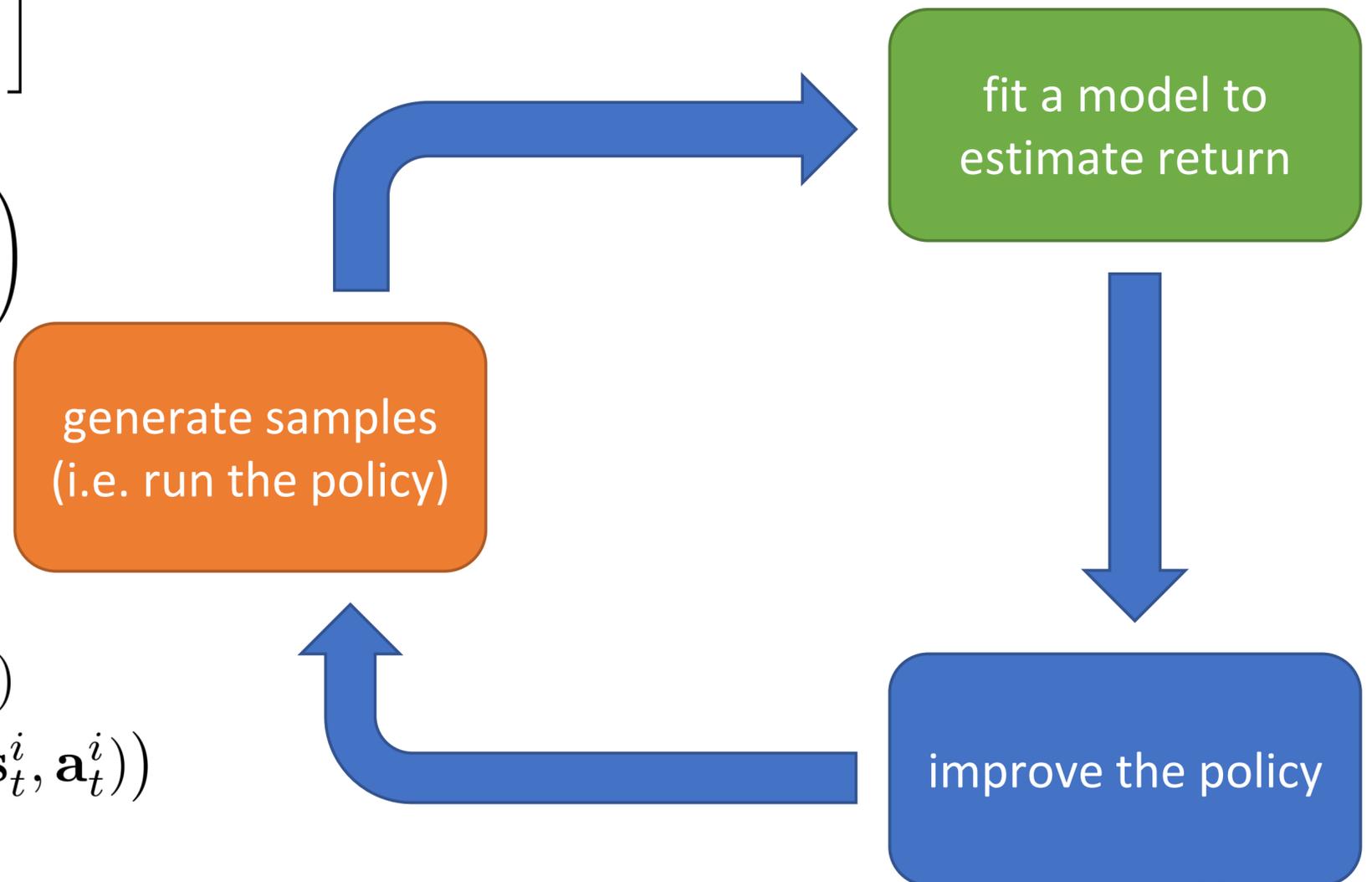
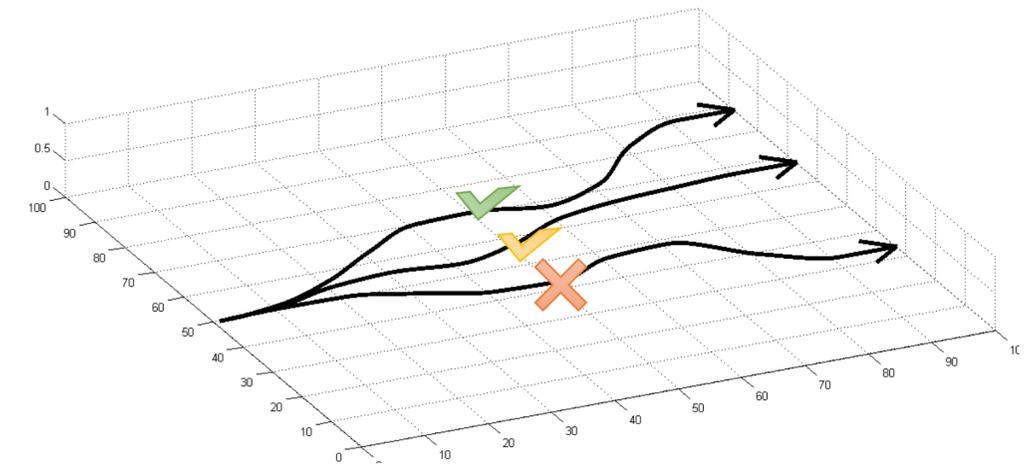
$$\nabla_{\theta} J(\theta) = E_{\tau \sim \pi_{\theta}(\tau)} \left[\left(\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right) \left(\sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) \right) \right]$$

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left(\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right) \left(\sum_{t=1}^T r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right)$$

$$\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$$

REINFORCE algorithm:

1. sample $\{\tau^i\}$ from $\pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)$ (run the policy)
2. $\nabla_{\theta} J(\theta) \approx \sum_i \left(\sum_t \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t^i | \mathbf{s}_t^i) \right) \left(\sum_t r(\mathbf{s}_t^i, \mathbf{a}_t^i) \right)$
3. $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$

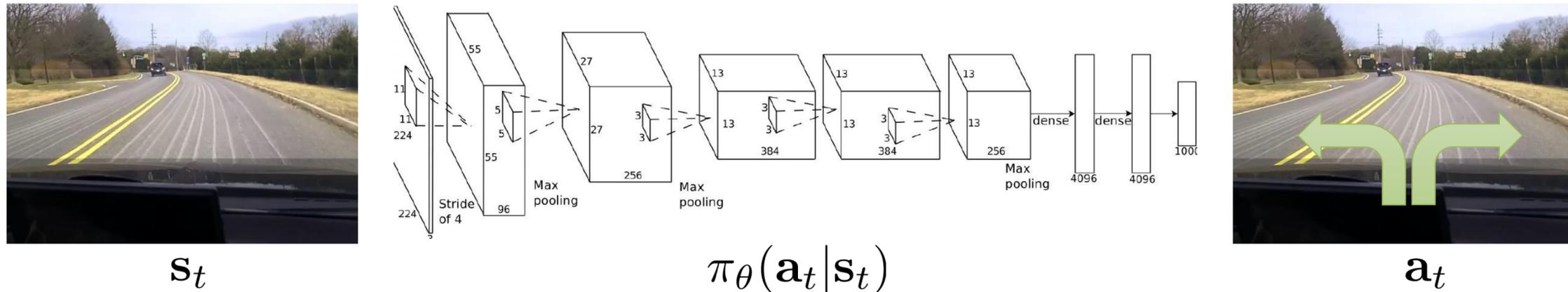


Comparison to maximum likelihood

policy gradient:
$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left(\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right) \left(\sum_{t=1}^T r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right)$$

Multi-task learning algorithms can readily be applied!

maximum likelihood:
$$\nabla_{\theta} J_{\text{ML}}(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left(\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right)$$



What did we just do?

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left(\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right) \left(\sum_{t=1}^T r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right)$$

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \underbrace{\nabla_{\theta} \log \pi_{\theta}(\tau_i)}_T r(\tau_i) \sum_{t=1}^T \nabla_{\theta} \log_{\theta} \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t})$$

maximum likelihood:

$$\nabla_{\theta} J_{\text{ML}}(\theta) \approx \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} \log \pi_{\theta}(\tau_i)$$

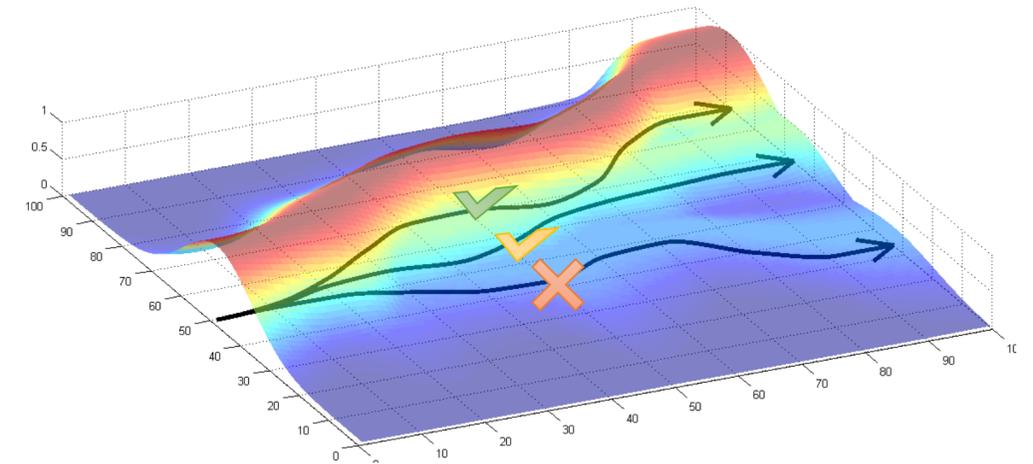
good stuff is made more likely

bad stuff is made less likely

simply formalizes the notion of “trial and error”!

REINFORCE algorithm:

1. sample $\{\tau^i\}$ from $\pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)$ (run it on the robot)
2. $\nabla_{\theta} J(\theta) \approx \sum_i \left(\sum_t \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t^i | \mathbf{s}_t^i) \right) \left(\sum_t r(\mathbf{s}_t^i, \mathbf{a}_t^i) \right)$
3. $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$



Policy Gradients

policy gradient: $\nabla_{\theta} J(\theta) = \underline{E_{\tau \sim \pi_{\theta}(\tau)}} \left[\left(\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right) \left(\sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) \right) \right]$

Pros:

+ Simple

+ Easy to combine with existing multi-task & meta-learning algorithms

Cons:

- Produces a **high-variance** gradient

- Can be mitigated with **baselines** (used by all algorithms in practice), trust regions

- Requires **on-policy** data

- Cannot reuse existing experience to estimate the gradient!

- Importance weights can help, but also high variance

The Plan

Multi-task reinforcement learning problem

Policy gradients & their multi-task/meta counterparts

Q-learning

Multi-task Q-learning

Value-Based RL: Definitions

Value function: $V^\pi(\mathbf{s}_t) = \sum_{t'=t}^T \mathbb{E}_\pi [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \mid \mathbf{s}_t]$ total reward starting from \mathbf{s} and following π
"how good is a state"

Q function: $Q^\pi(\mathbf{s}_t, \mathbf{a}_t) = \sum_{t'=t}^T \mathbb{E}_\pi [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \mid \mathbf{s}_t, \mathbf{a}_t]$ total reward starting from \mathbf{s} , taking \mathbf{a} ,
and then following π
"how good is a state-action pair"

They're related: $V^\pi(\mathbf{s}_t) = \mathbb{E}_{\mathbf{a}_t \sim \pi(\cdot \mid \mathbf{s}_t)} [Q^\pi(\mathbf{s}_t, \mathbf{a}_t)]$

If you know Q^π , you can use it to **improve** π .

Set $\pi(\mathbf{a} \mid \mathbf{s}) \leftarrow 1$ for $\mathbf{a} = \operatorname{argmax}_{\bar{\mathbf{a}}} Q^\pi(\mathbf{s}, \bar{\mathbf{a}})$ New policy is at least as good as old policy.

Value-Based RL: Definitions

Value function: $V^\pi(\mathbf{s}_t) = \sum_{t'=t}^T \mathbb{E}_\pi [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \mid \mathbf{s}_t]$ total reward starting from \mathbf{s} and following π
"how good is a state"

Q function: $Q^\pi(\mathbf{s}_t, \mathbf{a}_t) = \sum_{t'=t}^T \mathbb{E}_\pi [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \mid \mathbf{s}_t, \mathbf{a}_t]$ total reward starting from \mathbf{s} , taking \mathbf{a} ,
and then following π
"how good is a state-action pair"

For the optimal policy π^* : $Q^*(\mathbf{s}_t, \mathbf{a}_t) = \mathbb{E}_{\mathbf{s}' \sim p(\cdot \mid \mathbf{s}, \mathbf{a})} \left[r(\mathbf{s}, \mathbf{a}) + \gamma \max_{\mathbf{a}'} Q^*(\mathbf{s}', \mathbf{a}') \right]$

Bellman equation

Value-Based RL

Value function: $V^\pi(\mathbf{s}_t) = ?$

Q function: $Q^\pi(\mathbf{s}_t, \mathbf{a}_t) = ?$

Q* function: $Q^*(\mathbf{s}_t, \mathbf{a}_t) = ?$

Value* function: $V^*(\mathbf{s}_t) = ?$

Reward = 1 if I can play it in a month, 0 otherwise

IMPROVISATION TEST EXAMPLES AND IDEAS FOR ROCKSCHOOL GRADE 1 DRUMS EXAM
Written by Theo Lawrence / TL Music Lessons

$\text{♩} = 70$

Exercise 1 - Rock

Exercise 2 - Rock

Exercise 3 - Rock

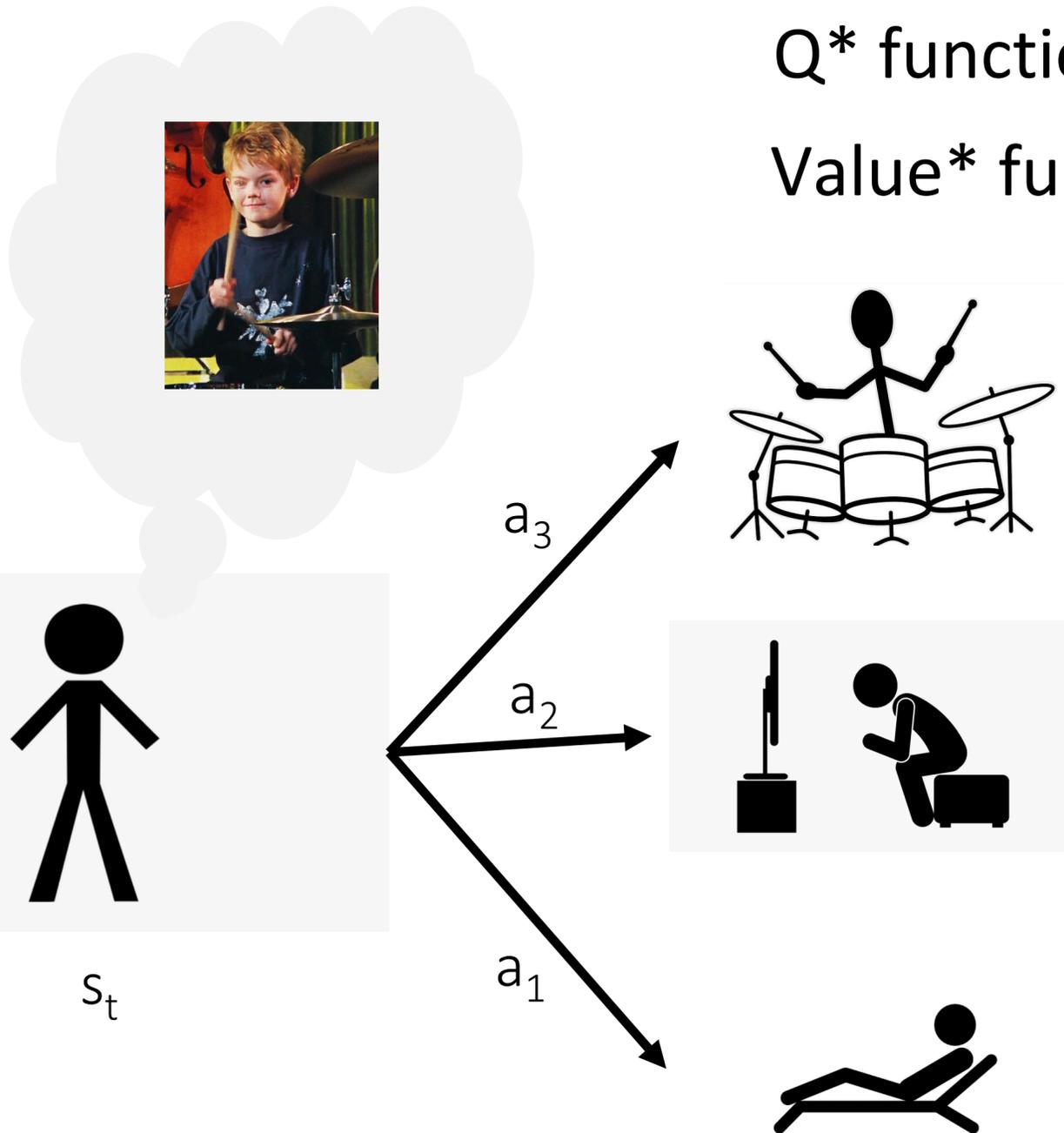
Exercise 4 - Rock

Exercise 5 - Funk Rock

Exercise 6 - Rock

Exercise 7 - Blues

Exercise 8 - Blues



Set $\pi(\mathbf{a}_1 | \mathbf{s}) \leftarrow \pi(\mathbf{a}_1 | \mathbf{s})$ for $\mathbf{a} = \operatorname{argmax}_{\bar{\mathbf{a}}} Q^\pi(\mathbf{s}, \bar{\mathbf{a}})$ New policy is at least as good as old policy.

Fitted Q-iteration Algorithm

full fitted Q-iteration algorithm:

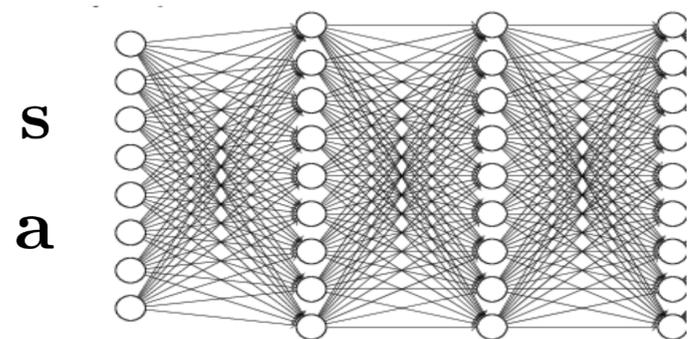
1. collect dataset $\{(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i)\}$ using some policy
2. set $\mathbf{y}_i \leftarrow r(\mathbf{s}_i, \mathbf{a}_i) + \gamma \max_{\mathbf{a}'_i} Q_\phi(\mathbf{s}'_i, \mathbf{a}'_i)$
3. set $\phi \leftarrow \arg \min_{\phi} \frac{1}{2} \sum_i \|Q_\phi(\mathbf{s}_i, \mathbf{a}_i) - \mathbf{y}_i\|^2$

Algorithm hyperparameters

dataset size N , collection policy

iterations K

gradient steps S



$Q_\phi(\mathbf{s}, \mathbf{a})$
parameters ϕ

Result: get a policy $\pi(\mathbf{a}|\mathbf{s})$ from $\arg \max_{\mathbf{a}} Q_\phi(\mathbf{s}, \mathbf{a})$

Important notes:

We can **reuse data** from previous policies!
an **off-policy** algorithm using replay buffers

This is not a gradient descent algorithm!

Example: Q-learning Applied to Robotics

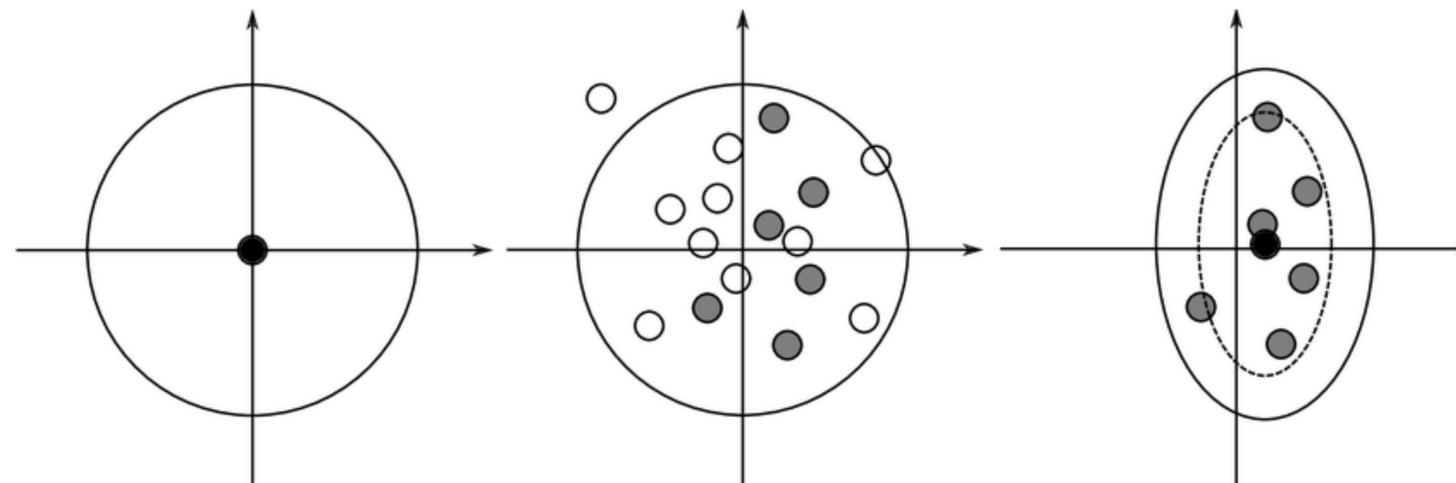
1. collect dataset $\{(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i)\}$ using some policy

2. set $\mathbf{y}_i \leftarrow r(\mathbf{s}_i, \mathbf{a}_i) + \gamma \max_{\mathbf{a}'_i} Q_\phi(\mathbf{s}'_i, \mathbf{a}'_i)$

3. set $\phi \leftarrow \arg \min_{\phi} \frac{1}{2} \sum_i \|Q_\phi(\mathbf{s}_i, \mathbf{a}_i) - \mathbf{y}_i\|^2$

Continuous action space?

Simple optimization algorithm ->
Cross Entropy Method (CEM)

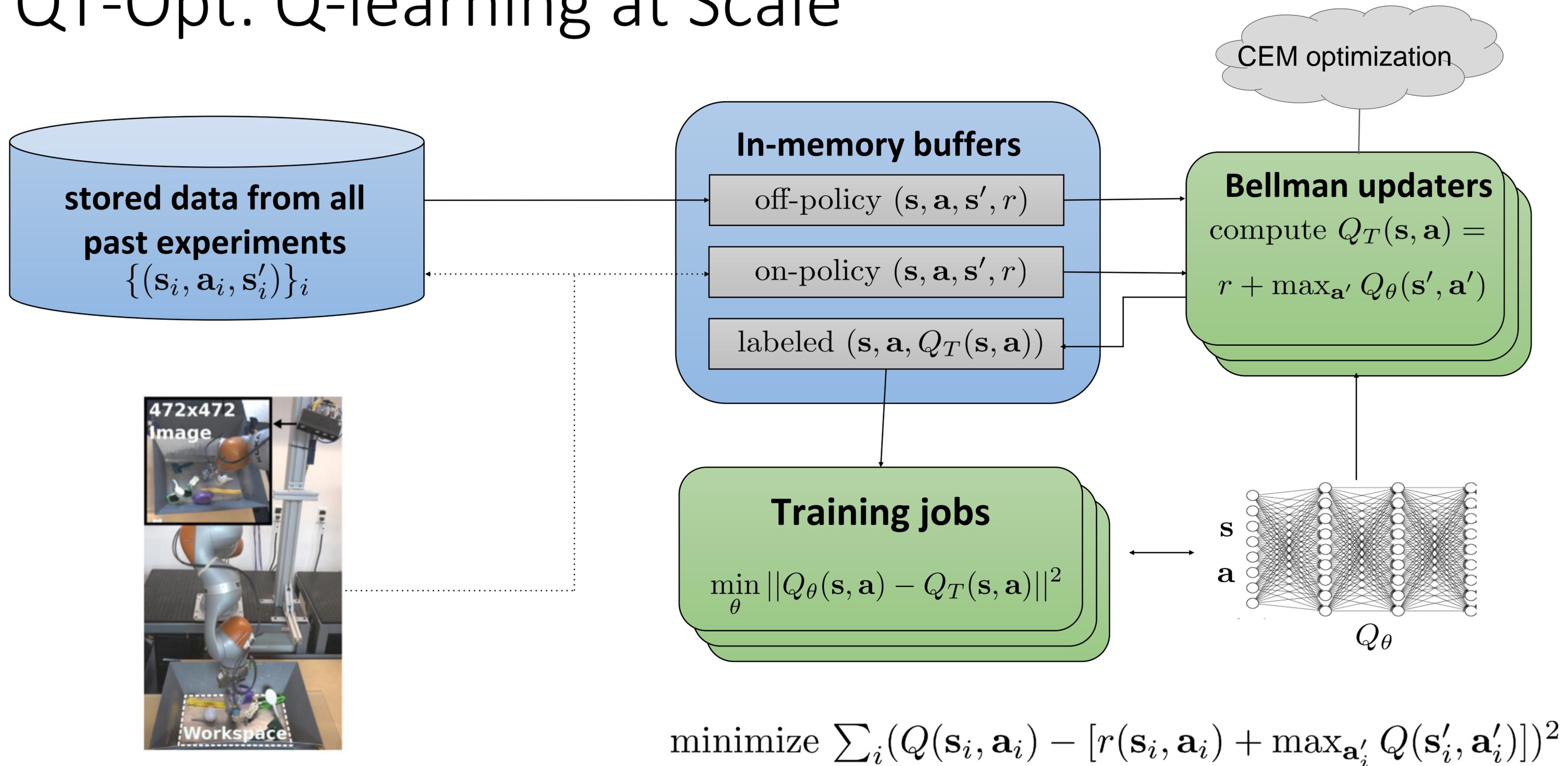


1. Start with the normal distribution $N(\mu, \sigma^2)$.

2. Evaluate some parameters from this distribution and select the best (in grey)

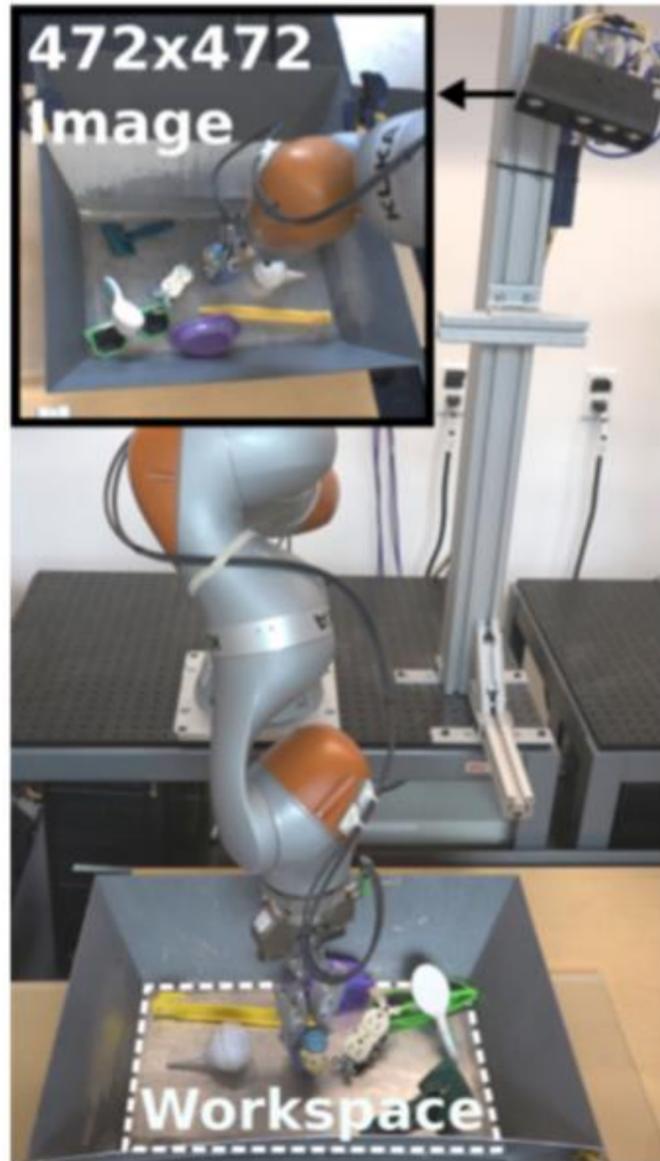
3. Compute the mean and std.dev. of the best, add some noise and goto to 1

QT-Opt: Q-learning at Scale



$$\text{minimize } \sum_i (Q(s_i, a_i) - [r(s_i, a_i) + \max_{a'_i} Q(s'_i, a'_i)])^2$$

QT-Opt: MDP Definition for Grasping



State: over the shoulder RGB camera image, no depth

Action: 4DOF pose change in Cartesian space + gripper control

Reward: binary reward at the end, if the object was lifted. Sparse. No shaping

Automatic success detection:



Q-learning

Bellman equation: $Q^*(\mathbf{s}_t, \mathbf{a}_t) = \mathbb{E}_{\mathbf{s}' \sim p(\cdot | \mathbf{s}, \mathbf{a})} \left[r(\mathbf{s}, \mathbf{a}) + \gamma \max_{\mathbf{a}'} Q^*(\mathbf{s}', \mathbf{a}') \right]$

Pros:

- + More sample efficient than on-policy methods
- + Can incorporate off-policy data (including a fully offline setting)
- + Can update the policy even without seeing the reward
- + Relatively easy to parallelize

Cons:

- Harder to apply standard meta-learning algorithms (DP algorithm)
- Lots of “tricks” to make it work
- Potentially could be harder to learn than just a policy

The Plan

Multi-task reinforcement learning problem

Policy gradients & their multi-task/meta counterparts

Q-learning

Multi-task Q-learning

Multi-Task RL Algorithms

Policy: $\pi_{\theta}(\mathbf{a}|\bar{\mathbf{s}}) \rightarrow \pi_{\theta}(\mathbf{a}|\bar{\mathbf{s}}, \mathbf{z}_i)$

Q-function: $Q_{\phi}(\bar{\mathbf{s}}, \mathbf{a}) \rightarrow Q_{\phi}(\bar{\mathbf{s}}, \mathbf{a}, \mathbf{z}_i)$

Analogous to multi-task supervised learning: stratified sampling, soft/hard weight sharing, etc.

What is different about **reinforcement learning**?

The data distribution is
controlled by the agent!

Why mention it now?

Should we share **data** in addition to sharing **weights**?

An example

Task 1: passing



Task 2: shooting goals



What if you accidentally perform a good pass when trying to shoot a goal?

Store experience as normal. *and* Relabel experience with task 2 ID & reward and store.

“hindsight relabeling” “hindsight experience replay” (HER)

Goal-conditioned RL with hindsight relabeling

1. Collect data $\mathcal{D}_k = \{(\mathbf{s}_{1:T}, \mathbf{a}_{1:T}, \mathbf{s}_g, r_{1:T})\}$ using some policy

2. Store data in replay buffer $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_k$

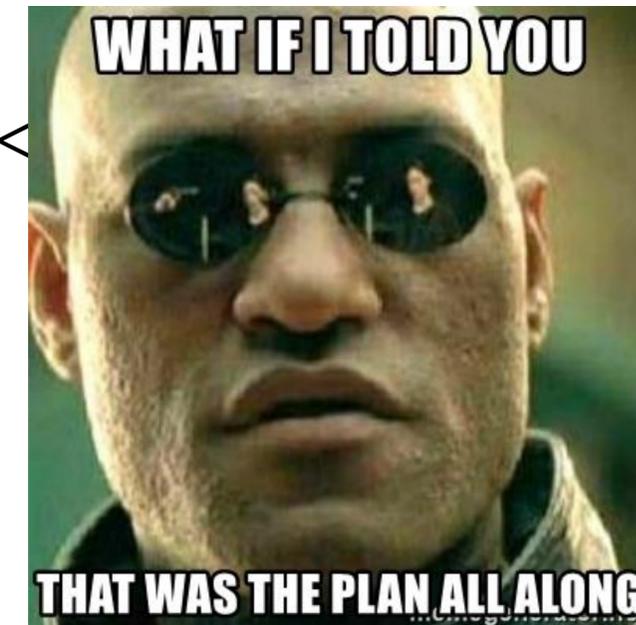
3. Perform **hindsight relabeling**:

a. Relabel experience in \mathcal{D}_k using last state as goal:

$$\mathcal{D}'_k = \{(\mathbf{s}_{1:T}, \mathbf{a}_{1:T}, \mathbf{s}_T, r'_{1:T}) \text{ where } r'_t = -d(\mathbf{s}_t, \mathbf{s}_T)\}$$

b. Store relabeled data in replay buffer $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}'_k$

4. Update policy using replay buffer \mathcal{D}



Result: exploration challenges alleviated

Why mention it now?

Task 1: close a drawer



Task 2: open a drawer



Can we use episodes from drawer opening task for drawer closing task?

How does that answer change for Q-learning vs Policy Gradient?

Multi-task RL with relabeling

1. Collect data $\mathcal{D}_k = \{(\mathbf{s}_{1:T}, \mathbf{a}_{1:T}, \mathbf{z}_i, r_{1:T})\}$ using some policy

2. Store data in replay buffer $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_k$

3. Perform **hindsight relabeling**:

a. Relabel experience in \mathcal{D}_k for task \mathcal{T}_j :

$$\mathcal{D}'_k = \{(\mathbf{s}_{1:T}, \mathbf{a}_{1:T}, \mathbf{z}_j, r'_{1:T}) \text{ where } r'_t = r_j(\mathbf{s}_t)\}$$

b. Store relabeled data in replay buffer $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}'_k$

4. Update policy using replay buffer \mathcal{D}

\leftarrow Which task \mathcal{T}_j to choose?

- randomly
- task(s) in which the trajectory gets high reward

Eysenbach et al. Rewriting History with Inverse RL

Li et al. Generalized Hindsight for RL

When can we apply relabeling?

- reward function form is known, evaluatable
- dynamics consistent across goals/tasks
- using an off-policy algorithm*

Hindsight relabeling for goal-conditioned RL

Example: goal-conditioned RL, simulated robot manipulation

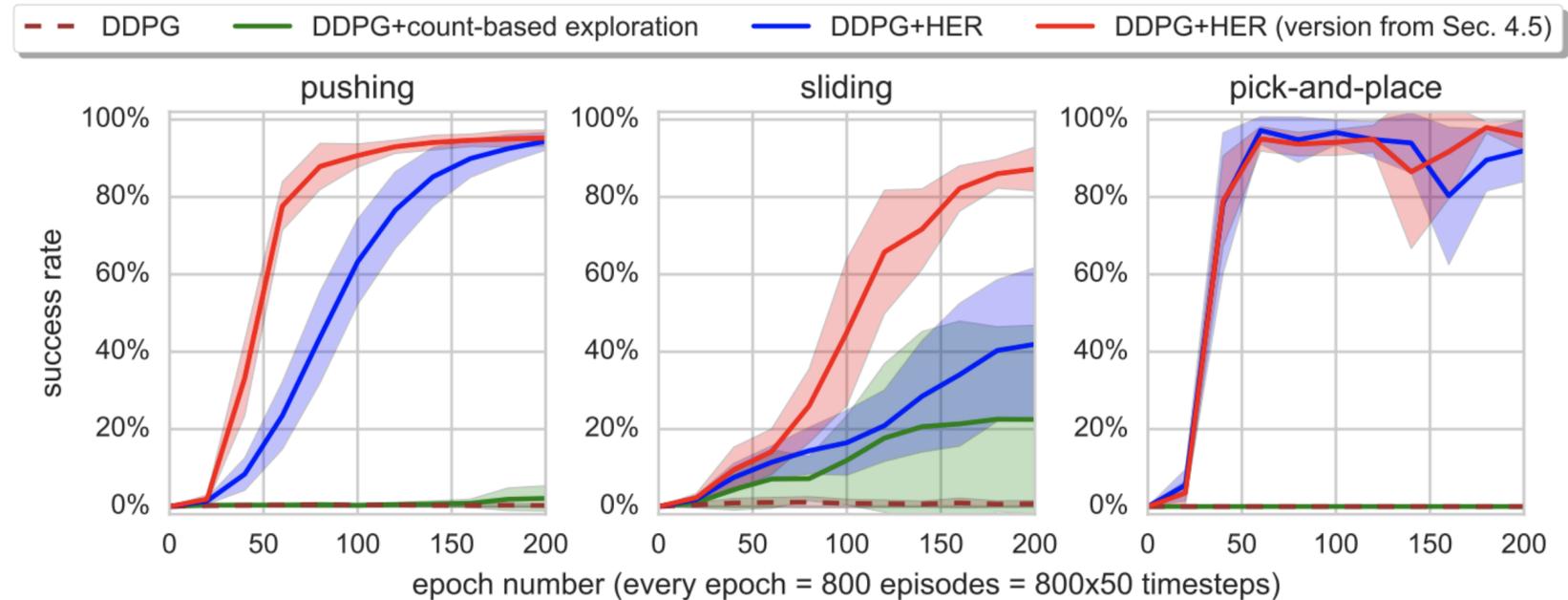
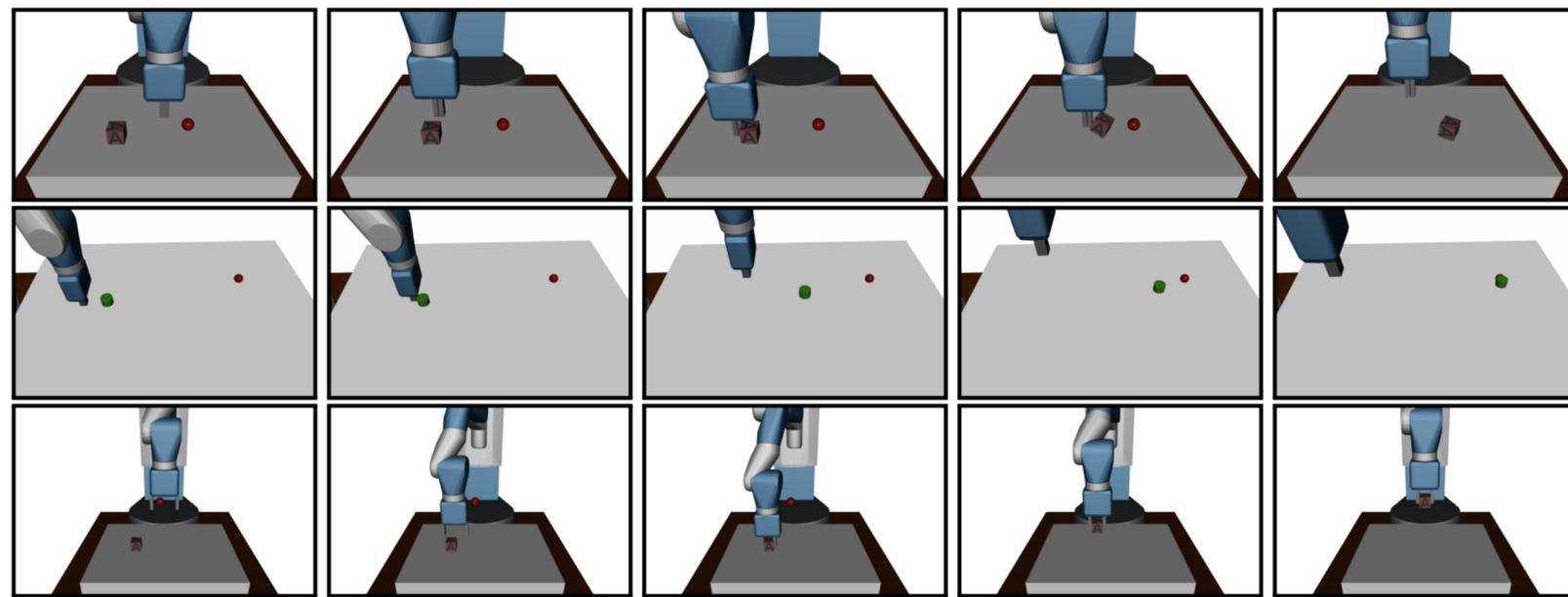


Figure 2: Different tasks: *pushing* (top row), *sliding* (middle row) and *pick-and-place* (bottom row). The red ball denotes the goal position.

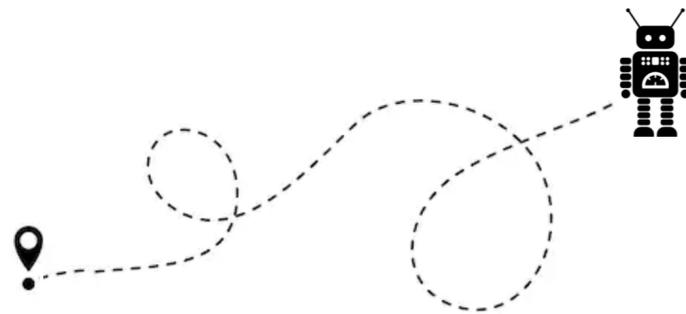
Time Permitting: What about image observations?

Recall: need a distance function between current and goal state!

$$r'_t = -d(\mathbf{s}_t, \mathbf{s}_T)$$

Use binary 0/1 reward? Sparse, but accurate.

Random, unlabeled interaction is *optimal* under the 0/1 reward of reaching the last state.

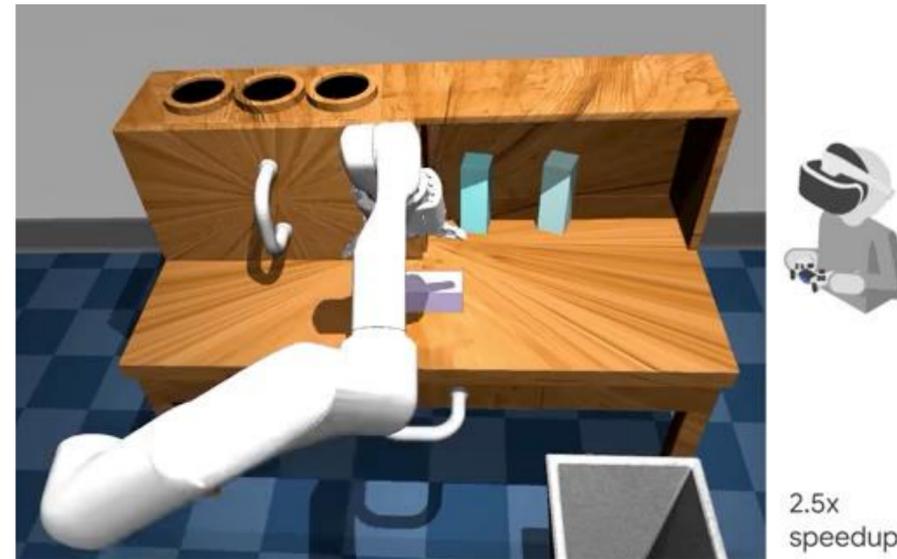


Can we use this insight for **better learning**?

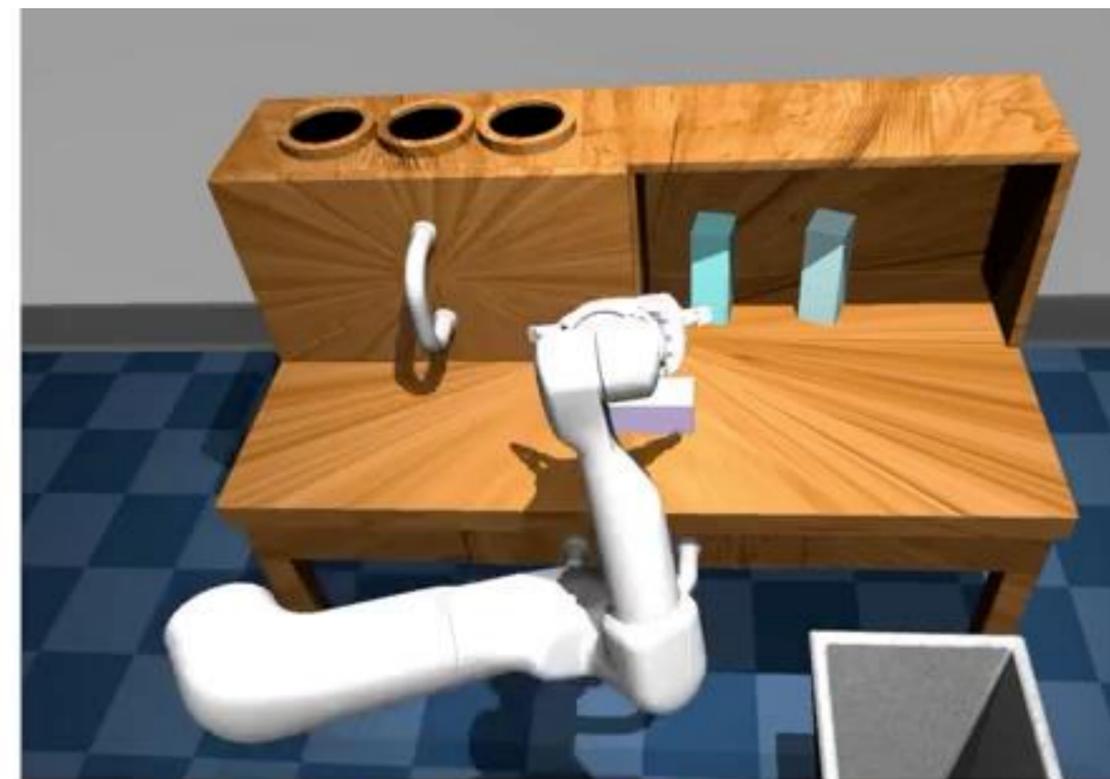
If the data is **optimal**, can we use **supervised imitation learning**?

1. Collect data $\mathcal{D}_k = \{(\mathbf{s}_{1:T}, \mathbf{a}_{1:T})\}$ using some policy
2. Perform **hindsight relabeling**:
 - a. Relabel experience in \mathcal{D}_k using last state as goal:
 $\mathcal{D}'_k = \{(\mathbf{s}_{1:T}, \mathbf{a}_{1:T}, \mathbf{s}_T, r'_{1:T})\}$ where $r'_t = -d(\mathbf{s}_t, \mathbf{s}_T)$
 - b. Store relabeled data in replay buffer $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}'_k$
3. Update policy using **supervised imitation on** replay buffer \mathcal{D}

Collect data from "human play", perform goal-conditioned imitation.



Goal



Single Play-LMP policy

The Plan

Multi-task reinforcement learning problem

Policy gradients & their multi-task/meta counterparts

Q-learning

Multi-task Q-learning

How **data** can be **shared** across tasks.

Many Remaining Questions: The Next Three Weeks

How can we use a **model** in multi-task RL?

Model-based RL - **Oct 19**

What about **meta-RL** algorithms?

Meta-RL - **Oct 21**

Can we learn **exploration strategies** across tasks?

Meta-RL: Learning to explore - **Oct 26**

What about **hierarchies** of tasks?

Hierarchical RL - **Nov 2**

Additional RL Resources

Stanford CS234: Reinforcement Learning

UCL Course from David Silver: Reinforcement Learning

Berkeley CS285: Deep Reinforcement Learning