

Computational Sensorimotor Learning (Spring'21)

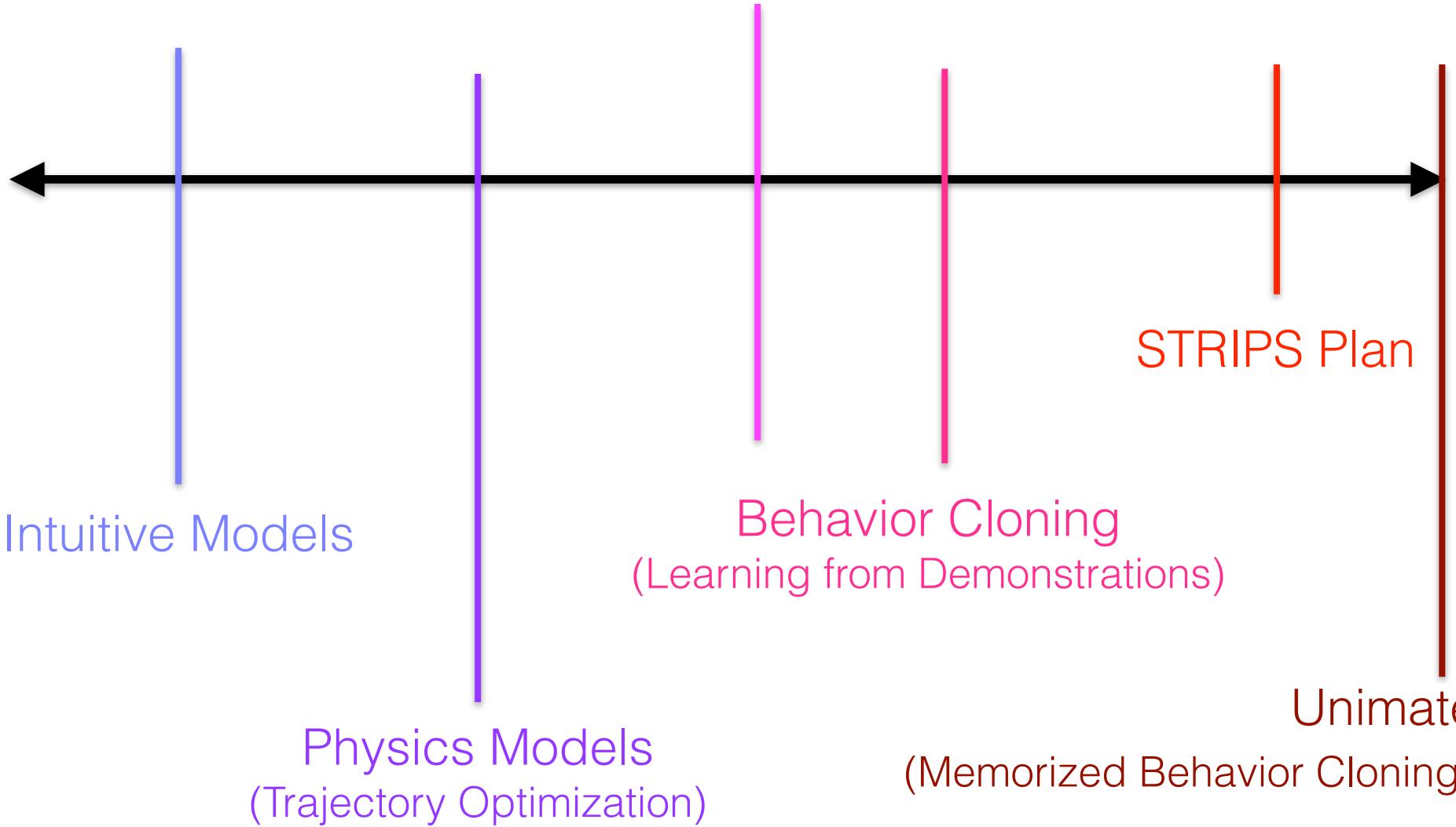
Pulkit Agrawal

Lecture 2
Feb 18 2021

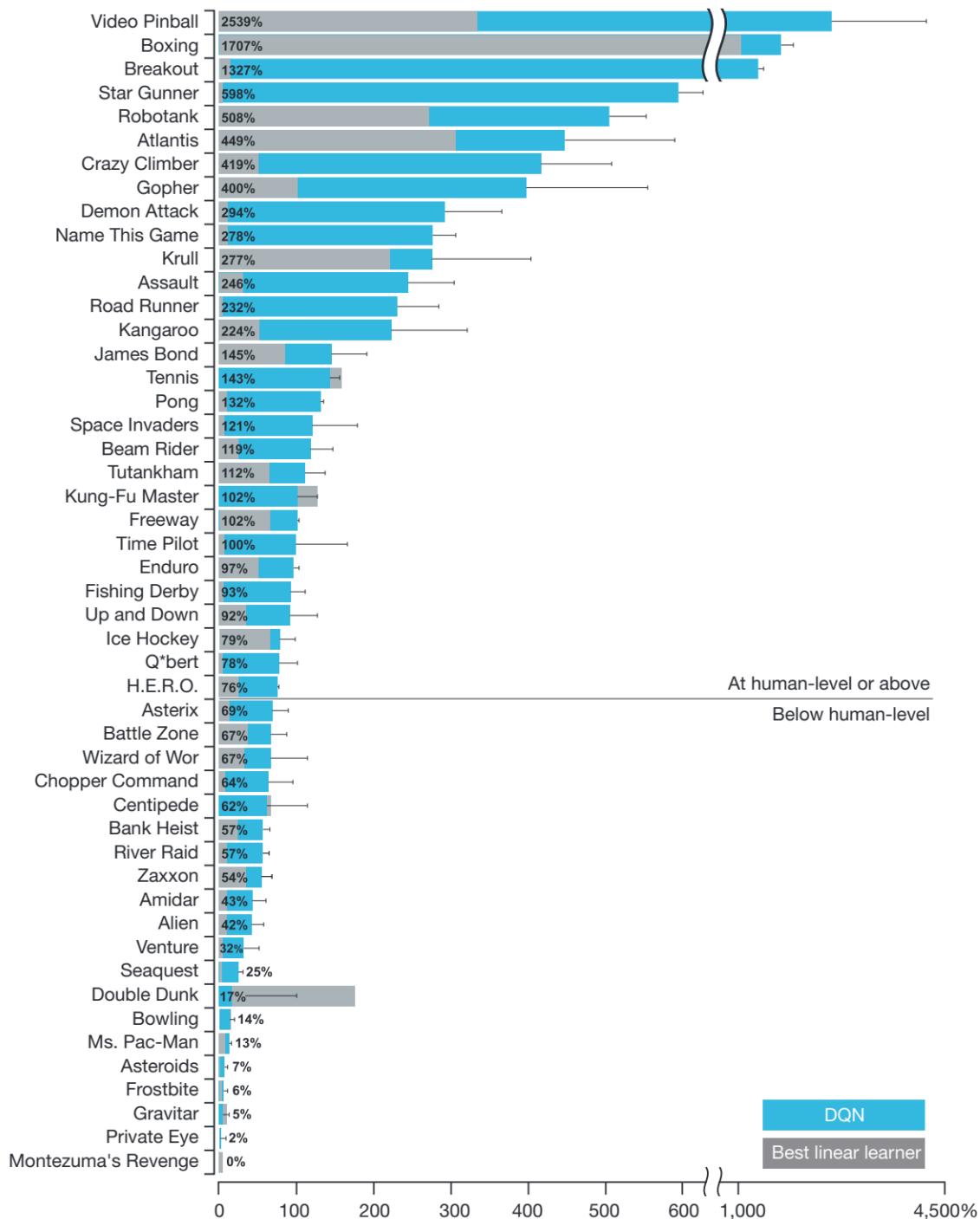
Self-Learnt
Behavior

Hard-Coded
Behavior

Imitation Learning
(Learning by Observing)



Hard to achieve Super-Human Performance with behavior cloning



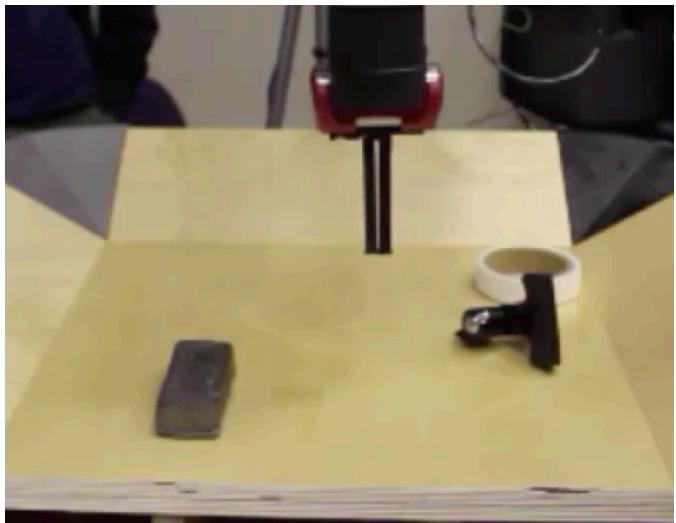
Hard to achieve

Super-Human
Performance

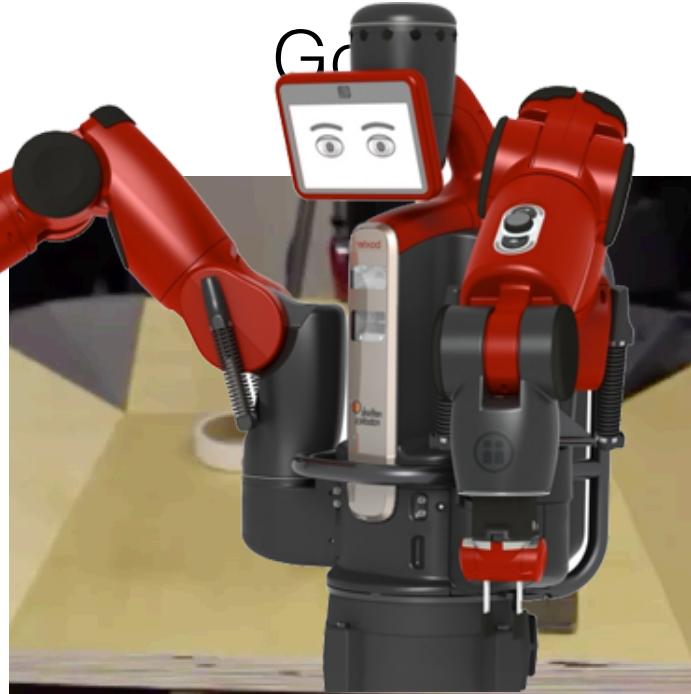
with
behavior cloning

Let the machine
automatically figure out
decision making rules!

Current Observation



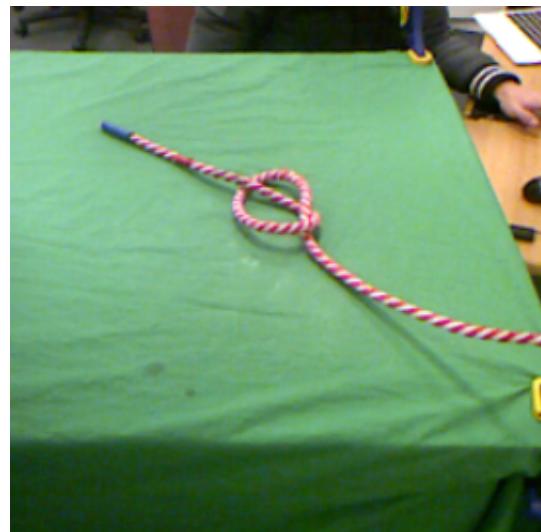
Actions?
→



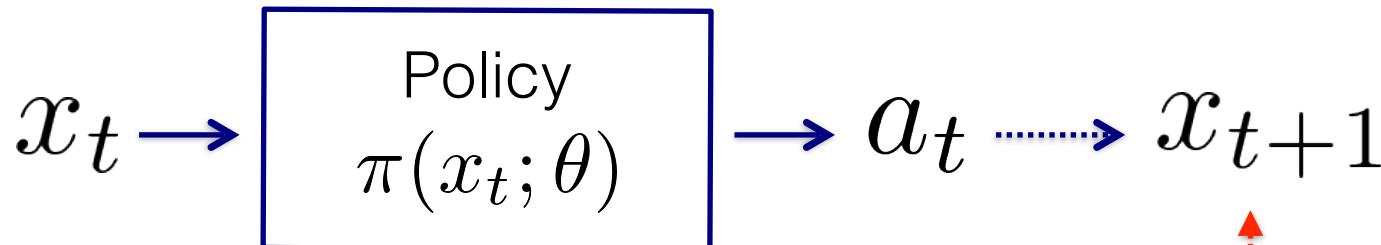
Goal



Actions?
→



Brief Overview of Reinforcement Learning



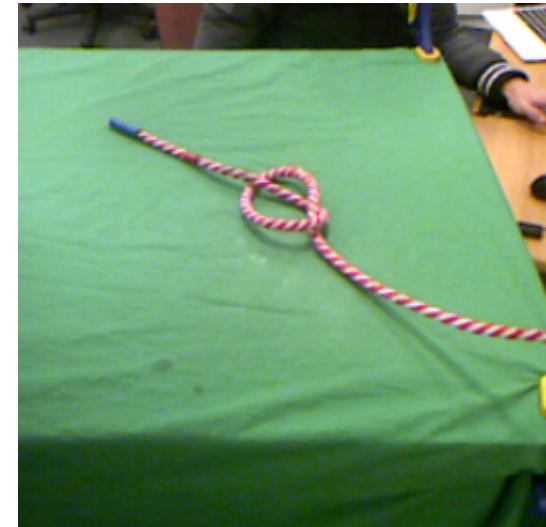
initially random
(how to learn this?)

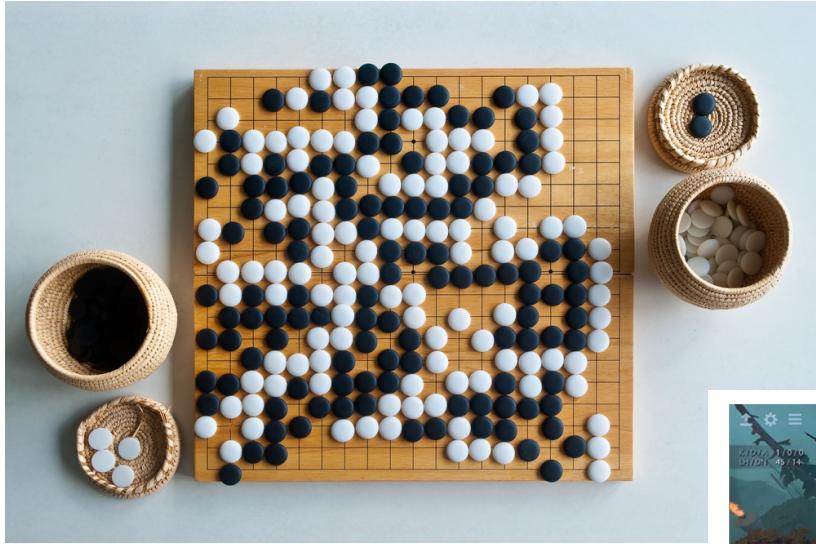
$$\max_{\theta} \mathbb{E} \left(\sum_{t=1}^T r_t \right)$$

x_G



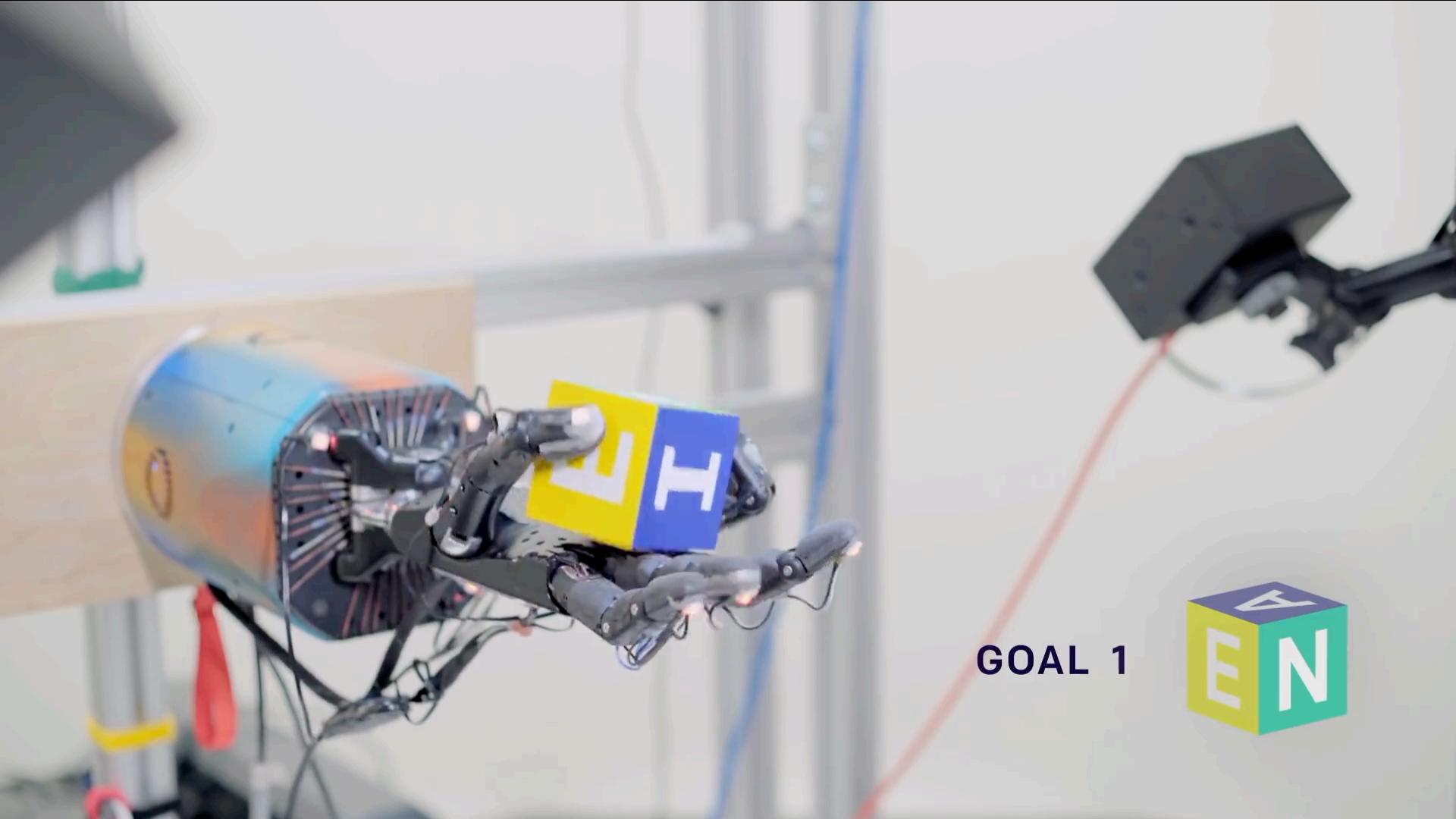
Actions? $\xrightarrow{\hspace{1cm}}$



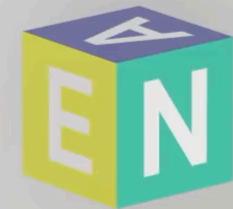


Open AI Five playing DOTA

Learning Dexterity



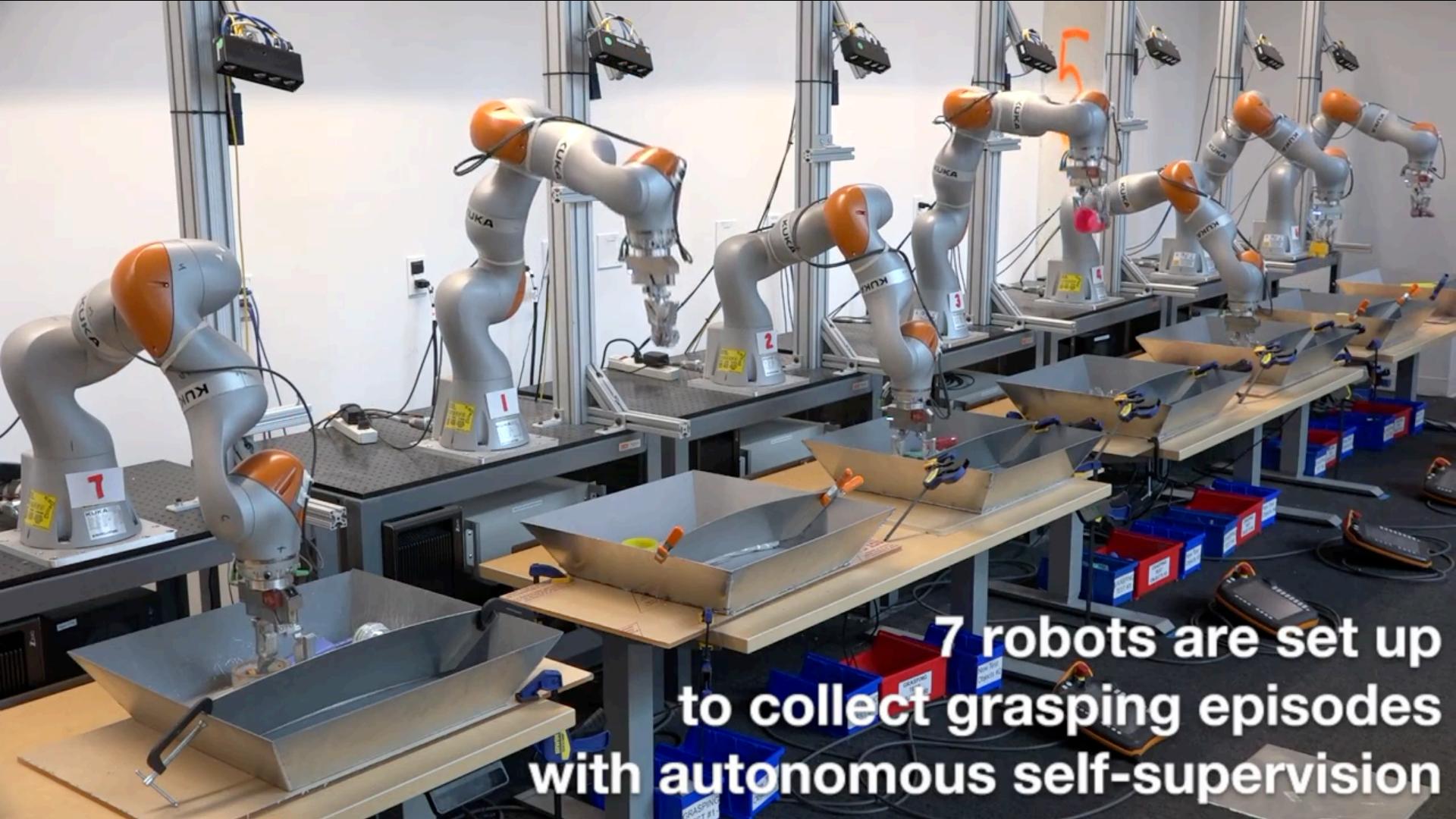
GOAL 1



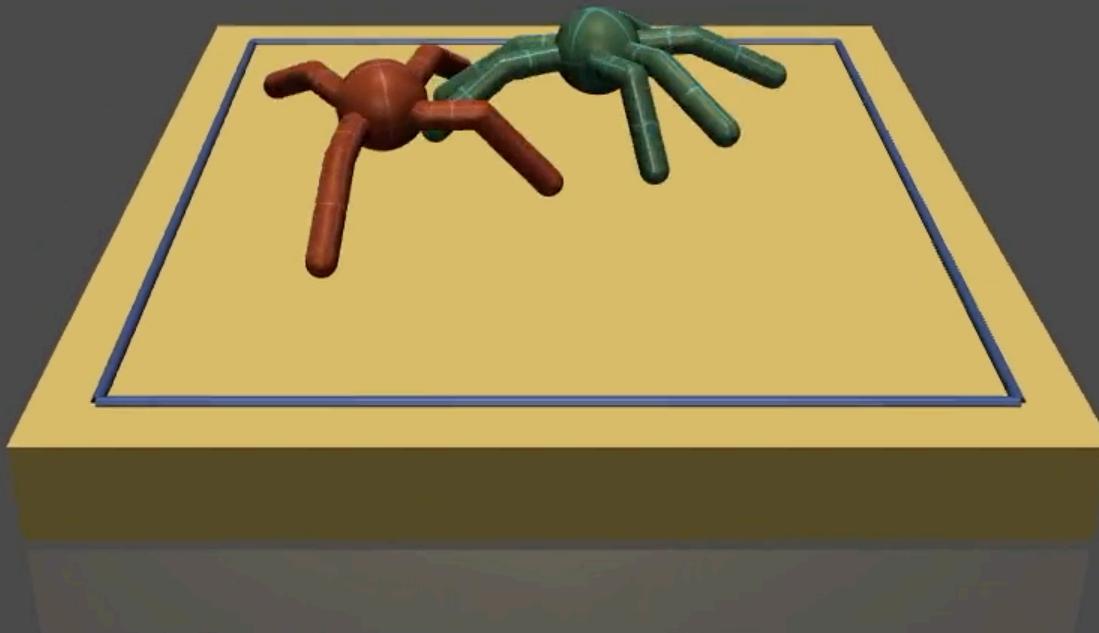
Locomotion Strategies



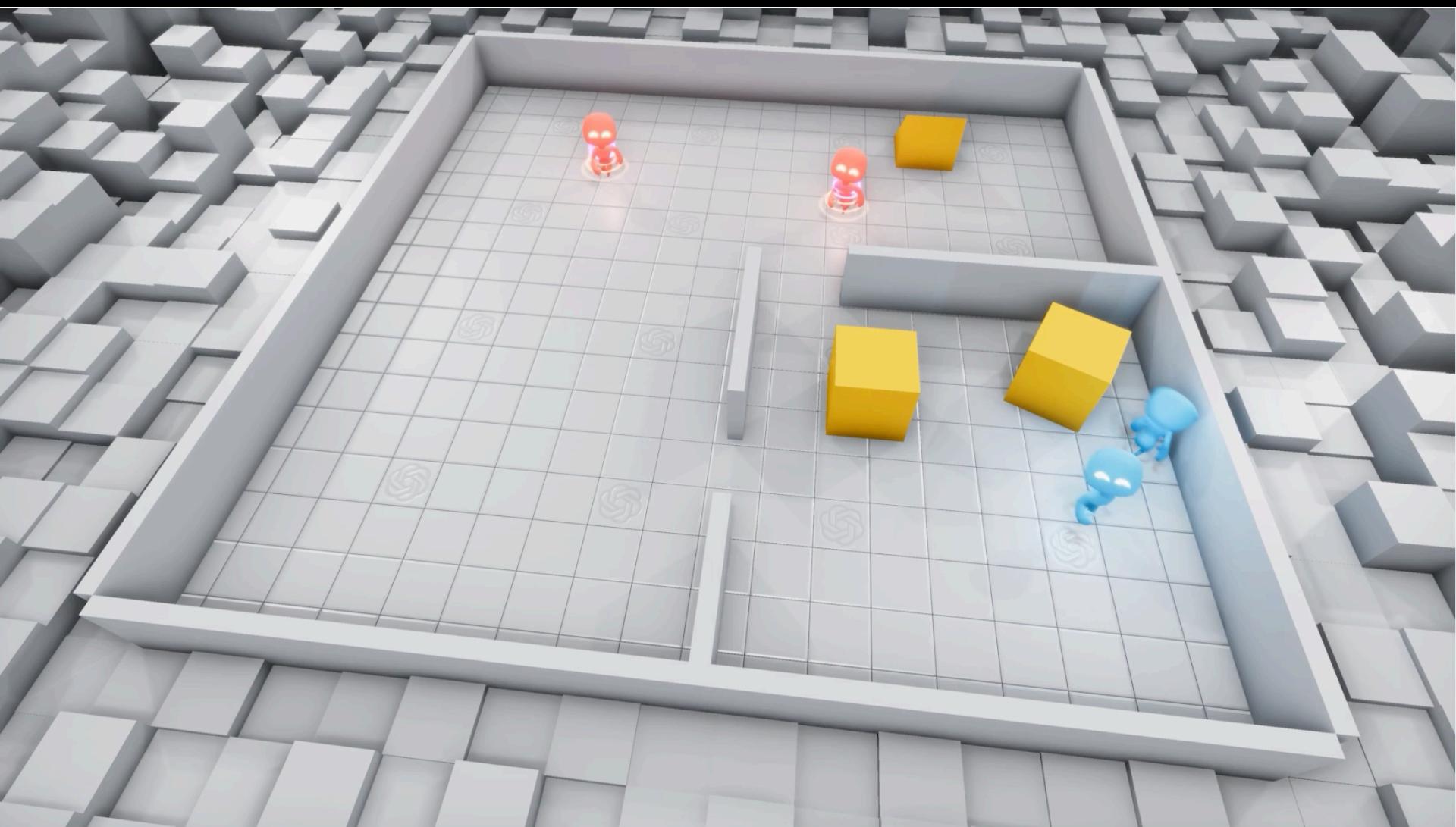
Robots Learning to Grasp Objects



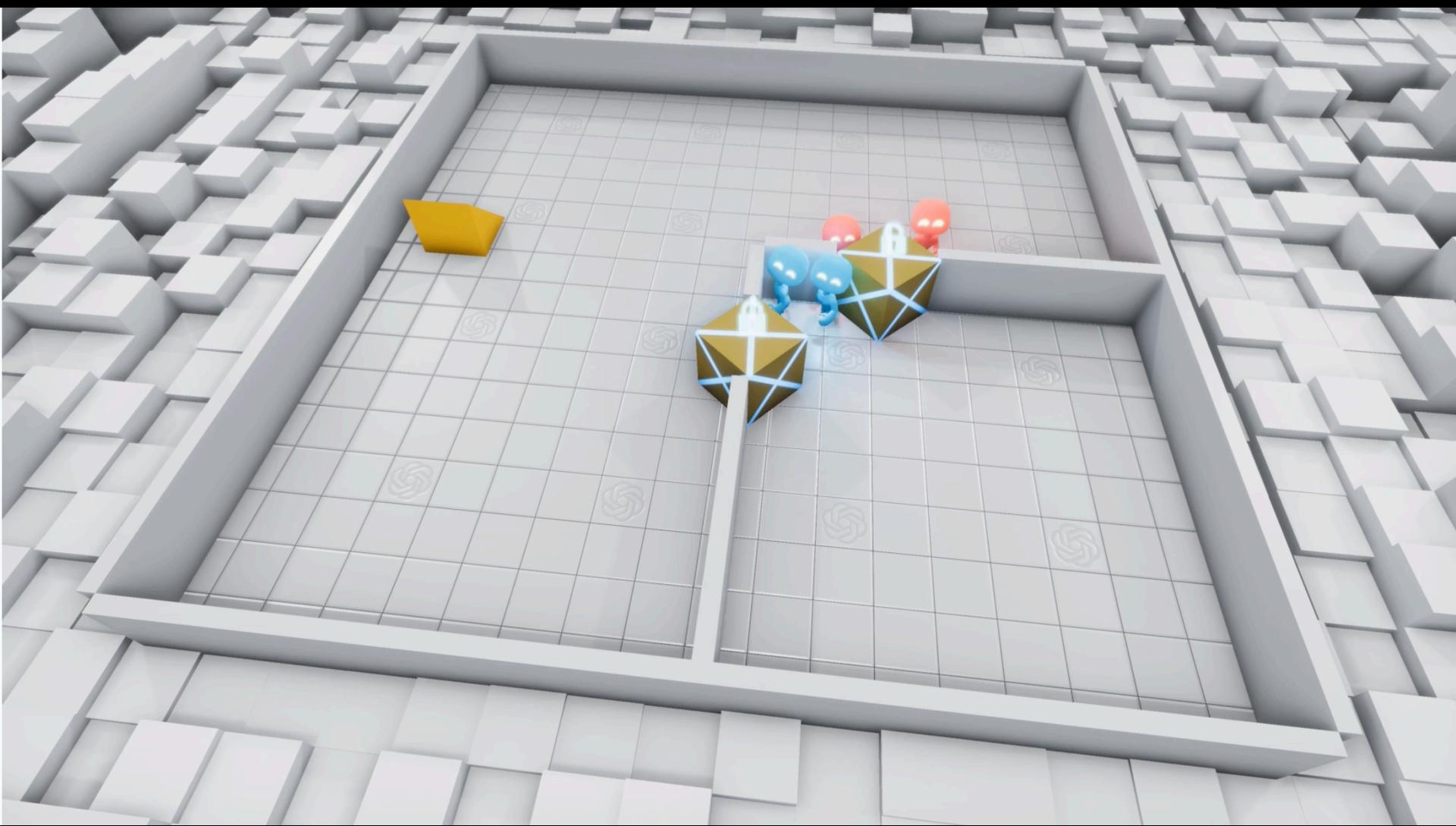
Automatic Discovery of Skills



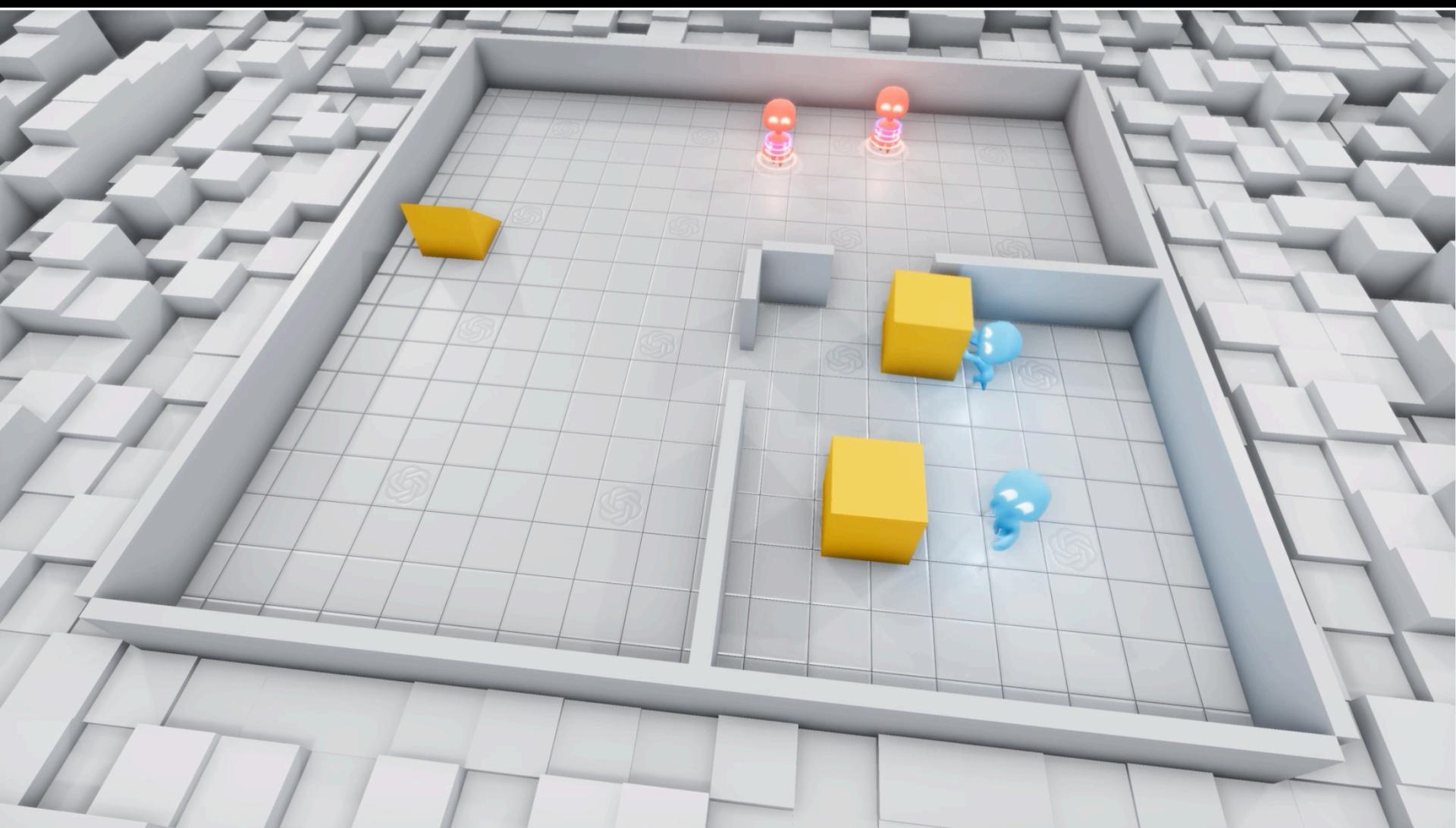
Initial Random Exploration



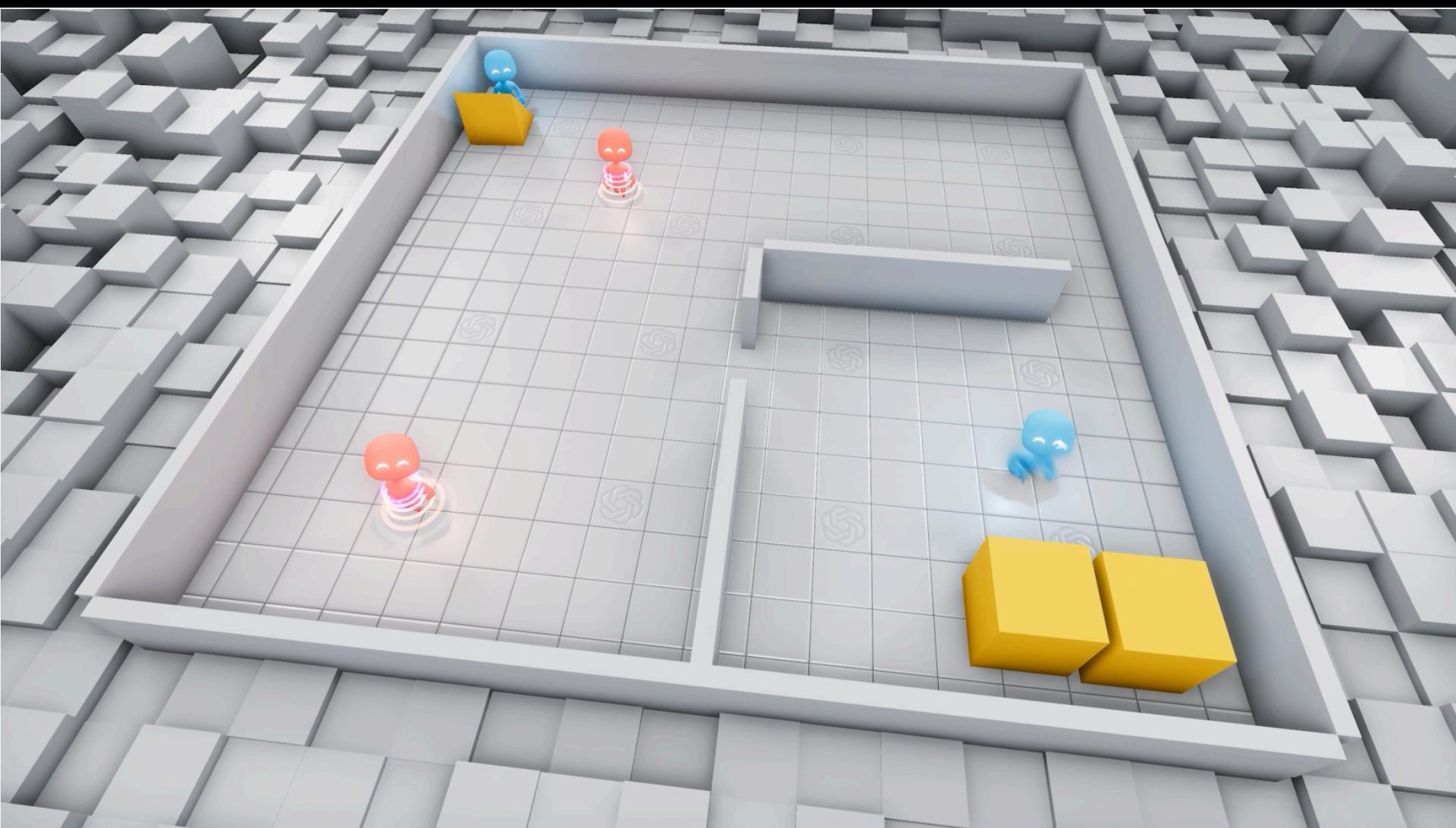
Use of blockers



Using the Ramp



Hiding the Ramp



Internet to remote places: Balloon stabilization



And more problems

Market Summary > Tesla Inc
NASDAQ: TSLA

457.92 USD +38.3
Sep 15, 3:19 PM EDT · Disclaim

1 day 5 days

CHIP COOLING
Heating power Today's chips dissipate 10 times the heat of a typical hotplate. For optimal operation, chips must be cooled below 85°C.
Source: IBM Zurich Research Laboratory

Processors should not reach more than 85°C.

+ Follow

1. MICRO CHANNELS
High performance micro-channel coolers are attached directly to the backside of the processor. In the cooler, water is distributed by a network of very fine channels for efficient heat removal.

2. HEAT EXCHANGER
The heat removed from the data center is delivered to a second circuit.

3. DIRECT REUSE OF WASTE HEAT
The heat removed from the data center can directly be repurposed for a second usage, e.g. for heating of buildings.

How does this all work?

Reinforcement Learning

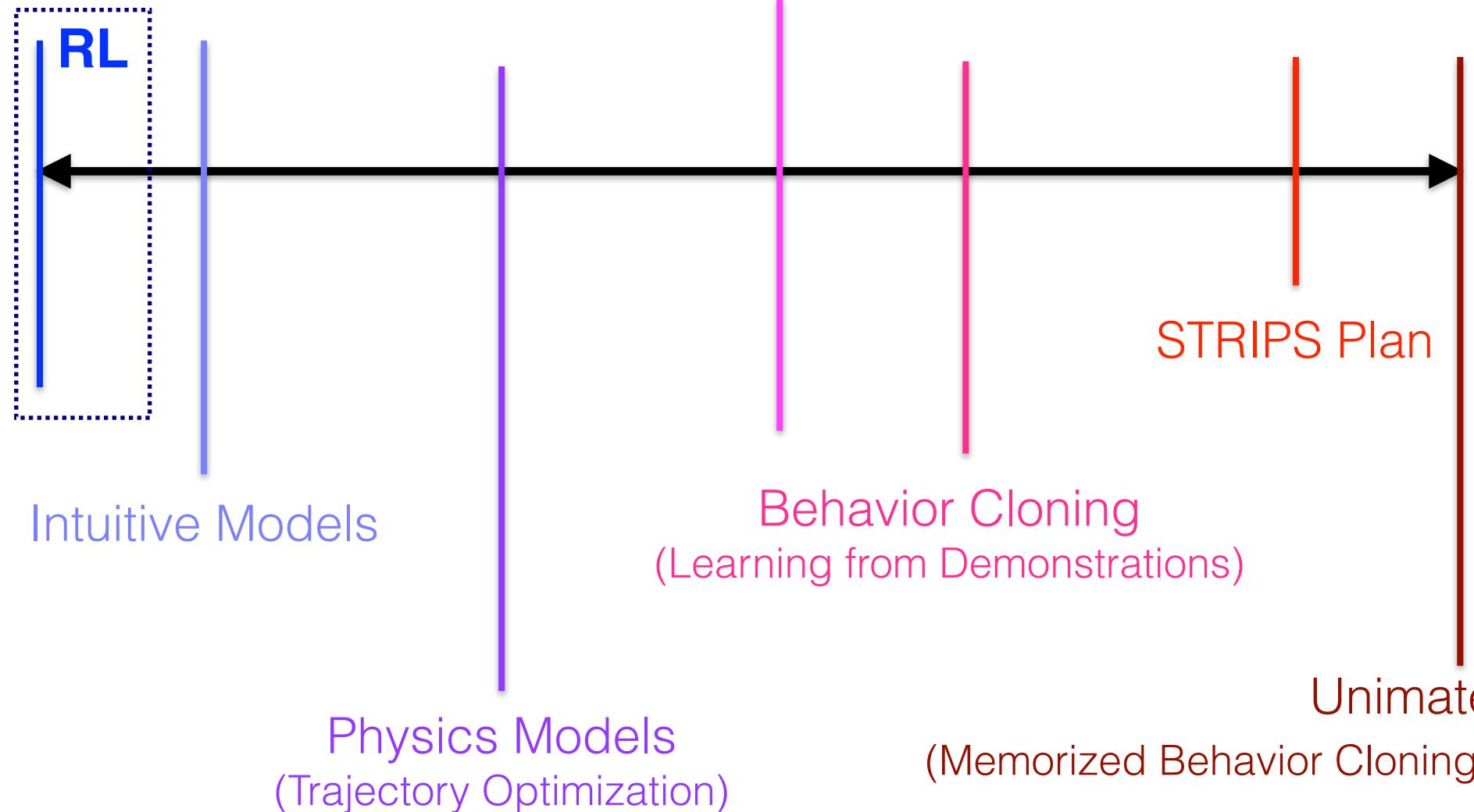
Bandits

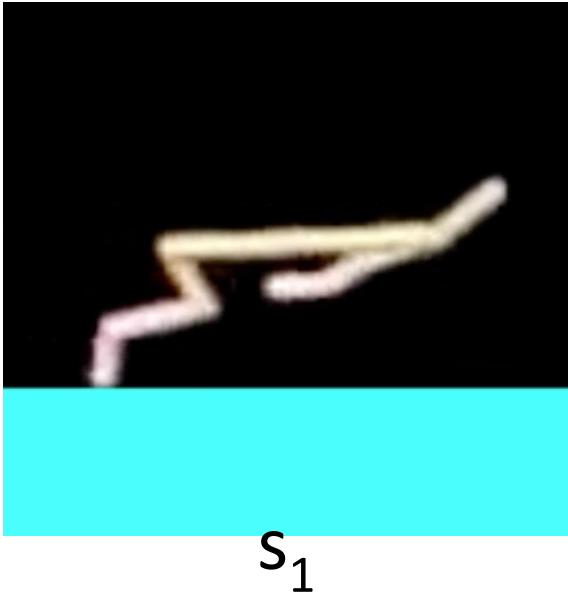
Contextual Bandits

Self-Learnt
Behavior

Hard-Coded
Behavior

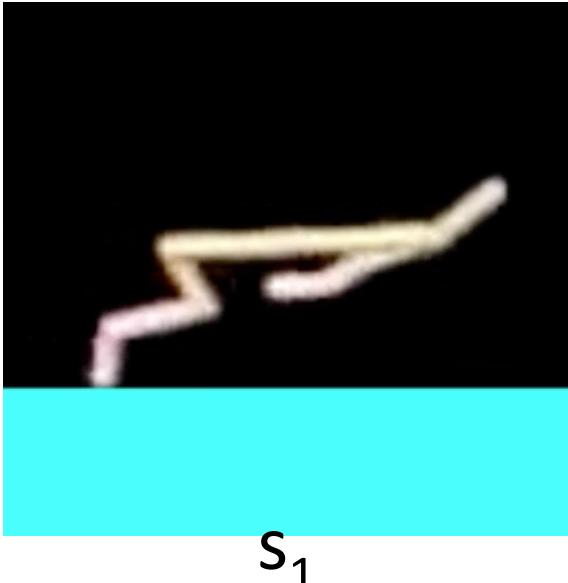
Imitation Learning
(Learning by Observing)





State ($s_1, s_2 \dots$)

- Location/rotation
of joints
- Or, the image
- Or, both



State ($s_1, s_2 \dots$)

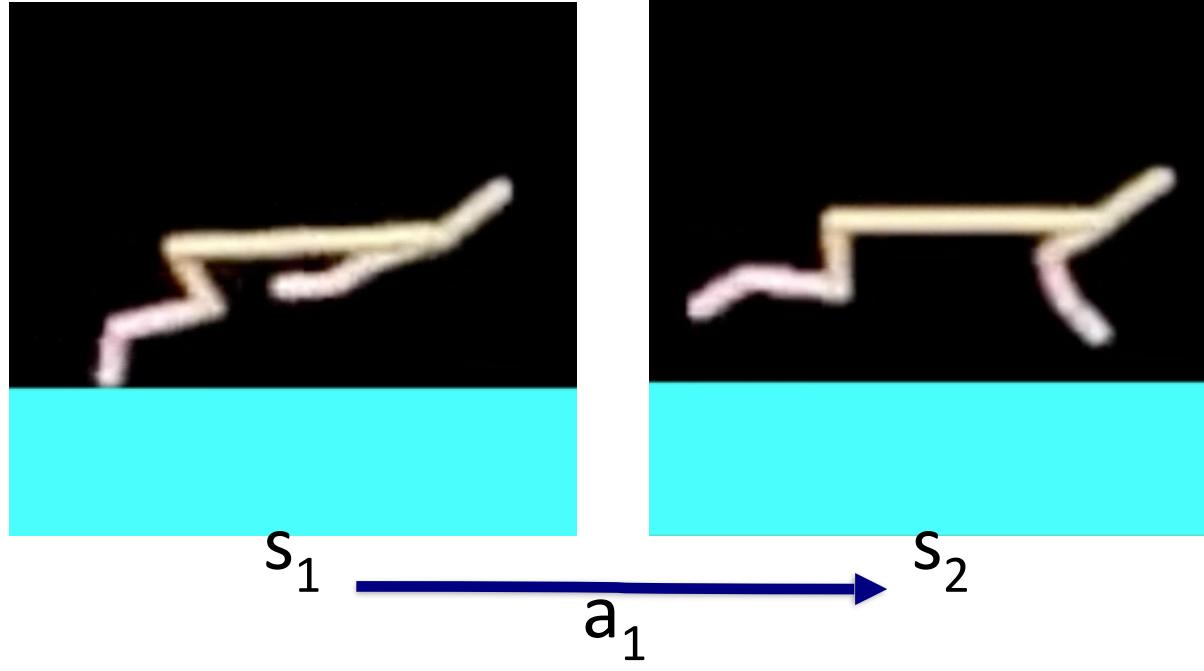
- Location/rotation
of joints

- Or, the image

- Or, both

Action ($a_1, a_2 \dots$)

Torques on each joint

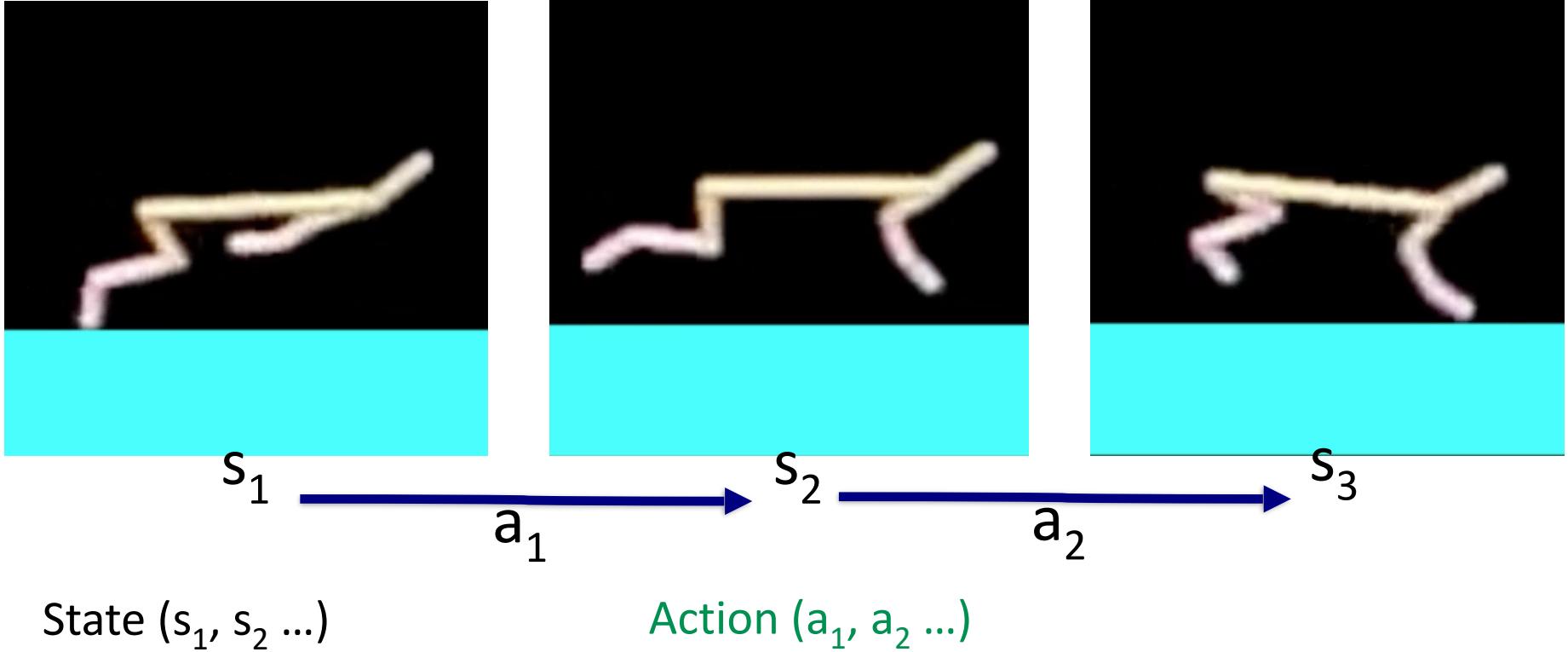


State ($s_1, s_2 \dots$)

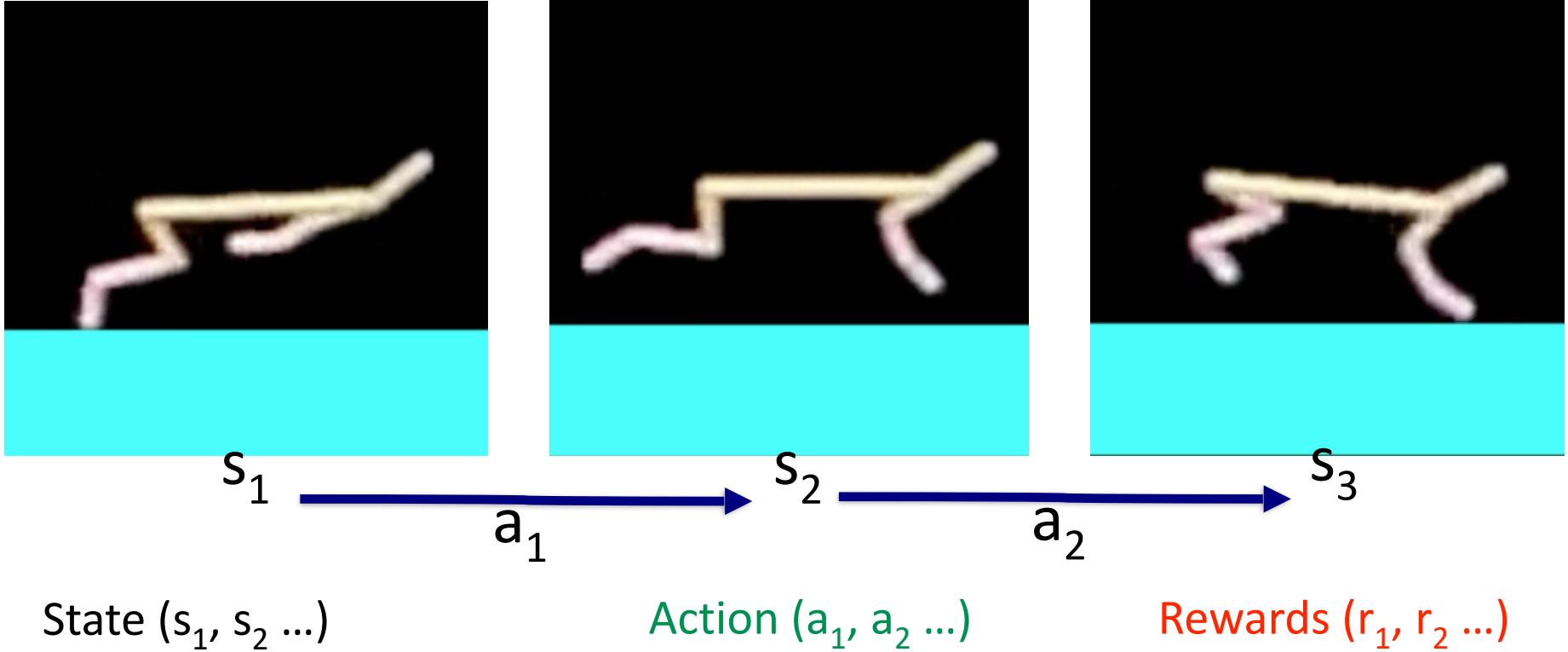
Action ($a_1, a_2 \dots$)

- Location/rotation
of joints
 - Or, the image
 - Or, both

Torques on each joint

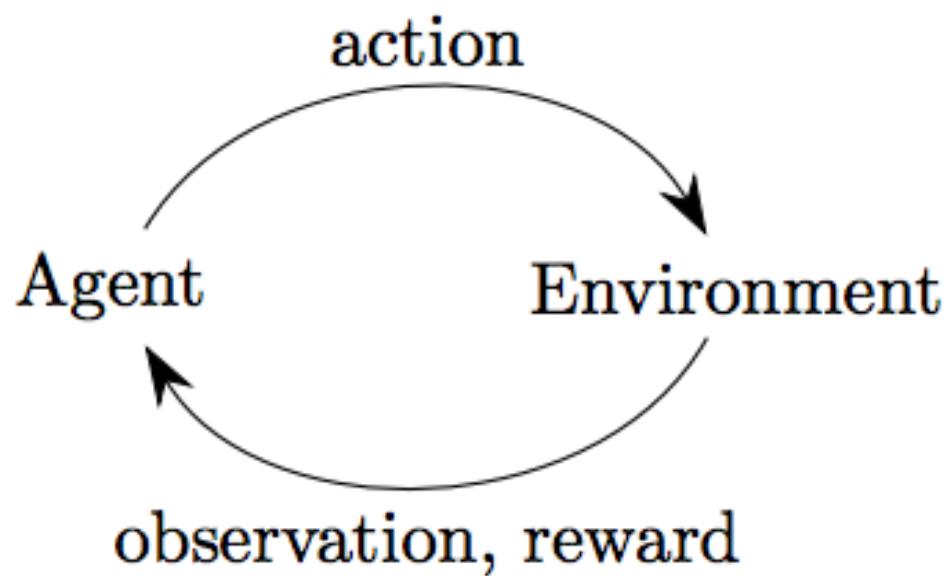


- Location/rotation
of joints Torques on each joint
- Or, the image
- Or, both



- Location/rotation of joints
 - Torques on each joint
 - Speed of the Cheetah
- Or, the image
- Or, both

Problem Formulation



Problem Formulation

Do Actions: $a_1, a_2, a_3, \dots, a_T$

Problem Formulation

Do Actions: $a_1, a_2, a_3, \dots, a_T$

Get Rewards: $r_1, r_2, r_3, \dots, r_T$

Problem Formulation

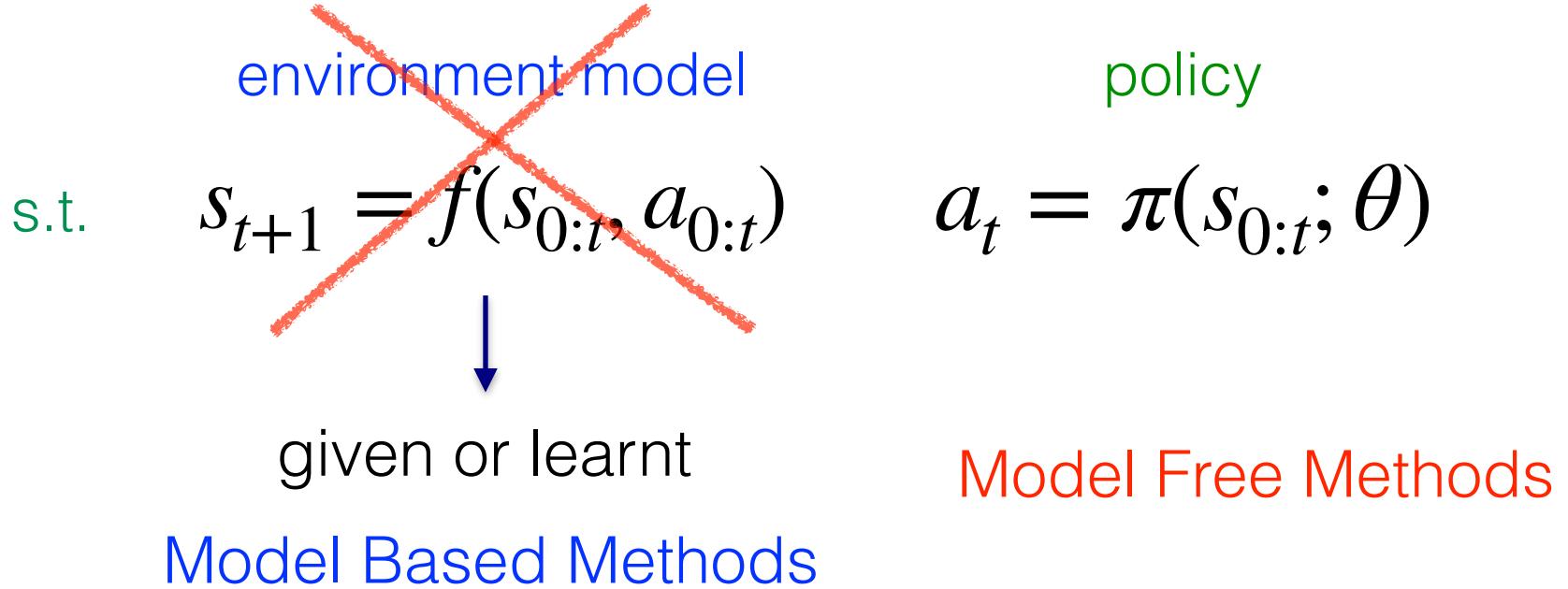
Do Actions: $a_1, a_2, a_3, \dots, a_T$

Get Rewards: $r_1, r_2, r_3, \dots, r_T$

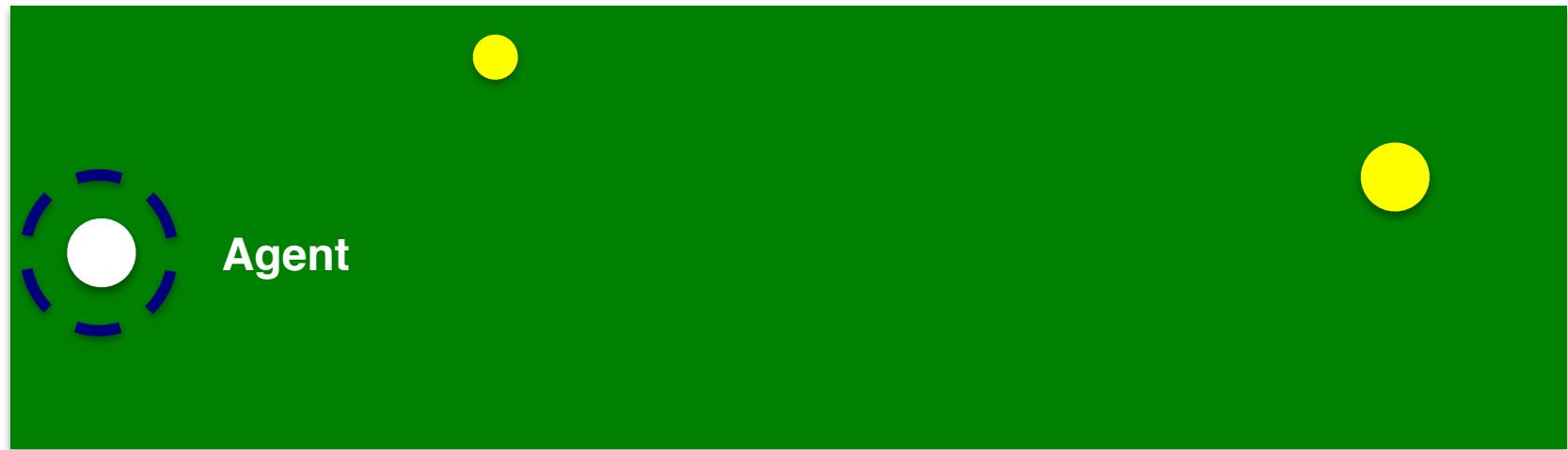
$$\max \sum_{t=1}^T r_t$$

Problem Formulation

$$\max \sum_{t=1}^T r_t$$



The RL Problem



The RL Problem



The RL Problem

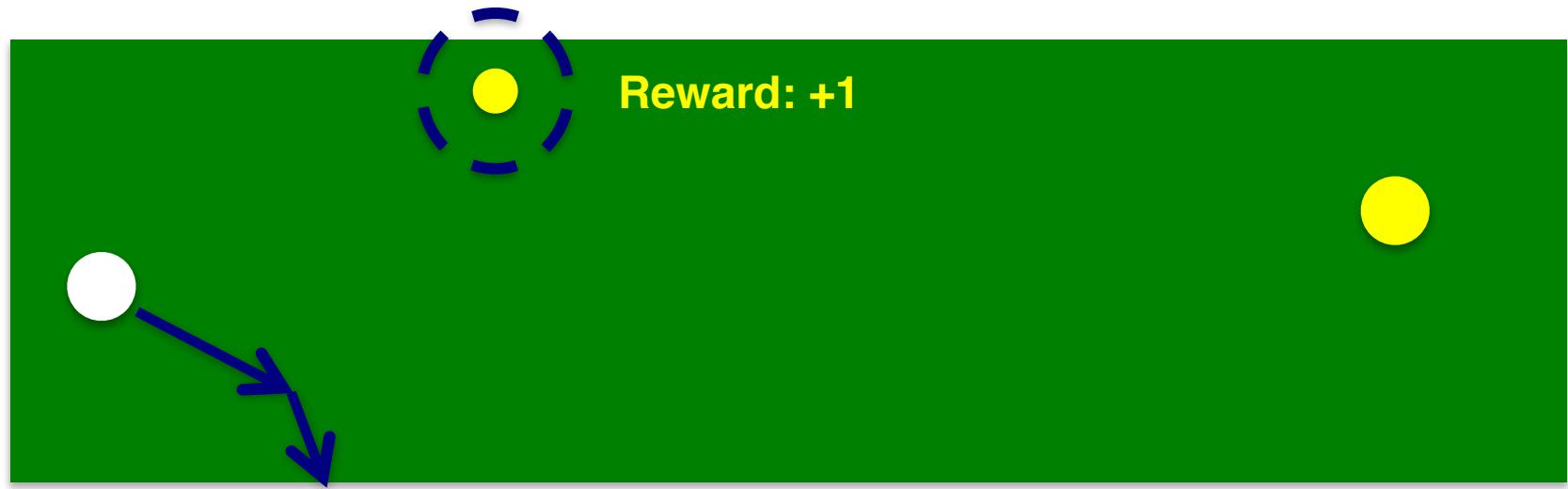


The RL Problem



Reward: -1

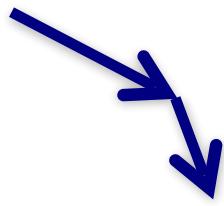
Solving the MDP



The RL Problem



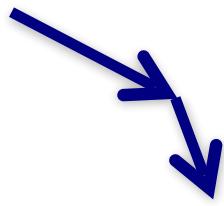
Another Attempt



Reward: -1

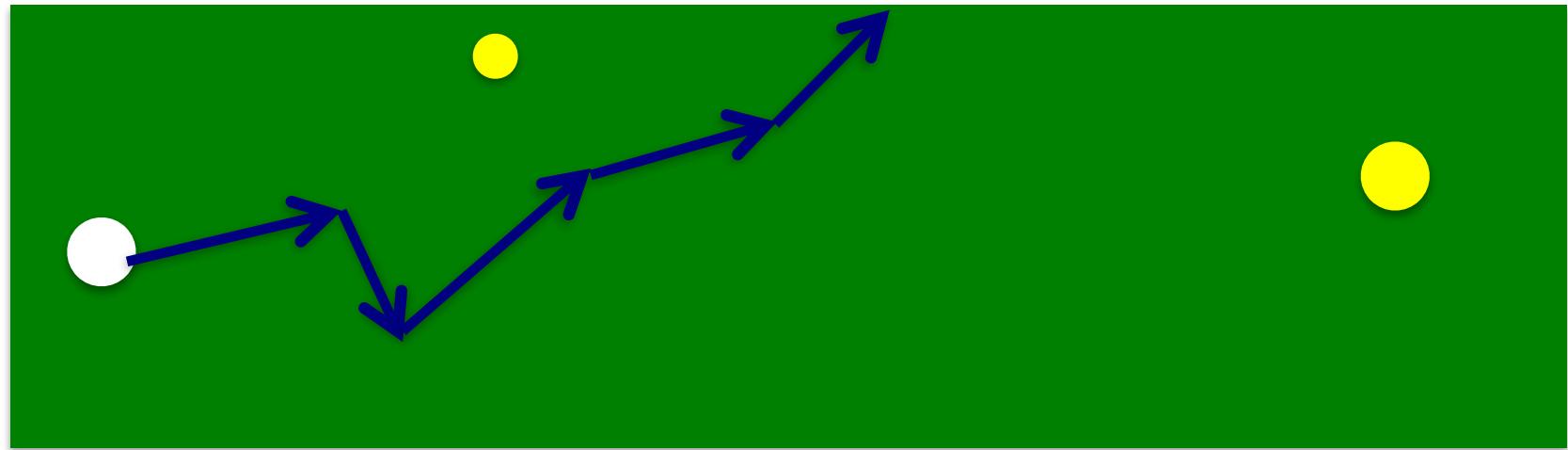


Another Attempt

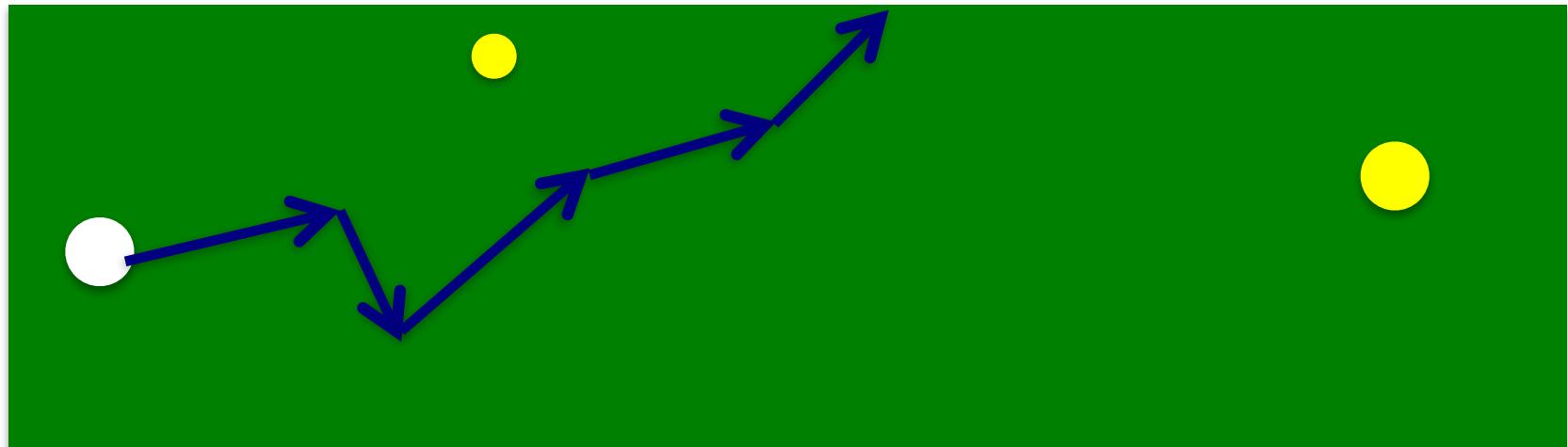
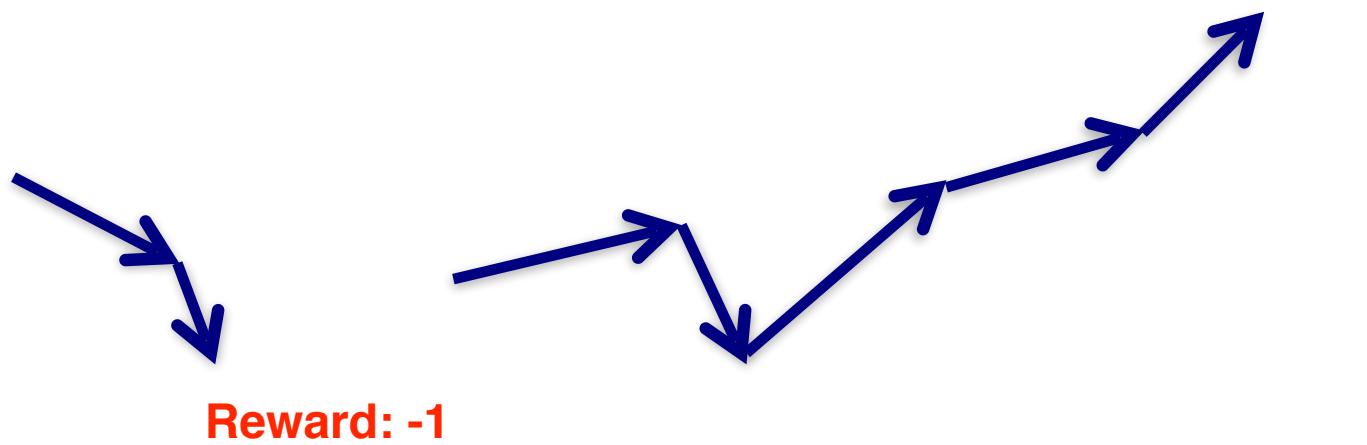


Reward: -1

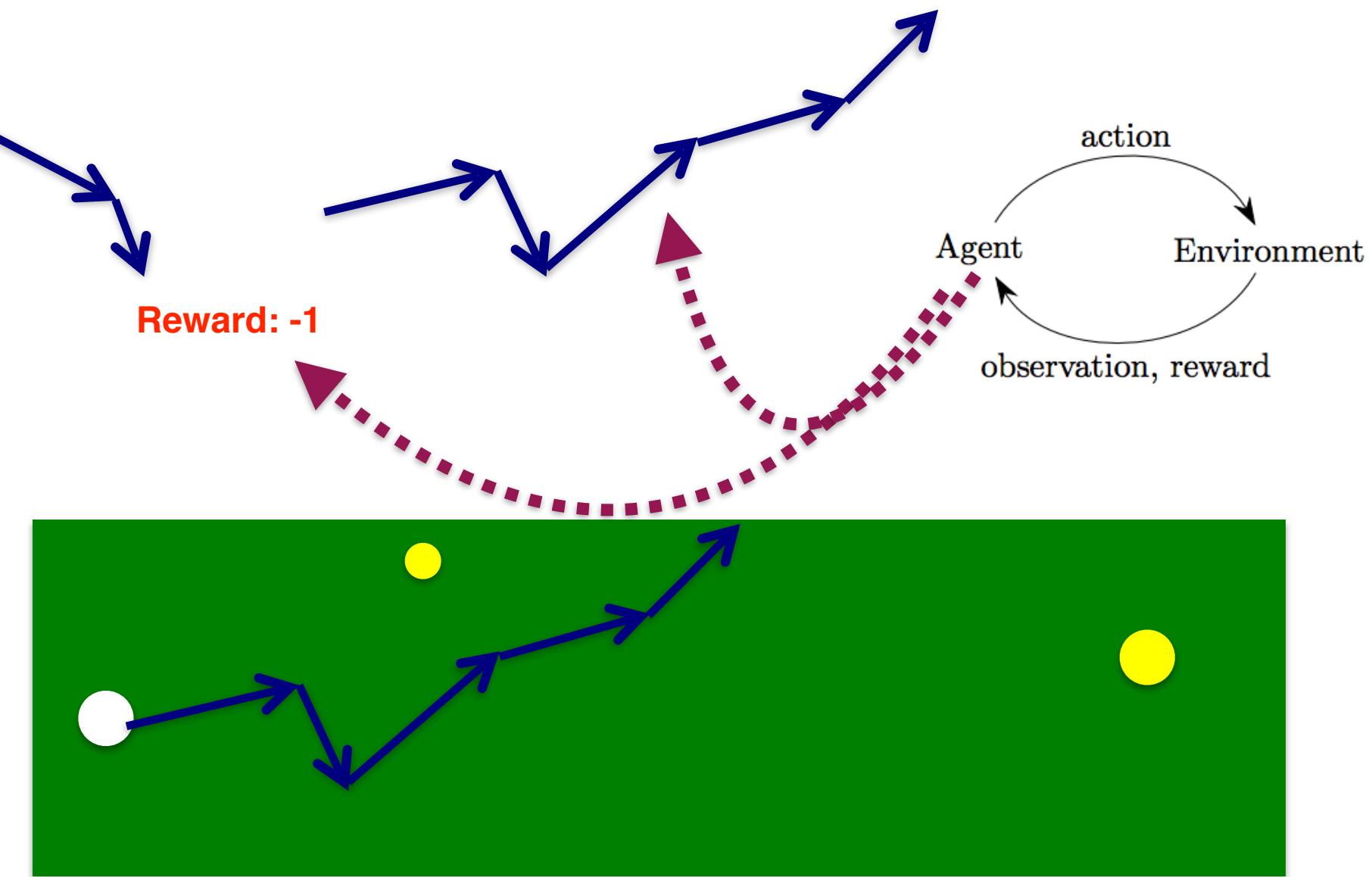
Reward: -1



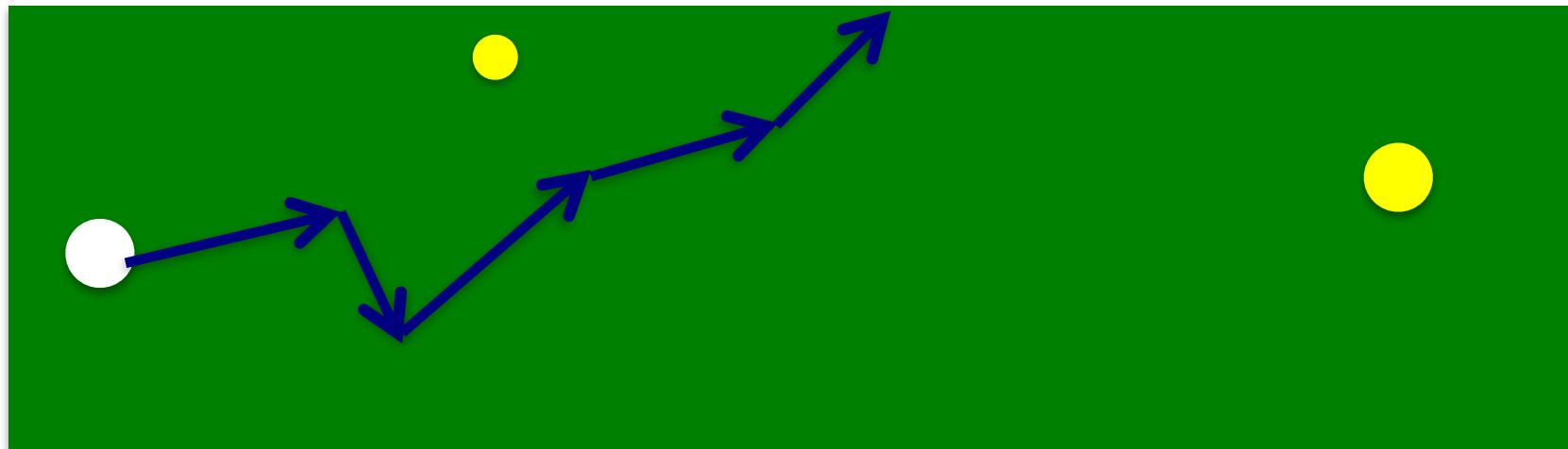
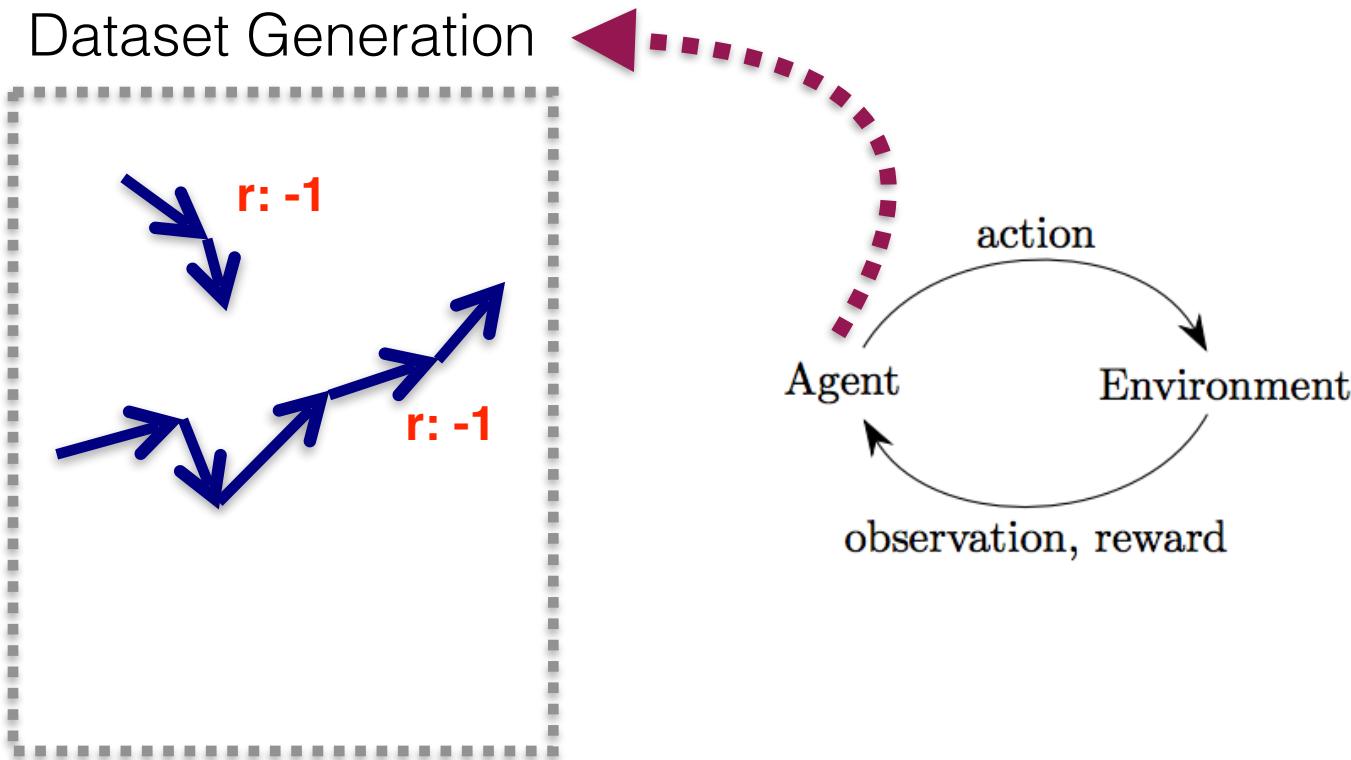
Another Attempt



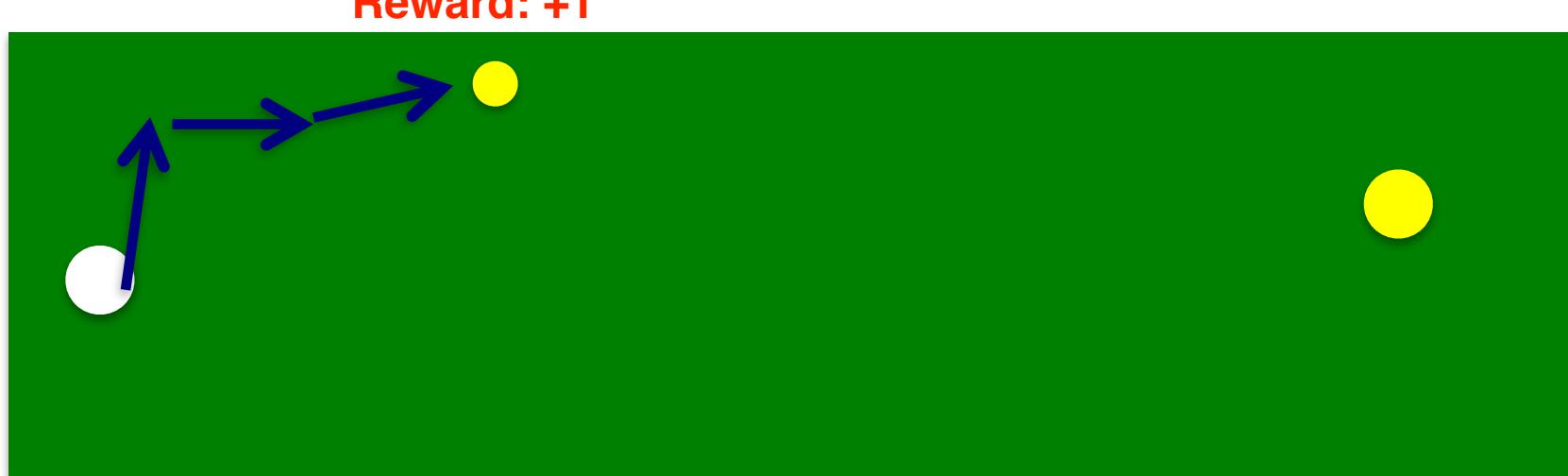
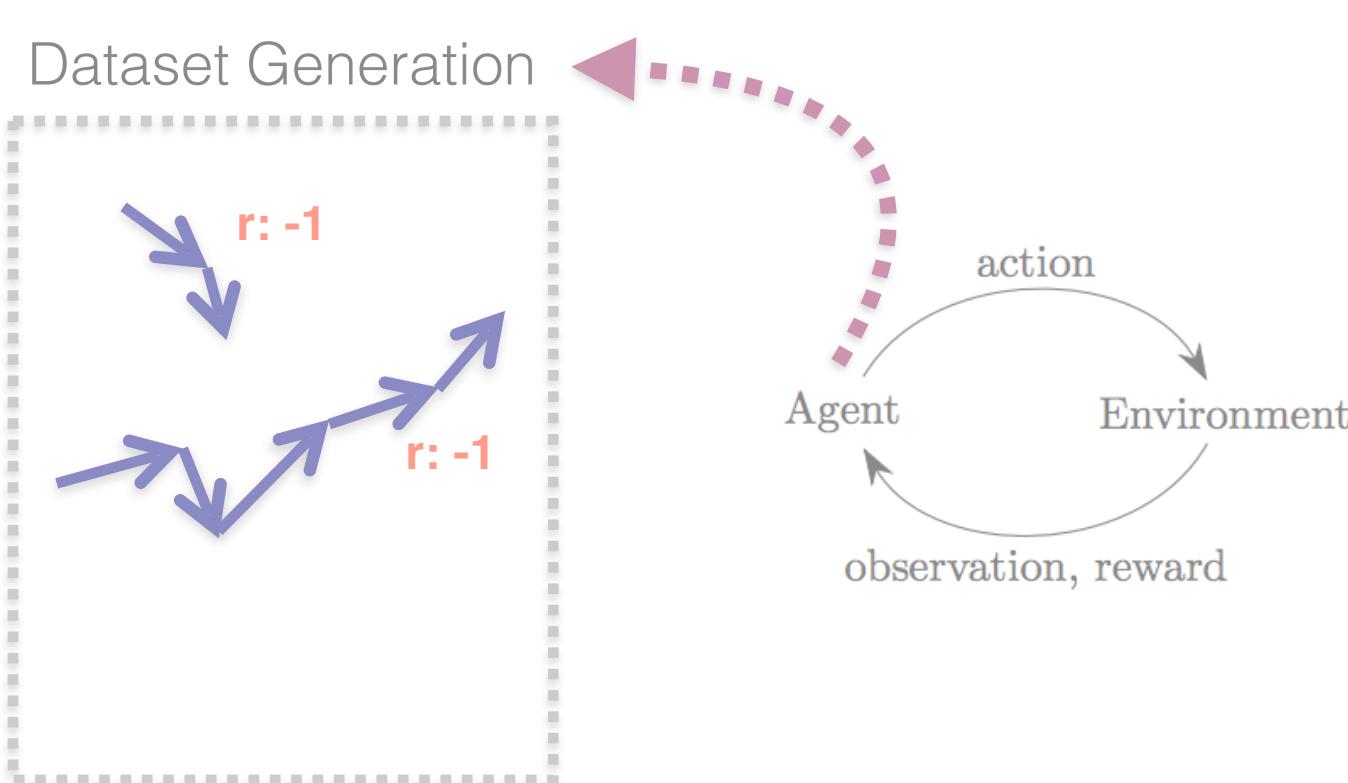
Another Attempt



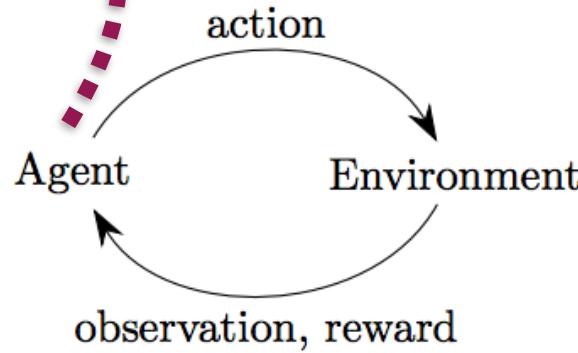
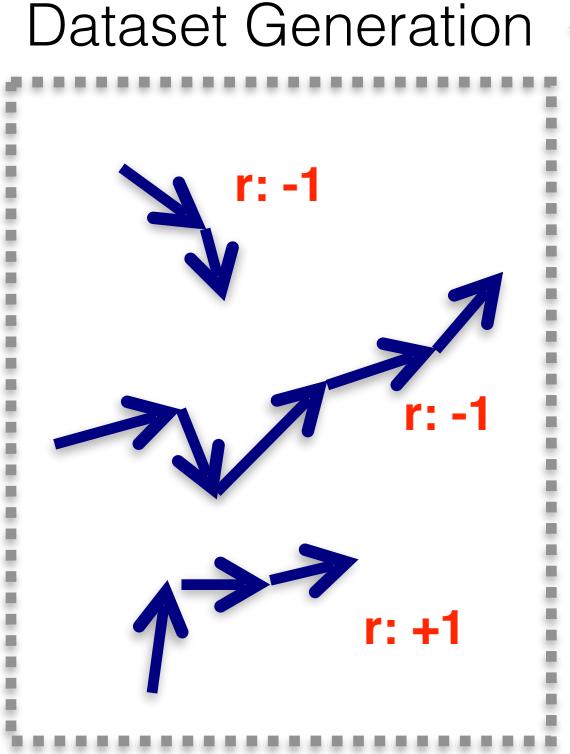
Dataset Generation



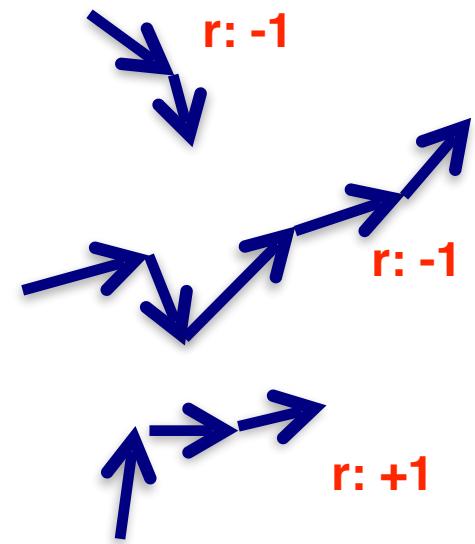
Dataset Generation



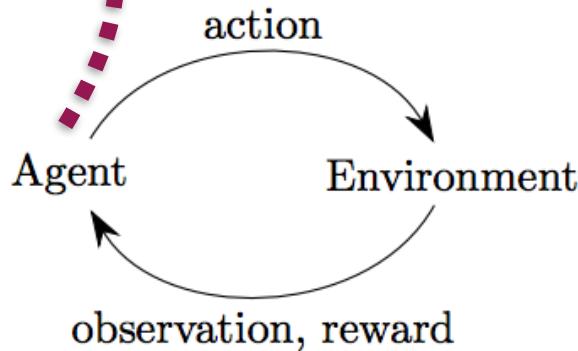
Dataset Generation



Dataset Generation

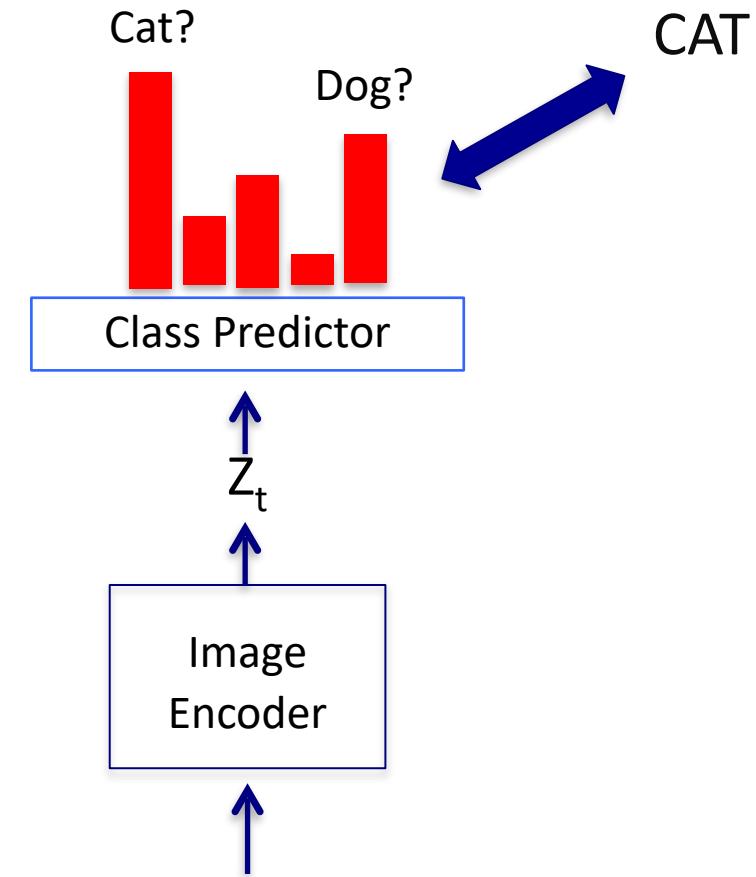
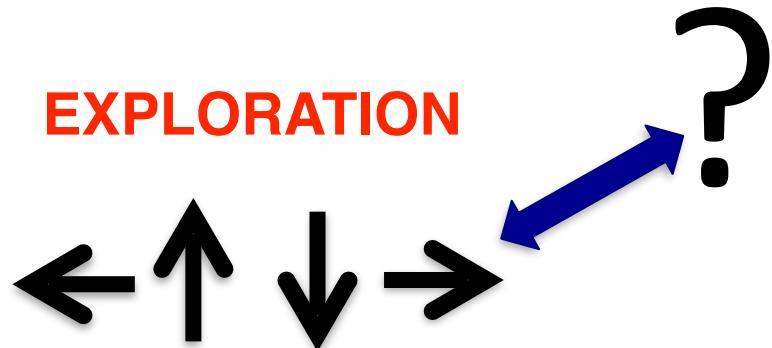


Supervised Learning



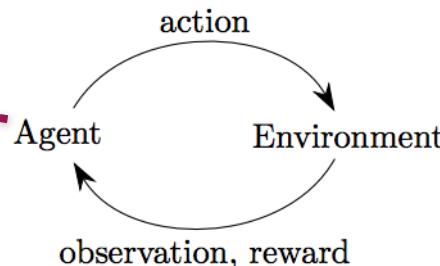
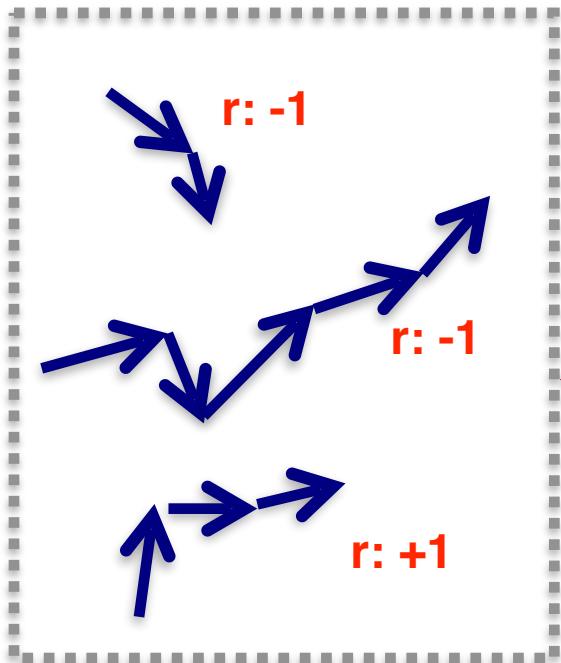
In supervised learning, dataset is GIVEN

In reinforcement learning, the agent collects its own data
(i.e, it needs to explore)

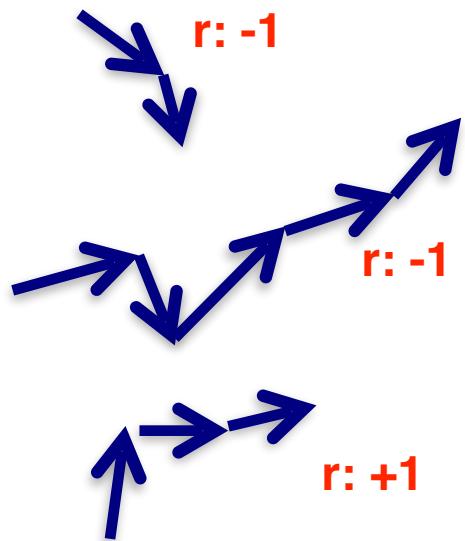


Dataset

Is exploration a problem?

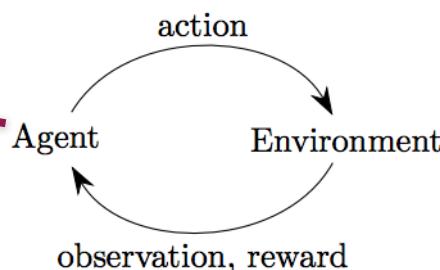


Dataset



Is exploration a problem?

Learn

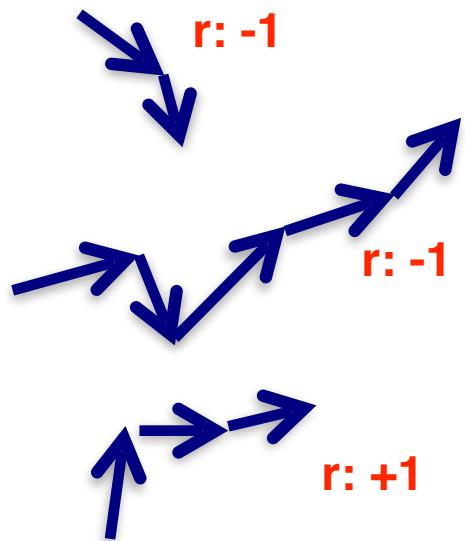


Goal

$$a_t = \pi(s_{0:t}; \theta)$$

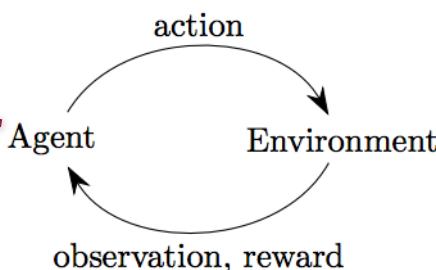


Dataset



Is exploration a problem?

Learn



Goal

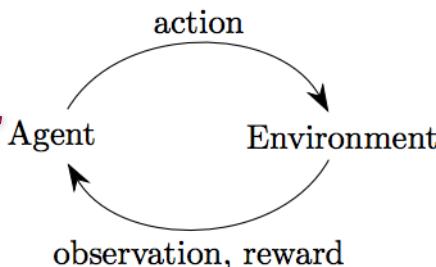
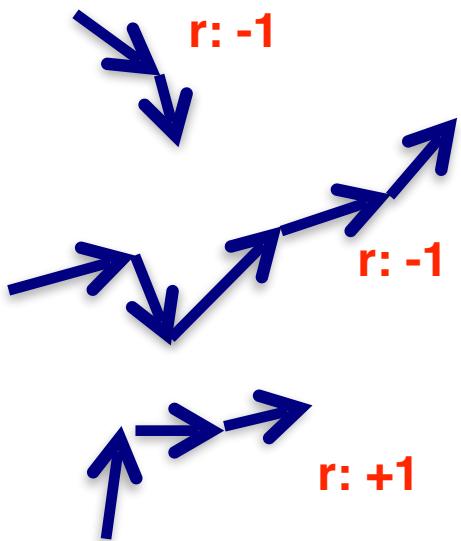
$$a_t = \pi(s_{0:t}; \theta)$$

Looks good!
(is there a problem?)



Dataset

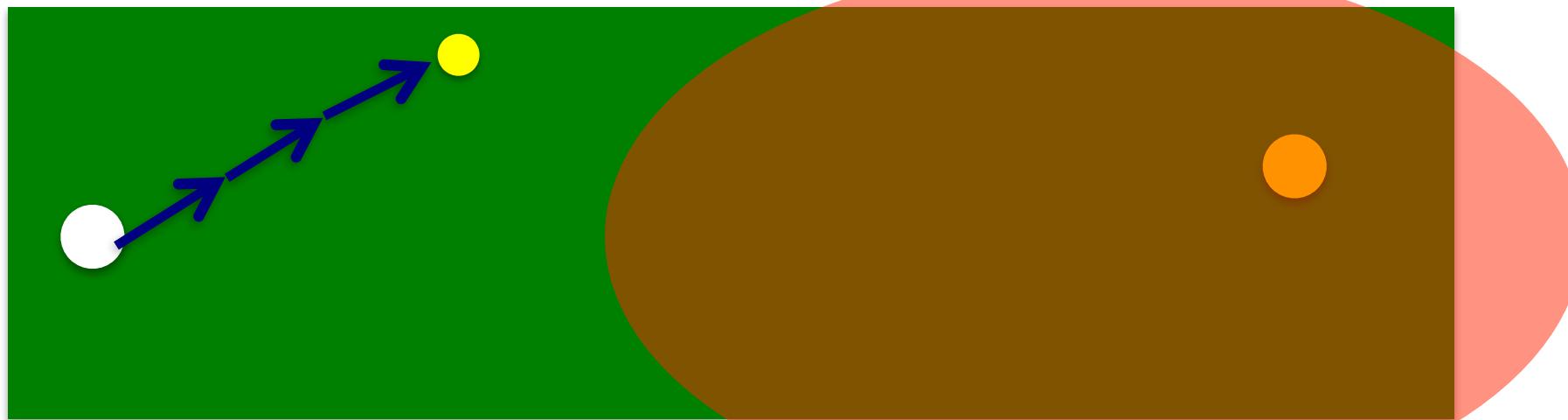
Is exploration a problem?



Goal

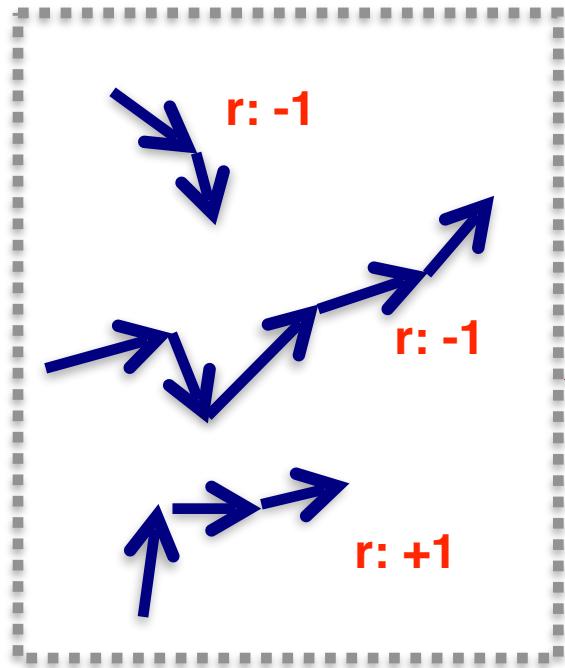
$$a_t = \pi(s_{0:t}; \theta)$$

Looks good!
(is there a problem?)



Might not explore the state space!

Dataset

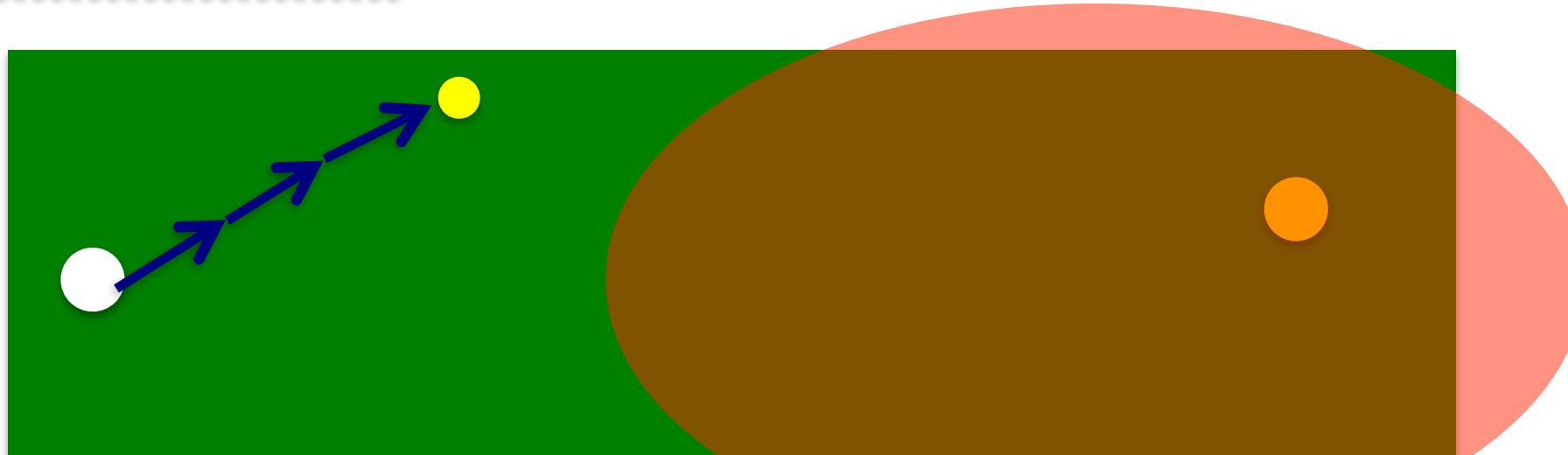
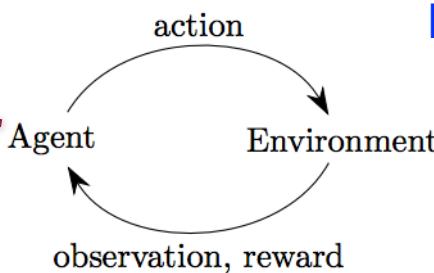


Is exploration a problem?

Yes!

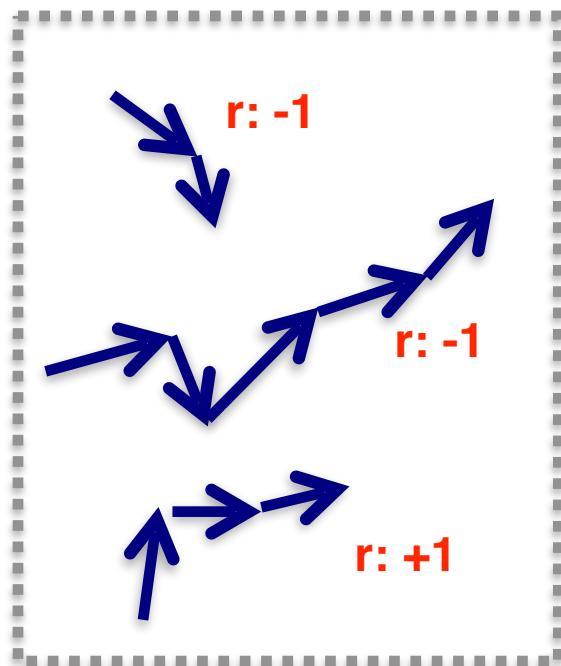
(Can learn sub-optimal behavior)

Exploration-Exploitation Dilemma



Might not explore the state space!

Dataset



Is exploration a problem?

Yes!

(Can learn sub-optimal behavior)

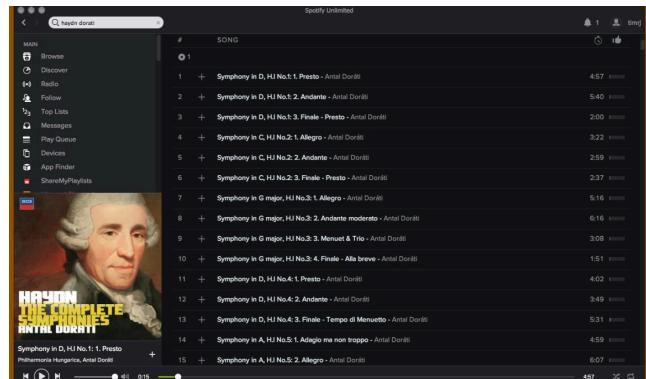
Exploration-Exploitation Dilemma

Exploration can be quite hard!

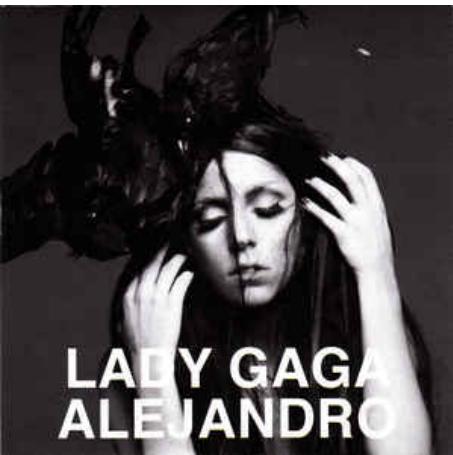
When?



Imagine your favorite playlist



(they want you hooked)



Explore by
Suggesting other music

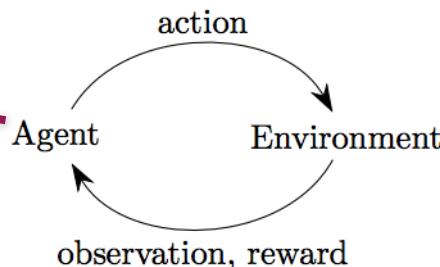
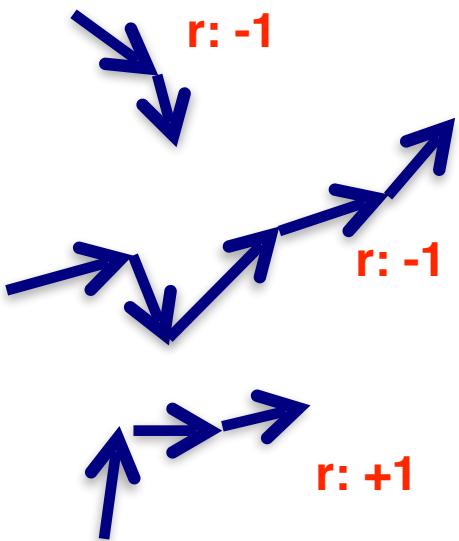
Imagine your favorite playlist

Sometimes Exploration can be very costly!



Dataset

What questions might be of interest?



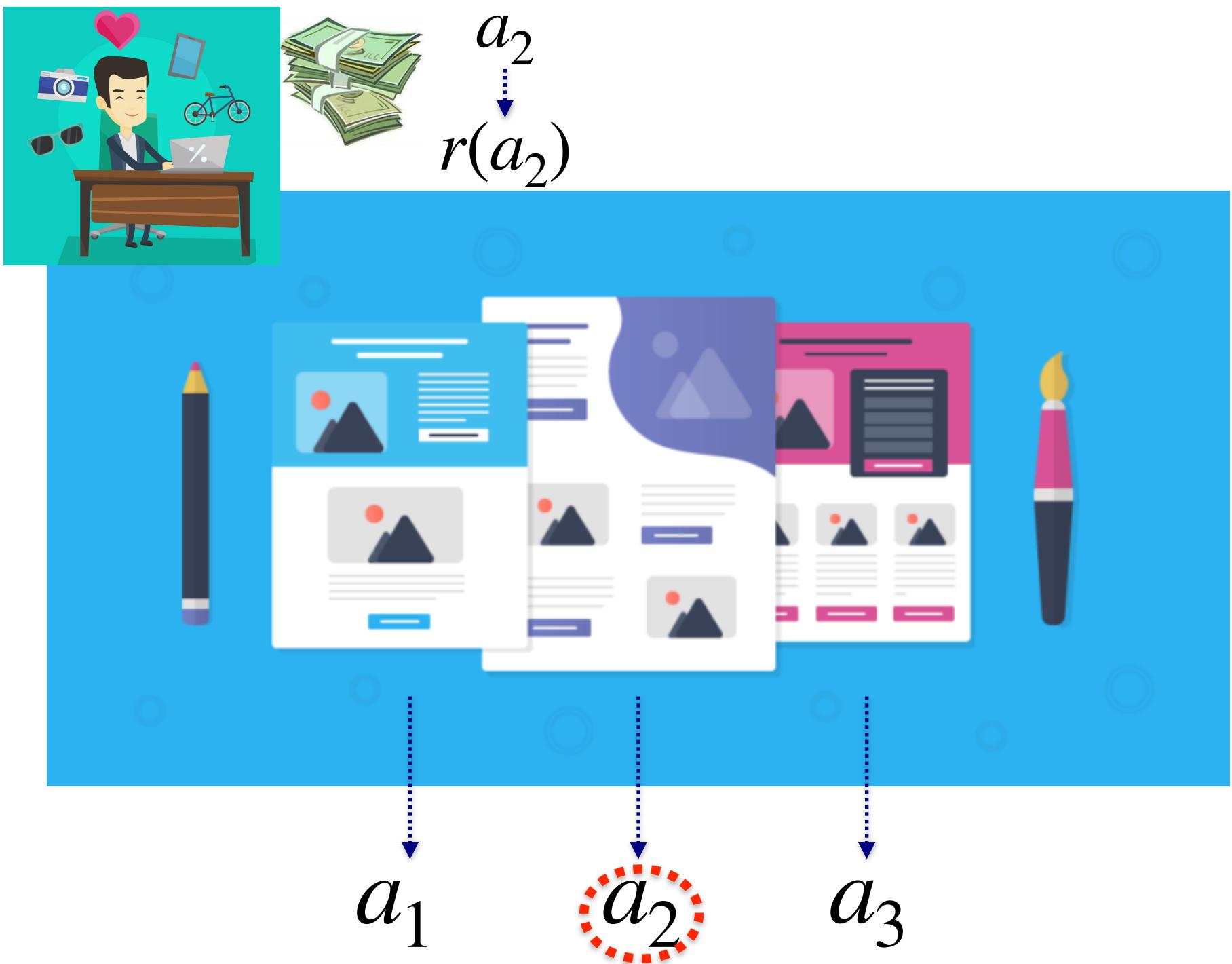
Is there a method that will achieve the highest reward?

How fast will reach it?
(i.e., what is the overall regret)



Consider landing pages for your website

 a_1 a_2 a_3





a_2 a_3

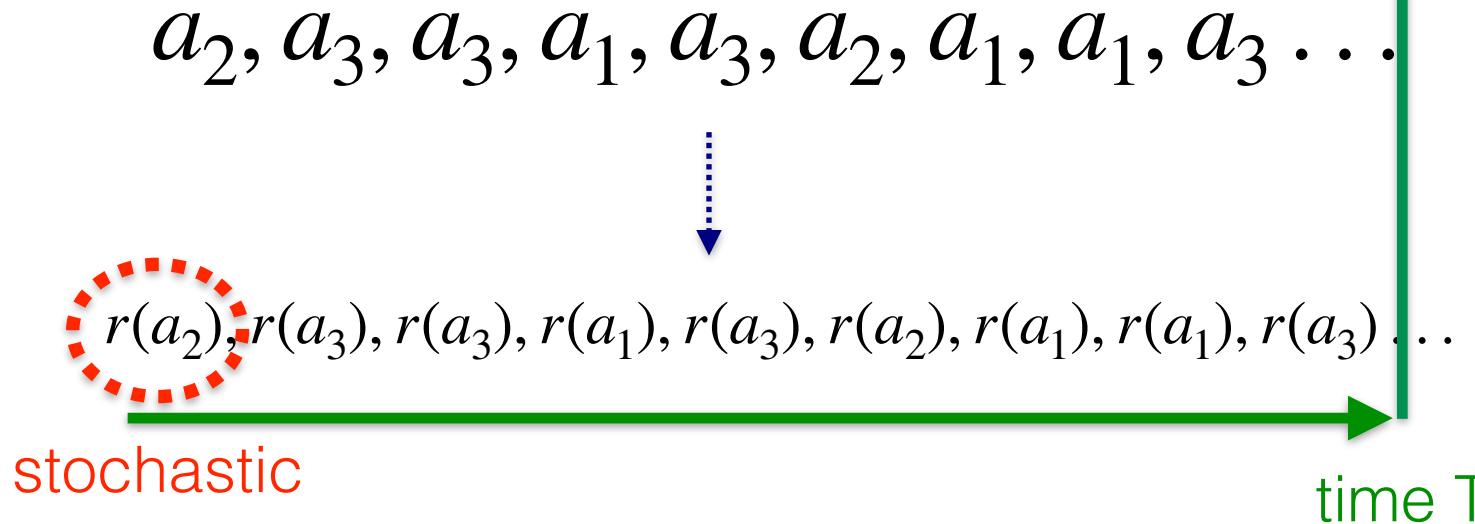
$r(a_2)$ $r(a_3)$



a_1 a_2

a_3

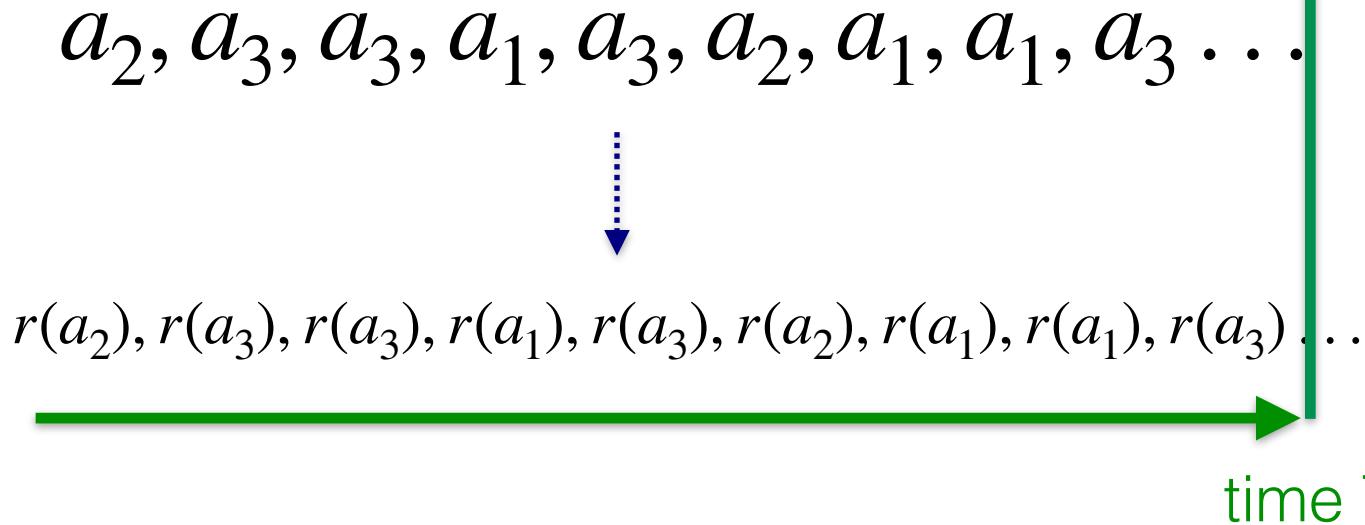
Multi-Arm Bandit Problem



$$\sum_{t=1}^T r(a_i^t) \quad i \in [1, N]$$



Multi-Arm Bandit Problem



Goal: maximize
reward over time

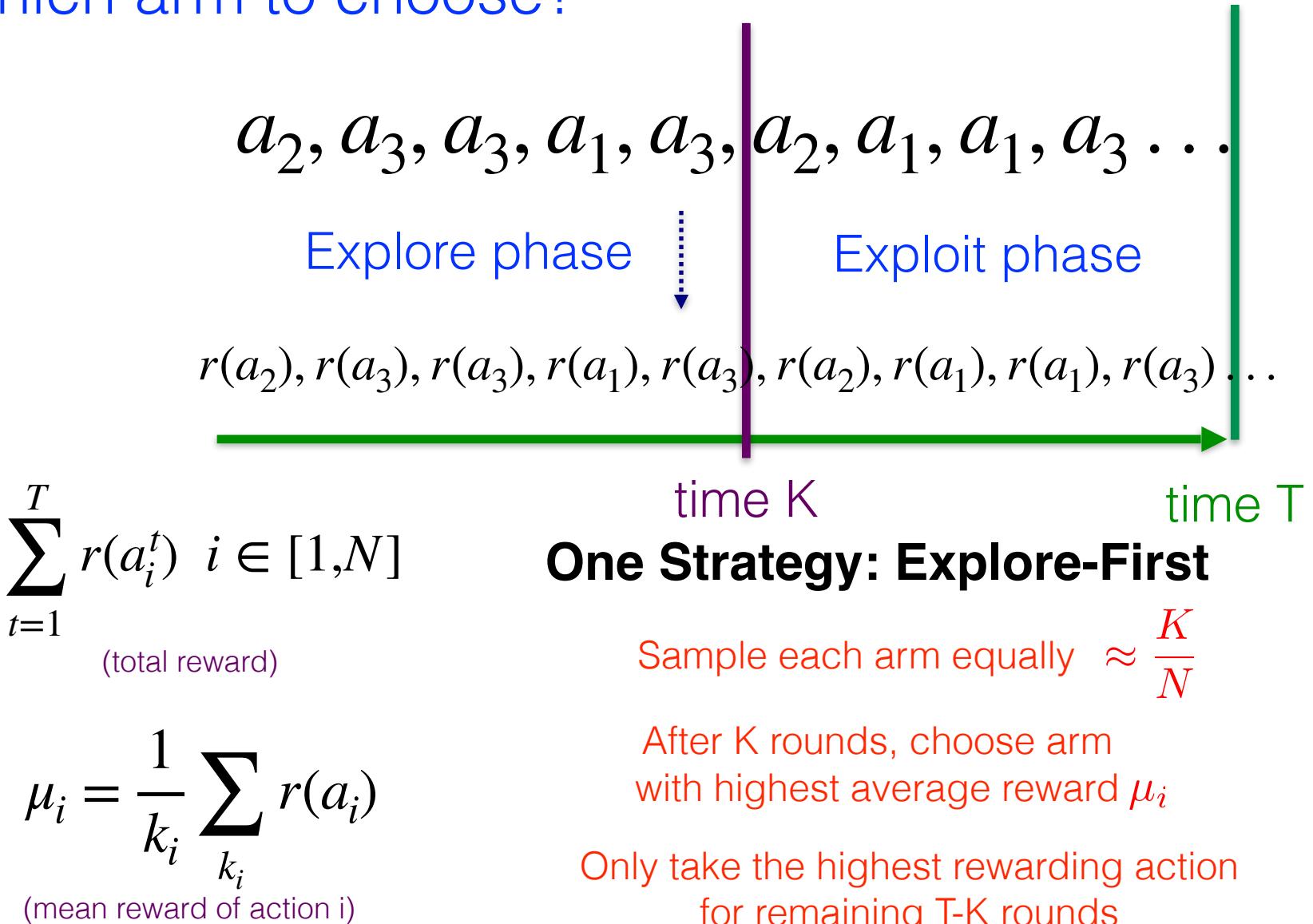
$$\sum_{t=1}^T r(a_i^t) \quad i \in [1, N]$$

Return of i^{th} arm: $\mu_i = \frac{1}{k_i} \sum_{k_i} r(a_i)$

(i^{th} arm is pulled k_i times)

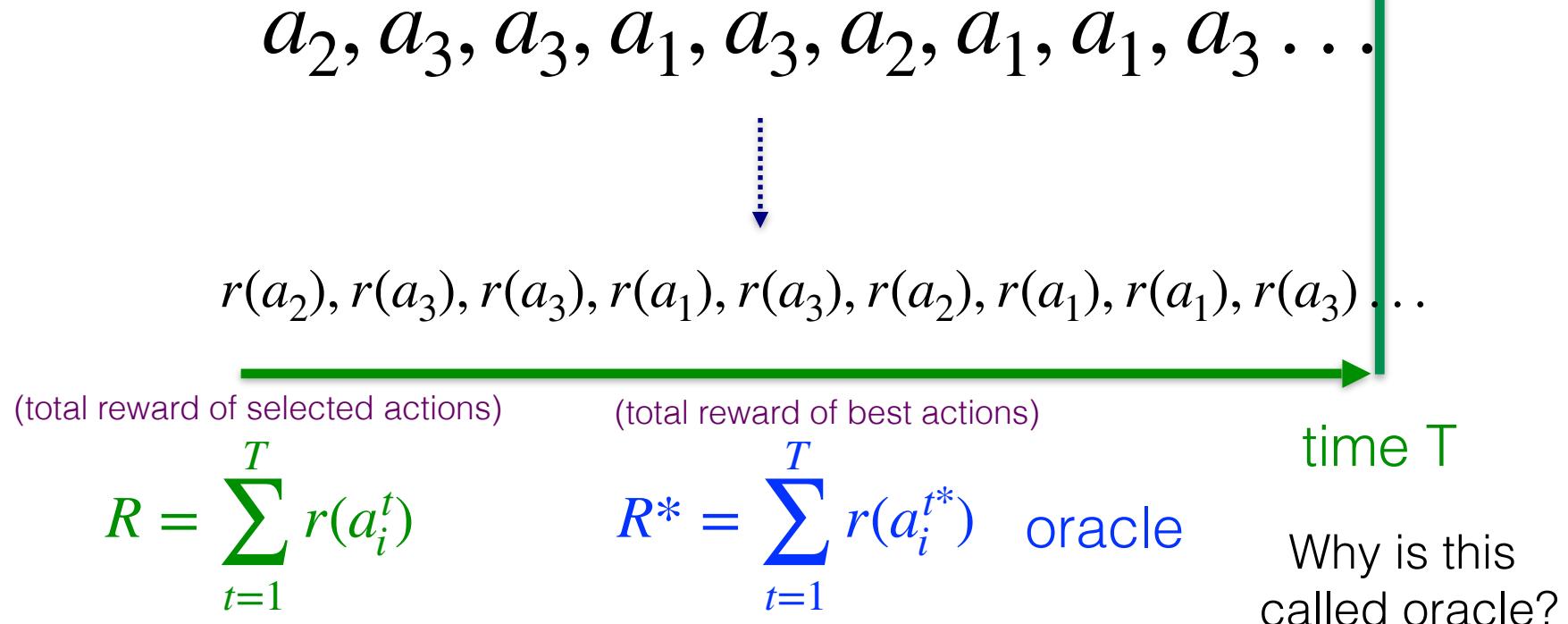
If rewards
are deterministic
can we find
the best solution?

Which arm to choose?



Is this the best we can do?

What do we mean by best?



regret $\|R^* - R\|$

As in life, goal is to minimize regret

Assume $r \in [0, 1]$

Worst that we can do: T

Explore-First: $T^{2/3} \times O(N \log T)^{1/3}$

Not just the asymptotic performance, but how fast we reach!

What do we mean by best?

$a_2, a_3, a_3, a_1, a_3, a_2, a_1, a_1, a_3 \dots$

↓

$r(a_2), r(a_3), r(a_3), r(a_1), r(a_3), r(a_2), r(a_1), r(a_1), r(a_3) \dots$

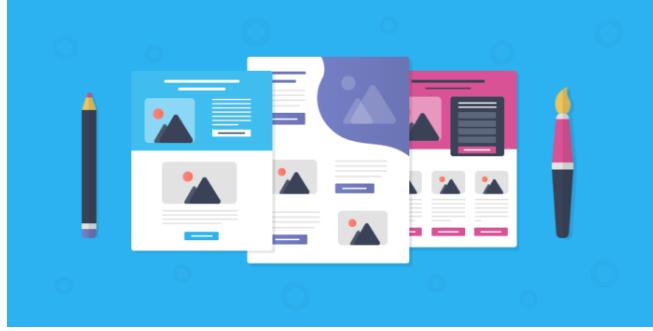
Does there exist an optimal algorithm?

$$R = \sum_{t=1}^T r(a_i^t)$$
$$R^* = \sum_{t=1}^T r(a_i^{t*}) \text{ oracle}$$

time T

regret $\|R^* - R\|$

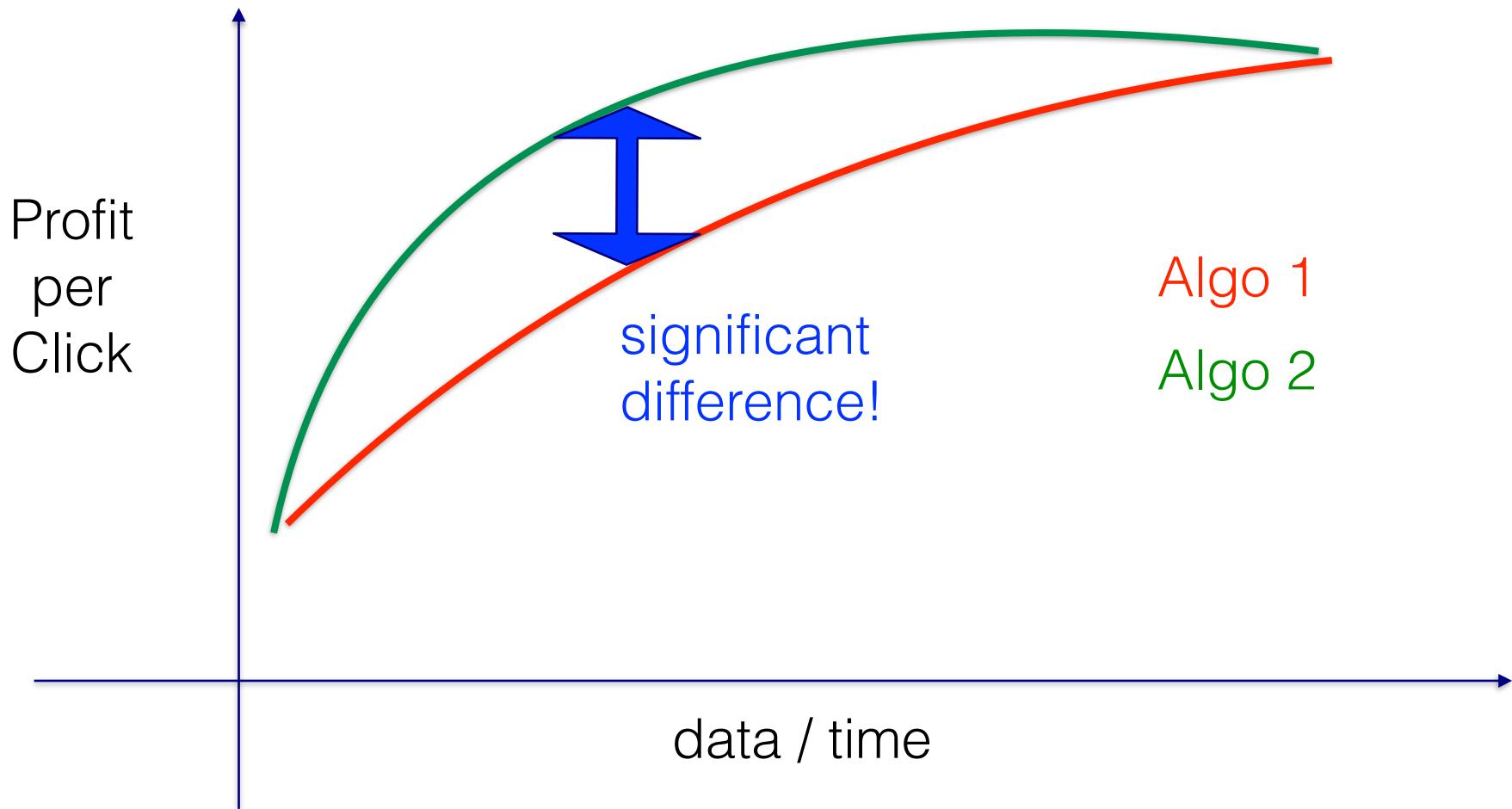
Not just the asymptotic performance, but how fast we reach!



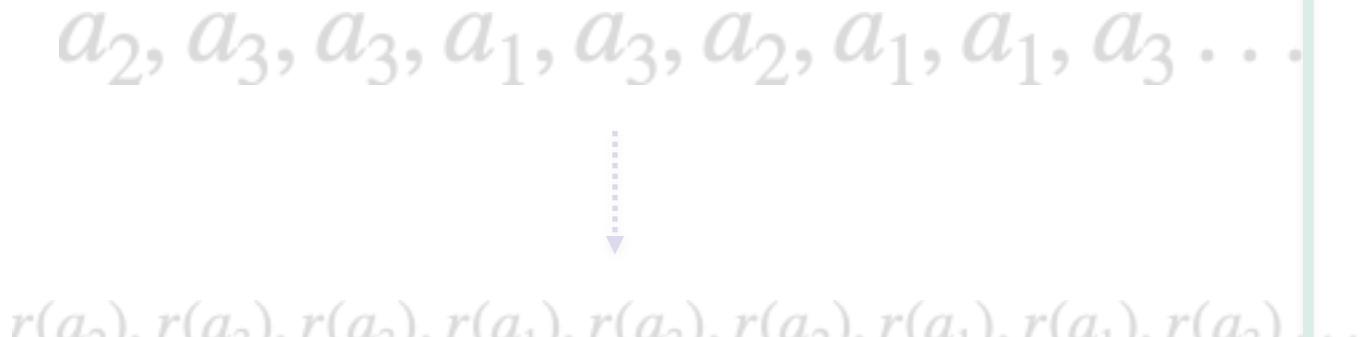
Online Decisions!

What's the difference
from supervised learning?

(i.e., can't wait to collect data first)



What do we mean by best?



Does there exist an optimal algorithm?

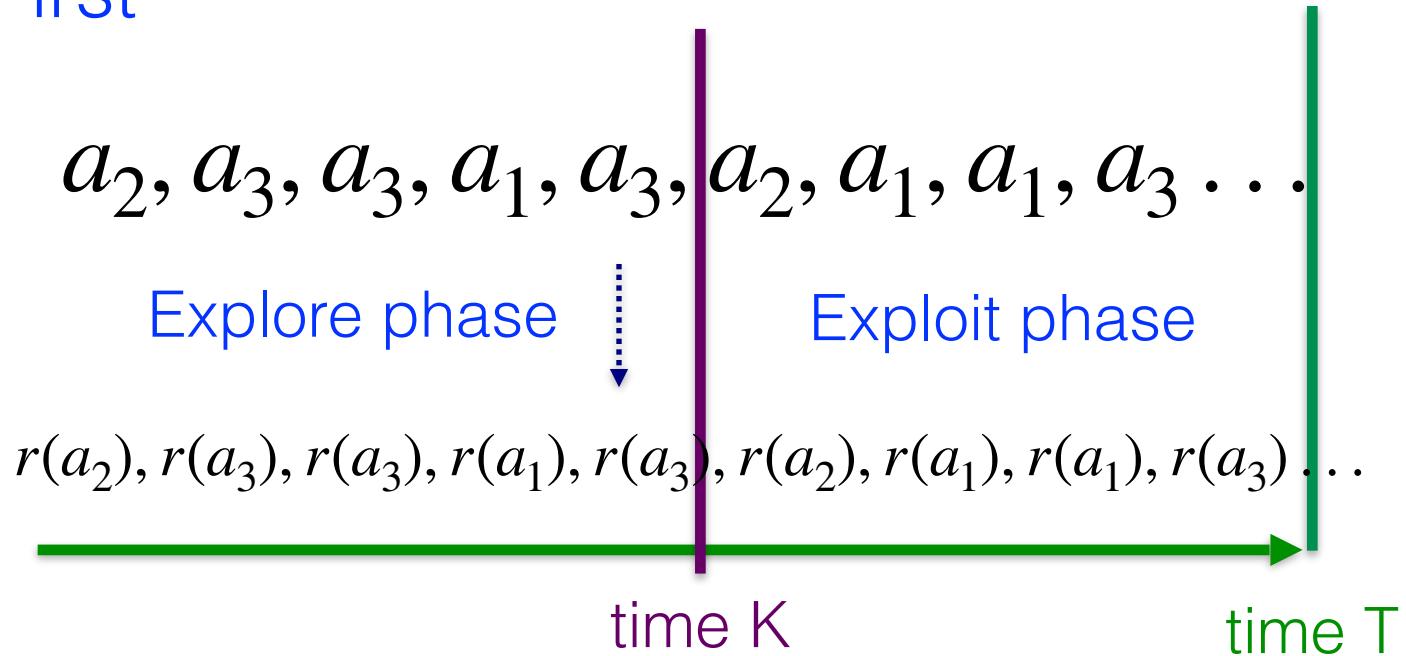
Upper Confidence Bound (UCB) Algorithm

$$R = \sum_{t=1}^T r(a_i^t)$$
$$R^* = \sum_{t=1}^T r(a_i^{t*}) \text{ oracle}$$

regret $\|R^* - R\|$

Not just the asymptotic performance, but how fast we reach!

Explore-First



Non-adaptive exploration

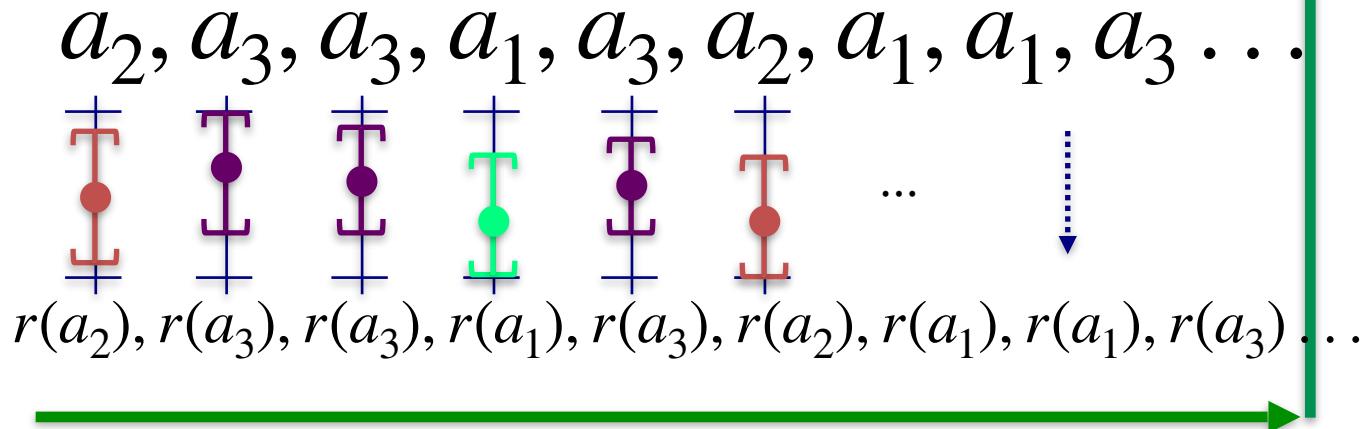
Explore + exploit separately

vs

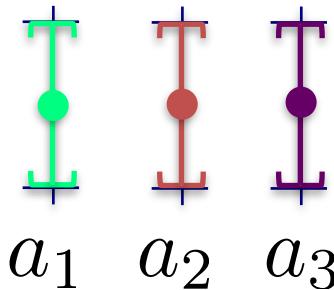
Adaptive exploration

Explore + Exploit simultaneously

Upper Confidence Bound (UCB) Algorithm



Initial confidence intervals:



$$\mu_i = \frac{1}{k_i} \sum_{k_i} r(a_i)$$

(mean reward of action i)

Optimism in face of uncertainty

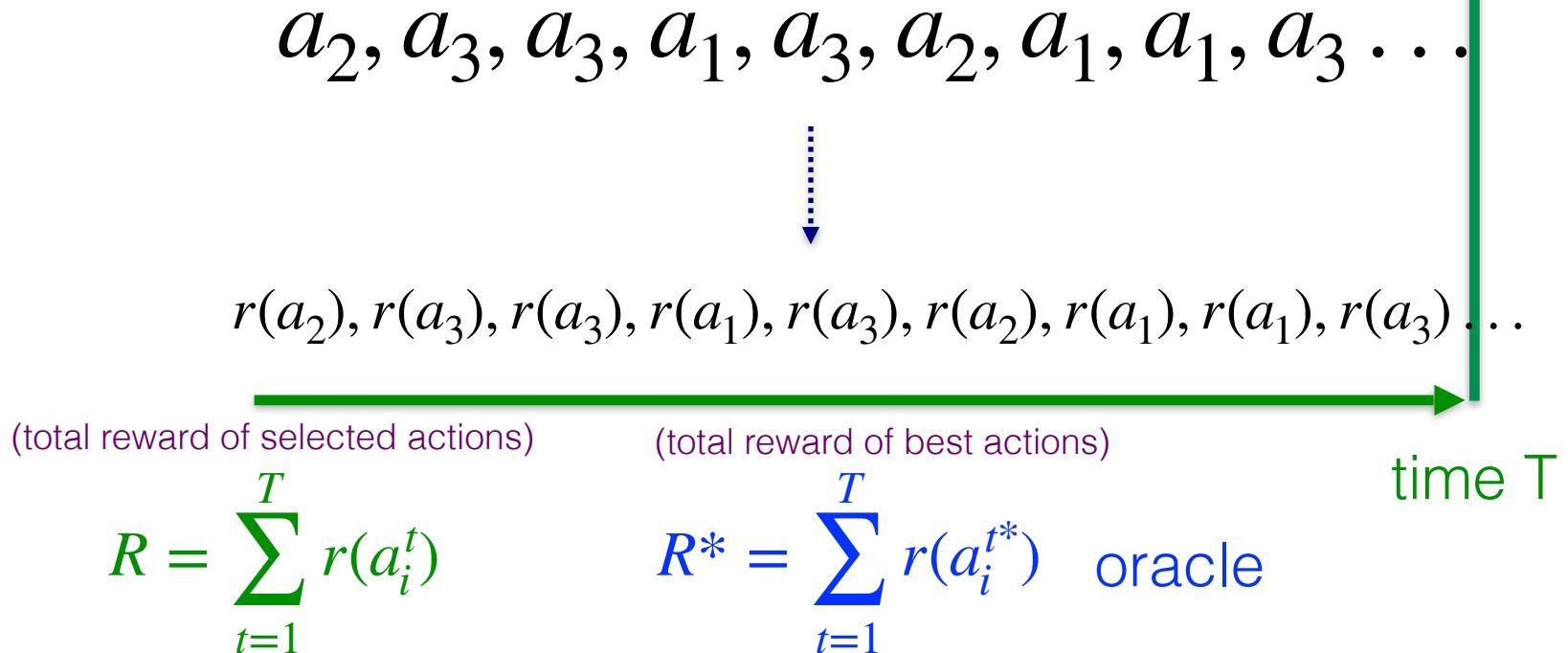
Exploitation

$$a_{t+1} = \arg \max_i \mu_i(t) +$$

$$\sqrt{\frac{4 \log t}{k_i}}$$

Exploration bonus
for rare actions
(optimism)

How good is UCB?



regret $\|R^* - R\|$

As in life, goal is to
minimize regret

Optimal! (up to log factors)

Assume $r \in [0, 1]$

Worst that we can do: T

Explore-First: $T^{2/3} \times (N \log T)^{1/3}$

UCB: $(NT \log T)^{1/2}$

Upper Confidence Bound (UCB) Algorithm

$$\arg \max_i \hat{\mu}_i(t) + \sqrt{\frac{4 \log t}{k_i}}$$

Where this come from?

We want

$$\arg \max_i \mu_i$$

We have

$$\arg \max_i \hat{\mu}_i$$

Construct

$$\arg \max_i \hat{\mu}'_i$$

Principle of optimism: find $\hat{\mu}'_i \geq \mu_i$

$$\hat{\mu}_i = \frac{1}{k_i} \sum_{k_i} r(a_i)$$

Empirical estimate of μ_i (unknown)

Initially (with few samples)
This estimate is going to be bad

$$p(\mu_i \geq \hat{\mu}'_i) \leq \delta$$

If $r(a_i)$ are 1-subgaussian and if

$$\hat{\mu}'_i : \hat{\mu}_i + \sqrt{\frac{2 \log \frac{1}{\delta}}{k_i}}$$

Is true!

Upper Confidence Bound (UCB) Algorithm

$$\arg \max_i \mu_i(t) + \sqrt{\frac{4 \log t}{k_i}} \quad \mu_i = \frac{1}{k_i} \sum_{k_i} r(a_i)$$

Upper Bound on average number of
sub-optimal actions

$$\frac{16 |A| \log T}{\Delta^2} + O(1)$$

$$\Delta = \mu_{best} - \mu_{second_best}$$

$|A|$: number of actions

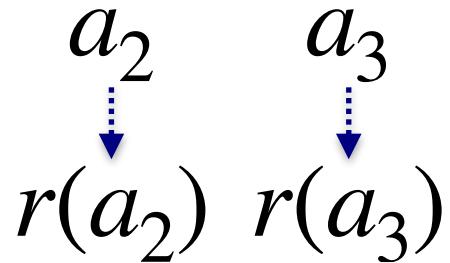

$$a_2 \quad a_3$$

$r(a_2) \quad r(a_3)$


$$a_1 \quad a_2 \quad a_3$$



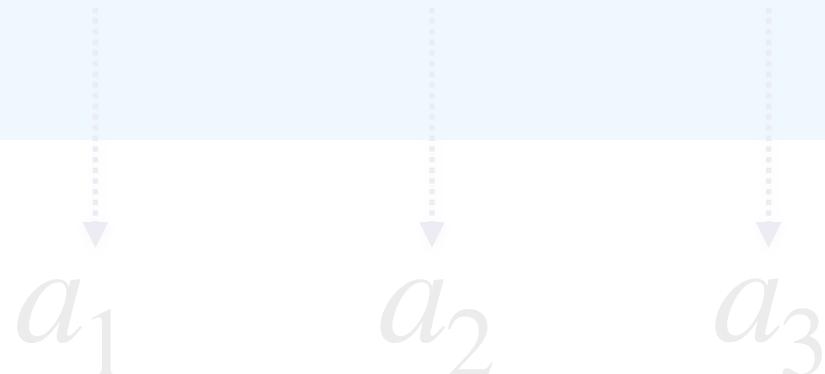
(male, 30s, computer-savvy)

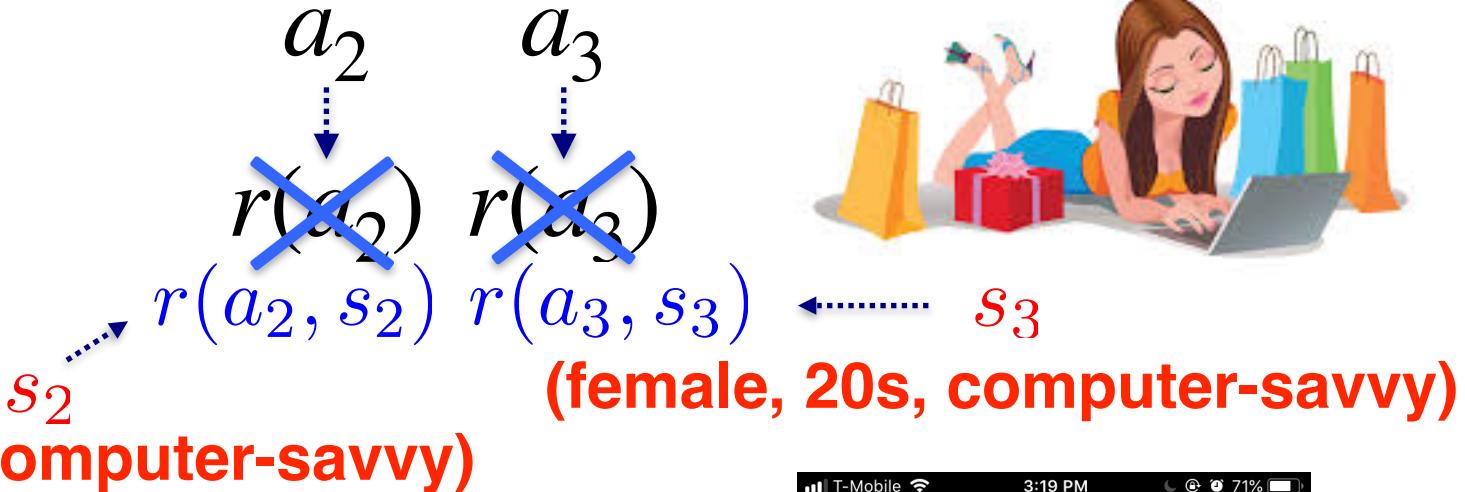


(female, 20s, computer-savvy)

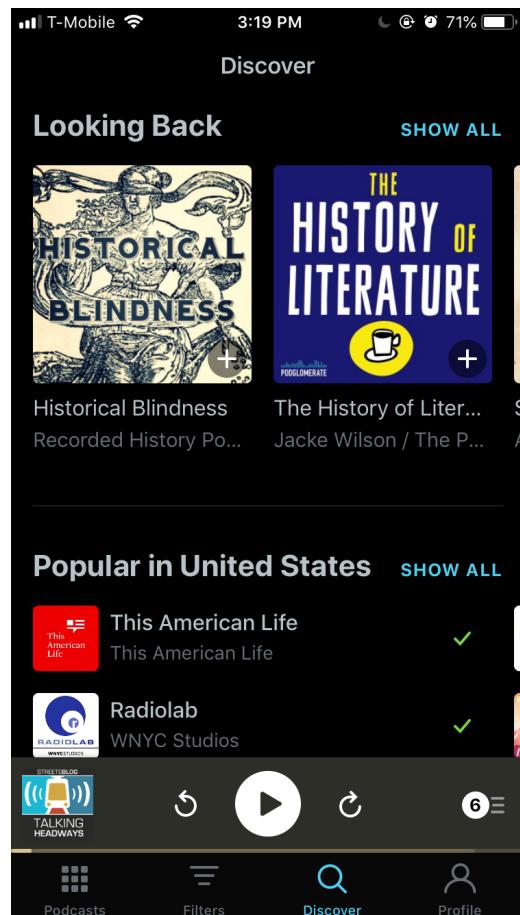
How to use these “features” in decision making?

Contextual Bandits

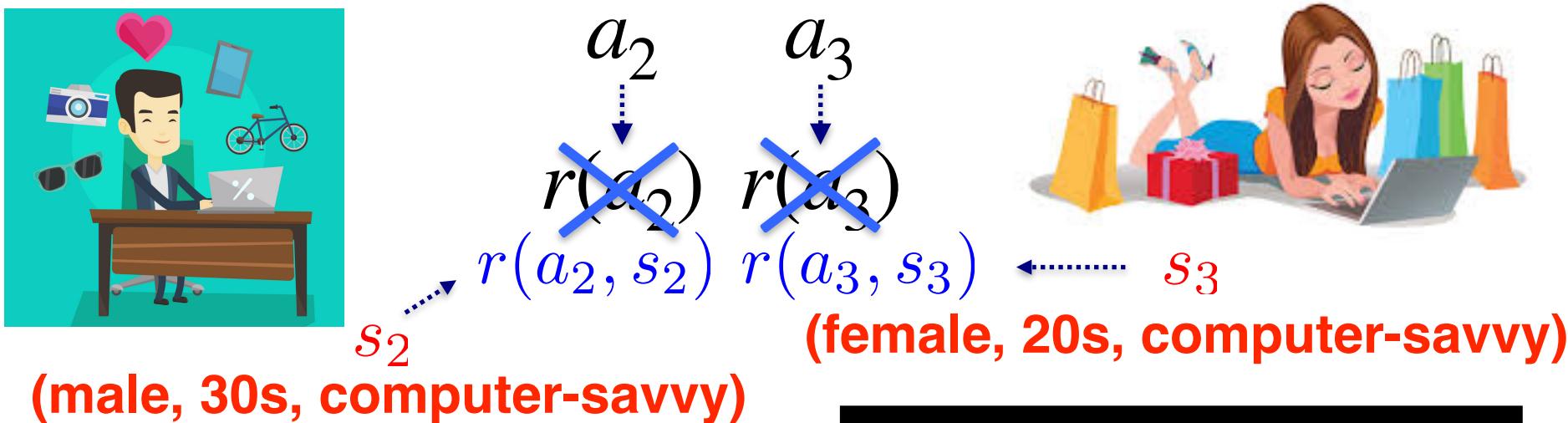




Example:
Podcast
recommendations



(slide co-designed with Cathy Wu)



Example types of contexts

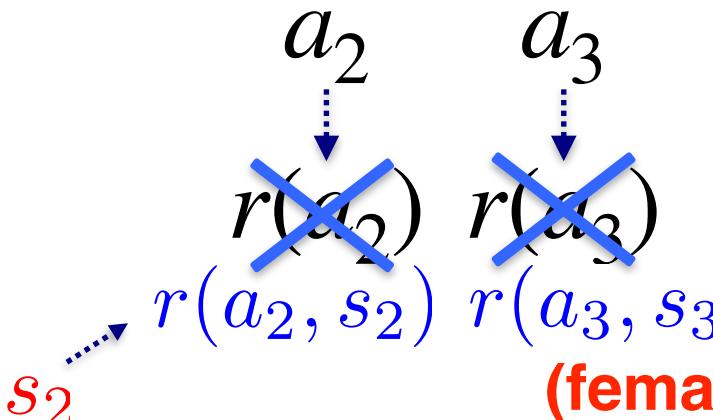
- User demographics
 - Discrete vs continuous

Switzerland

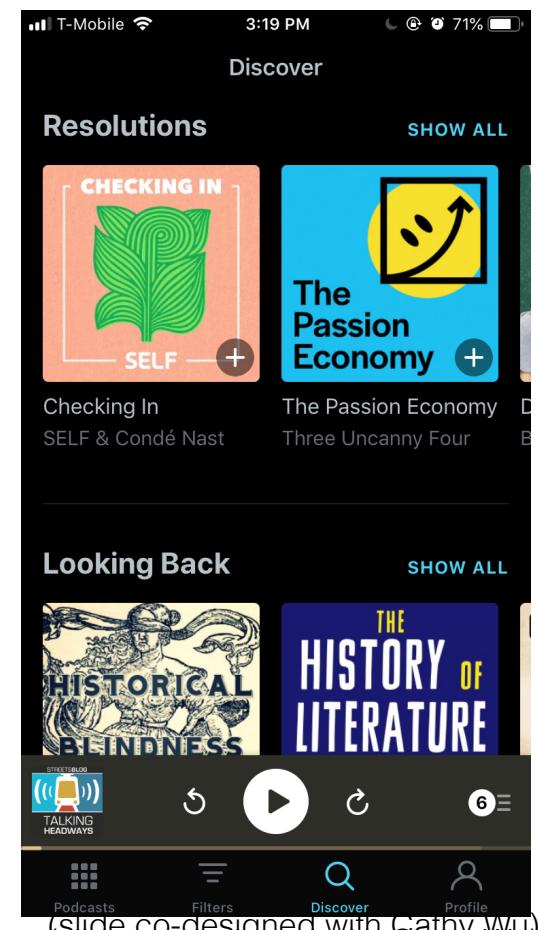
VS

(lat, long) coordinates

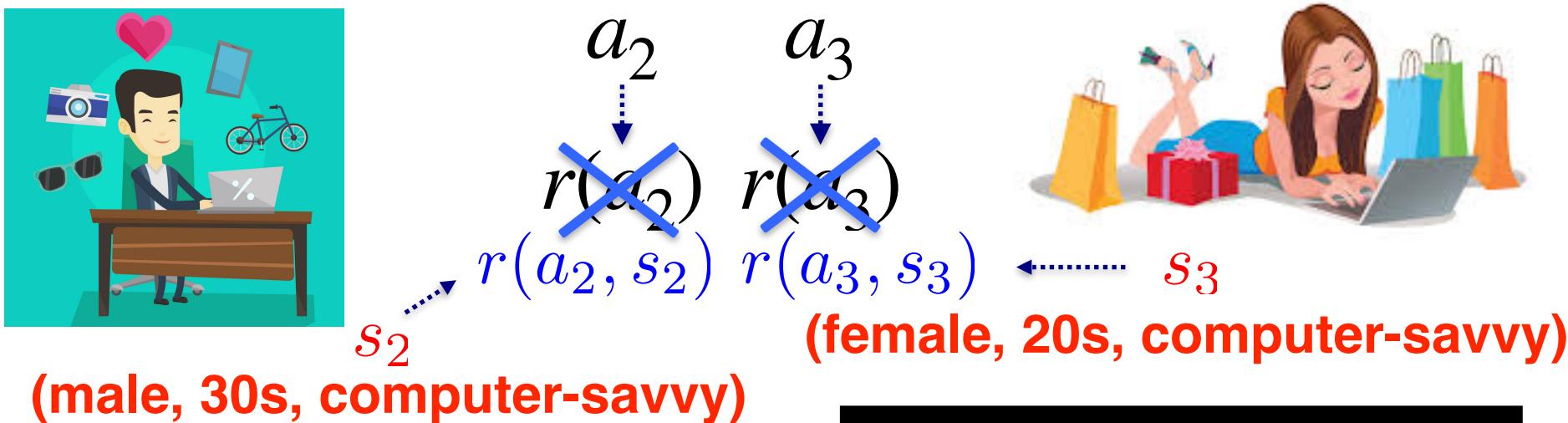




(male, 30s, computer-savvy)



Happy 2021!
New years resolutions?



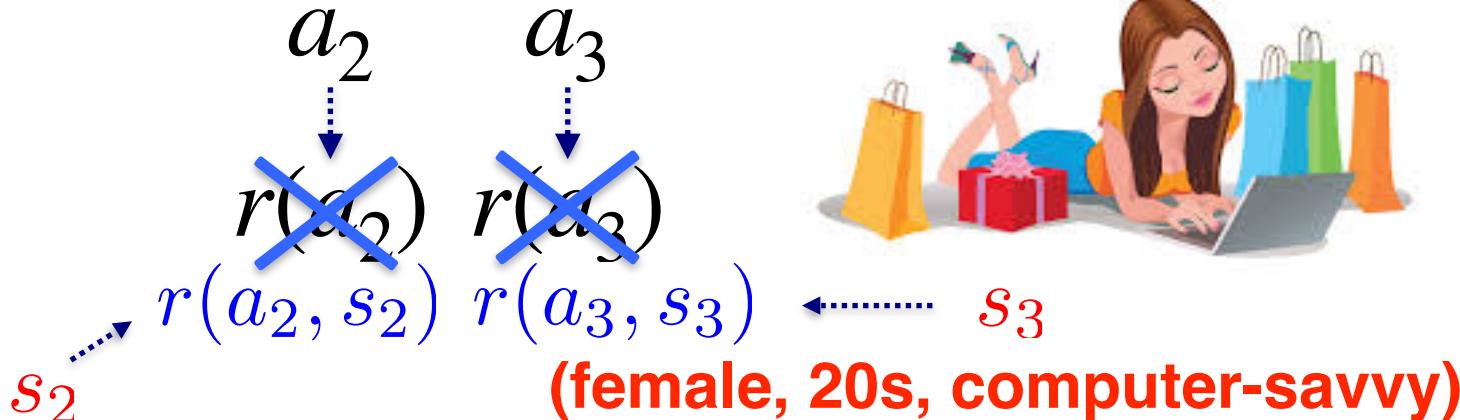
Naïve approach:
independent bandit problems
(one for each context)

Challenge:
may not handle continuous
contexts well





(male, 30s, computer-savvy)



Better approach:
LinUCB

Strategy:
select arm with highest

$$\text{UCB}_t(a|s_t) = \max_{\theta \in C_t} s_{t,a} \cdot \theta_a$$

Confidence interval
(shrinks with each visit)

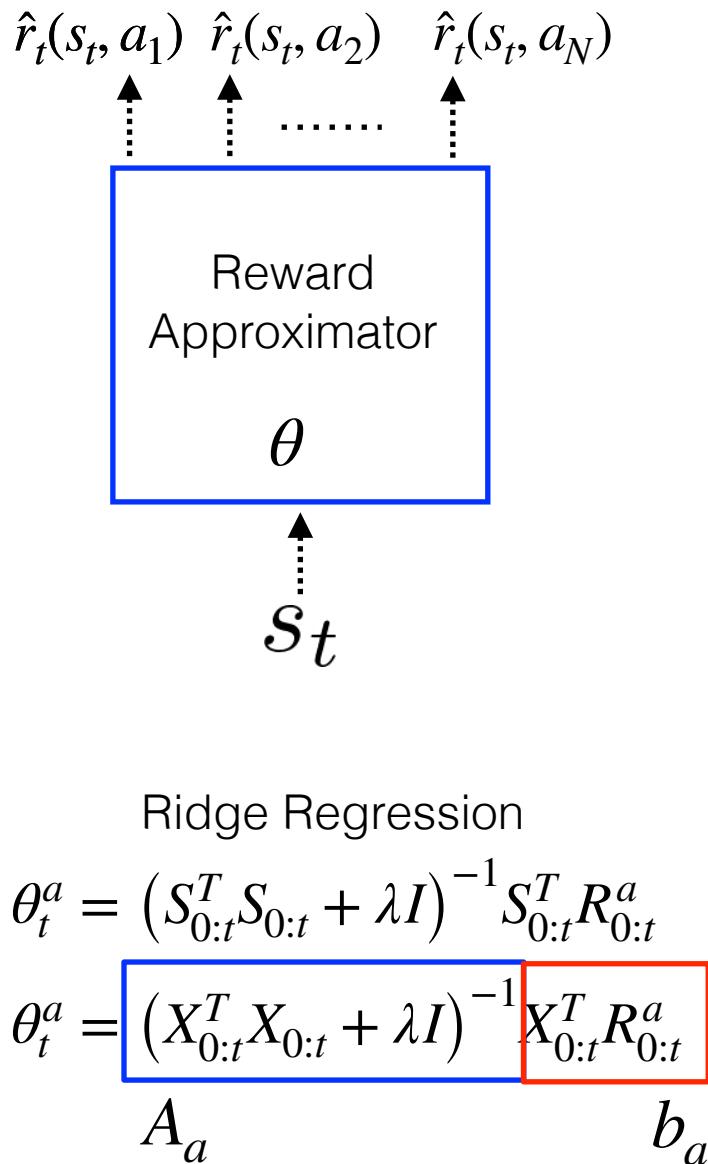
Assume:
expected rewards are linear
in context

$$\mu(a|s) = s_a \cdot \theta_a$$

θ_a fixed but
unknown

Optimal! (up to log factors)

Expanding on LinUCB



Estimate reward for each action

$$\hat{r}_t^a = s_t \theta_t^a$$

Choose the best one (or sample)

$$a_t \leftarrow \max_a \hat{r}_t^a$$

All time steps until t

$$\hat{R}_{0:t}^a = S_{0:t} \theta^a$$

Solve for the parameters

$$\min_{\theta_a} \|R_{0:t}^a - \hat{R}_{0:t}^a\|_2^2$$

(X: features of S)

Need to solve
Online!

$$\theta_t^a = \left(X_{0:t}^T X_{0:t} + \lambda I \right)^{-1} X_{0:t}^T R_{0:t}^a$$

$$A_a \qquad \qquad b_a$$

Algorithm 1 LinUCB with disjoint linear models.

```

0: Inputs:  $\alpha \in \mathbb{R}_+$ 
1: for  $t = 1, 2, 3, \dots, T$  do
2:   Observe features of all arms  $a \in \mathcal{A}_t$ :  $\mathbf{x}_{t,a} \in \mathbb{R}^d$ 
3:   for all  $a \in \mathcal{A}_t$  do
4:
5:
6:
7:
8:      $\boldsymbol{\theta}_a \leftarrow \mathbf{A}_a^{-1} \mathbf{b}_a$  . Exploration Bonus
9:
10:    end for
11:
12:
13:    . Online Update
14:  end for

```

$$\theta_t^a = \left(X_{0:t}^T X_{0:t} + \lambda I \right)^{-1} X_{0:t}^T R_{0:t}^a$$

$$A_a \qquad \qquad b_a$$

Algorithm 1 LinUCB with disjoint linear models.

```

0: Inputs:  $\alpha \in \mathbb{R}_+$ 
1: for  $t = 1, 2, 3, \dots, T$  do
2:   Observe features of all arms  $a \in \mathcal{A}_t$ :  $\mathbf{x}_{t,a} \in \mathbb{R}^d$ 
3:   for all  $a \in \mathcal{A}_t$  do
4:
5:
6:
7:
8:      $\boldsymbol{\theta}_a \leftarrow \mathbf{A}_a^{-1} \mathbf{b}_a$ 
9:      $p_{t,a} \leftarrow \hat{\boldsymbol{\theta}}_a^\top \mathbf{x}_{t,a} + \alpha \sqrt{\mathbf{x}_{t,a}^\top \mathbf{A}_a^{-1} \mathbf{x}_{t,a}}$  Exploration Bonus
10:    end for
11:    Choose arm  $a_t = \arg \max_{a \in \mathcal{A}_t} p_{t,a}$  with ties broken arbitrarily, and observe a real-valued payoff  $r_t$ 
12:     $\mathbf{A}_{a_t} \leftarrow \mathbf{A}_{a_t} + \mathbf{x}_{t,a_t} \mathbf{x}_{t,a_t}^\top$ 
13:     $\mathbf{b}_{a_t} \leftarrow \mathbf{b}_{a_t} + r_t \mathbf{x}_{t,a_t}$  Online Update
14:  end for

```

Pros and Cons of “disjoint models” (separate $\boldsymbol{\theta}_a$ for each action) ?

$$\theta_t^a = \left(X_{0:t}^T X_{0:t} + \lambda I \right)^{-1} X_{0:t}^T R_{0:t}^a$$

$$A_a \qquad \qquad b_a$$

Algorithm 1 LinUCB with disjoint linear models.

```

0: Inputs:  $\alpha \in \mathbb{R}_+$ 
1: for  $t = 1, 2, 3, \dots, T$  do
2:   Observe features of all arms  $a \in \mathcal{A}_t$ :  $\mathbf{x}_{t,a} \in \mathbb{R}^d$ 
3:   for all  $a \in \mathcal{A}_t$  do
4:
5:
6:
7:
8:      $\boldsymbol{\theta}_a \leftarrow \mathbf{A}_a^{-1} \mathbf{b}_a$ 
9:      $p_{t,a} \leftarrow \hat{\boldsymbol{\theta}}_a^\top \mathbf{x}_{t,a} + \alpha \sqrt{\mathbf{x}_{t,a}^\top \mathbf{A}_a^{-1} \mathbf{x}_{t,a}}$  Exploration Bonus
10:   end for
11:   Choose arm  $a_t = \arg \max_{a \in \mathcal{A}_t} p_{t,a}$  with ties broken arbitrarily, and observe a real-valued payoff  $r_t$ 
12:    $\mathbf{A}_{a_t} \leftarrow \mathbf{A}_{a_t} + \mathbf{x}_{t,a_t} \mathbf{x}_{t,a_t}^\top$ 
13:    $\mathbf{b}_{a_t} \leftarrow \mathbf{b}_{a_t} + r_t \mathbf{x}_{t,a_t}$  Online Update
14: end for

```

What if there are new news articles?

$$\theta_t^a = \left(X_{0:t}^T X_{0:t} + \lambda I \right)^{-1} X_{0:t}^T R_{0:t}^a$$

$$A_a \qquad \qquad b_a$$

Algorithm 1 LinUCB with disjoint linear models.

```

0: Inputs:  $\alpha \in \mathbb{R}_+$ 
1: for  $t = 1, 2, 3, \dots, T$  do
2:   Observe features of all arms  $a \in \mathcal{A}_t$ :  $\mathbf{x}_{t,a} \in \mathbb{R}^d$ 
3:   for all  $a \in \mathcal{A}_t$  do
4:     if  $a$  is new then
5:        $\mathbf{A}_a \leftarrow \mathbf{I}_d$  ( $d$ -dimensional identity matrix)
6:        $\mathbf{b}_a \leftarrow \mathbf{0}_{d \times 1}$  ( $d$ -dimensional zero vector)
7:     end if
8:      $\hat{\theta}_a \leftarrow \mathbf{A}_a^{-1} \mathbf{b}_a$ 
9:      $p_{t,a} \leftarrow \hat{\theta}_a^\top \mathbf{x}_{t,a} + \alpha \sqrt{\mathbf{x}_{t,a}^\top \mathbf{A}_a^{-1} \mathbf{x}_{t,a}}$  Exploration Bonus
10:   end for
11:   Choose arm  $a_t = \arg \max_{a \in \mathcal{A}_t} p_{t,a}$  with ties broken arbitrarily, and observe a real-valued payoff  $r_t$ 
12:    $\mathbf{A}_{a_t} \leftarrow \mathbf{A}_{a_t} + \mathbf{x}_{t,a_t} \mathbf{x}_{t,a_t}^\top$ 
13:    $\mathbf{b}_{a_t} \leftarrow \mathbf{b}_{a_t} + r_t \mathbf{x}_{t,a_t}$  Online Update
14: end for

```

Context: News articles + User Features