

Multi-armed Bandits

Matt Barnes

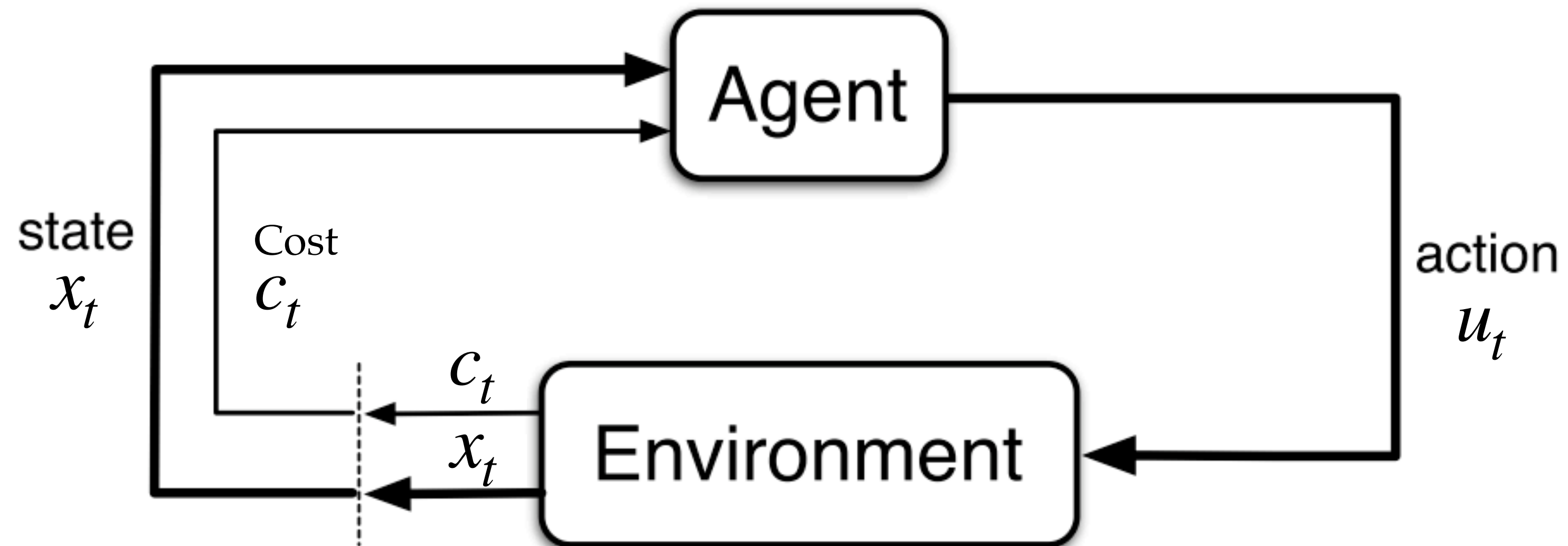
TAs: Matthew Rockett, Gilwoo Lee, Matt Schmittle

Housekeeping

- Remember to fill out course evaluations

Recap

- Markov Decision Processes are a very general class of models, which encompass planning and reinforcement learning.



Recap

“Markov” means that _____ captures all information about the history
 x_1, x_2, \dots, x_t

The most recent state x_t

Recap

- The difference between planning and reinforcement learning is whether the
_____ are known

transition model / dynamics / environment

Recap

- The three general methods for reinforcement learning are...

(1) Model-based

(2) Approximate dynamic programming

(3) Policy gradient

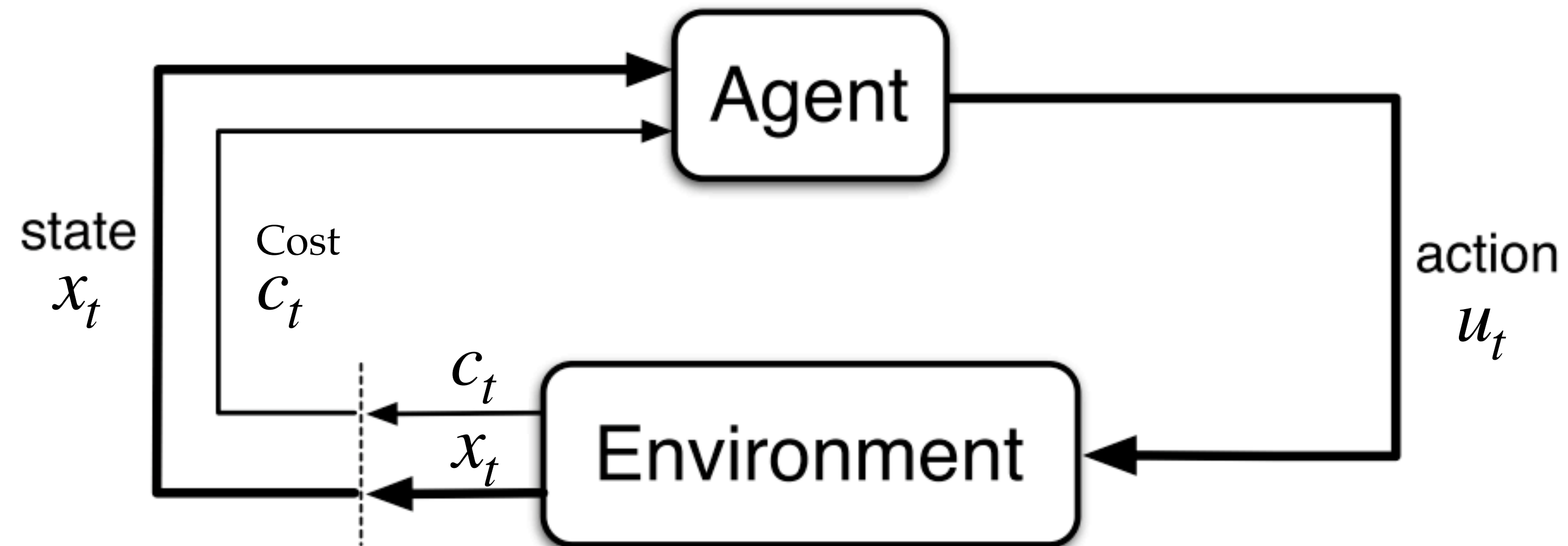
- (2) and (3) are both _____ methods

Model-free

What if the MDP
only has a single state?

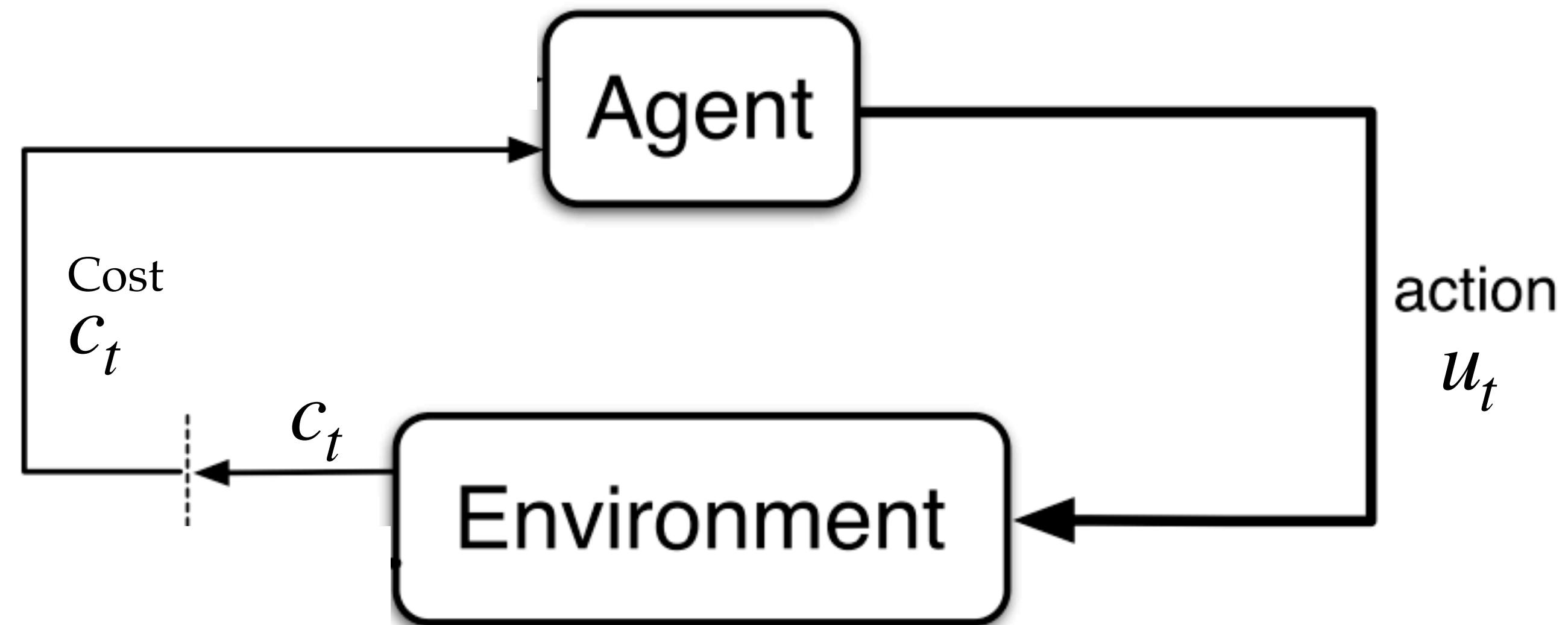
The MDP view of bandits

Reinforcement Learning

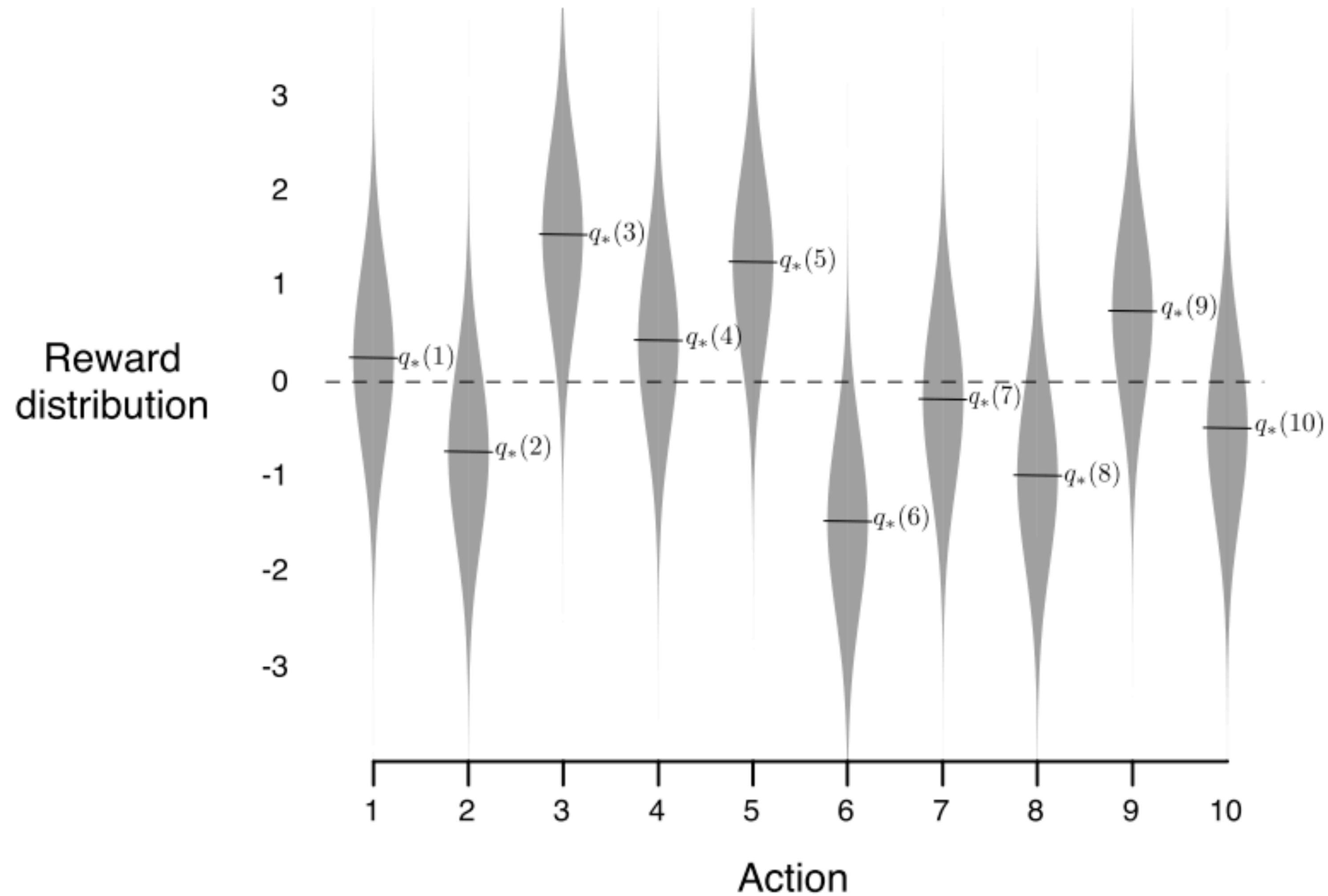


The MDP view of bandits

Multi-armed Bandits



Another view of the bandit problem



The original MAB problem

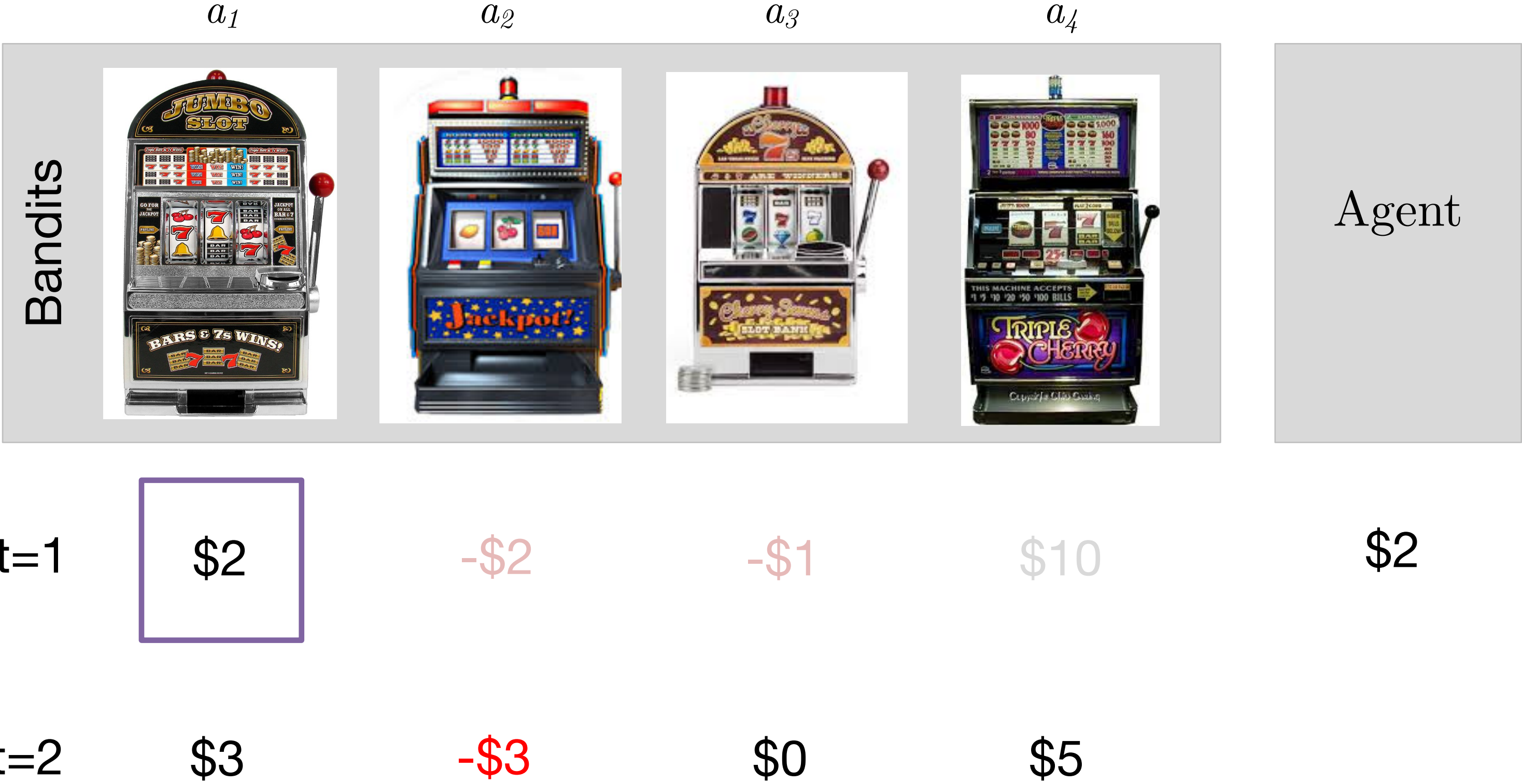


Side note: Throughout this lecture I use “reward” and “cost” interchangeably. You can think about reward as negative cost.

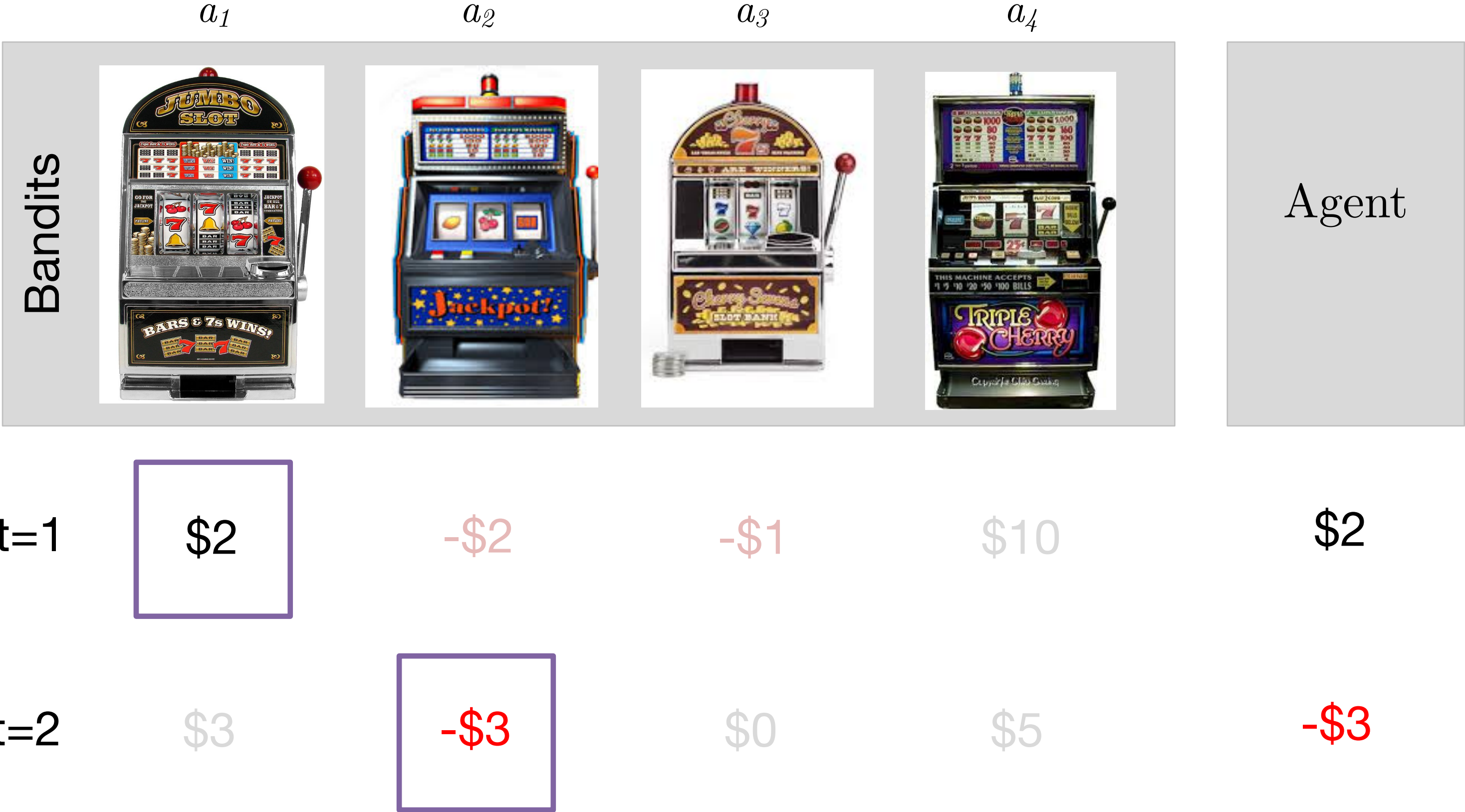
The original MAB problem







The original MAB problem







The original MAB problem



The original MAB problem

	a_1	a_2	a_3	a_4	
Bandits					Agent
t=1	\$2	-\$2	-\$1	\$10	\$2
t=2	\$3	-\$3	\$0	\$5	-\$3
t=3	\$1	-\$4	\$1	\$1000000	\$1

The original MAB problem

	a_1	a_2	a_3	a_4	
Bandits					Agent
t=1	\$2	-\$2	-\$1	\$10	\$2
t=2	\$3	-\$3	\$0	\$5	-\$3
t=3	\$1	-\$4	\$1	\$1000000	\$1
t=4	\$3	-\$4	-\$5	\$3	-\$5

Real world bandit successes

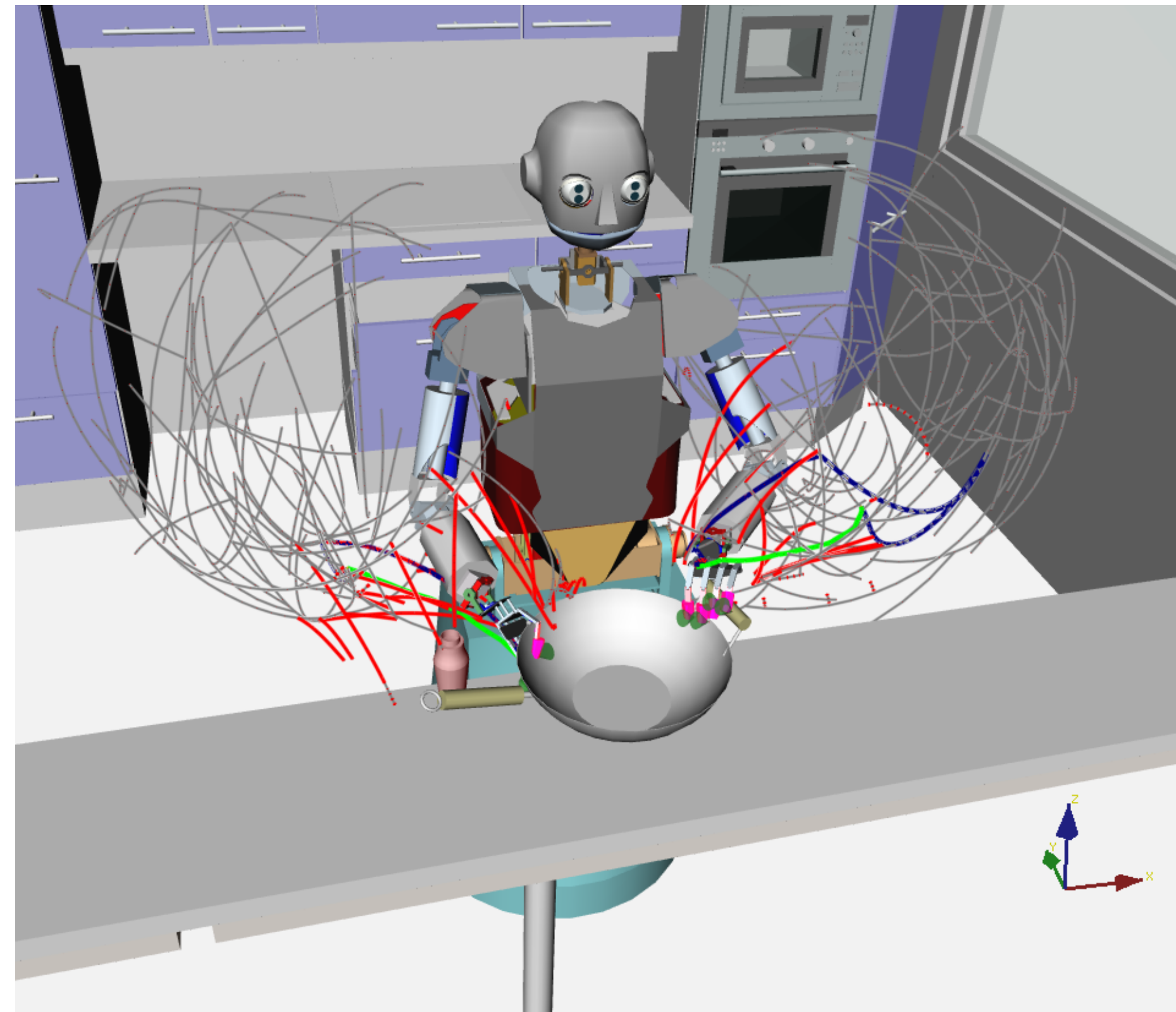
Which advertisement to display?

Reward: user clicks on the selected ad

The screenshot shows the International New York Times website. At the top, there are navigation links for U.S., INTERNATIONAL, 中文, and ESPAÑOL. Below this is a header with the 'Times JOURNEYS' logo, the text 'Small Group Tours to Unconventional Destinations', and a 'VIEW BROCHURE' button. The main headline is 'International New York Times' with the date 'Tuesday, August 30, 2016' and weather information '85°F' and 'Nasdaq -0.23% ↓'. Below the headline is a navigation bar with links for World, U.S., Politics, N.Y., Business, Opinion, Tech, Science, Health, Sports, Arts, Style, Food, Travel, Magazine, T Magazine, Real Estate, and ALL. The main content area features a large advertisement for 'outthink wind' with a blue background and a white airplane silhouette. Below the advertisement is a Google search bar with the text 'machine learning' and a search button. The search results show 'About 22,100,000 results (0.82 seconds)' and a top result for 'Machine Learning Andrew Ng - Get Certified in Machine Learning' with a link to 'www.coursera.org/machine-learning'. Below the link is a snippet: 'Earn a Stanford Certificate! 72% of Coursera participants surveyed reported career benefits – HBR'. At the bottom, there are four links: 'How Coursera Works', 'Coursera Specializations', 'Coursera - Join for Free', and 'Big Data Specialization'.

Other exciting applications

Which grasp? Reward: robot picks up object



Other exciting applications

Which treatment? Reward: patient gets healthy



Other exciting applications

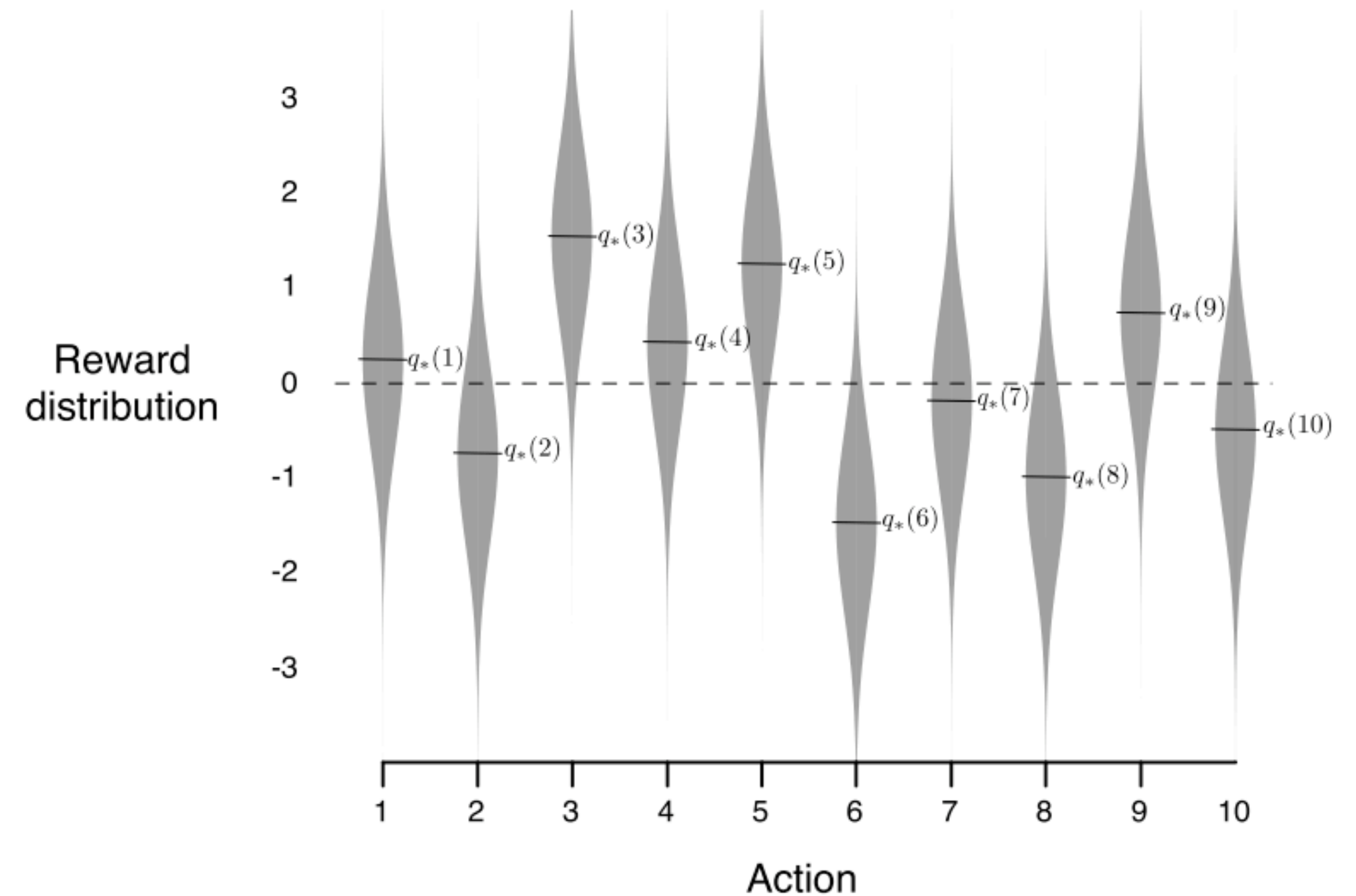
Packet routing: Which path to send data along?



The *stochastic* bandit setting

At each time-step $t = 1, 2, \dots$

- Choose action a_k
- Receive stochastic reward $r_t \sim \mathbb{P}_k(c)$



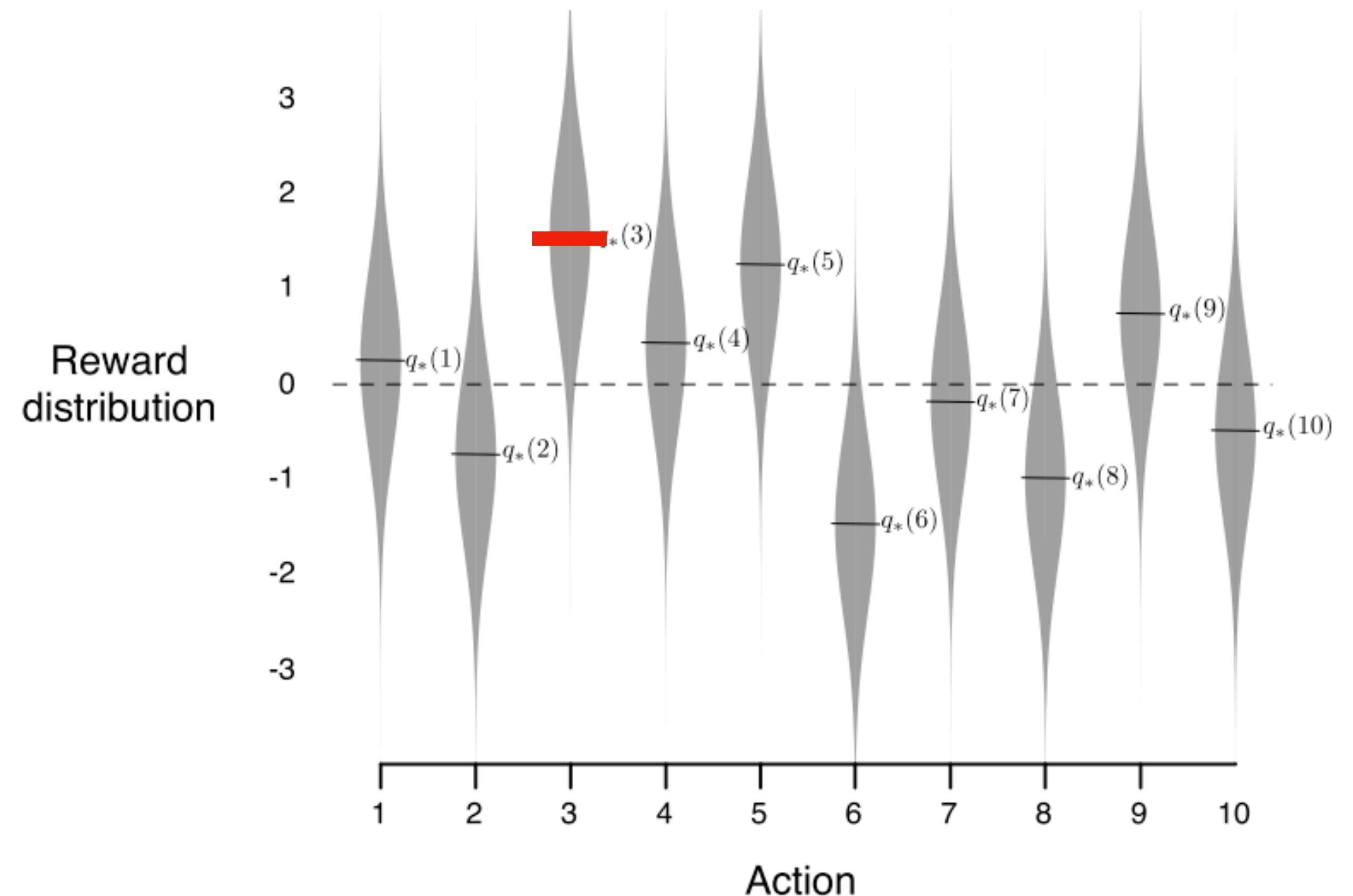
The *stochastic* bandit setting

- What is the best action?

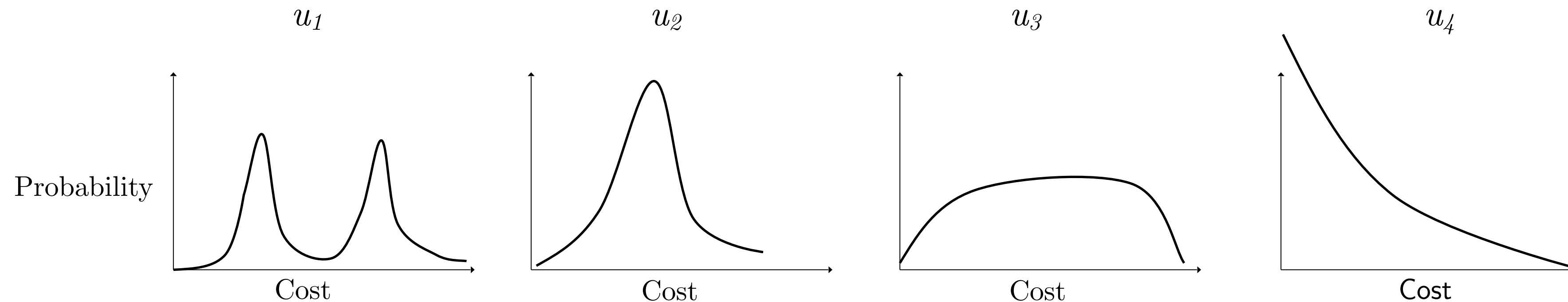
\mathcal{U}_3

- Why?

Optimal action is one with highest expected reward



The *stochastic* bandit setting

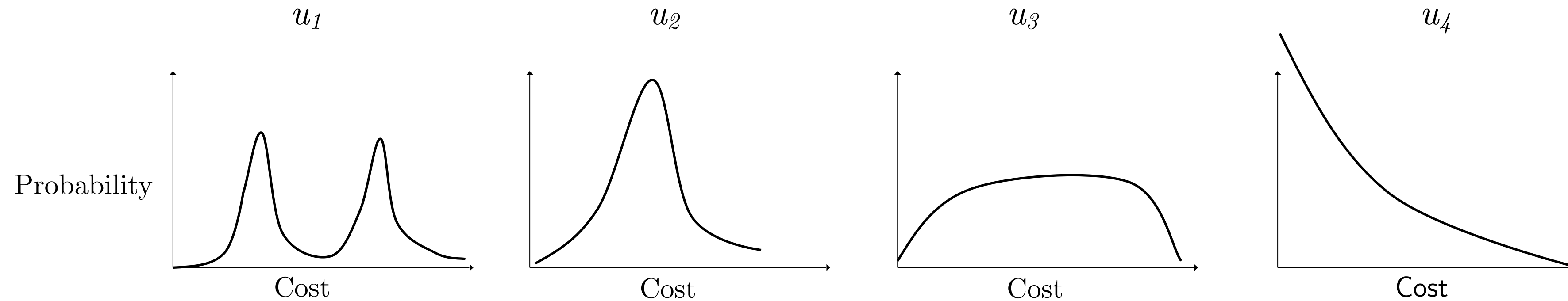


Key challenge: How do we maximize our gambling earnings?

Requires both

- (a) Playing arms we don't know much about (**exploring**)
- (b) Earning money on arms we know will pay off well (**exploiting**)

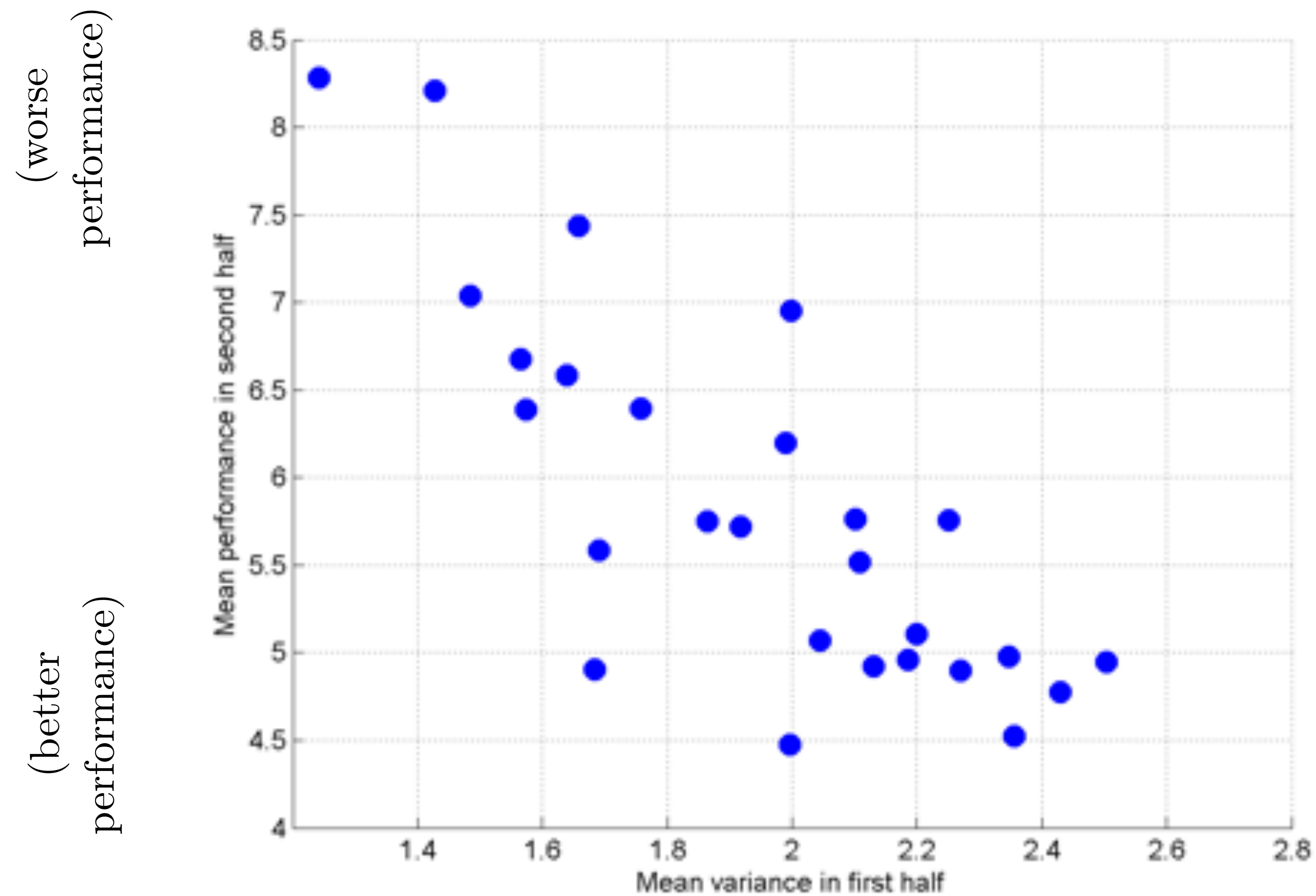
A dumb algorithm



Step 1) **Explore:** Play each arm n times

Step 2) **Exploit:** Choose the arm with the best estimated reward forever

“Exploration vs. Exploitation” is a fundamental tradeoff



“Exploration vs. Exploitation” is a fundamental tradeoff

Key takeaway: Algorithm needs to transition
from *exploring* to *exploiting*

A slightly less dumb algorithm: The ε -greedy algorithm

Maintain *estimates* of each action's expected cost. At each time step, choose action a_t according to

- **Exploit:** With probability $(1 - \varepsilon)$ choose the action with the lowest estimated cost
- **Explore:** With probability ε , choose a random action

Update estimates.

The ε -greedy algorithm

Say we chose action k at time steps $t = t_{k_1}, \dots, t_{k_{n_k}}$

The estimator of the expected cost is the empirical mean of the observed costs for this action:

$$\hat{c}_k = \frac{1}{n_k} \sum_{i=1}^{n_k} c_{t_{k_i}}$$

Why can't we just choose the action with the lowest estimated cost??

If \hat{c} of the optimal action is less than the true expected cost of some other action, then we will never choose the optimal action!

The ε -greedy algorithm

The downside of ε -greedy

- It's pretty dumb about how it chooses to explore or exploit
- It will keep exploring forever!

A Smart Algorithm: Upper Confidence Bound (UCB) Algorithm

- Instead of maintaining an estimate of the expected reward, greedily choose actions with the largest *upper confidence bound* on the expected reward

The diagram shows the UCB formula: $\hat{c}_k + \sqrt{\frac{2}{n_k} \log \left(\frac{1}{\delta} \right)}$. A red box highlights the square root term. Three red arrows point from text labels to parts of the formula: one from the left to \hat{c}_k , one from below to n_k , and one from the right to the box. The text labels are in red.

$\hat{c}_k + \sqrt{\frac{2}{n_k} \log \left(\frac{1}{\delta} \right)}$ Bonus which encourages exploration

The estimate of the mean
(same as ε -greedy)

Number of times we
have tried action u_k

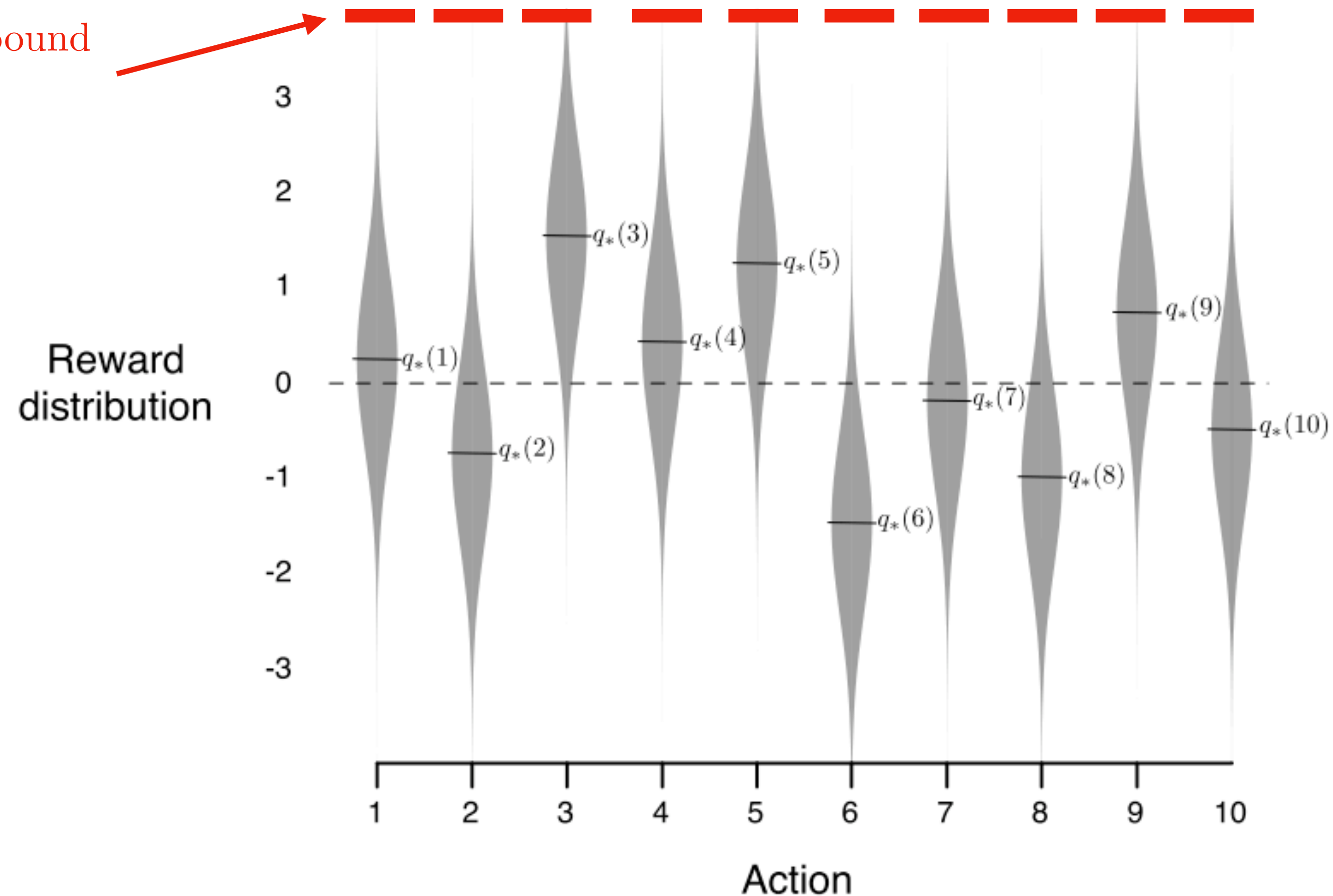
Constant, trades off
exploration and exploitation

- “Optimism in the face of uncertainty.” Naturally trades off choosing actions:
 - With a high estimated reward
 - We are uncertain about (have not tried many times)

A Smart Algorithm: Upper Confidence Bound (UCB) Algorithm

At $t=0$, $n_k = 0$ for all k and the upper confidence bounds are infinite.

The upper confidence bound



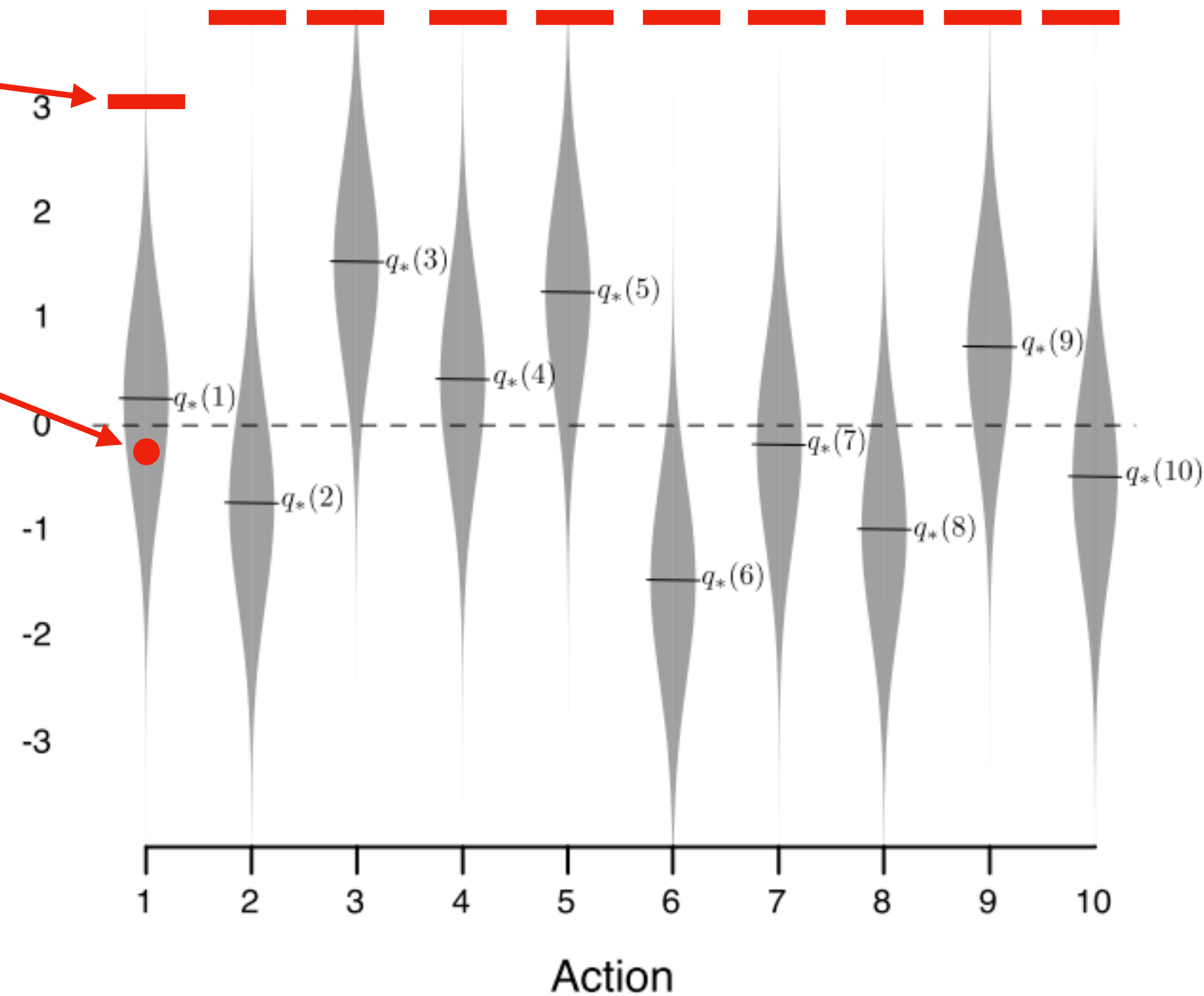
A Smart Algorithm: Upper Confidence Bound (UCB) Algorithm

We pull the first arm

The upper confidence bound

The sample r

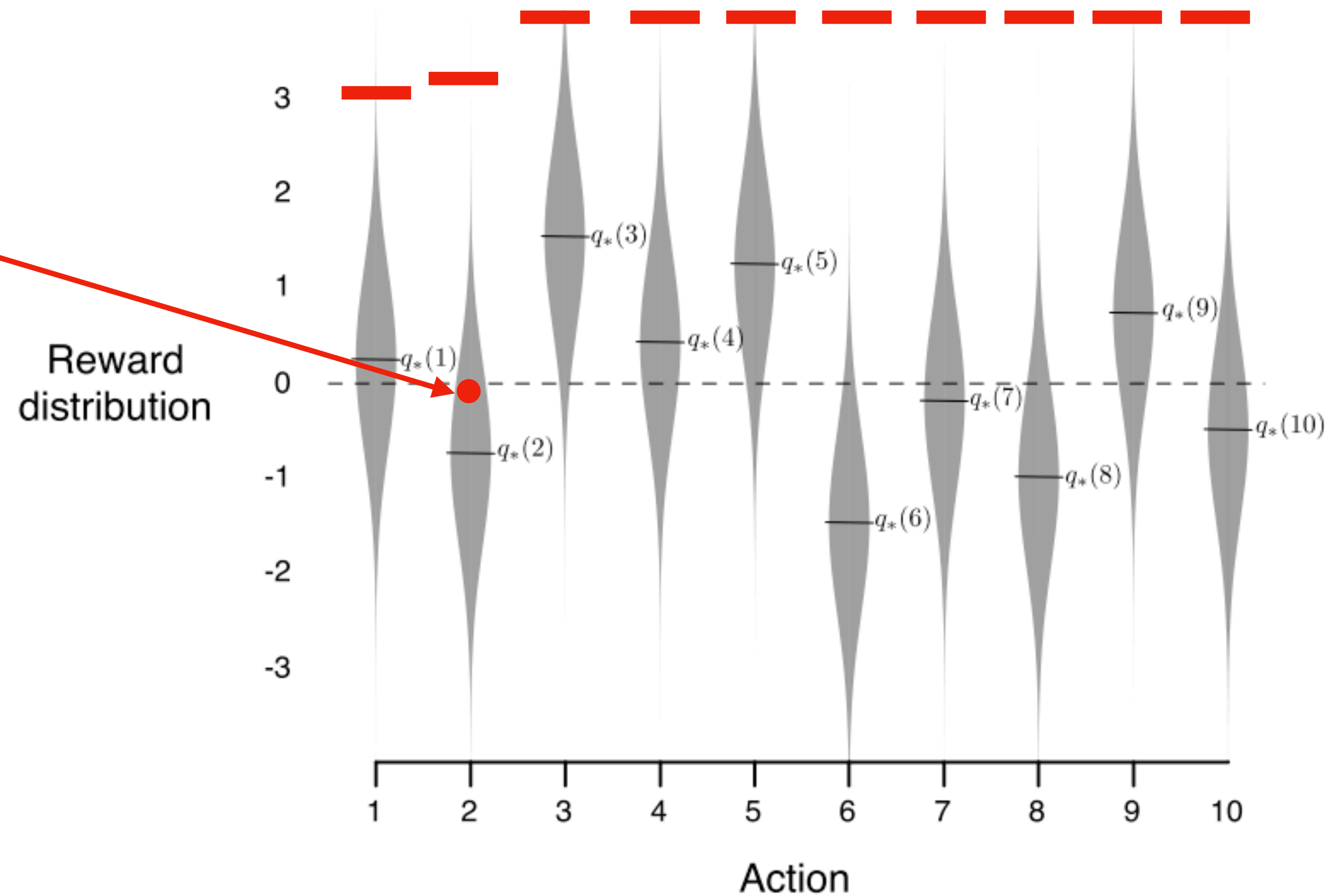
Reward
distribution



A Smart Algorithm: Upper Confidence Bound (UCB) Algorithm

We pull the second arm

The sample r

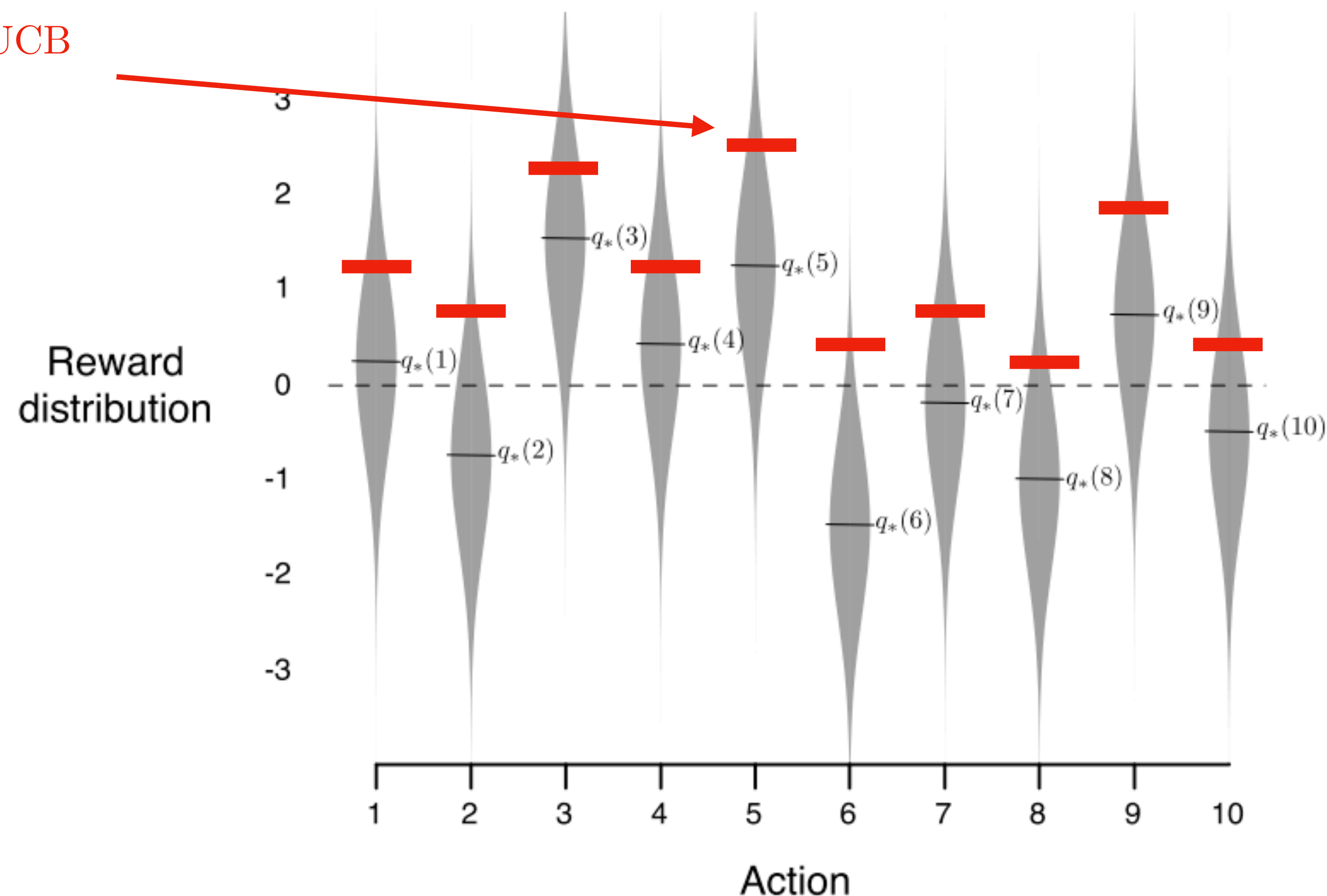


A Smart Algorithm:

Upper Confidence Bound (UCB) Algorithm

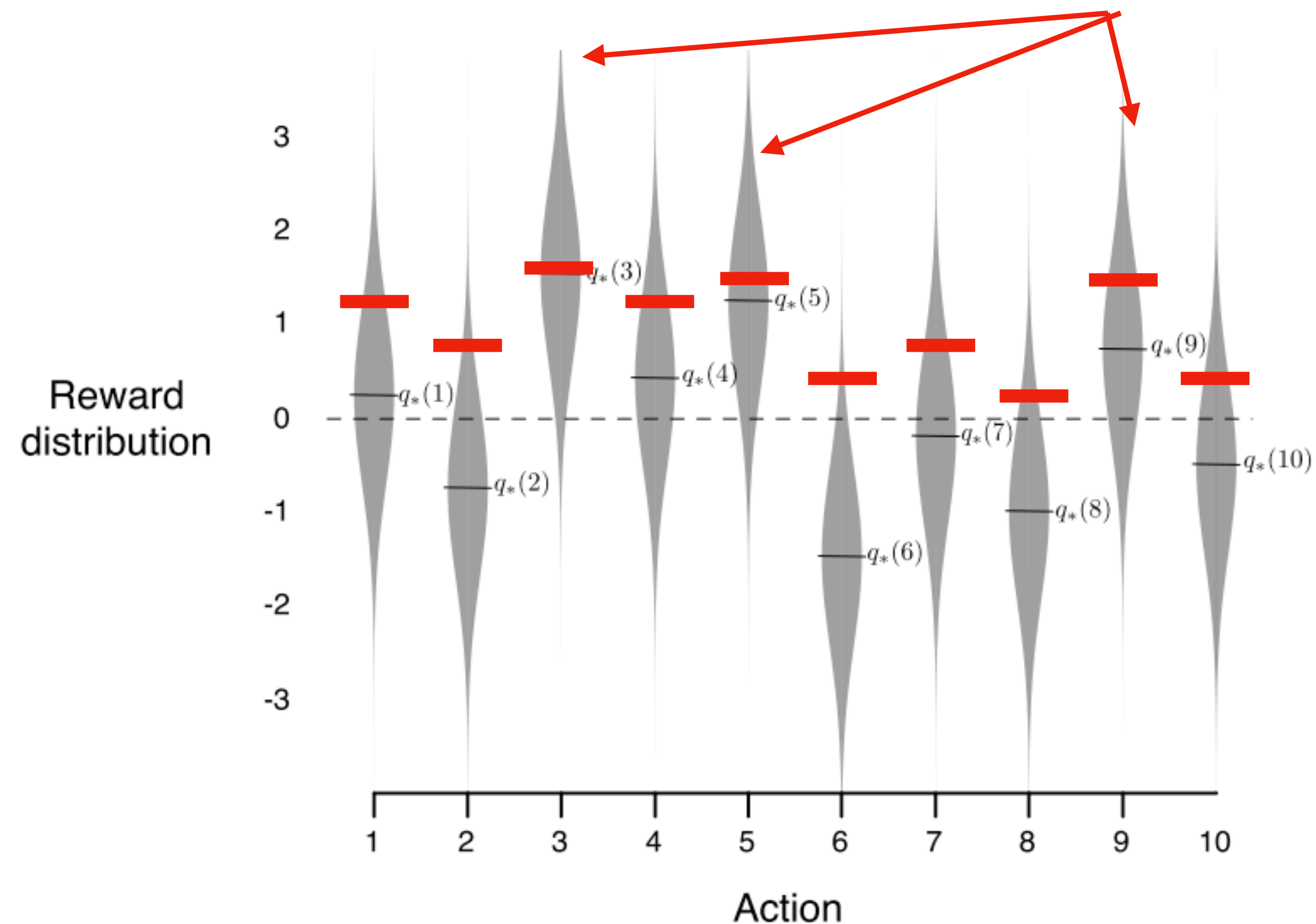
- By now, each arm has been pulled a couple times, and we have a more reasonable UCB.
- Note that these UCB's are *optimistic*, and allow us to focus our actions on promising actions in the same way as planning heuristics (e.g. A^*).

Action with the best UCB



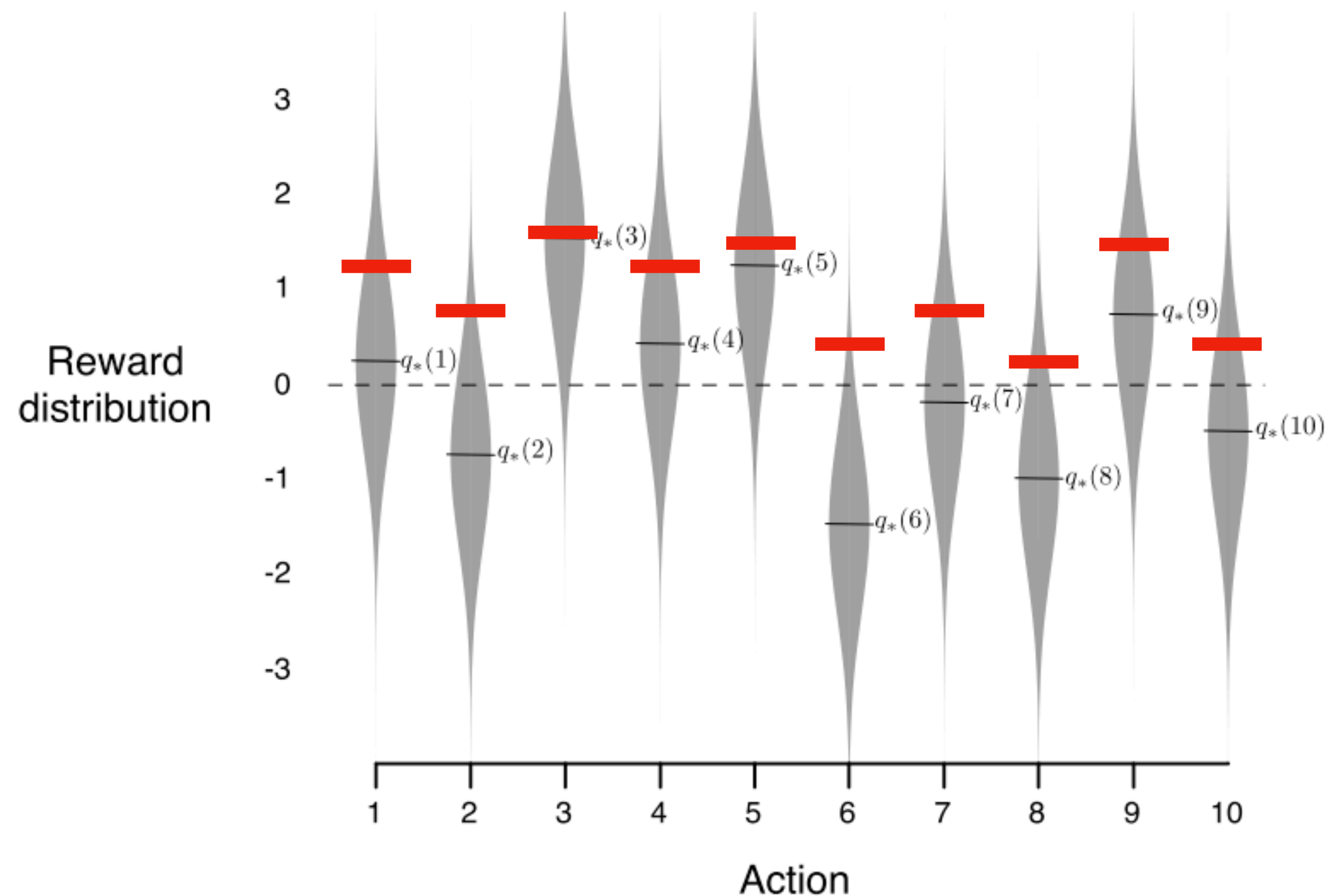
A Smart Algorithm: Upper Confidence Bound (UCB) Algorithm

Only need to pull these three arms!
Will effectively stop exploring once all UCB's are less than r_3

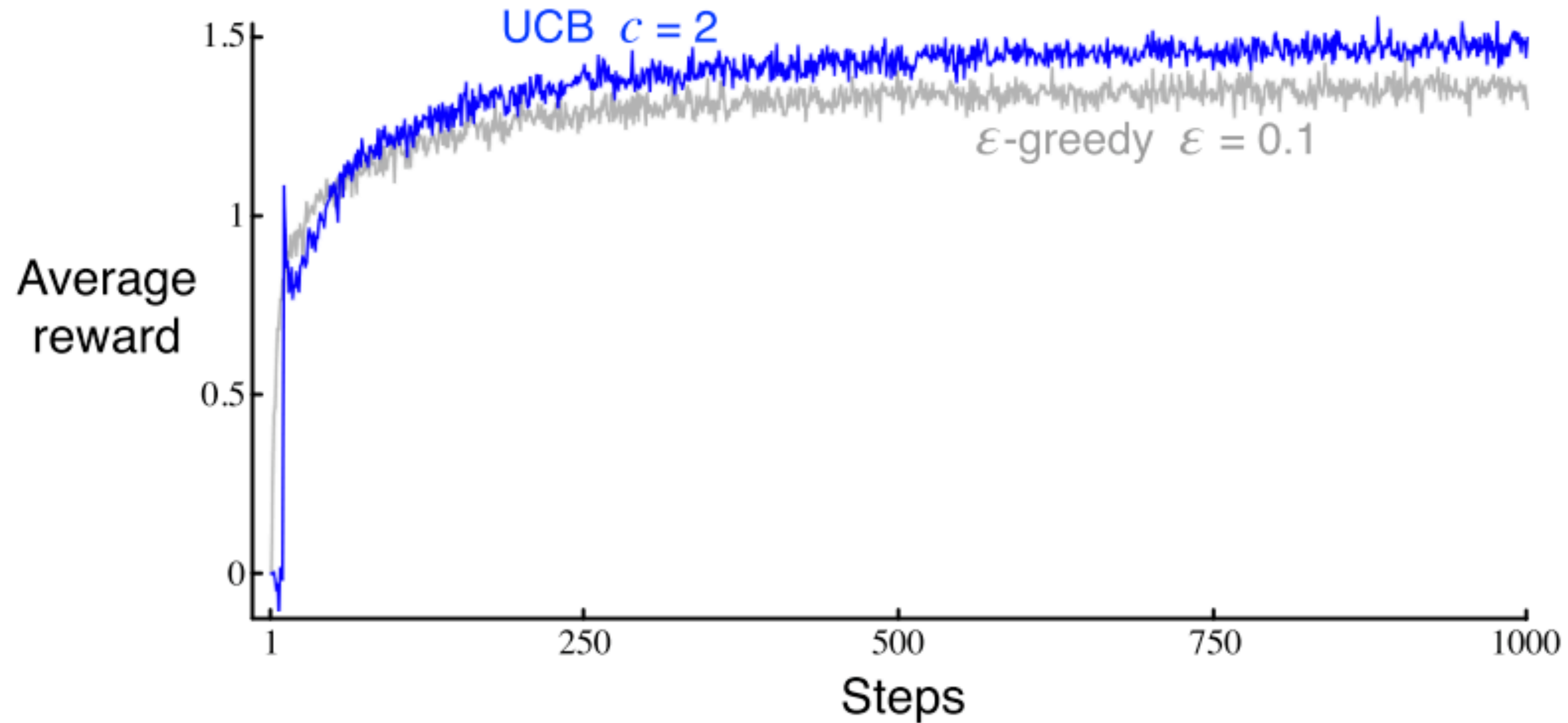


A Smart Algorithm: Upper Confidence Bound (UCB) Algorithm

- Focus **exploration** on more promising arms, evident by the better arms having tighter UCB'
- Naturally transitions from exploration to **exploitation**



UCB outperforms ϵ -greedy in practice (and in theory)



Contextual bandits are a powerful variation of the classic bandit problem

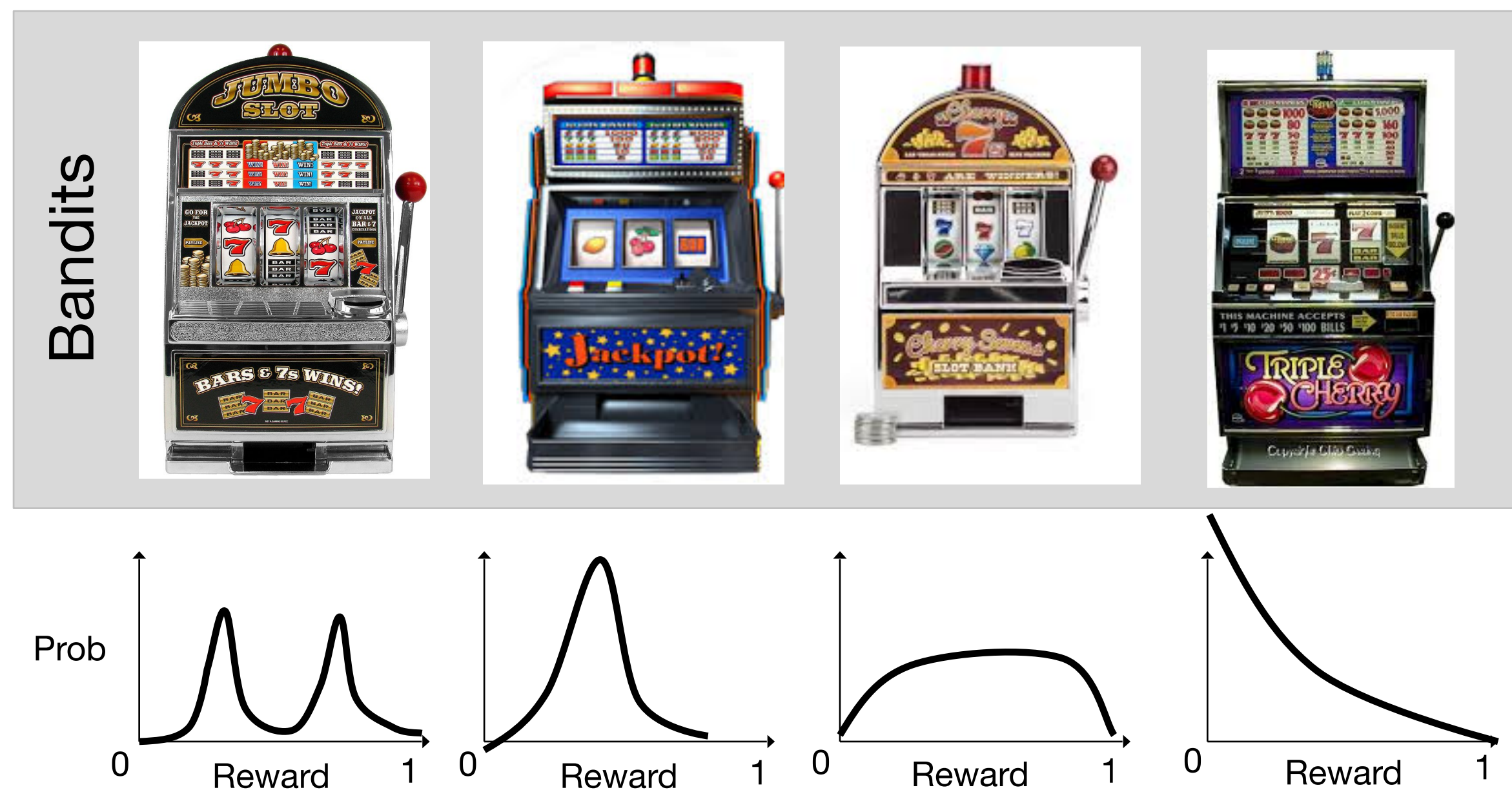
Example: Advertising – which ad to show user?



Additional context: User history, type of ad, IP address, time of day
Reward: Does user click the ad?

Another variation: Adversarial rewards

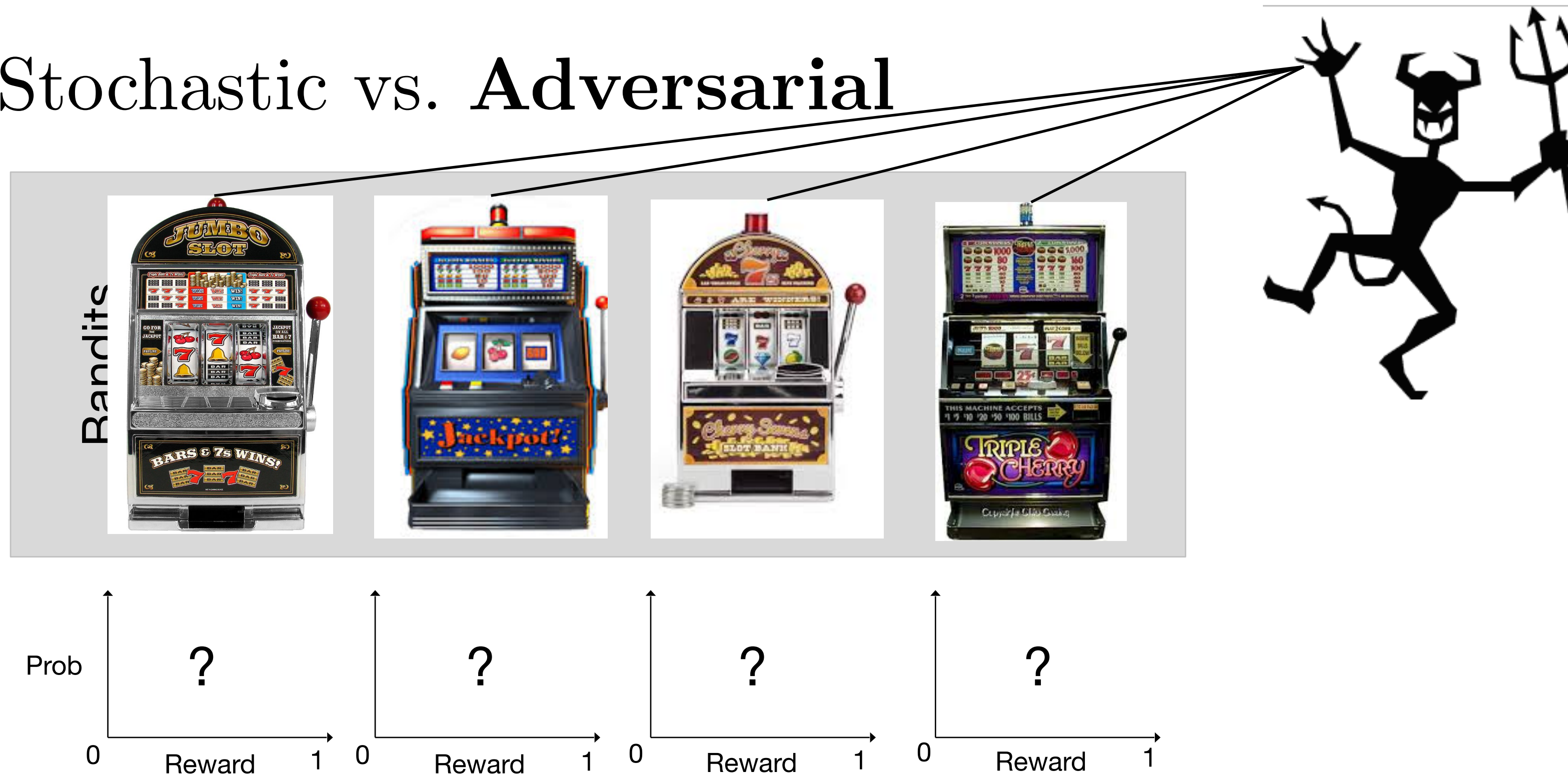
Stochastic vs. Adversarial



- Stochastic bandits assume reward drawn i.i.d. from some bounded distribution

Another variation: Adversarial rewards

Stochastic vs. **Adversarial**



- Adversarial bandits only assume reward is bounded, samples not drawn IID from some distribution
- Surprisingly, we can prove strong guarantees in this setting, too (EXP3 algorithm)

Recap

The multi-armed bandit problem is reinforcement learning with __ state(s)

1

Recap

What is the fundamental trade-off in bandit (and reinforcement learning) problems?

Exploration vs. exploitation

Recap

Reward is equivalent to _____

Negative cost

Recap

The ε -greedy algorithm randomly explores with probability

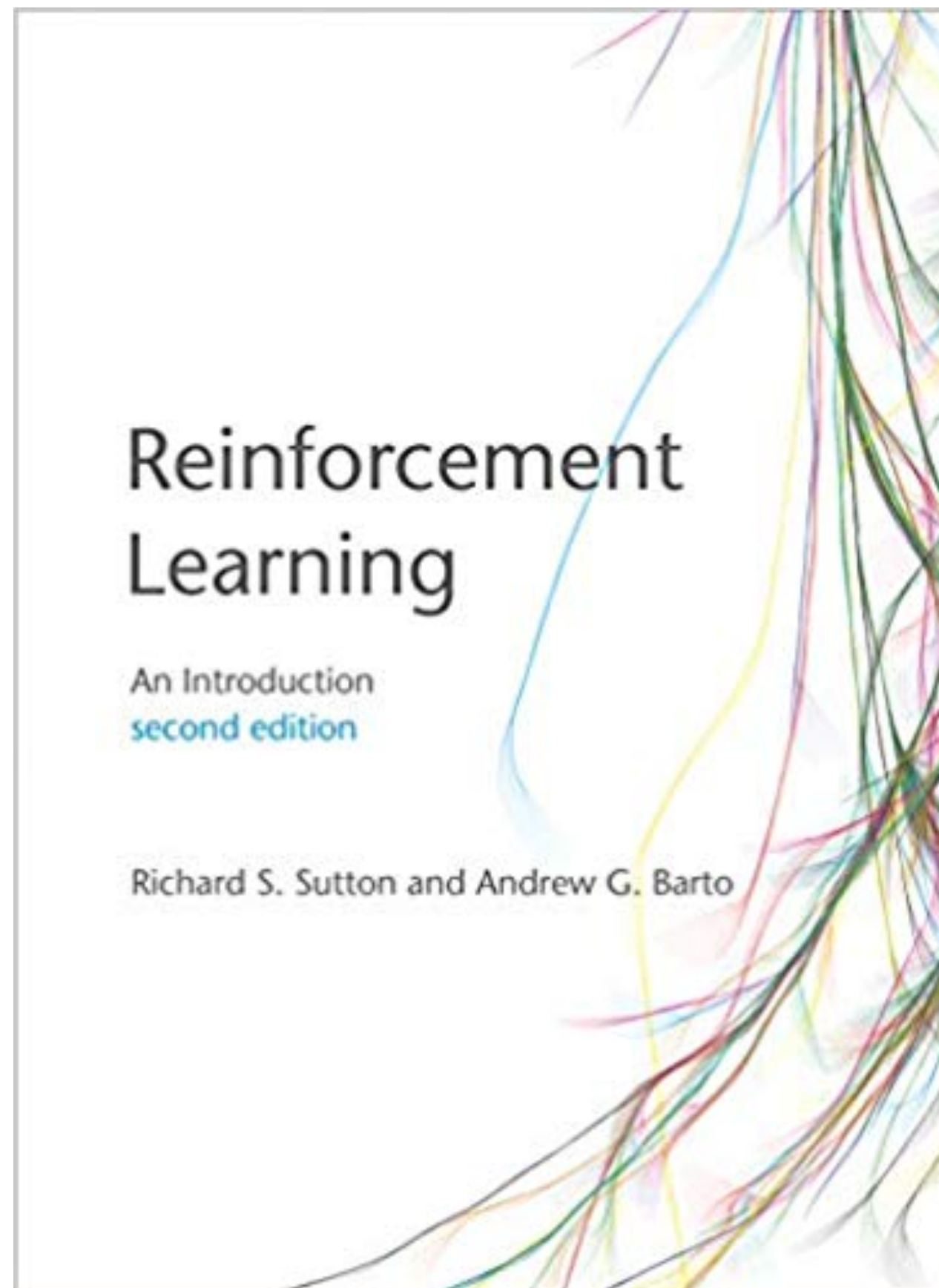
ε

Recap

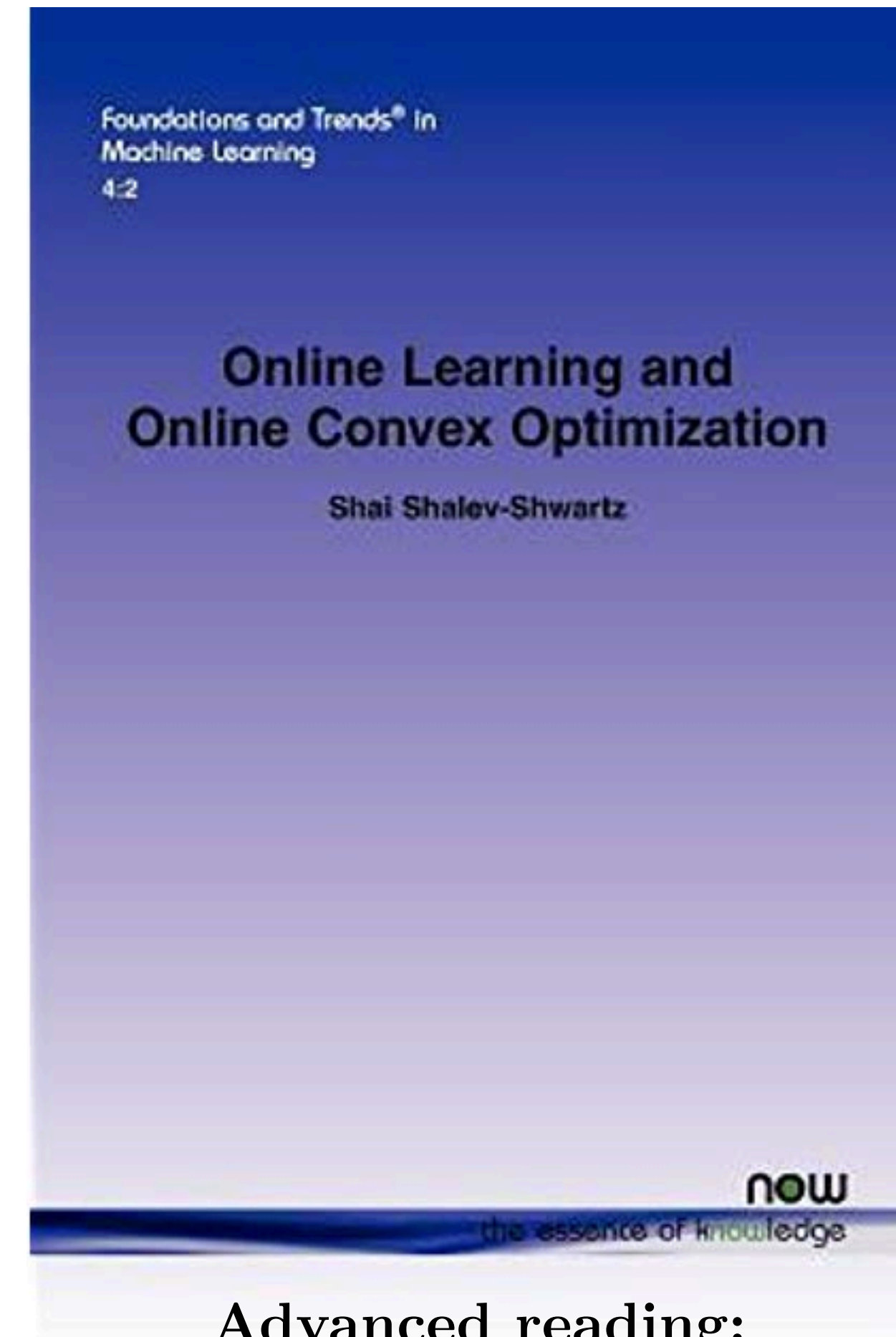
The UCB algorithm chooses actions according to the estimated reward, plus a bonus (i.e. confidence interval) which decreases with respect to _____

The number of times we try that action.

Further Reading



Chapter 2: Multi-armed bandits



**Advanced reading:
Connection to online convex optimization**

Preview of the next lecture

- **Imitation learning** — what if we have help from a (human) expert?

