

### Motivation: RL and Bandits

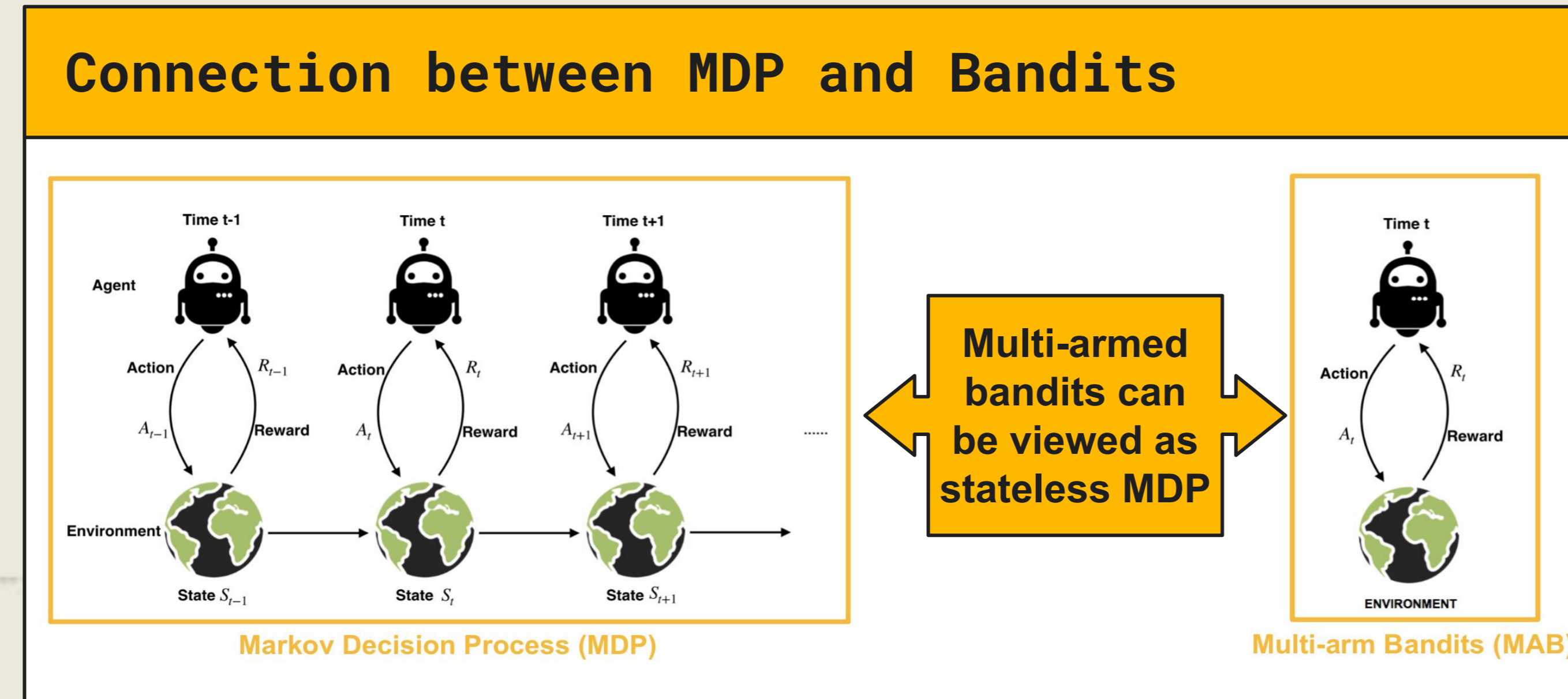
Reinforcement Learning (RL)

- Self-driving Cars
- Robotics
- Gaming

Multi-Armed Bandits (MAB)

- Online Advertising
- Clinical Trials
- Recommendation

- Real-world environments are **complex and uncertain**
- Training data are **expensive**
- Our novel algorithms allows an agent to **converge to optimal policies with proven sample efficiency**



### Check out our work at UAI 2023 (Oral presentation)!!!

### Challenge: Exploitation vs Exploration Trade-Off

**Exploitation**

Take actions with high empirical reward to gain pay-off

**Exploration**

Take less observed actions to gather information

### Stochastic Multi-Armed Bandits (MAB)

A stochastic MAB instance  $\Theta := ([K]; \mu_1, \mu_2, \dots, \mu_K)$   
 In every round  $t = 1, 2, \dots, T$

- Environment generates a reward vector  $(X_1(t), \dots, X_j(t), \dots, X_K(t))$  where  $X_j(t) \sim \text{Ber}(\mu_j)$
- Simultaneously, Learner pulls an arm  $J_t \in [K]$
- Environment reveals  $X_{J_t}(t)$ ; Learner observes and obtains  $X_{J_t}(t)$

Goal: pull arms sequentially to maximize cumulative reward  
 Regret:  $\mathcal{R}(T; \Theta) = \mathbb{E} \left[ \sum_{t=1}^T \left( \max_{j \in [K]} \mu_j - \mu_{J_t} \right) \right]$

## Optimistic Thompson Sampling balances exploration-exploitation trade-off in RL

### Our algorithms enjoy elegant analyses and tight regret bounds

### Episodic Markov Decision Processes (MDP)

An MDP instance  $M := (T, H, [S], [A], \{\mu\}_{[S] \times [A] \times [H]}, \{\bar{P}\}_{[S] \times [A] \times [H]}, p_0)$

- Number of episodes:  $T$
- Number of rounds in each episode:  $H$
- Mean reward function:  $\{\mu_{s,a,t}\}$
- Transition probability distribution function:  $\{\bar{P}_{s,a,t}\}$
- Deterministic initial state distribution:  $p_0$

Policy:  $\pi = (\pi(\cdot, 1), \pi(\cdot, 2), \dots, \pi(\cdot, H))$  with each  $\pi(\cdot, t) : S \rightarrow A$  taking a state  $s_t$  as input and outputs an action  $a_t$  that will played in that state

Goal: play a sequence of policies  $\pi_1, \pi_2, \dots, \pi_k, \dots, \pi_K$  to accumulate as much reward as possible

Regret:  $\mathcal{R}(T; M) = \mathbb{E} \left[ \sum_{k=1}^K (V_1^{\pi_*}(s_1^k) - V_1^{\pi_k}(s_1^k)) \right]$ , where  $V_t^{\pi}(s)$  is the value function and  $\pi_*$  is the optimal policy

### UPPER Confidence BOUND (UCB) vs Thompson Sampling (TS) in Bandits

Unknown parameters:  $(\mu_1, \mu_2, \dots, \mu_K)$   
 Empirical parameters:  $(\hat{\mu}_{1, O_1(t-1)}, \hat{\mu}_{2, O_2(t-1)}, \dots, \hat{\mu}_{K, O_K(t-1)})$

UPPER confidence BOUND (UCB):  $\bar{\mu}_{j,t} = \hat{\mu}_{j, O_j(t-1)} + \sqrt{\frac{2 \log(t)}{O_j(t-1)}}$  Pull arm  $J_t = \arg \max \bar{\mu}_{j,t}$

Thompson Sampling (TS):  $\theta_{j,t} \sim \mathcal{N}(\hat{\mu}_{j, O_j(t-1)}, \frac{1}{O_j(t-1)})$

Optimistic TS (O-TS):  $\theta_{j,t} \sim \mathcal{N}^+(\hat{\mu}_{j, O_j(t-1)}, \frac{1}{O_j(t-1)})$

Optimistic TS+ (O-TS+):  $\theta_{j,t} \sim \mathcal{N}^+(\hat{\mu}_{j, O_j(t-1)}, \frac{1}{O_j(t-1)})$

**Regret UPPER BOUND**

$O(\sqrt{KT \ln(T)})$

**More Optimistic Distributions!!**

**Key idea: Sampled parameters are always better than empirical parameters!**  
 O-TS for bandits was originally proposed and empirically evaluated in Chapelle and Li [2011], May et al. [2012].

### 0-TS-MDP vs 0-TS-MDP+ in MDPs

Unknown parameters:  $\mu_{s,a,t}, \bar{P}_{s,a,t}$   
 Empirical parameters:  $\hat{\mu}_{s,a,t}^{k-1}, \hat{P}_{s,a,t}^{k-1}$

Model-based: Construct a model  $M^k$  in each episode  $k$ . Find the best policy  $\pi_k$  for  $M^k$

UCB-VI:  $M^k = \{[S], [A], H, \bar{\mu}^k, \hat{P}^{k-1}\}$ , where  $\bar{\mu}_{s,a,t}^k = \hat{\mu}_{s,a,t}^{k-1} + \tilde{O}\left(\sqrt{\frac{H^2}{O_{s,a,t}^{k-1}}}\right)$

O-TS-MDP:  $M^k = \{[S], [A], H, \theta^k, \hat{P}^{k-1}\}$ , where:  $\theta_{s,a,t}^k \sim \mathcal{N}^+(\hat{\mu}_{s,a,t}^{k-1}, \tilde{O}\left(\frac{H^3 s}{O_{s,a,t}^{k-1}}\right))$

O-TS-MDP+:  $M^k = \{[S], [A], H, \theta^k, \hat{P}^{k-1}\}$ , where:  $\theta_{s,a,t}^k \sim \mathcal{N}^+(\hat{\mu}_{s,a,t}^{k-1}, \tilde{O}\left(\frac{H^2}{O_{s,a,t}^{k-1}}\right))$

**Regret UPPER BOUND**

O-TS :  $\tilde{O}(\sqrt{AS^2 H^4 T})$

O-TS+ :  $\tilde{O}(\sqrt{ASH^3 T})$

O-TS-MDP enjoys an elegant theoretical analysis, avoiding bounding the absolute value of approximation error. O-TS-MDP+ has the same regret bound as UCB-VI [Azar et al., 2017] and can be viewed as a randomized version of UCB-VI.

### Acknowledgement

This work was supported by Alberta Machine Intelligence Institute (Amii), the Canada CIFAR AI Program and the Natural Sciences and Engineering Research Council of Canada (NSERC) Discovery Grants.