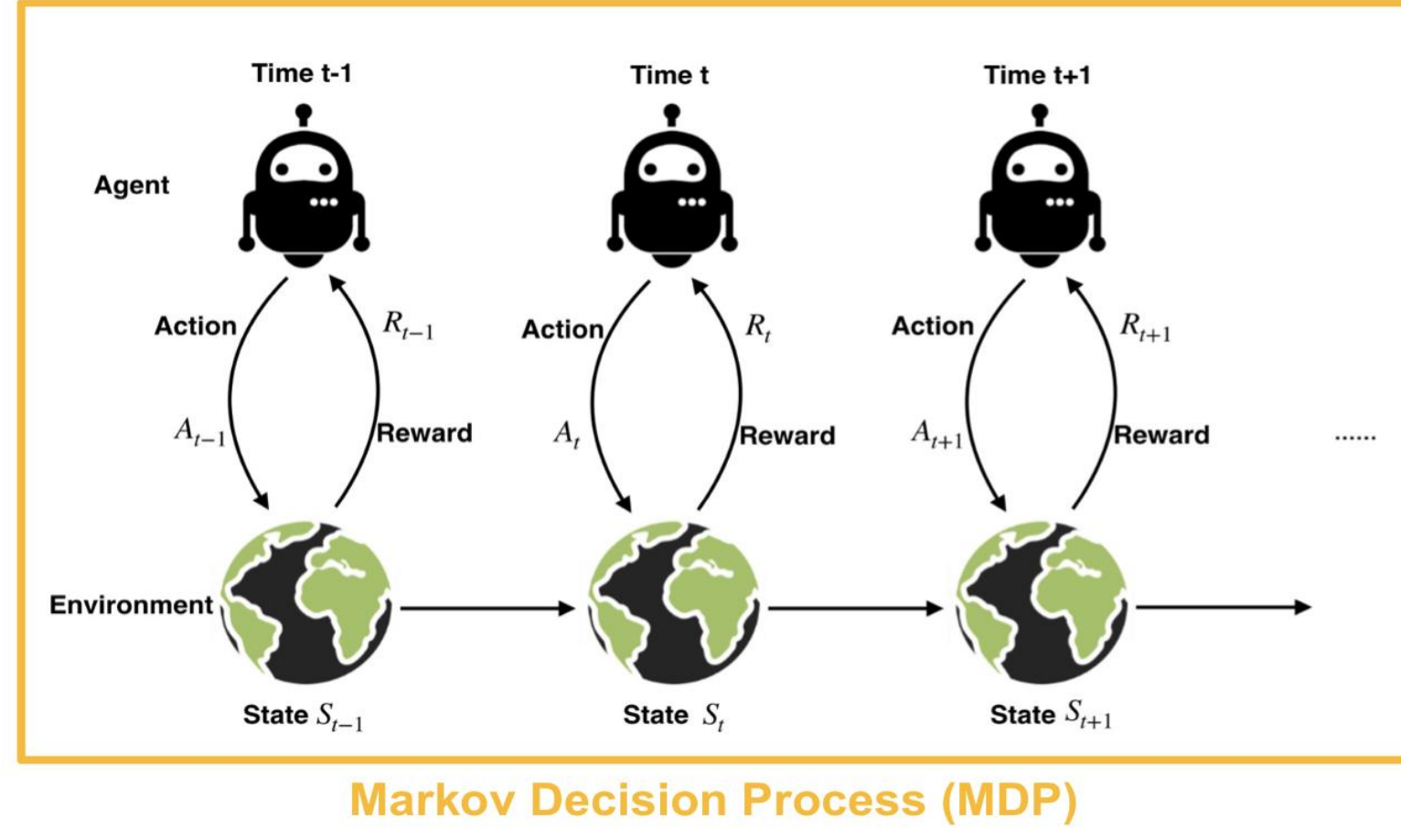
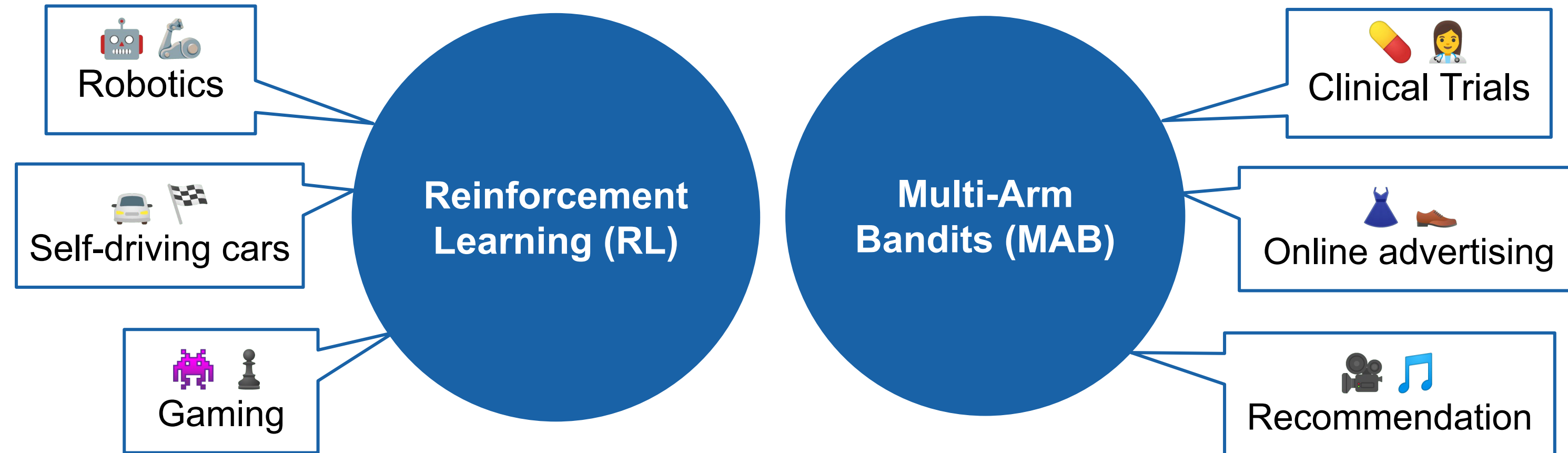
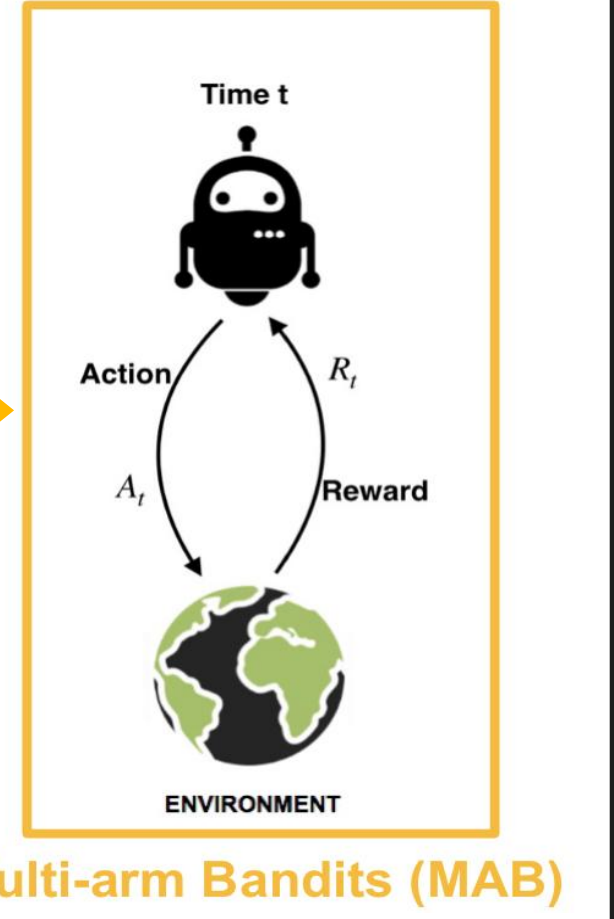


Motivation: RL and Bandits



Multi-arm bandits can be viewed as stateless MDP

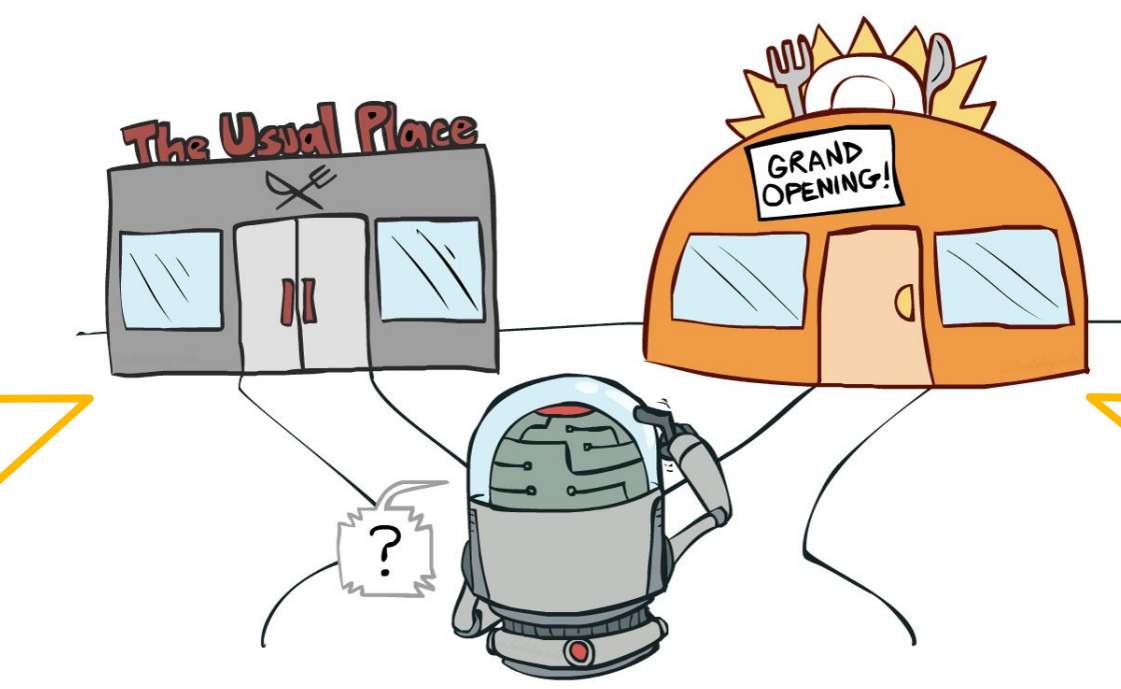


Check out our work at UAI 2023 (Oral presentation)!!!

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

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

Stochastic Multi-Armed Bandits (MAB)

Goal: pull arms sequentially to maximize cumulative reward

$$\text{Regret: } \mathcal{R}(T; \Theta) = \mathbb{E} \left[\sum_{t=1}^T \left(\max_{j \in [K]} \mu_j - \mu_{J_t} \right) \right]$$

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 $\pi(\cdot, t) : S \rightarrow \mathcal{A}$ taking a state s_t as input and outputs an action a_t that will be 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

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

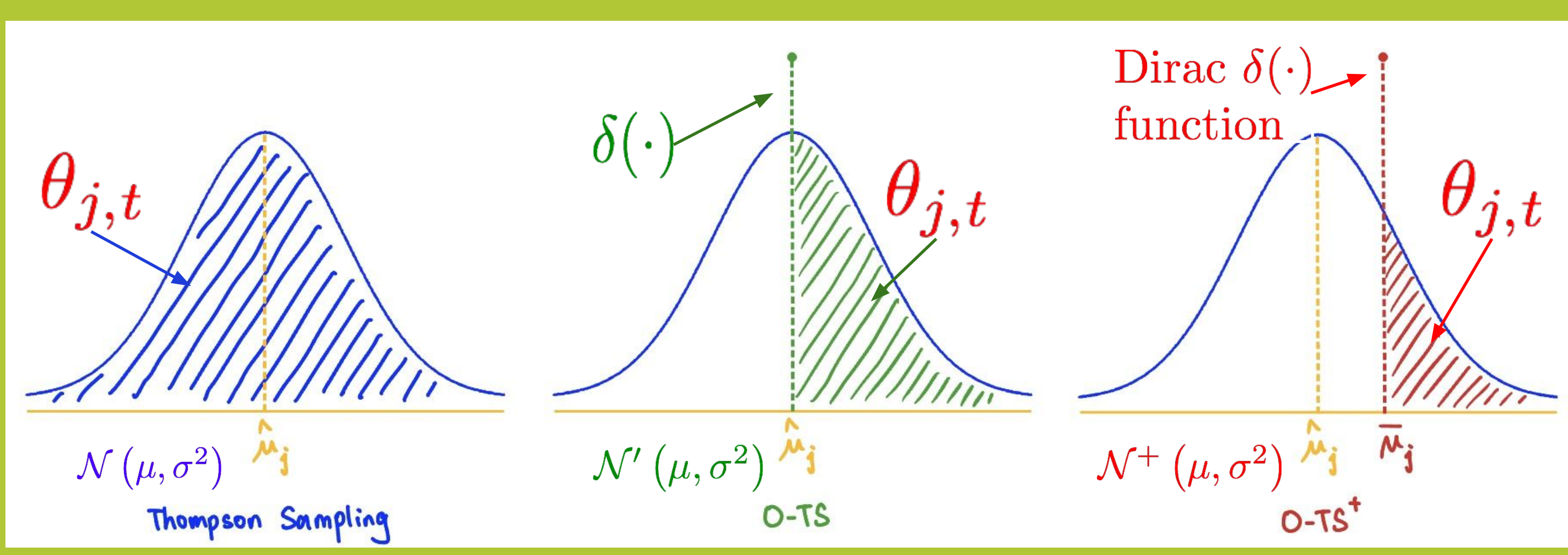
UPPER Confidence BOUND (UCB) vs Thompson Sampling (TS) in Bandits

Unknown parameters: $(\mu_1, \mu_2, \dots, \mu_K)$
Empirical parameters: $(\hat{\mu}_1, \hat{\mu}_2, \dots, \hat{\mu}_K)$

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

$$J_t = \arg \max \theta_{j,t}$$

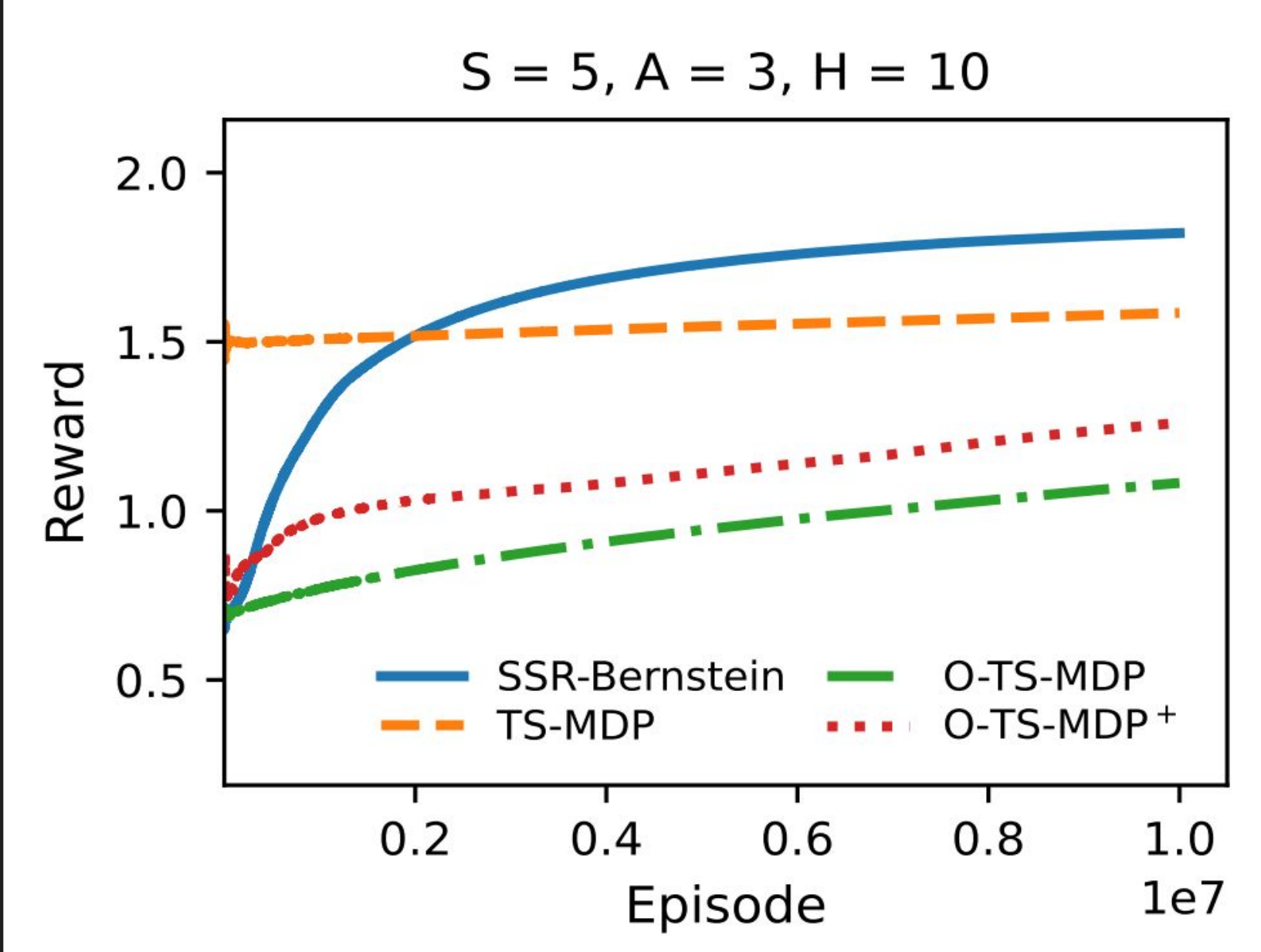


Regret
UPPER BOUND
 $O(\sqrt{KT \ln(T)})$

More Optimistic Distributions!!

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

Experiments for MDP



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 π_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: $\theta_{s,a,t}^k \sim \mathcal{N}'\left(\hat{\mu}_{s,a,t}^{k-1}, \tilde{O}\left(\frac{H^3 S}{O_{s,a,t}^{k-1}}\right)\right)$
 O-TS-MDP⁺: $M^k = \{[S], [A], H, \theta^k, \hat{P}^{k-1}\}$, where: $\theta_{s,a,t}^k \sim \mathcal{N}^+\left(\hat{\mu}_{s,a,t}^{k-1}, \tilde{O}\left(\frac{H^2}{O_{s,a,t}^{k-1}}\right)\right)$

Regret UPPER BOUND
 O-TS : $\tilde{O}(\sqrt{AS^2 H^4 T})$
 O-TS⁺ : $\tilde{O}(\sqrt{ASH^3 T})$

- 0-TS-MDP enjoys an elegant theoretical analysis, avoiding bounding the absolute value of approximation error.
- 0-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) and the Canada CIFAR AI Program.