

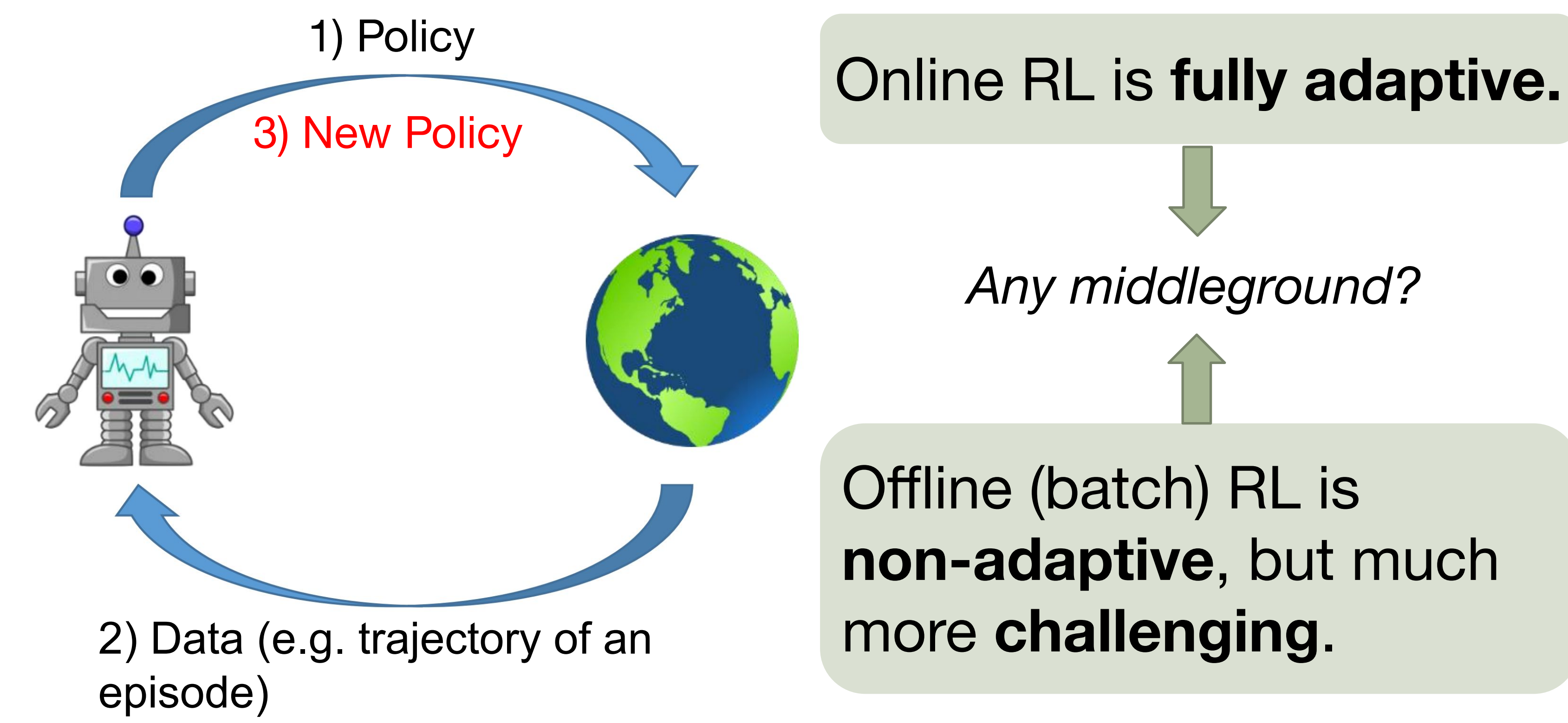
Provably Efficient Q-Learning with Low Switching Cost

Yu Bai^{1,2}, Tengyang Xie³, Nan Jiang³, Yu-Xiang Wang⁴
¹Stanford University ²Salesforce Research ³UIUC ⁴UC Santa Barbara



Motivation: RL with limited adaptivity?

- In many domains (recommendation, medical, ...), **deploying a new policy** is more prohibitive than gathering data with the existing policy.



Proposed framework: low switching cost RL

Setup: Episodic MDP with horizon H . RL algorithm plays K episodes ($T = K \cdot H$ steps.) Measure PAC/Regret.

Definition: the **switching cost** between two (deterministic) policies (π, π') is number of different actions they would take, (summed) for all (h, s) :

$$n_{\text{switch}}(\pi, \pi') := \#\{(h, s) \in [H][S] : \pi_h(s) \neq \pi'_h(s)\}$$

Definition: the **switching cost** of an RL algorithm that plays with policies π^1, \dots, π^K is

$$N_{\text{switch}} := \sum_{k=1}^{K-1} n_{\text{switch}}(\pi^k, \pi^{k+1})$$

Goal: fast exploration with low switching cost

Prior work: Q-Learning with UCB exploration: $\tilde{O}(\sqrt{H^{4.3}SAT})$ regret, but $N_{\text{switch}} = \Theta(HSK)$ **linear in K** 😞 [Jin et al. 2018]

Any low regret algorithm such that N_{switch} **sublinear in K** ?

Recap: UCB2 for bandits

Algorithm (UCB2): Repeat

- Select the arm that maximizes the UCB
- If this is the r -th time it's selected, play the arm exactly $\tau(r) - \tau(r-1)$ times, where $\tau(r) = (1 + \alpha)^r$

Theorem [Auer et al. 2002]: UCB2 achieves same regret as UCB, and only **log(K) policy switches**: $N_{\text{switch}} = O(A \log(K/A))$ 😊

Idea: Integrate UCB2 into Q-Learning!

Our Algorithm: Q-Learning with UCB2 scheduling

Key idea: update the policy only when Q has been updated $\tau(r) = (1 + \alpha)^r$ times.

Definition: The *triggering sequence* $\{t_n\}$ with parameter (α, r_*) is

$$\{t_n\}_{n \geq 1} = \{1, 2, \dots, \tau(r_*)\} \cup \{\tau(r_*) + 1, \tau(r_*) + 2, \dots\}$$

Algorithm 2 Q-learning with UCB2-Hoeffding (UCB2H) Exploration

input Parameter $\eta \in (0, 1)$, $r_* \in \mathbb{Z}_{>0}$, and $c > 0$.

Initialize: $\tilde{Q}_h(x, a) \leftarrow H$, $Q_h \leftarrow \tilde{Q}_h$, $N_h(x, a) \leftarrow 0$ for all $(x, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]$. // Two sets of Q: Running estimate \tilde{Q} Policy network Q

for episode $k = 1, \dots, K$ **do**

Receive x_1 .

for step $h = 1, \dots, H$ **do**

Take action $a_h \leftarrow \arg \max_{a'} Q_h(x_h, a')$, and observe x_{h+1} . // Take action according to Q

$t = N_h(x_h, a_h) \leftarrow N_h(x_h, a_h) + 1$;

$b_t = c\sqrt{H^3 \ell/t}$ (Hoeffding-type bonus);

$\tilde{Q}_h(x_h, a_h) \leftarrow (1 - \alpha_t)\tilde{Q}_h(x_h, a_h) + \alpha_t[r_h(x_h, a_h) + \tilde{V}_{h+1}(x_{h+1}) + b_t]$. // Update \tilde{Q} via Q-Learning

$\tilde{V}_h(x_h) \leftarrow \min \{H, \max_{a' \in \mathcal{A}} \tilde{Q}_h(x_h, a')\}$. // $\{t_n\}$ is the triggering sequence above

if $t \in \{t_n\}_{n \geq 1}$ (where t_n is defined in (1)) **then**

(Update policy) $Q_h(x_h, \cdot) \leftarrow \tilde{Q}_h(x_h, \cdot)$. // Set Q to be \tilde{Q} occasionally according to UCB2 scheduling

end if

end for

end for

Theoretical Result

Theorem 1: Our Q-Learning with UCB2-{Hoeffding, Bernstein} exploration achieves $\tilde{O}(\sqrt{H^{4.3}SAT})$ regret and **logarithmic switching cost**:

$$N_{\text{switch}} \leq O(H^3 SA \log(K/A))$$

Proof highlight: analysis of error propagation under *delayed* Q updates.

Application: concurrent /parallel RL

Setup: M agents play an episode in parallel, and can only *communicate after each episode*.



Idea: policy not scheduled to switch in M episodes \Rightarrow can parallelize to M non-communicating agents

Theorem 2 (Nearly linear speedup in PAC concurrent RL): There exists concurrent versions of our algorithm, s.t. given M machines it can find ϵ optimal policy in $\tilde{O}\left(H^3 SA + \frac{H^{\{5,4\}} SA}{\epsilon^2 M}\right)$ rounds.

\rightarrow Also improves upon prior work [Guo et al. 2015] in (H, S, ϵ) dependence.

Lower bound on switching cost

Simple Observation: you “need” to switch $HS(A-1)$ times to at least try out all the possible actions to take.

Theorem 3 (lower bound): Any algorithm that has switching cost $N_{\text{switch}} \leq HSA/2$ has to suffer from **linear (trivial) worst-case regret**:

$$\sup_{\text{MDP}} \text{Regret}(K) \geq KH/4$$

Remark: Our algorithm achieves $N_{\text{switch}} = \tilde{O}(H^3 SA)$, so still an H^2 gap between the lower and upper bounds.

Discussion & future work

- Close the gap on the switching cost.
- Alternative notions of limited adaptivity:
 - Hard constraint on switching cost.
 - RL with only $O(1)$ rounds of adaptivity.
- Connections to fully offline/batch RL.