

Near Optimal Exploration-Exploitation in Non-Communicating Markov Decision Processes

Ronan Fruit (INRIA), Matteo Pirotta (INRIA), Alessandro Lazaric (FAIR)



Motivations

- Learning in an **unknown** environment means to balance
 - Exploration
 - Exploitation
- All theoretically-grounded approaches requires **prior knowledge**
 - This information is hard to get!
 - Limit their applicability
- This is about learning without prior **knowledge!!**

Online Learning in MDPs

- Markov Decision Process** $M = \{S, A, r, p\}$
 - states: $S = S^c \cup S^T$ → **communicating** set: $\forall s, s' \in S^c, \exists \pi : \mathbb{P}^\pi(s \rightarrow s') > 0$
 - actions: $A = (A_s)_{s \in S}$
 - mean rewards: $r(s, a)$ → **transient** set: $S^c \cap S^T = \emptyset$
 - transition probabilities: $p(s'|s, a)$
 - Possible next states: $\Gamma^S = \max_{s \in S, a \in A_s} \|p(\cdot|s, a)\|_0$

Optimality criterion: long-term average reward

For any policy $\pi \in \Pi^{\text{SR}}(M)$ starting from $s \in S$:

$$\text{GAIN: } g_M^\pi(s) := \lim_{T \rightarrow +\infty} \mathbb{E} \left[\frac{1}{T} \sum_{t=1}^T r(s_t, a_t) \right]$$

$$\text{BIAS: } h_M^\pi(s) := C\text{-}\lim_{T \rightarrow +\infty} \mathbb{E} \left[\sum_{t=1}^T (r(s_t, a_t) - g_M^\pi(s_t)) \right]$$

In **weakly communicating** MDPs: any **optimal policy** $\pi^* \in \arg \max_{\pi} \{g^\pi(s)\}$ has **constant** gain

Learning problem: cumulative regret minimization

The true M^* is unknown, thus it is g^* $\Delta(\mathcal{A}, T) = Tg^* - \sum_{t=1}^T r_t(s_t, a_t)$

Asm. 1 The initial state $s_1 \in S^c$

Diameter and Span: [Jaksch et al. 2010; Bartlett and Tewari, 2009]

$$D^S = \max_{s, s' \in S} \left\{ \min_{\pi: S \rightarrow \mathcal{P}(A)} \mathbb{E}_\pi [T(s'|s)] \right\}$$

$$sp_S \{h^*\} = \max_{s \in S} \{h^*(s)\} - \min_{s \in S} \{h^*(s)\}$$

- D depends on **all** policies (*global property*)
- $sp_S \{h^*\}$ on **only** π^*
- $sp_S \{h^*\} \leq D$ (always)

In weakly communicating MDPs $D = \infty$ but $sp_S \{h^*\} \leq \infty$

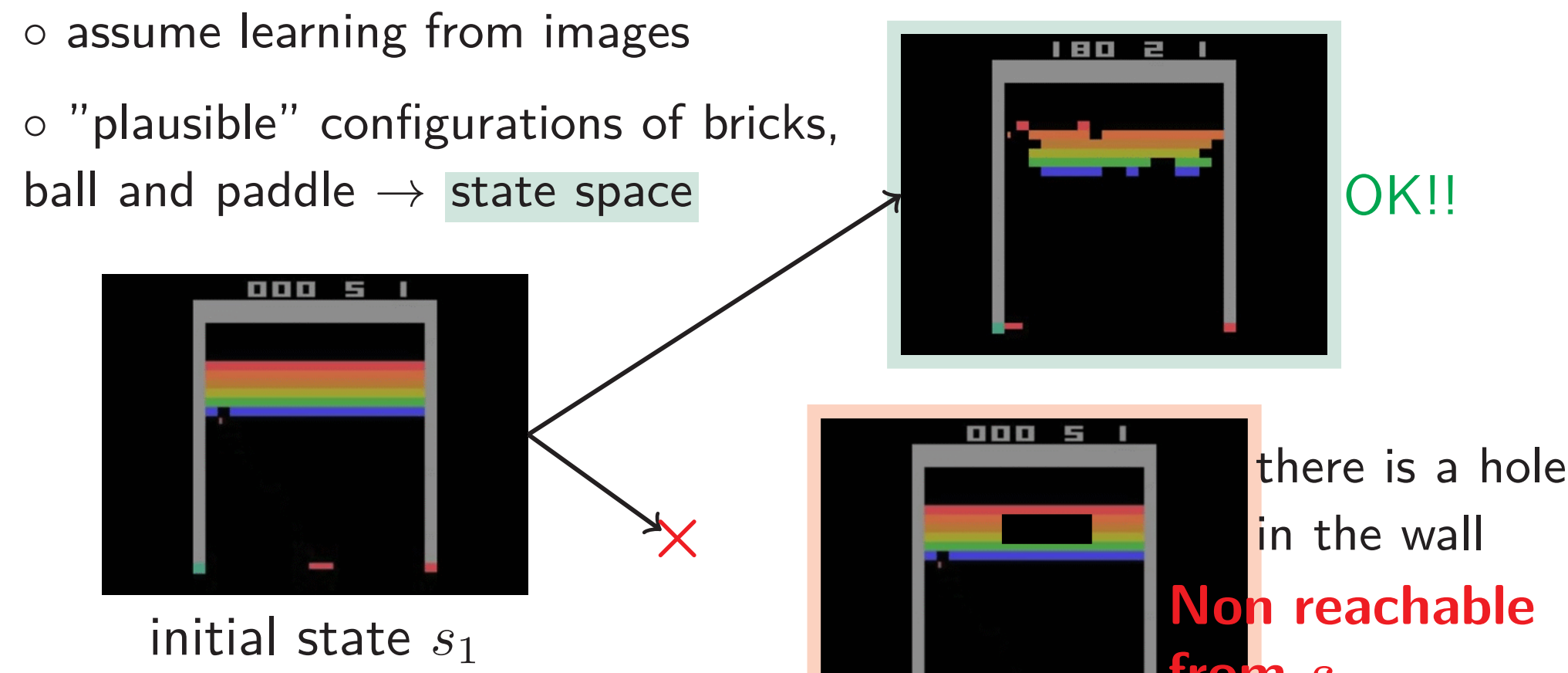
References

- Jaksch, Ortner, and Auer. Near-optimal regret bounds for reinforcement learning. Journal of Machine Learning Research, 2010.
- Bartlett and Tewari. REGAL: A regularization based algorithm for reinforcement learning in weakly communicating MDPs. In UAI 2009.
- Agrawal and Jia. Optimistic posterior sampling for reinforcement learning: worst-case regret bounds. In NIPS 2017, 2017.
- Fruit, Pirotta, Lazaric and Ortner. Efficient Bias-Span-Constrained Exploration-Exploitation in Reinforcement Learning. In ICML 2018, 2018.
- Dietterich. Hierarchical reinforcement learning with the MAXQ value function decomposition. J. Artif. Intell. Res., 13:227-303, 2000.

Prior Knowledge & Misspecified states

- UCRL and OPT-PSRL assume **communicating** MDPs
 - \Rightarrow all states are reachable
 - $\Rightarrow D < +\infty$ **regret:** $\tilde{O}(D^S \sqrt{\Gamma^S S A T}) / \tilde{O}(D^S \sqrt{S A T})$
- reasonable! but rarely verified in practice

BREAKOUT EXAMPLE



Problem: misspecified state!

- In weakly communicating or misspecified problems ($D = +\infty$) UCRL and OPT-PSRL \rightsquigarrow linear regret
- REGAL.C and SCAL exploits the knowledge $sp_S \{h_{M^*}^*\} \leq H$ implicitly “removes” non-reachable states able to learn in weakly comm. MDPs **regret:** $\tilde{O}(H \sqrt{\Gamma^S S A T})$
- knowing H not easier than designing well-specified states

*similar assumptions in Bayesian regret

Truncated Upper-Confidence for Reinforcement Learning (TUCRL)

Plain OFU (e.g., UCRL) executed on

- S and $S^T \neq \emptyset$ is over-exploring (*linear regret*)
- $S_k^c \subseteq S^c$ may under-exploration (*linear regret*)

$$S_k^c := \left\{ s \in S \mid \sum_{a \in A_s} N_k(s, a) > 0 \right\} \cup \{s_{t_k}\}$$

Empirical estimate of S^c

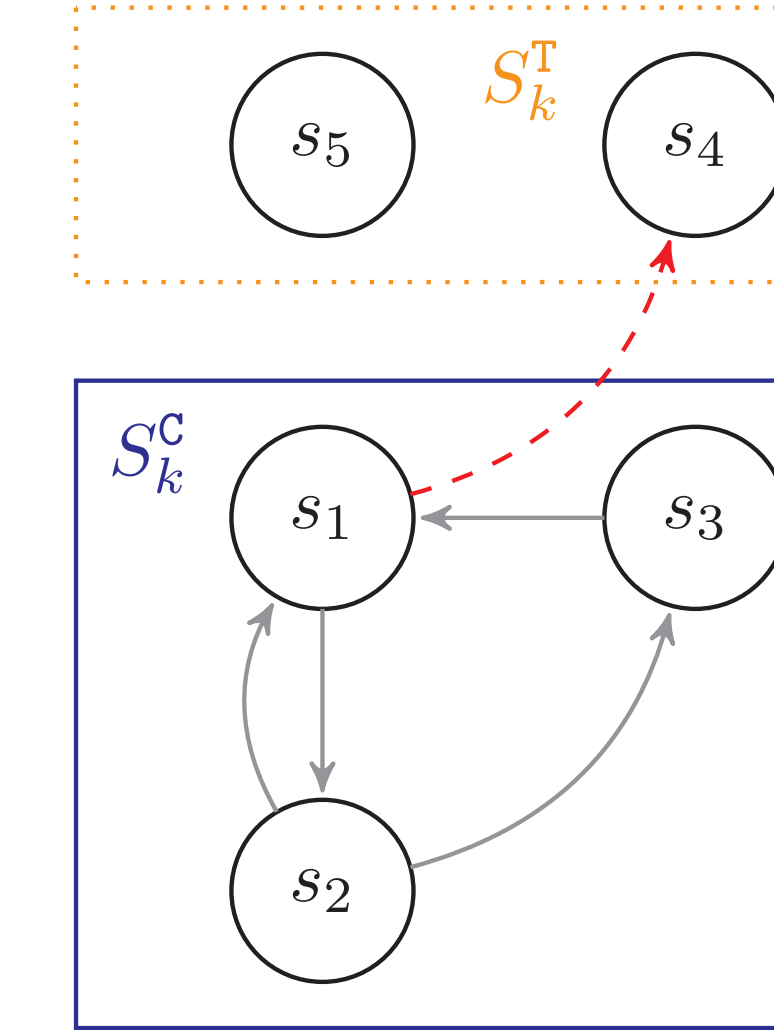
s_{t_k} : starting state of episode k

★ TUCRL learns in weakly communicating or misspecified problems **WITHOUT PRIOR KNOWLEDGE**

WHY DOES TUCRL WORK?

- reconsiders transitions periodically (i.e., decrease ρ_k) \Rightarrow avoid *under-exploration*
- exploits Bernstein confidence interval:

$$\beta_{p,k}^{sas'} = \sqrt{\frac{\alpha \sigma_k^2(s'|s, a)}{N_k(s, a)}} + \frac{\beta}{N_k(s, a)}$$
 - $(s, a) \rightarrow s'$ not observed $\Rightarrow \sigma_k^2 = 0$ and $\hat{p}_k = 0$ $\Rightarrow \beta/N_k(s, a) < \rho_k$
 - fast shrinking $\approx O(1/N_k(s, a))$
- we set $\rho_k = O(SA/t_k)$ equivalent to removing transitions s.t.: $N_k(s, a) > \sqrt{t_k/SA} \Rightarrow \tilde{p}(s'|s, a) = 0, \forall s' \in S_k^T$



Does this transition $s_1 \rightarrow s_4$ exists? Is $s_4 \in S_k^T$? Should be enabled it at episode k ? Explore or not?

At episode k we know:

$$\hat{p}_k(s'|s, a) \text{ and } |\tilde{p}(s'|s, a) - \hat{p}(s'|s, a)| \leq \beta_{p,k}^{sas'}$$

empirical mean confidence interval

TUCRL idea:

“guess” a lower bound ρ_k to the trans. probabilities $\tilde{p}(s'|s, a) + \beta_{p,k}^{sas'} < \rho_k \Rightarrow s \rightarrow s'$ **FORBIDDEN**

Maximum probability given conf. interval

⚠ $(\rho_k)_k$ should be non-increasing!!

TUCRL ALGORITHM

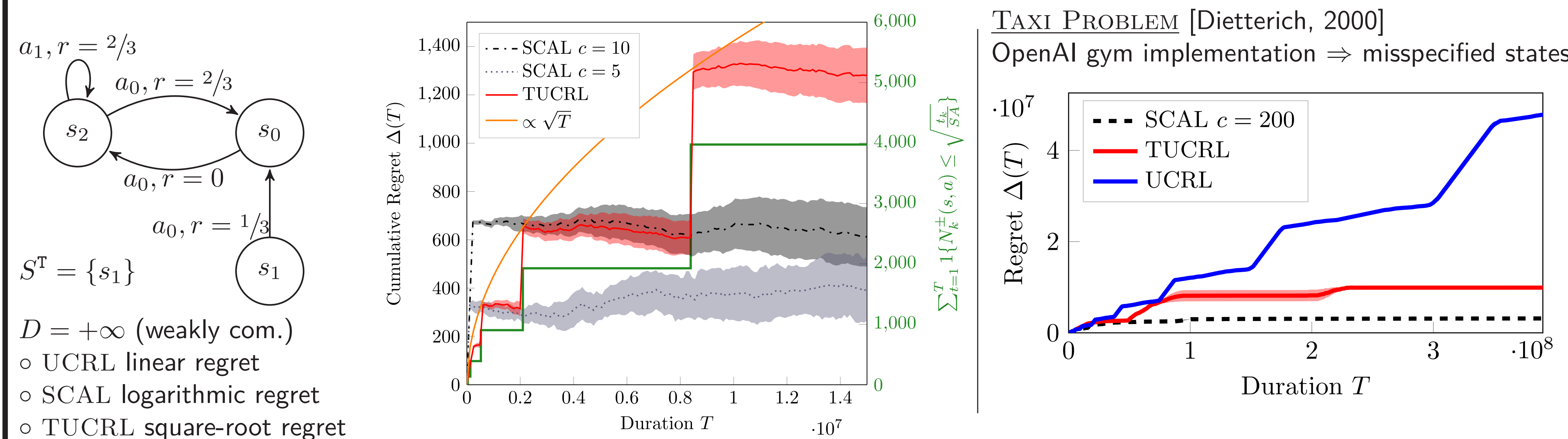
For episode $k = 1, 2, \dots$

- Confidence set:** TUCRL builds $\mathcal{M}_k = \{M = (S, A, \tilde{r}, \tilde{p}) : \tilde{r}(s, a) \in B_{r,k}(s, a), \tilde{p}(\cdot|s, a) \in B_{p,k}(s, a)\}$ with:

$$B_{p,k}(s, a) = \{\tilde{p}(\cdot|s, a) : \|\tilde{p}(\cdot|s, a) - \hat{p}(\cdot|s, a)\|_1 \leq \sum_{s'} \beta_{p,k}^{sas'}\}$$

$$\cap \{\tilde{p}(\cdot|s, a) : N_k(s, a) > \sqrt{t_k/SA} \Rightarrow \forall s' \in S_k^T, \tilde{p}(s'|s, a) = 0\}$$
- Planning:** TUCRL solves $(\tilde{M}_k, \tilde{\pi}_k) = \arg \max_{M \in \mathcal{M}_k, \pi} \{g_M^\pi\}$
- Execution:** of policy $\tilde{\pi}_k$ in the true MDP

Numerical Experiments



Regret of TUCRL

$$\Delta(\text{TUCRL}, T) = \tilde{O} \left(D^c \sqrt{\Gamma^c S^c A T} + \left(D^c \right)^2 S^3 A \right)$$

- Adaptability** to communicating part $D^c := D^{S^c}, \Gamma^c := \Gamma^{S^c}$ and S^c
- Regret due to the **early stage** where TUCRL suffers linear regret

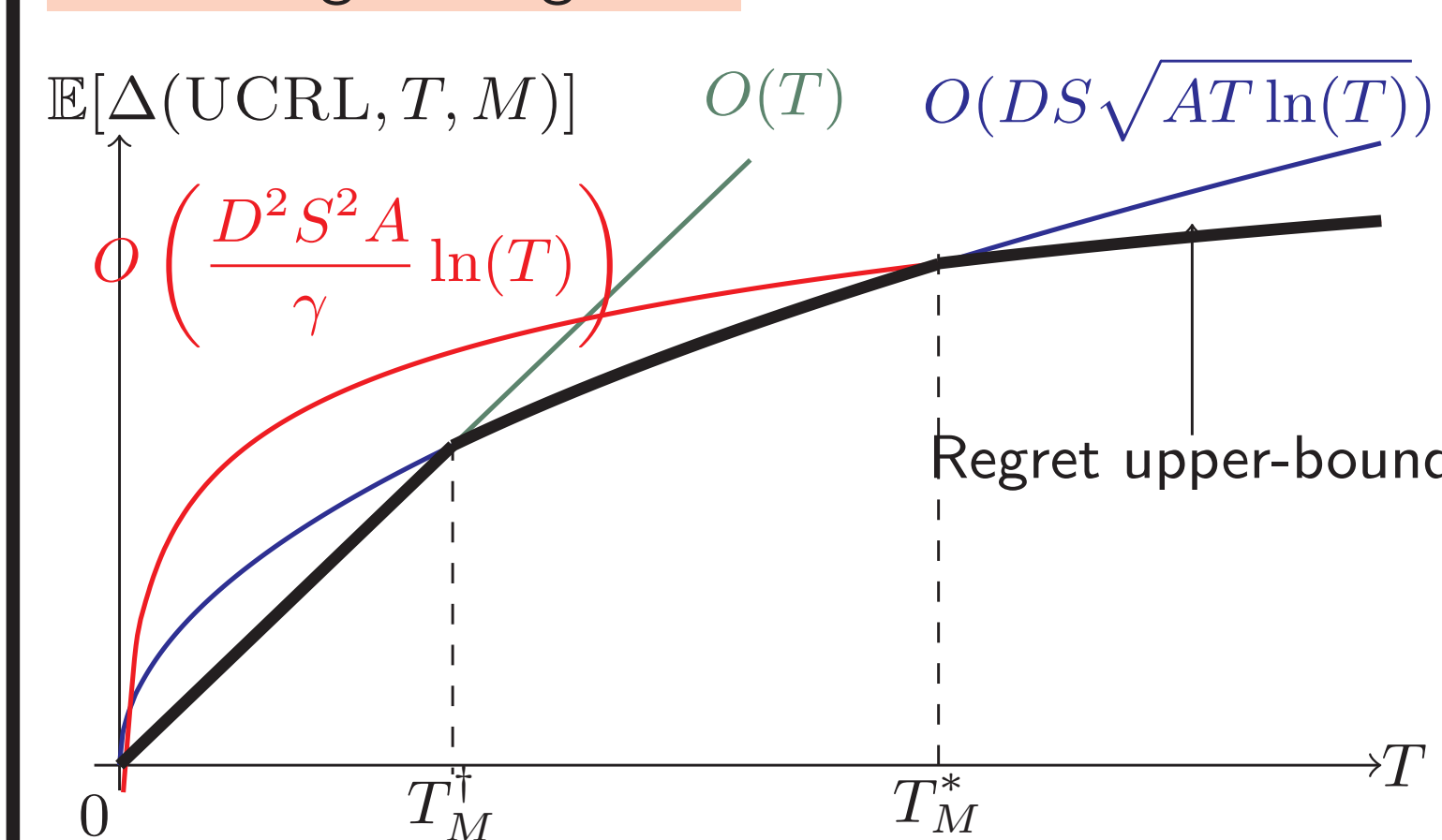
If M^* is **communicating**

First term: same as UCRL

Second term: bigger than UCRL by a factor S^c/Γ^c

Exploration-Exploitation Dilemma with infinite diameter

UCRL regret “regimes”



efficient algorithm: time T_M^\dagger to achieve sub-linear regret is *polynomial* in the parameters of the MDP

In communicating MDPs, UCRL achieves sublinear regret in $T_M^\dagger = \tilde{O}((D^S)^2 \Gamma^S S A) \Rightarrow$ *efficient algorithm*

Impossibility result

Without prior knowledge, any efficient learning algorithm must satisfy $T_M^* = +\infty$ when M has *infinite diameter* (i.e., it cannot achieve logarithmic regret)

No logarithmic regret without prior knowledge!

