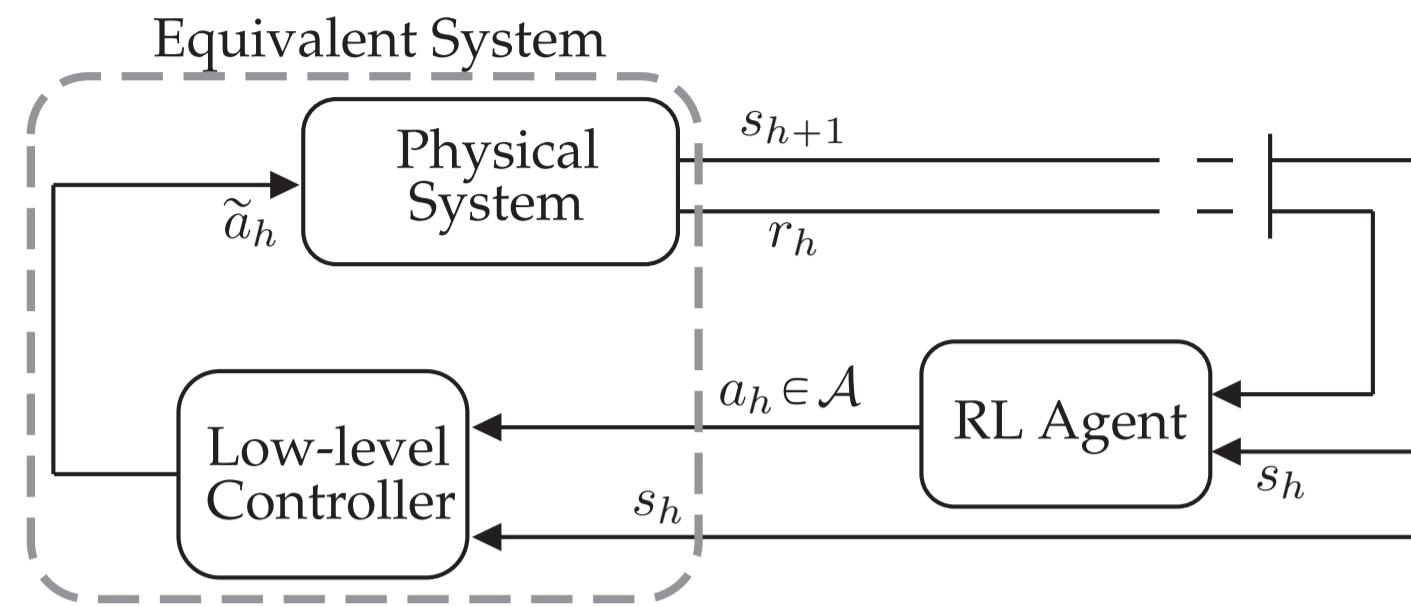
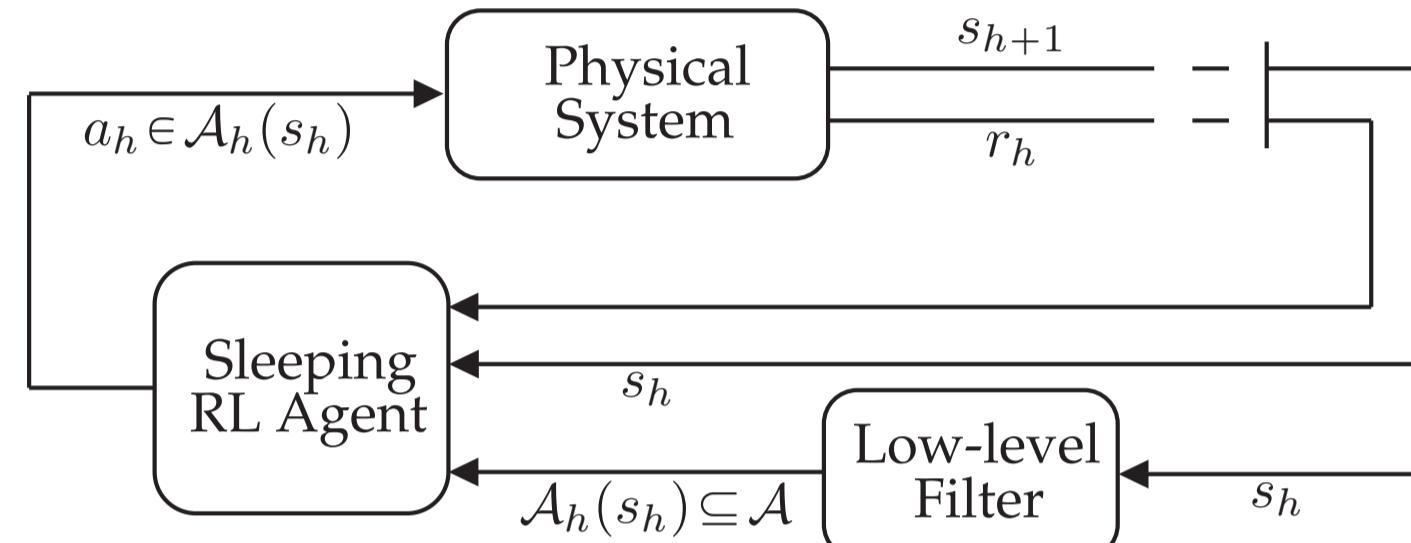
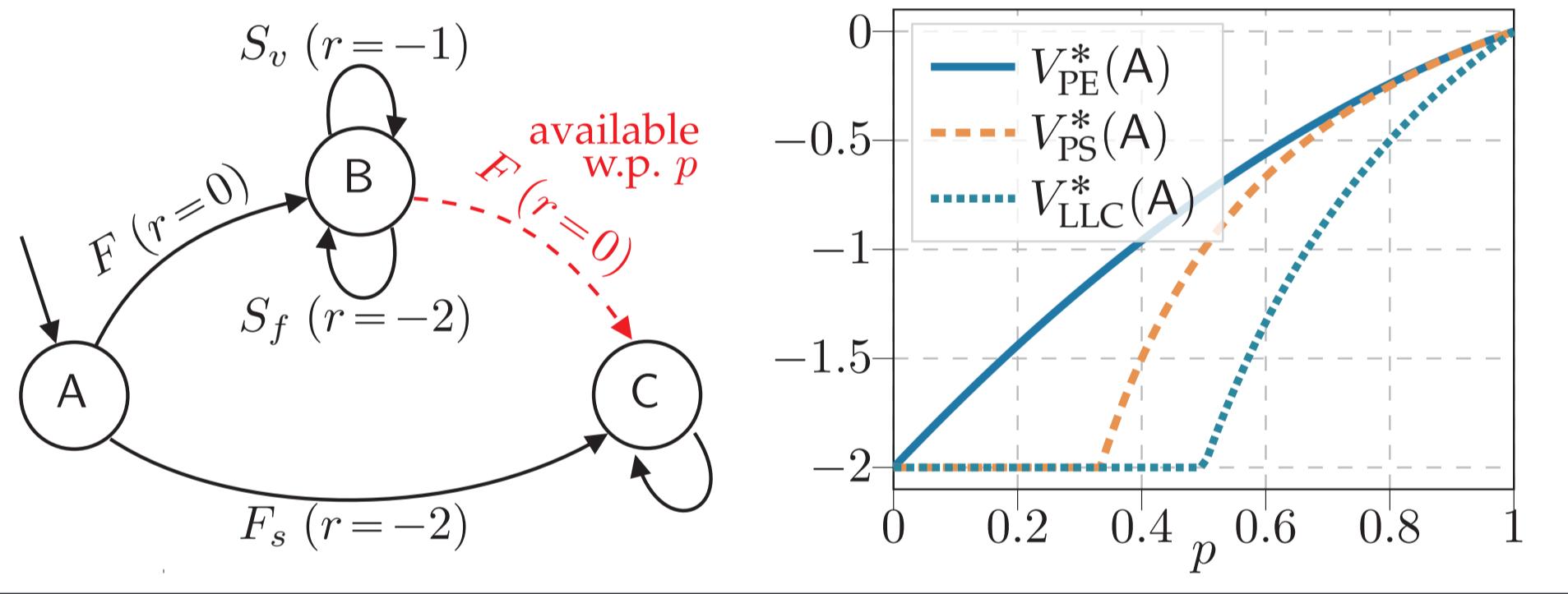


**MOTIVATION**
**CLASSIC REINFORCEMENT LEARNING**

**SLEEPING REINFORCEMENT LEARNING**

**EXAMPLE**

**SETTING – SLEEPING MDPs**

$$\mathcal{M} := (\underbrace{\mathcal{S}, \mathcal{A}, P, R, H, \bar{s}}_{\text{Standard Finite States and Actions MDP}}, \underbrace{C}_{\text{Stochastic Action Availability Function}})$$

We consider an **episodic** scenario and we call  $K$  the number of episodes

**Available action set**  $\mathcal{A}_{k,h}(s) \subseteq \mathcal{A}$  is selected according to function  $C$

**ACTION DISCLOSURE MODELS**

- **Per-episode disclosure:** at the beginning of the current episode
- **Per-stage disclosure:** when the agent is asked to choose an action

**REFERENCES**

- M. G. Azar, I. Osband, and R. Munos. Minimax regret bounds for reinforcement learning. In *International Conference on Machine Learning*, 2017.
- O. D. Domingues, P. Ménard, E. Kaufmann, and M. Valko. Episodic reinforcement learning in finite mdps: Minimax lower bounds revisited. In *Algorithmic Learning Theory*, 2021.
- R. D. Kleinberg, A. Niculescu-Mizil, and Y. Sharma. Regret bounds for sleeping experts and bandits. In *Conference on Learning Theory*, 2008.

**PER-EPIISODE DISCLOSURE**
**PER-EPIISODE REGRET**

Given algorithm  $\mathfrak{A}$  running  $(\pi_k)_{k=1}^K$ , we evaluate its performances in terms of **per-episode (PE) regret** w.r.t. the **optimal policy**  $\pi_k^*$  for each episode:

$$R_{\text{PE}}(\mathfrak{A}, K) := \sum_{k \in [K]} (V_{k,1}^{\pi_k^*}(\bar{s}) - V_1^{\pi_k}(\bar{s}))$$

**LOWER BOUND**

Same as for MDPs ( $|\mathcal{S}| = S$ ,  $|\mathcal{A}| = A$ ):

$$\mathbb{E}[R_{\text{PE}}(\mathfrak{A}, K)] \geq \Omega(H^{3/2}\sqrt{SAK})$$

Regardless of  $C$ .

**ALGORITHM (AR-UCBVI)**

We extend UCBVI to handle the action availabilities. AR-UCBVI is the **optimistic counterpart** of performing an **action-restricted value iteration**. This solution matches the LB (for every  $C$ ):

$$\mathbb{E}[R_{\text{PE}}(\text{AR-UCBVI}, K)] \leq \tilde{\mathcal{O}}(H^{3/2}\sqrt{SAK})$$

**PER-STAGE DISCLOSURE**
**PER-STAGE REGRET**

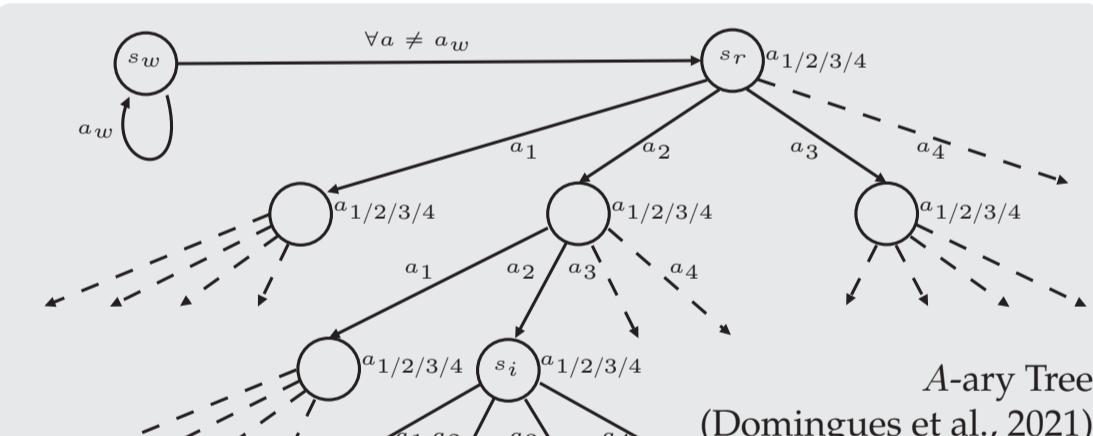
Given algorithm  $\mathfrak{A}$  running  $(\pi_k)_{k=1}^K$ , we evaluate its performances in terms of **per-stage (PS) regret** w.r.t. the **best constant policy**  $\pi^*$ :

$$R_{\text{PS}}(\mathfrak{A}, K) = KV_1^{\pi^*}(\bar{s}) - \sum_{k \in [K]} V_1^{\pi_k}(\bar{s})$$

**MARKOVIAN AVAILABILITY ( $C = C^{\text{Markov}}$ )**

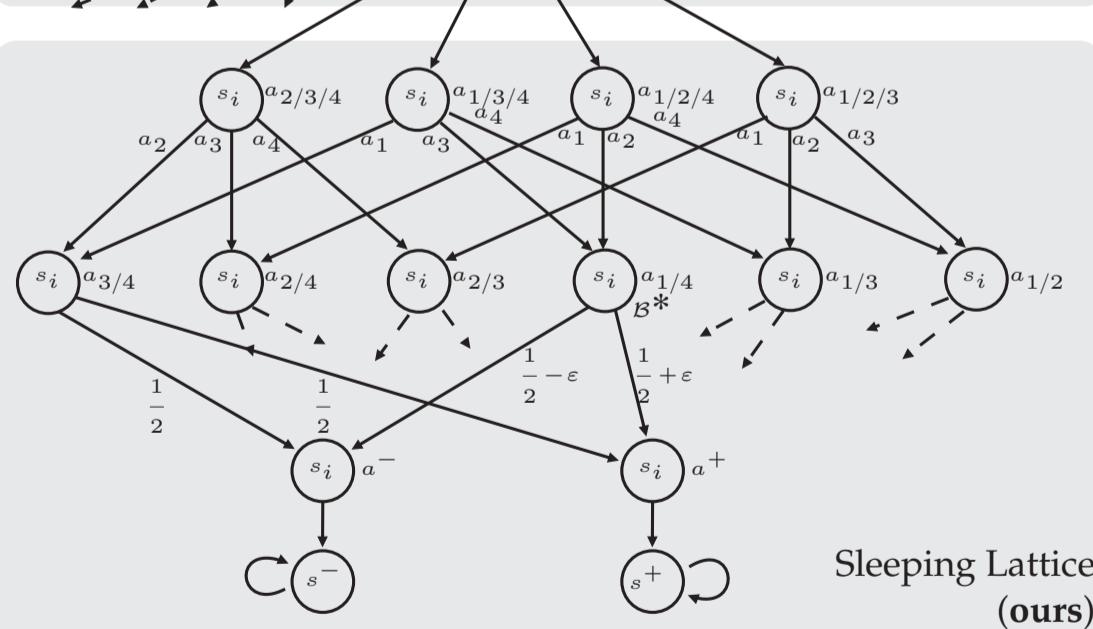
At every stage  $h$ , the **available actions**  $\mathcal{A}_h(s_h)$  depend on:

- current state  $s_h$
- previous state  $s_{h-1}$  and action  $a_{h-1}$
- previous available actions  $\mathcal{A}_{h-1}(s_{h-1})$

**LOWER BOUND**


The regret  $\mathbb{E}[R_{\text{PS}}(\mathfrak{A}, K)]$  that every algorithm suffer is:

$$\Omega(H^{3/2}\sqrt{SAK}2^{A/2})$$


**HOW TO LEARN?**

We can map the Sleeping MDP with an **Augmented MDP** in which we encode action availability sets in the state representation. This MDP will present an **augmented state space**  $\tilde{\mathcal{S}} := \mathcal{S} \times \mathcal{P}(\mathcal{A})$  with cardinality  $S2^A$ .

Given this mapping, we can run UCBVI with Bernstein-Freedman bonus and get:

$$\mathbb{E}[R_{\text{PS}}(\text{UCBVI}, K)] \leq \tilde{\mathcal{O}}(H^{3/2}\sqrt{SAK}2^A)$$

**INDEPENDENT AVAILABILITY ( $C = C^{\text{ind}}$ )**

At every stage  $h$ , the **available actions**  $\mathcal{A}_h(s_h)$  depend only on **state  $s_h$**

**LOWER BOUND**

$$\mathbb{E}[R_{\text{PS}}(\mathfrak{A}, K)] \geq \Omega(H^{3/2}\sqrt{SAK})$$

**ALGORITHM (S-UCBVI)**

S-UCBVI extends UCBVI to learn action availabilities  $C_h^{\text{ind}}(\cdot|s)$ . It performs **optimism** on the estimates of both **state-action**  $\hat{Q}_{k,h}(s, a)$  and **state**  $\hat{V}_{k,h}(s)$  **value functions**:

$$\hat{Q}_{k,h}(s, a) = R_h(s, a) + \sum_{s' \in \mathcal{S}} \hat{P}_{k,h}(s'|s, a) \hat{V}_{k,h+1}(s') + b_{k,h}^Q(s, a)$$

$$\hat{V}_{k,h}(s) = \sum_{\mathcal{B} \in \mathcal{P}(\mathcal{A})} \hat{C}_{k,h}^{\text{ind}}(\mathcal{B}|s) \max_{a \in \mathcal{B}} \hat{Q}_{k,h}(s, a) + b_{k,h}^V(s)$$

where  $b_{k,h}^Q(s, a)$  and  $b_{k,h}^V(s)$  are **Bernstein-Freedman** bonuses to account for the **uncertainty** on  $P$  and  $C$ . S-UCBVI enjoys a regret of:

$$\mathbb{E}[R_{\text{PS}}(\text{S-UCBVI}, K)] \leq \tilde{\mathcal{O}}(H^{3/2}\sqrt{SAK})$$

**EXPERIMENTS - STOCHASTIC FROZEN LAKE**

Frozen Lake with holes appearing during the episode with probability  $p$

