

Regret Minimization in MDPs with Options without Prior Knowledge

Ronan Fruin, Matteo Pirotta, Alessandro Lazaric, Emma Brunskill

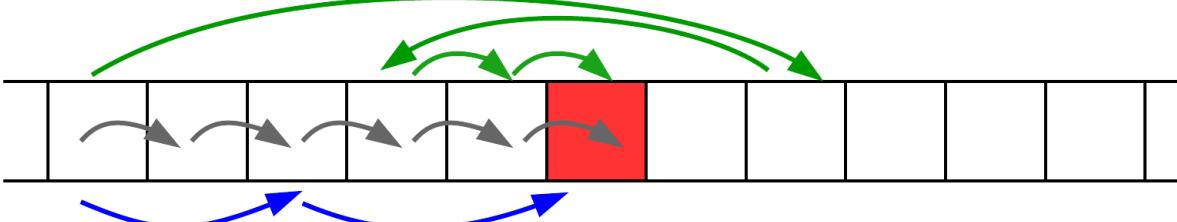


Motivations

- "Flat" RL: difficult to learn **complex behaviours** (eg, sequence of subgoals) \Rightarrow Humans abstract from **low-level actions**
- Hierarchical RL: decompose large problems into smaller ones by imposing constraints on value function or **policy**
- Possible implementation: **options** [Sutton et al., 1999]
- Empirical observations: introducing options in an MDP can **speed up** learning but can also be **harmful** [Jong et al., 2008].
 \Rightarrow Lack of theoretical motivation and understanding of options

Theoretical analysis of learning with options

- Adding options does not just reduce the space of stationary policies, the **exploration** is also greatly affected



- Navigability and **temporal abstraction** can be preserved in an UCRL-like algorithm when **prior knowledge** about options is given [Fruit and Lazaric, 2017]

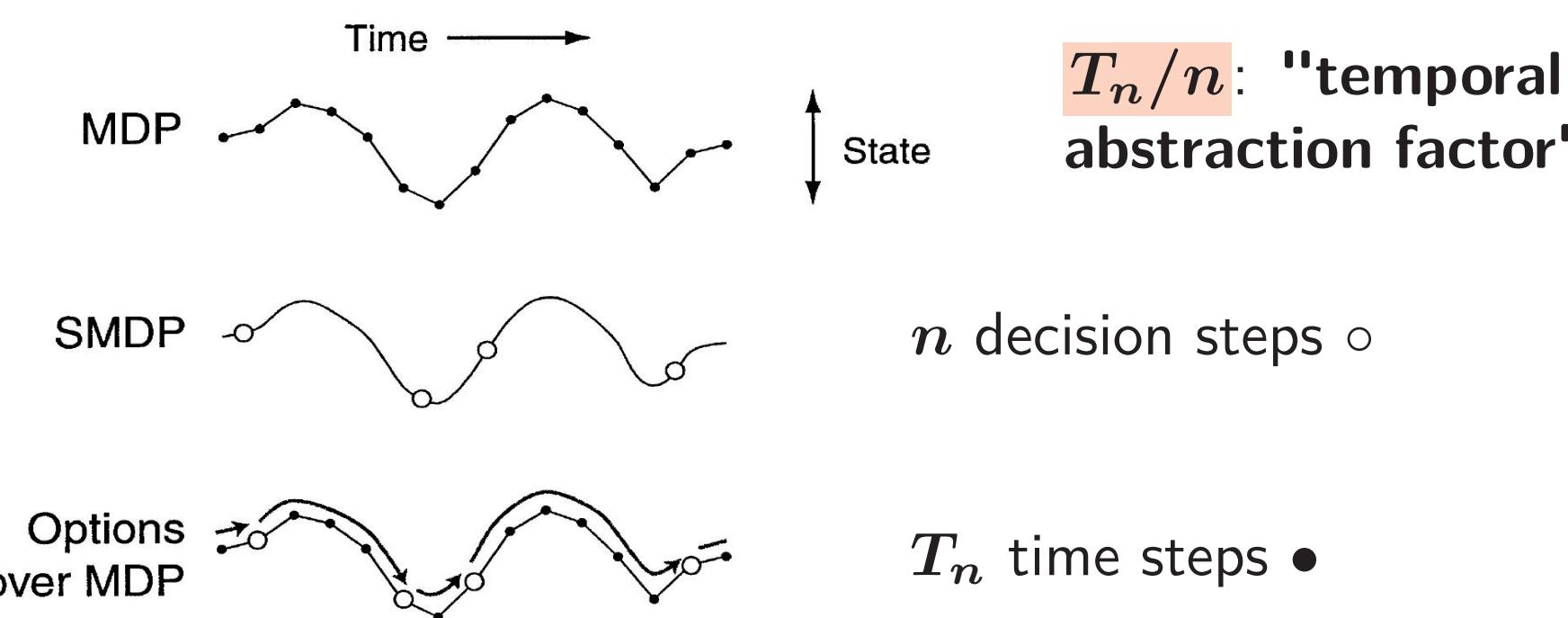
Can we preserve temporal abstraction **without** prior knowledge?

Temporal Abstraction in MDPs

Markov Decision Process

$$M = \{S, A, r, p\}$$

- **S**: states
- **A** = $(A_s)_{s \in S}$: actions
- $r(s, a)$: rewards
- $p(s'|s, a)$: transition probabilities



Proposition 1. A set of options \mathcal{O} defined on an MDP M induces a semi-MDP $M' = \{S_{\mathcal{O}}, \mathcal{O}, R_{\mathcal{O}}, p_{\mathcal{O}}, \tau_{\mathcal{O}}\}$ where the cumulative reward $R_{\mathcal{O}}$ and holding time $\tau_{\mathcal{O}}$ of every option $o \in \mathcal{O}$ are sub-exponential random variables with parameters $\sigma_R(o)$ and $\sigma_{\tau}(o)$.

Online Learning for SMDPs

Optimality criterion

An optimal policy π^* maximizes the **long-term average reward**

$$\rho^*(M') = \max_{\pi} \left\{ \lim_{n \rightarrow +\infty} \mathbb{E}^{\pi} \left[\frac{R_n}{T_n} \right] \right\}; \quad T_n = \sum_{i=1}^n \tau_i; \quad R_n = \sum_{i=1}^n r_i$$

Performance Measure

The learning agent aims at **minimizing the cumulative regret**

$$\Delta(M', n) = \rho^*(M')T_n - R_n$$

Diameter of an SMDP

$$D(M') = \max_{s, s' \in S} \left\{ \min_{d \in D_M^{MD}} \left\{ \mathbb{E}^{d \infty} [T(s') | s_0 = s] \right\} \right\}$$

$$\text{where } T(s') = \inf \left\{ \sum_{i=1}^n \tau_i : n \in \mathbb{N}, s_n = s' \right\}$$

FREE-PARAMETER SMDP-UCRL

SMDP optimistic optimality equation:

$$\tilde{\rho}_{SUCRL}^* = \max_{o \in \mathcal{O}_s} \left\{ \max_{\tilde{R}, \tau} \left\{ \frac{\tilde{R}(s, o)}{\tau(s, o)} + \frac{1}{\tau(s, o)} \left(\max_p \left\{ \sum_{s' \in S} \tilde{p}(s'|s, o) \tilde{u}^*(s') \right\} - \tilde{u}^*(s) \right) \right\} \right\} \quad \text{Exploits } \sigma_R \text{ and } \sigma_{\tau} \text{ to compute } \tilde{R} \text{ and } \tilde{\tau}$$

An option can be seen as an **absorbing Markov Chain** by merging initial and absorbing states \Rightarrow Irreducible MC

Irreducible Transformation τ : first-hitting-time of $\{s_3, s_4\}$
 $R = \sum_{t=1}^{\tau} r(s_t, \pi_o(s_t))$

$$\tau: \text{First-return-time in } s_0 \Rightarrow \begin{cases} \bar{\tau} = 1/\mu(s_0) \\ \bar{R}/\bar{\tau} = r^{\tau}\mu \end{cases}$$

Irreducible optimistic optimality equation:

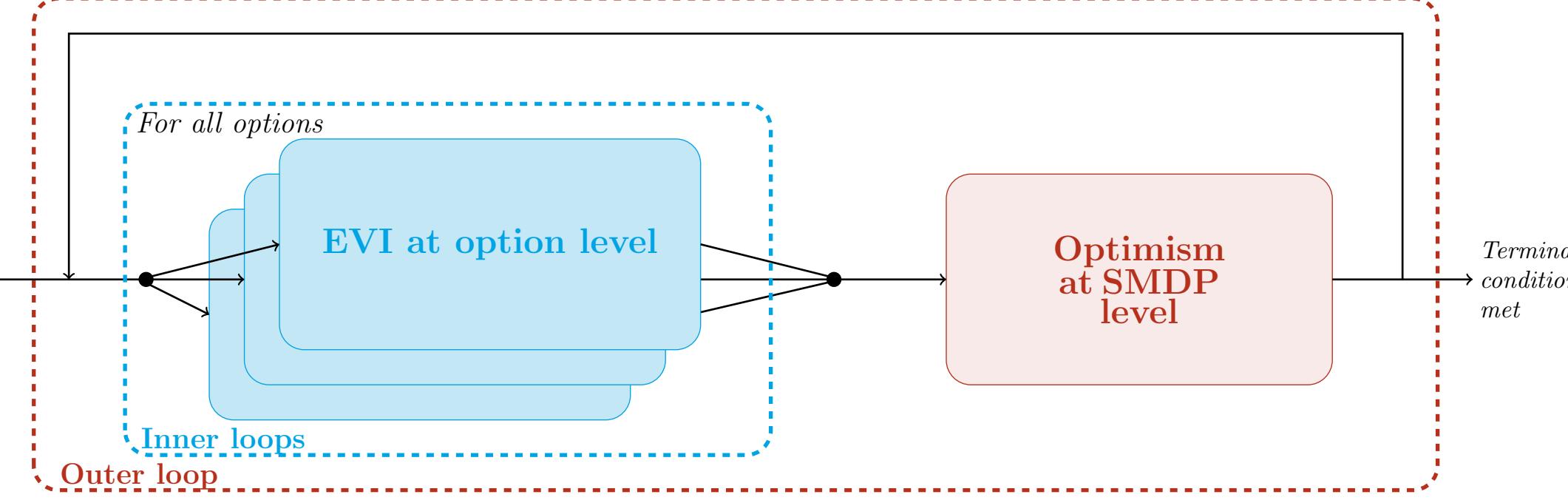
$$\tilde{\rho}_{FSUCRL}^* = \max_{o \in \mathcal{O}_s} \left\{ \max_{\tilde{\mu}_o, \tilde{r}_o} \left\{ \sum_{s' \in S_o} \tilde{\rho}_o(s') \tilde{\mu}_o(s') + \tilde{\mu}_o(s) \left(\max_{\tilde{b}_o} \{\tilde{b}_o \tilde{u}^*\} - \tilde{u}^*(s) \right) \right\} \right\} \quad \text{Exploits knowledge acquired over time}$$

HOW DO WE SOLVE THE PROBLEM?

FSUCRL

- 1: **for** episodes $k = 1, 2, \dots$ **do**
- 2: Compute ε_k -approximation of the optimal policy $\tilde{\pi}_k$ giving $\tilde{\rho}_{FSUCRL}^*$
 - Use Extended Value Iteration for each option
 - Use optimism at option level
- 3: Execute $\tilde{\pi}_k$ to acquire new experience
- 4: **end for**

NESTED EXTENDED VALUE ITERATION



Regret bound in an SMDP [Fruit and Lazaric, 2017]

$$\Delta(M', T_n) = O \left((D' \sqrt{S'} + C(M', n, \delta)) \sqrt{S' A' n \log \left(\frac{n}{\delta} \right)} \right)$$

where $C(M', n, \delta)$ depends on knowing $\sigma_R(o)$ and $\sigma_{\tau}(o)$.

Advantages of temporal abstraction: $n \ll T_n \Rightarrow \Delta(M', T_n) \leq \Delta(M, T_n)$ if $D' \approx D$

Special case: Bound for MDPs [Jaksch et al., 2010]:

$$\Delta(M, T_n) = O \left(DS \sqrt{AT_n \log \left(\frac{T_n}{\delta} \right)} \right)$$

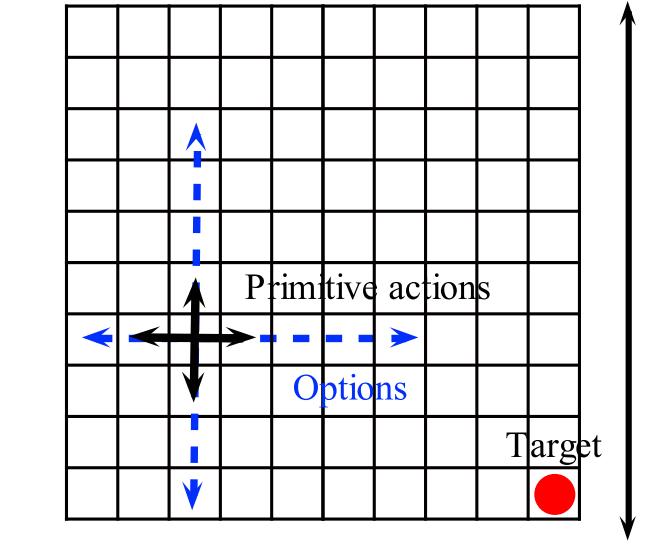
To assess the potential benefit of options, we study the ratio of the regrets SMDP/MDP:
 $\mathcal{R} = \frac{\Delta(M, \text{SMDP-UCRL}, T_n)}{\Delta(M, \text{MDP-UCRL}, T_n)} \leq 1$

Proxy for \mathcal{R} : ratio of the upper-bounds

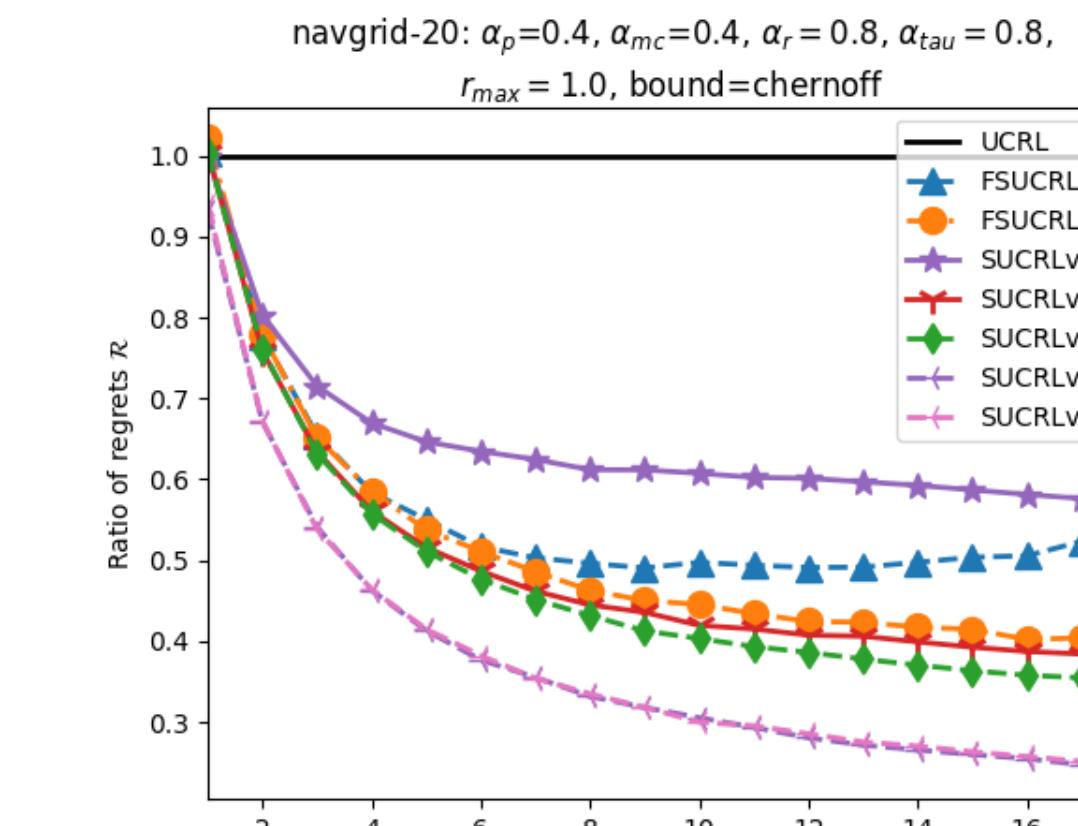
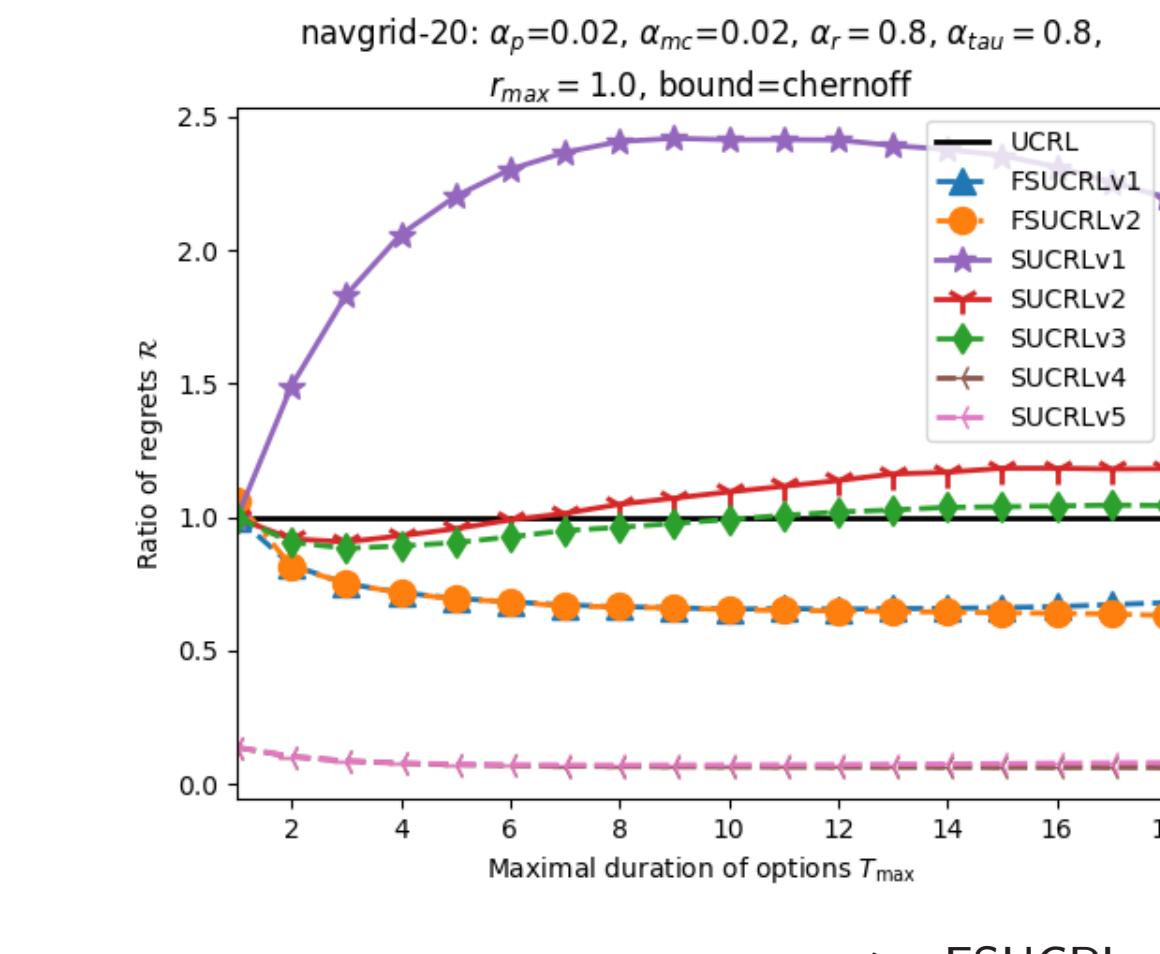
$$\tilde{\mathcal{R}} \sim \frac{D'}{D} \sqrt{\frac{O}{A}} \sqrt{\frac{n}{T_n}}$$

Experimental Results

Gridworld domain

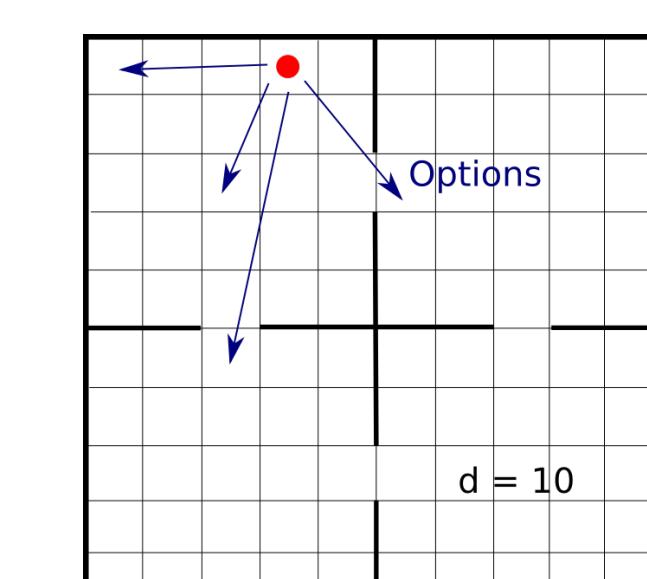


- The options are a.s. bounded $\tau \leq T_{\max}$

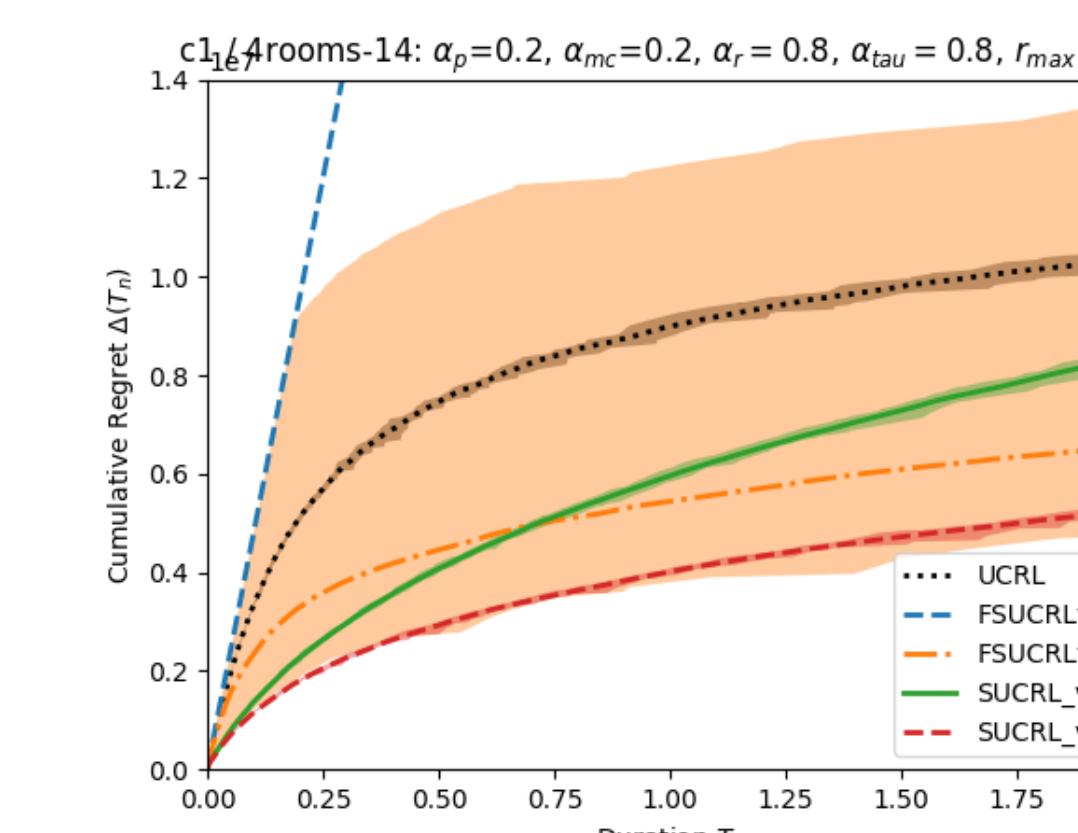
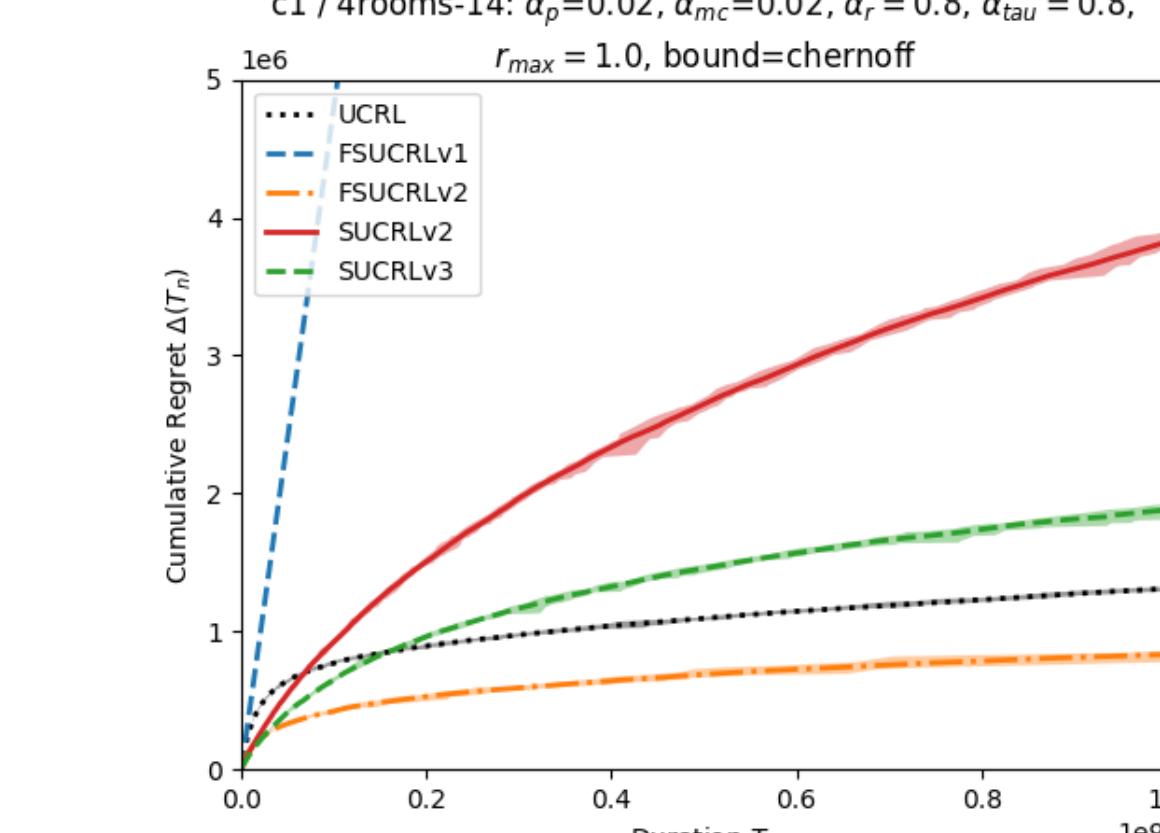


► FSUCRL always outperforms UCRL

Four Rooms Domain



- Size of state/action space is preserved



Conclusion

Take away message

- Temporal abstraction can be achieved without prior knowledge
 - better regret due to exploitation of correlations
 - but navigability within an option should not be compromised too much

Future work

- Can we leverage on these insights to design better options?
 - single task and transfer settings

References

- [Brunskill et al. (2014)] Brunskill, E. and Li, L. Pac-inspired option discovery in lifelong RL. In ICML, 2014.
- [Jaksch et al. (2010)] Jaksch, T., Ortner, R., and Auer, P. Near-optimal regret bounds for RL. In JMLR, 2010.
- [Jong et al. (2008)] Jong, N. K., Hester, T., and Stone, P. The utility of temporal abstraction in RL. In AAMAS, 2008.
- [Sutton et al. (1999)] Sutton, R., Precup, D., and Singh, S. Between mdps and semi-mdps: A framework for temporal abstraction in RL. In AI, 1999.
- [Fruit and Lazaric. (2017)] Fruit, R., and Lazaric, A., Exploration-Exploitation in MDPs with Options. In AISTATS, 2017.