**Lemma 1** *For any policy $\pi$ and any recurrent class $R_\pi^{\otimes i}$ in the Markov chain $MC_\pi^\otimes$, $MC_\pi^\otimes$ satisfies one of the following conditions.*

1. $\delta_{\pi,i}^\otimes \cap \bar{F}_j^\otimes \neq \emptyset$ , $\forall j \in \{1, \ldots, n\}$,

2. $\delta_{\pi,i}^\otimes \cap \bar{F}_j^\otimes = \emptyset$ , $\forall j \in \{1, \ldots, n\}$.

Let $\mathcal{SV}^*$ be the set of optimal supervisors. Let $D_{SV^*}^\otimes$ denote the product DES $D^\otimes$ controlled by the optimal supervisor $SV^*$.

For a Markov chain $MC_{SV}^\otimes$ induced by a product MDP $D^\otimes$ with a supervisor $SV$, let $S_{SV}^\otimes = T_{SV}^\otimes \sqcup R_{SV}^{\otimes 1} \sqcup \ldots \sqcup R_{SV}^{\otimes h}$ be the set of states in $MC_{SV}^\otimes$, where $T_{SV}^\otimes$ is the set of transient states and $R_{SV}^{\otimes i}$ is the recurrent class for each $i \in \{1, \ldots, h\}$, and let $R(MC_{SV}^\otimes)$ be the union of all recurrent classes in $MC_{SV}^\otimes$. Let $\delta_{SV,i}^\otimes$ be the set of transtions in a recurrent class $R_{SV}^{\otimes i}$, namely $\delta_{SV,i}^\otimes = \{(s^\otimes, e, s^{\otimes\prime}) \in \delta^\otimes; s^\otimes \in R_{SV}^{\otimes i}, P_T^\otimes(s^{\otimes\prime}|s^\otimes, e) > 0, P_E^\otimes(e|s^\otimes, SV(s^\otimes)) > 0\}$, and let $P_{SV}^\otimes : S_{SV}^\otimes \times S_{SV}^\otimes \to [0,1]$ such that $P_{SV}^\otimes = \sum_{e \in SV(s^\otimes)} P_T^\otimes(s^{\otimes\prime}|s^\otimes, e) P_E^\otimes(e|s^\otimes, SV(s^\otimes))$ be the transition probability under $SV$.

**Theorem 1** *Let $D^\otimes$ be the product DES corresponding to a DES $D$ and an LTL formuula $\varphi$. If there exists a supervisor $SV$ satisfying $\varphi$, then there exist a discount factor $\gamma^*$ and a positive reward $r_p$ such that any algorithm that maximizes the expected discounted reward with $\gamma > \gamma^*$ and $r_p > ||\mathcal{R}_1||_\infty$ will find a supervisor satisfying $\varphi$.*

**Proof 1** *Suppose that $SV^*$ be an optimal supervisor but does not satisfy the LTL formula $\varphi$. Then, for any recurrent class $R_{SV^*}^{\otimes i}$ in the Markov chain $MC_{SV^*}^\otimes$ and any accepting set $\bar{F}_j^\otimes$ of the product DES $D^\otimes$, $\delta_{SV^*,i}^\otimes \cap \bar{F}_j^\otimes = \emptyset$ holds by Lemma 1. Thus, the agent under the policy $\pi^*$ can obtain rewards only in the set of transient states. We consider the best scenario in the assumption. Let $p^k(s, s')$ be the probability of going to a state $s'$ in $k$ time steps after leaving the state $s$, and let $Post(T_{\pi^*}^\otimes)$ be the set of states in recurrent classes that can be transitioned from states in $T_{\pi^*}^\otimes$ by one event occurrence. For the initial state $s_{init}^\otimes$ in the set of transient states, it holds that*

$$V^{SV^*}(s_{init}^\otimes) = \sum_{k=0}^\infty \sum_{s^\otimes \in T_{\pi^*}^\otimes} \gamma^k p^k(s_{init}^\otimes, s^\otimes)$$
$$\sum_{s^{\otimes\prime} \in T_{\pi^*}^\otimes \cup Post(T_{\pi^*}^\otimes)} \sum_{e \in SV(s^\otimes)} P_T^\otimes(s^{\otimes\prime}|s^\otimes, e) P_E^\otimes(e|s^\otimes, SV(s^\otimes)) \mathcal{R}(s^\otimes, SV(s^\otimes), e, s^{\otimes\prime})$$
$$\leq r_p \sum_{k=0}^\infty \sum_{s^\otimes \in T_{\pi^*}^\otimes} \gamma^k p^k(s_{init}^\otimes, s^\otimes).$$

*By the property of the transient states, for any state $s^\otimes$ in $T_{\pi^*}^\otimes$, there exists a bounded positive value m such that $\sum_{k=0}^\infty \gamma^k p^k(s_{init}^\otimes, s^\otimes) \leq \sum_{k=0}^\infty p^k(s_{init}^\otimes, s^\otimes) <$*

$m$ [1]. Therefore, there exists a bounded positive value $\bar{m}$ such that $V^{\pi^*}(s_{init}^{\otimes}) < \bar{m}$. Let $\bar{SV}$ be a supervisor satisfying $\varphi$. We consider the following two cases.

1. Assume that the initial state $s_{init}^{\otimes}$ is in a recurrent class $R_{\bar{\pi}}^{\otimes i}$ for some $i \in \{1, \ldots, h\}$. For any accepting set $\bar{F}_j^{\otimes}$, $\delta_{\bar{\pi},i}^{\otimes} \cap \bar{F}_j^{\otimes} \neq \emptyset$ holds by the definition of $\bar{\pi}$. The expected discounted reward for $s_{init}^{\otimes}$ is given by

$$V^{\bar{\pi}}(s_{init}^{\otimes}) = \sum_{k=0}^{\infty} \sum_{s^{\otimes} \in R_{SV}^{\otimes i}} \gamma^k p^k(s_{init}^{\otimes}, s^{\otimes})$$
$$\sum_{s^{\otimes\prime} \in R_{SV}^{\otimes i}} \sum_{e \in SV(s^{\otimes})} P_T^{\otimes}(s^{\otimes\prime}|s^{\otimes}, e) P_E^{\otimes}(e|s^{\otimes}, SV(s^{\otimes})) \mathcal{R}(s^{\otimes}, \bar{SV}(s^{\otimes}), e, s^{\otimes\prime}).$$

Since $s_{init}^{\otimes}$ is in $R_{\bar{\pi}}^{\otimes i}$, there exists a set of positive numbers $K = \{k \; ; \; k \geq n, p^k(s_{init}^{\otimes}, s_{init}^{\otimes}) > 0\}$ [1]. We consider the worst scenario in this case. For the stopping time of first returning the initial state, it holds that

$$V^{\bar{\pi}}(s_{init}^{\otimes}) > \mathbb{E}[\gamma^k r_p - (1 + \ldots + \gamma^k)||\mathcal{R}_1||_\infty + \gamma^k V^{\bar{\pi}}(s_{init}^{\otimes})]$$
$$\geq \gamma^{\mathbb{E}[k]} r_p - (1 + \ldots + \gamma^{\mathbb{E}[k]})||\mathcal{R}_1||_\infty + \gamma^{\mathbb{E}[k]} V^{\bar{\pi}}(s_{init}^{\otimes})$$
$$= \frac{\gamma^{\mathbb{E}[k]} r_p - (1 + \ldots + \gamma^{\mathbb{E}[k]})||\mathcal{R}_1||_\infty}{1 - \gamma^{\mathbb{E}[k]}},$$

# References

[1] R. Durrett, *Essentials of Stochastic Processes*, 2nd Edition. ser. Springer texts in statistics. New York; London; Springer, 2012.

[2] L. Breuer, "Introduction to Stochastic Processes," [Online]. Available: https://www.kent.ac.uk/smsas/personal/lb209/files/sp07.pdf

[3] S.M. Ross, *Stochastic Processes*, 2nd Edition. University of California, Wiley, 1995.

[4] S. Singh, T. Jaakkola, M. L. Littman, and C. Szepesvári, "Convergence results for single-step on-policy reinforcement learning algorithms" *Machine Learning*, vol. 38, no. 3, pp, 287–308, 1998.

[5] J. Kretínský, T. Meggendorfer, S. Sickert, "Owl: A library for $\omega$-words, automata, and LTL," in *Proc. 16th International Symposium on Automated Technology for Verification and Analysis*, 2018, pp. 543–550.