

# Learning Optimal Strategies for Temporal Tasks in Stochastic Games

---

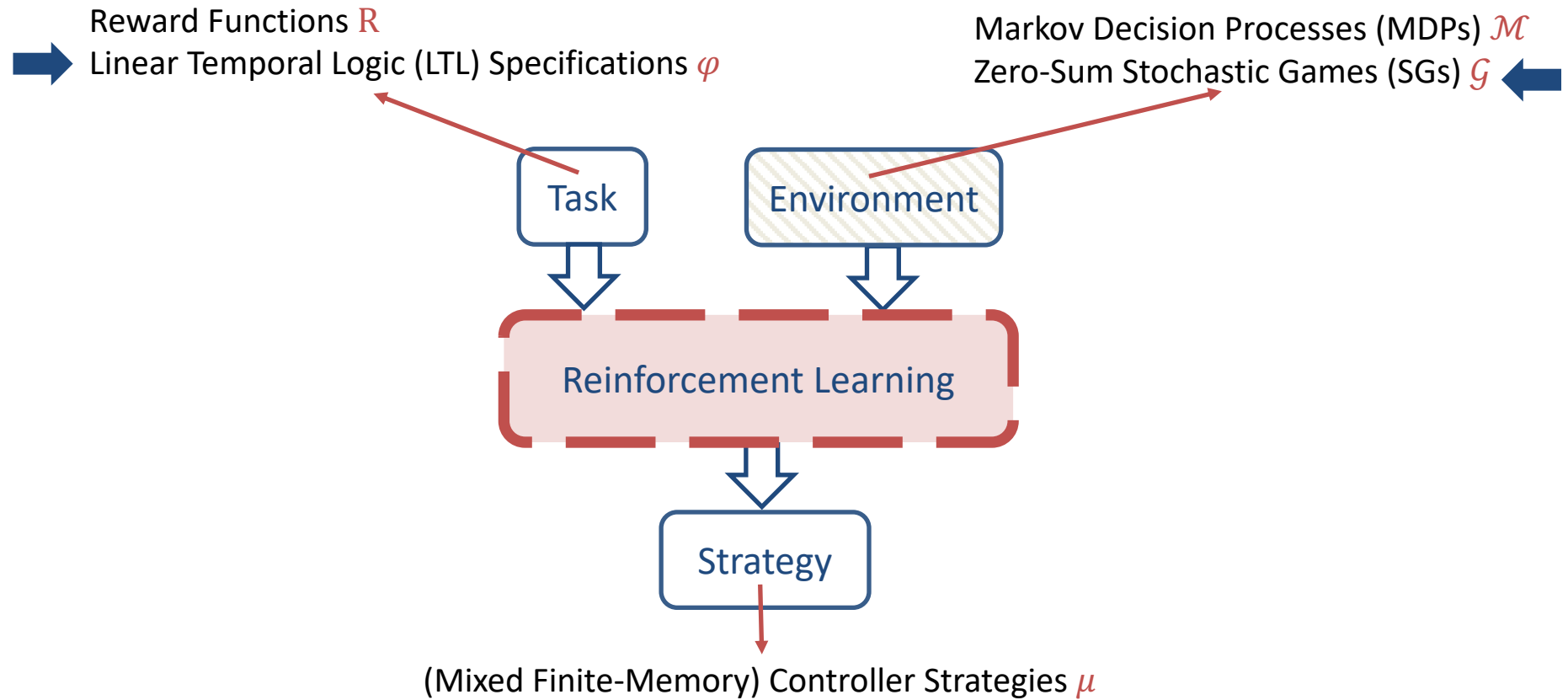
Miroslav Pajic

Department of Electrical and Computer Engineering

Duke University

April 30, 2021

- **Framework**



- **Problem**

Given an unknown  $\mathcal{G}$  and a specification  $\varphi$ ,  
learn a strategy in  $\operatorname{argmax}_{\mu} \min_{\nu} Pr_{\mu, \nu}(\mathcal{G} \models \varphi)$   
where  $\mu$  and  $\nu$  are controller and adversary strategies.

# Stochastic Games and Linear Temporal Logic

- Labeled Turn-Based Zero-Sum Stochastic Games**

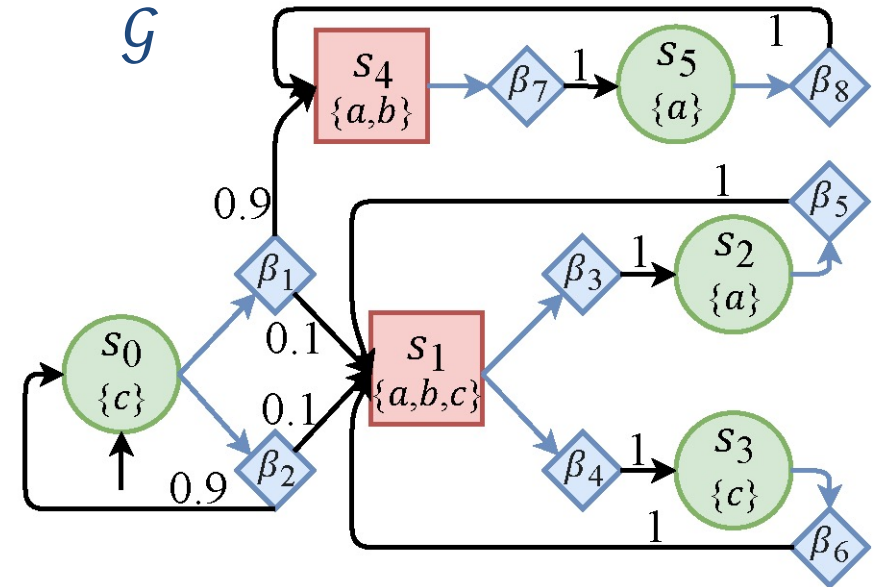
$$\mathcal{G} = (S, (S_\mu, S_\nu), s_0, A, P, AP, L)$$

- $S = S_\mu \cup S_\nu$  is a finite set of states;  $s_0$  is an initial state
- $S_\mu, S_\nu$  are the controller and the adversary states
- $A$  is a finite set of actions
- $P$  is the transition probability function (**unknown**)
- $AP$  is a set of labels/atomic propositions
- $L: S \rightarrow 2^{AP}$  is a labeling function

- LTL Grammar**

$$\varphi := \text{true} \mid a \mid \neg\varphi \mid \varphi_1 \wedge \varphi_2 \mid \bigcirc\varphi \mid \varphi_1 \text{ U } \varphi_2, \quad a \in AP$$

- $\varphi_1 \vee \varphi_2 := \neg(\neg\varphi_1 \wedge \neg\varphi_2)$ ;
- $\varphi_1 \rightarrow \varphi_2 := \neg\varphi_1 \vee \varphi_2$
- $\diamond\varphi := \text{true U } \varphi$
- $\square\varphi := \neg(\diamond\neg\varphi)$



- : Controller State
- : Adversary State
- ◆ : Actions

# RL Framework for LTL

- **Key Idea: Reduction**

- From the LTL objective

$$\operatorname{argmax}_{\mu} \min_{\nu} \Pr_{\mu, \nu}(\mathcal{G} \models \varphi)$$

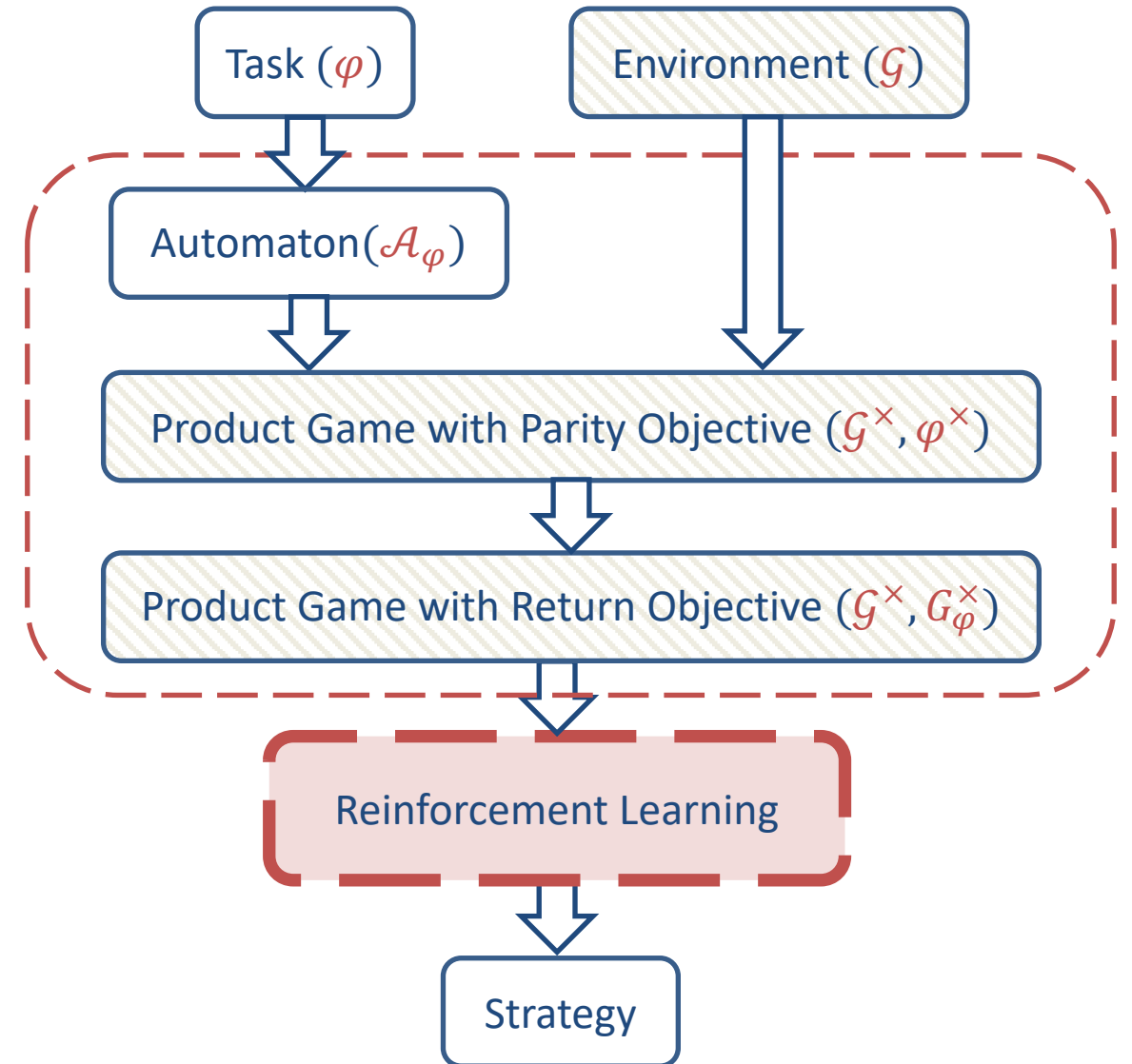
- To a return objective

$$\operatorname{argmax}_{\mu} \min_{\nu} \mathbb{E}_{\mu, \nu}[G_{\varphi}^{\times}]$$

$$\operatorname{argmax}_{\mu} \min_{\nu} \mathbb{E}_{\mu, \nu} \left[ \sum_{i=0}^{\infty} \gamma^i r_{(i)} \right]$$

- **Reduction Steps:**

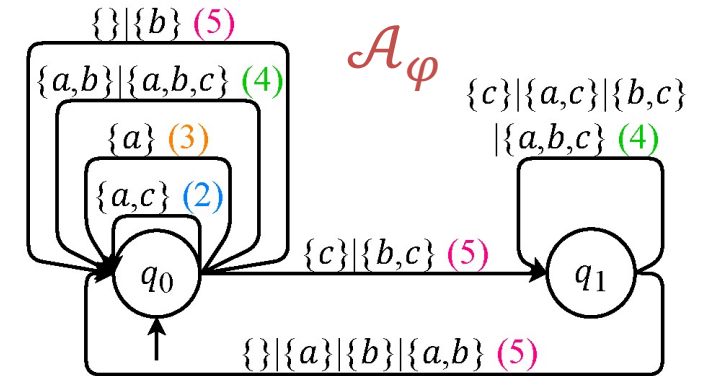
- LTL -> Automaton
- Product Game Construction
- Reduction from Parity to Return
- Model-free Learning



# Product Game Construction

- LTL to Deterministic Parity Automata (DPA) Translation**

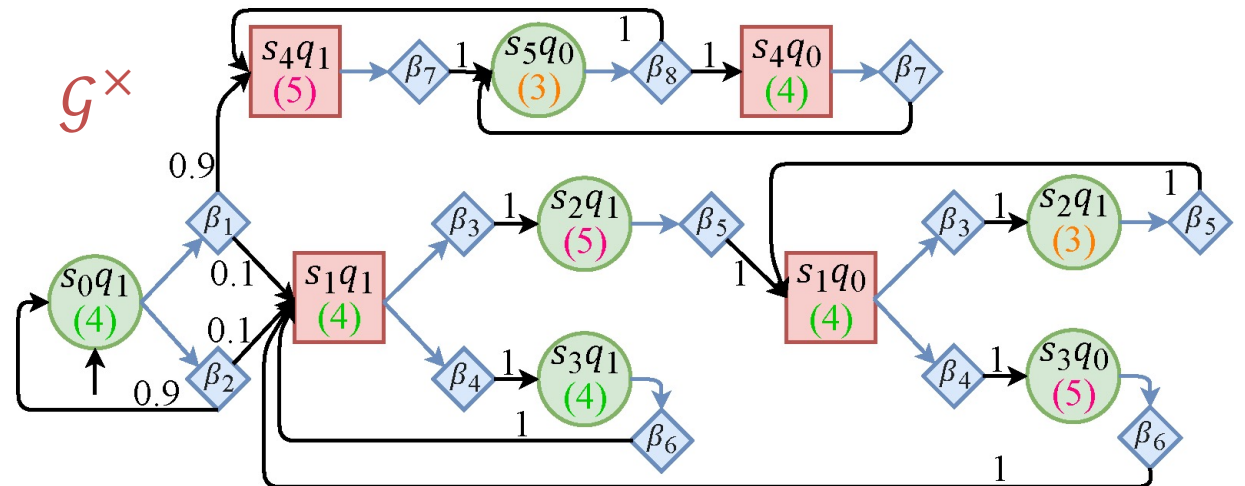
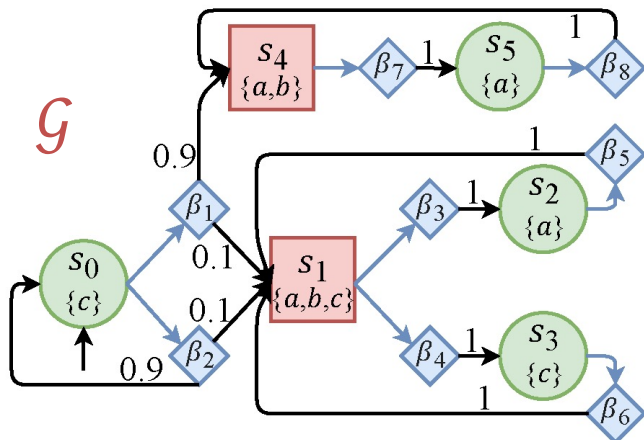
- The set of traces satisfying  $\varphi$  is an  $\omega$ -regular language
- A DPA  $\mathcal{A}_\varphi$  recognizing the language can be automatically constructed
- Example:  $\varphi = (\diamond \Box a \wedge \Box \diamond b) \vee \diamond \Box c$



- Product Game**

- Simultaneous execution of the SG  $\mathcal{G}$  and the DPA  $\mathcal{A}_\varphi$
- Does not have to be constructed explicitly
- Winning Condition: **Parity Objective**

$$\varphi^\times := \max\{Color(s^\times) \mid s^\times \in Inf(\pi^\times)\}$$



# Reduction I: Distinct Discount Factors - Büchi Conditions

- **Büchi Conditions**

- Two colors: Color 1 and Color 2
- Suffices for MDPs (with some additional structured nondeterminism)
- $\varphi^\times$ : Repeatedly visit states colored with 2
- Example:  $Pr((s_0^\times, \text{up}) \models \varphi^\times) = 0.9$  and  $Pr((s_0^\times, \text{down}) \models \varphi^\times) = 1$

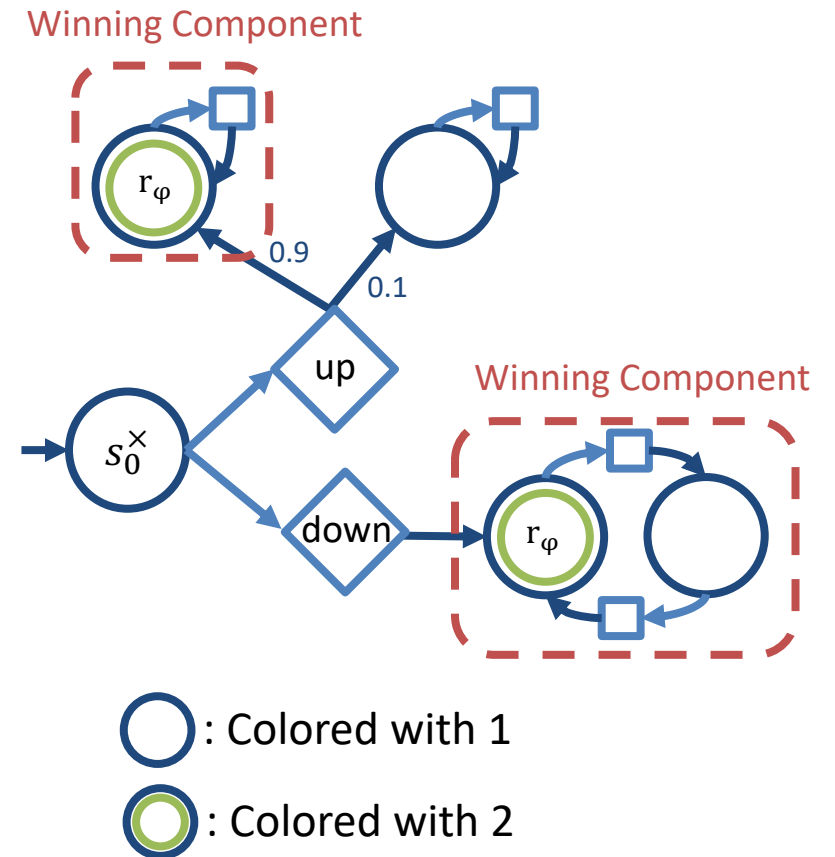
- **Reduction to Return Objectives**

- Reward Function:  $R_\varphi(s^\times) := \begin{cases} r_\varphi & \text{if } Color(s^\times) = 2 \\ 0 & \text{if } Color(s^\times) = 1 \end{cases}$
- Discount Function:  $\Gamma_\varphi(s^\times) := \begin{cases} 1 - r_\varphi & \text{if } Color(s^\times) = 2 \\ 1 - r_\varphi^2 & \text{if } Color(s^\times) = 1 \end{cases}$
- Example:  $q(s^\times, \text{up}) = 0.9$  and  $q(s^\times, \text{down}) = \frac{1}{1+r_\varphi(1-r_\varphi)}$

- **Theorem I [1,2]:**

- For a Büchi condition  $\varphi^\times$  and any strategy pair  $(\mu, \nu)$ ,

$$Pr_{\mu, \nu}(\mathcal{G}^\times \models \varphi^\times) = \lim_{r_\varphi \rightarrow 0^+} \mathbb{E}_{\mu, \nu} \left[ \sum_{i=0}^{\infty} \left( \prod_{j=1}^i \Gamma_\varphi(s_{(j)}^\times) \right) R_\varphi(s_{(i)}^\times) \right]$$



[1] A. K. Bozkurt, Y. Wang, M. M. Zavlanos, and M. Pajic. "Control Synthesis from Linear Temporal Logic Specifications using Model-Free Reinforcement Learning". ICRA, 2020.

[2] A. K. Bozkurt, Y. Wang, M. M. Zavlanos, and M. Pajic. "Model-Free Reinforcement Learning for Stochastic Games with Linear Temporal Logic Objectives". ICRA, 2021, accepted.

# Reduction I: Distinct Discount Factors - Generalization

- Parity Conditions with  $k > 2$

- Example:  $Pr((s_0^\times, \text{up}) \models \varphi^\times) = 0.9$  and  $Pr((s_0^\times, \text{down}) \models \varphi^\times) = 0$

- A distinct power for each color

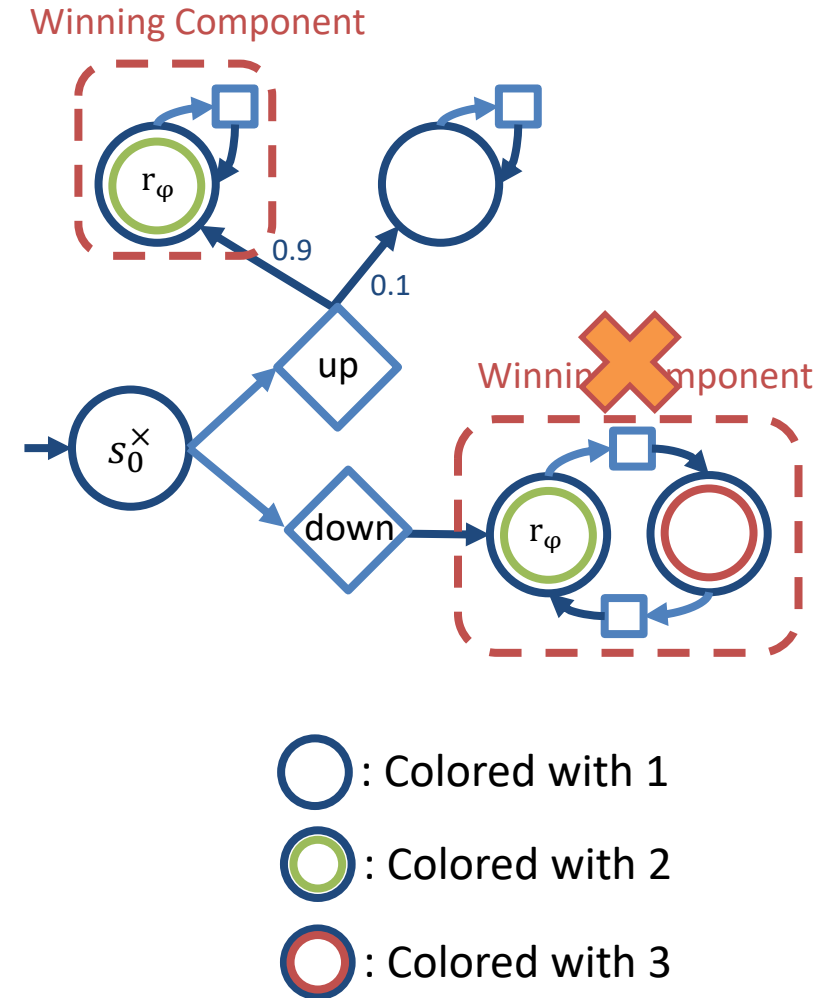
- Reward Function:  $R_\varphi(s^\times) := \begin{cases} r_\varphi^{k-Color(s^\times)} & \text{if } Color(s^\times) \text{ is even} \\ 0 & \text{if } Color(s^\times) \text{ is odd} \end{cases}$

- Discount Function:  $\Gamma_\varphi(s^\times) := 1 - r_\varphi^{k-Color(s^\times)}$

- Example:  $q(s^\times, \text{up}) = 0.9$  and  $q(s^\times, \text{down}) = \frac{1}{1 + \frac{1-r_\varphi^2}{r_\varphi}}$

- Distinct powers of rewards and discount factors captures the order
- Not Scalable

- An approximation is provided in [2]



# Case Study: Avoiding an Adversary

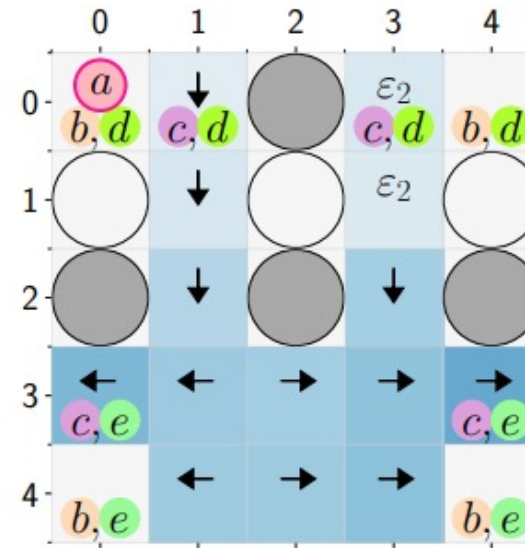
- Grid World**

- The agent and the adversary can take four actions: *North, South, East, West*
- The probability of moving in the **intended** direction: **0.8**
- The probability of moving in a direction **orthogonal** to the intended direction: **0.2**

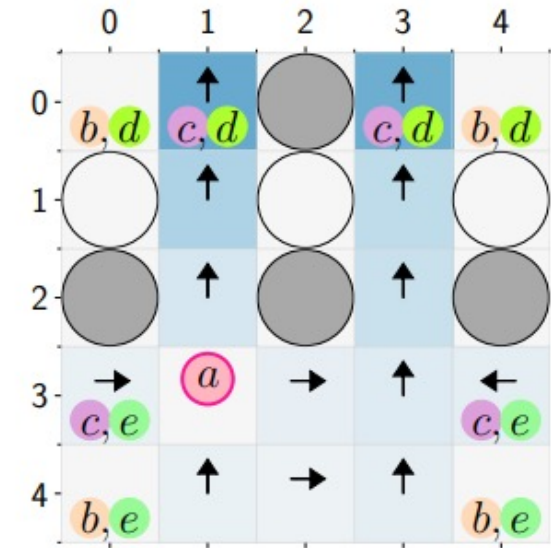
- Objective**

- Repeatedly visit a *b* and a *c* cell
- Reach a safe region labeled with *d* or *e* and do not leave
- Avoid the adversary (*a*) at all costs.

$$\varphi = \square \diamond b \wedge \square \diamond c \wedge (\diamond \square d \vee \diamond \square e) \wedge \square \neg a$$



(a) Adversary is at (0, 0) and  $i=1$



(b) Adversary is at (3, 1) and  $i=2$

The darker blue, the higher estimated satisfaction probability



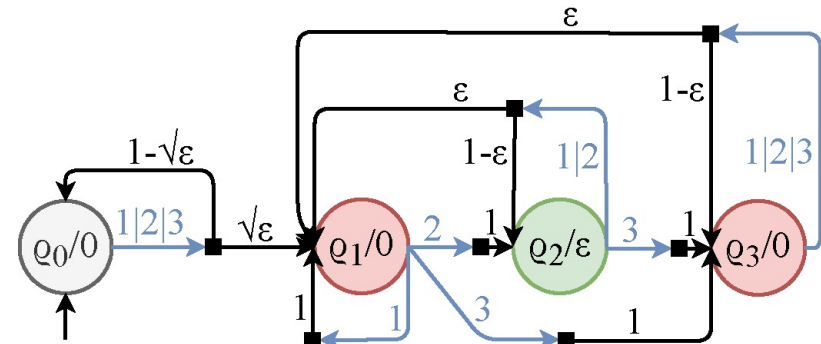
# Reduction II: Priority Reward Machines

- **Objective:**

- Design a reduction where the discount factors, rewards, transition probabilities **do not depend on the number of colors**

- **Priority Reward Machines (PRMs)**

- The Moore machines consisting of priority modes  $q_i$
- Output:  $R_\varphi^*(s^\times, \varrho) := \begin{cases} \varepsilon_\varphi & \text{if } \varrho > 0 \text{ and } \varrho \text{ is even} \\ 0, & \text{otherwise} \end{cases}$
- A priority mode  $q_i$  is overruled by  $q_j$  when Color  $j$  is consumed
- PRMs reset to  $q_1$  w.p.  $\varepsilon_\varphi$
- PRMs move from  $q_0$  to  $q_1$  w.p.  $\sqrt{\varepsilon_\varphi}$



- **Theorem II [3]:**

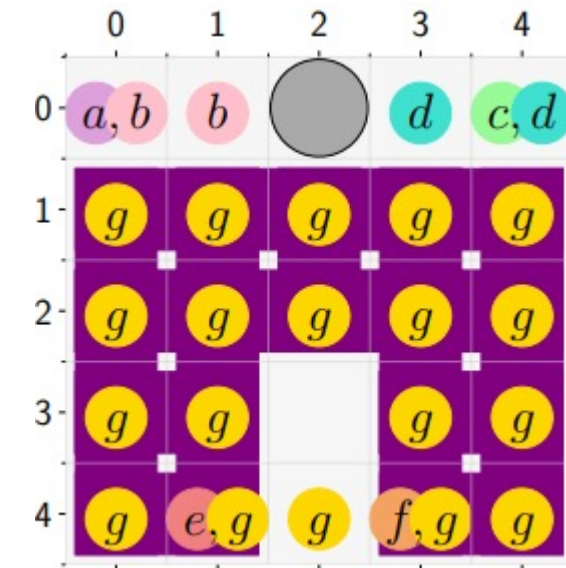
- For a parity condition  $\varphi^\times$  and any strategy pair  $(\mu, \nu)$ ,

$$Pr_{\mu, \nu}(\mathcal{G}^\times \models \varphi^\times) = \lim_{\varepsilon_\varphi \rightarrow 0^+} \mathbb{E}_{\mu, \nu} \left[ \sum_{i=0}^{\infty} (1 - \varepsilon_\varphi)^i R_\varphi^*(s_{(i)}^\times, \varrho_{(i)}) \right]$$

# Case Study: Surveillance

- **Grid World**
  - The agent can take four actions: *North, South, East, West*
  - The adversary can **disrupt** the movement so that the agent might move in a **perpendicular** direction
- **Objective**
  - Eventually perform one of the following surveillance tasks:
    - Repeatedly visit *a* without leaving the region *b*
    - Repeatedly visit *c* without leaving the region *d*
    - Repeatedly visit *e* and *f* without leaving the region *g*

$$\varphi = (\Box \diamond a \wedge \diamond \Box b) \vee (\Box \diamond c \wedge \diamond \Box d) \vee (\Box \diamond e \wedge \Box \diamond f \wedge \diamond \Box g)$$



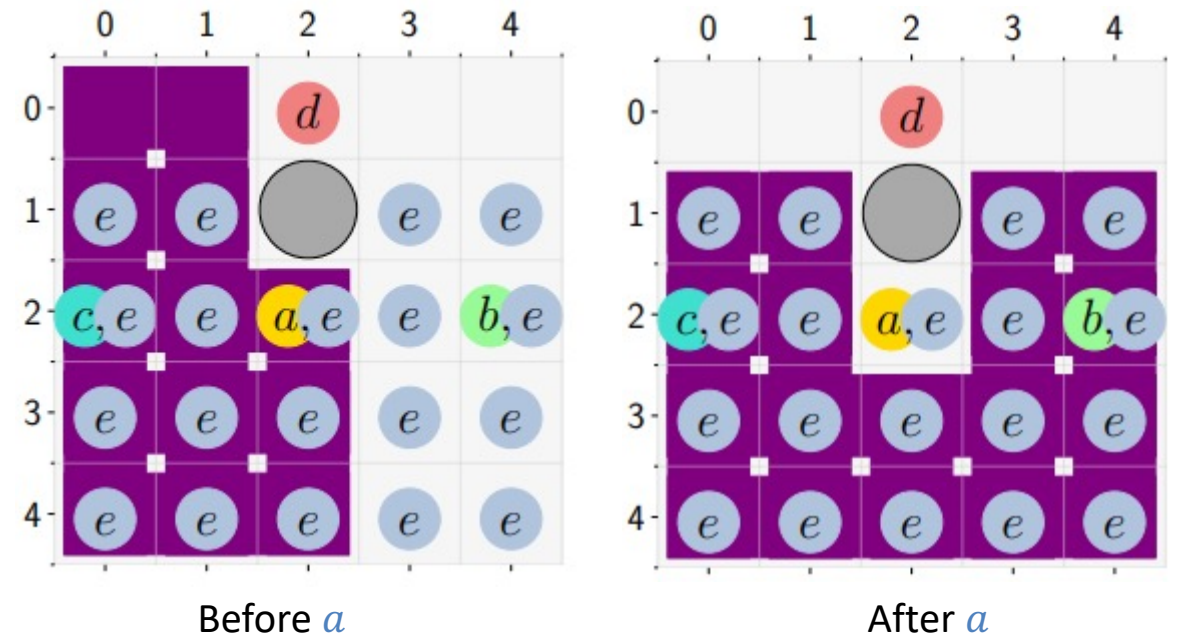
The cells visited under the optimal strategies are highlighted in purple.

# Case Study: Nursery Scenario

- **Objective**

- Start at  $(0,0)$
- Enter the region labeled with  $e$  and stay there
- Inform the adult  $a$  exactly once
- Repeatedly visit the baby  $b$  and the charger station  $c$
- Avoid the danger zone  $d$

$$\varphi = \diamond \square e \wedge \diamond a \wedge \square (a \rightarrow \bigcirc \square \neg a) \wedge \square \diamond b \wedge \square \neg d$$



The cells visited under the optimal strategies are highlighted in purple.

# Quality of Control Optimization under LTL Specifications

- **Multiple Objectives with Lexicographic Order:**
  - Priority 1: Safety  $\psi$ 
    - Safety LTL formula
    - Ensuring the safety is usually of utmost importance
  - Priority 2: LTL  $\varphi$ 
    - Important system specifications as an LTL formula other than safety
  - Priority 3: QoC  $R_{QoC}$ 
    - External Rewards
    - LTL cannot specify objectives including cost, yield or quality optimization
    - Example: Minimization of Energy Consumption

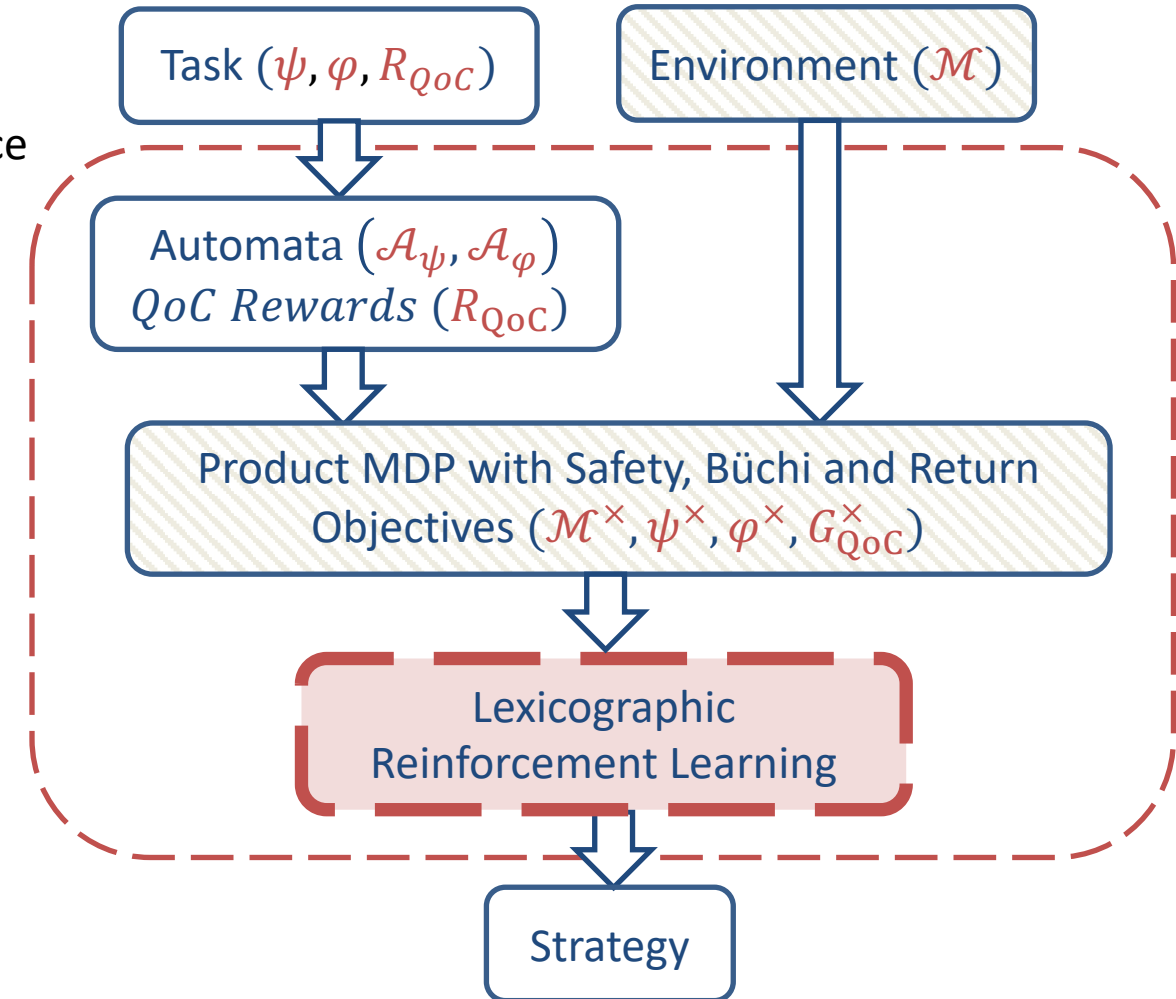
- **Problem:**

- For a given MDP  $\mathcal{M}$  learn a strategy  $\mu \in \Sigma_{\psi, \varphi}^{QoC}$  where

$$\Sigma_{\psi} := \operatorname{argmax}_{\mu} Pr_{\mu}(\mathcal{M} \models \psi)$$

$$\Sigma_{\psi, \varphi} := \operatorname{argmax}_{\mu \in \Sigma_{\psi}} Pr_{\mu}(\mathcal{M} \models \varphi)$$

$$\Sigma_{\psi, \varphi}^{QoC} := \operatorname{argmax}_{\mu \in \Sigma_{\psi, \varphi}} \mathbb{E}_{\mu} \left[ \sum_{i=0}^{\infty} \gamma^i r_{(i)}^{QoC} \right]$$



# QoC Optimization via Lexicographic RL

- **Algorithm I: QoC Optimization Under Safety and LTL Specifications**

- **Learning Actions Sets**

- $V_\psi(s^\times) \leftarrow \max_{a^\times} Q_\psi(s^\times, a^\times)$
- $\hat{A}_\psi^\times \leftarrow \{a^\times \mid V_\psi(s^\times) - Q_\psi(s^\times, a^\times) \leq \tau_\psi\}$
- $V_{\psi,\varphi}(s^\times) \leftarrow \max_{a^\times \in \hat{A}_\psi^\times} Q_{\psi,\varphi}(s^\times, a^\times)$
- $\hat{A}_{\psi,\varphi}^\times(s^\times) \leftarrow \{a^\times \in A_\psi^\times(s^\times) \mid V_{\psi,\varphi}(s^\times) - Q_{\psi,\varphi}(s^\times, a^\times) \leq \tau_\varphi\}$
- $\hat{A}_{\psi,\varphi}^{\times*} \leftarrow \underset{a^\times \in \hat{A}_{\psi,\varphi}^\times(s^\times)}{\operatorname{argmax}} Q_{\psi,\varphi}^{\text{QoC}}(s^\times, a^\times)$

- **Action Selection**

- Choose a random action w.p.  $e$  (during exploration)
- Choose a random action from  $\hat{A}_{\psi,\varphi}^\times(s^\times)$  w.p.  $v$  (for LTL)
- Choose an action in  $\hat{A}_{\psi,\varphi}^{\times*}$

- **Q-Value Updates**

- Use Q-learning to update  $Q_\psi$  and  $Q_{\psi,\varphi}$
- Use SARSA to update  $Q_{\psi,\varphi}^{\text{QoC}}$

- **Theorem III [4]:**

- **Algorithm I** learns a lexicographically near-optimal strategy  $\mu \in \Sigma_{\psi,\varphi}^{\text{QoC}^{-\nu}}$  for sufficiently small  $\tau_\psi, \tau_\varphi > 0$ , for the rewards provided in Reduction I.

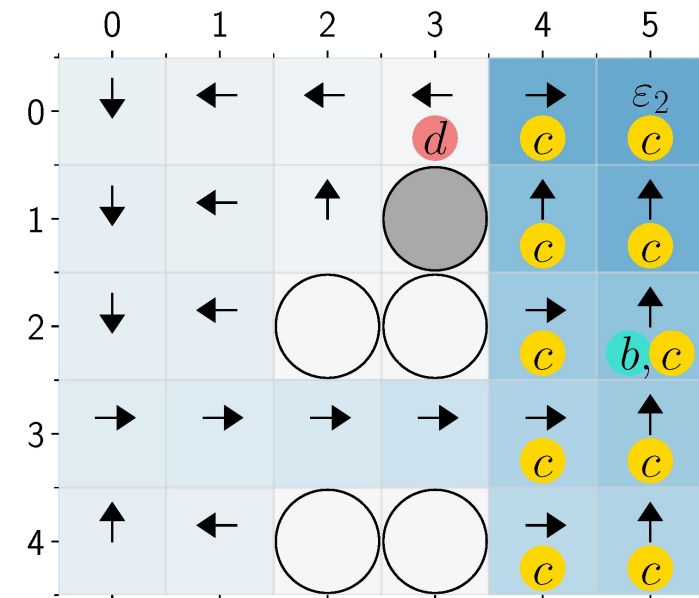
# Case Study: QoC Optimization

- Grid World**

- The agent can take four actions: *North, South, East, West*
- The agent moves in the *intended* direction w.p. **0.8**
- The agent moves in a direction *orthogonal* to the intended direction w.p. **0.2**

- QoC Optimization Example**

- Safety:
  - Avoid visiting a danger state *d* consecutively
  - $\psi = \Box \neg (d \wedge \bigcirc d)$
- LTL:
  - Occasionally visit a checkpoint *b*
  - $\varphi = \Box \diamond b$
- QoC:
  - Stay at the top-right corner as long as possible
  - $R^{\text{QoC}}(\langle 0,5 \rangle) = 1$



# Thank you

---



**Duke**  
UNIVERSITY

PRATT SCHOOL *of*  
**ENGINEERING**