

# Model-Free RL for Control Synthesis for MDPs and Stochastic Games

---

Alper Kamil Bozkurt, Yu Wang, Michael M. Zavlanos and  
Miroslav Pajic

Department of Electrical and Computer Engineering  
Department of Computer Science  
Pratt School of Engineering  
Duke University

# Preliminaries and Problem Statement

**Model:** (Labeled) Turn-Based Zero-Sum **Stochastic Games**

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

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

**Specification:** Linear Temporal Logic (**LTL**)

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

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

**Output:** Finite-Memory **Strategy**

$$\pi = (M, \Delta, \alpha, m_0)$$

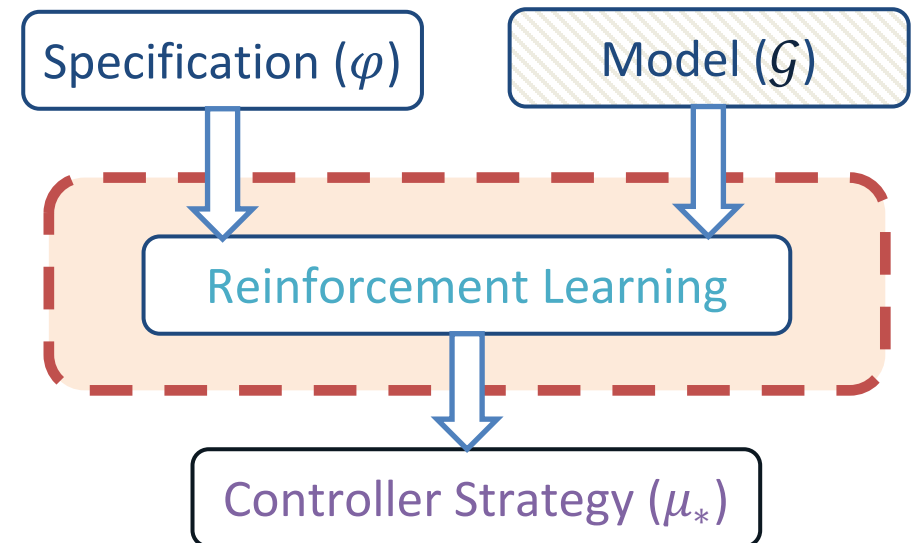
- $M$  is a finite set of modes;  $m_0$  is an initial state
- $\Delta: M \times S \rightarrow M$  is the transition function
- $\alpha: M \times S \rightarrow A$  maps the mode state pairs to actions

## Problem Statement

Given a stochastic game  $\mathcal{G}$  where the transition probabilities and the topology is unknown and an LTL specification  $\varphi$ , design a model-free RL algorithm that finds a finite-memory controller strategy  $\mu_*$  that satisfies

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

where  $\mu$  and  $\nu$  are controller and environment strategies

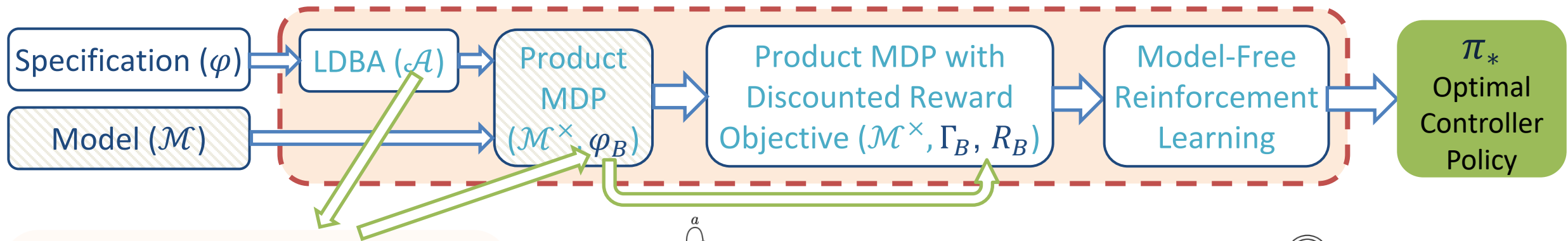


# Control Synthesis via RL for MDPs

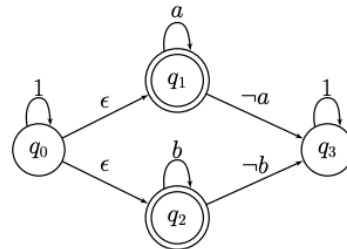
## Problem Statement for MDPs

Given an MDP  $\mathcal{M}$  where the transition probabilities and the topology are unknown and an LTL specification  $\varphi$ , design a model-free RL algorithm that finds a finite-memory objective policy  $\pi_*$  that satisfies

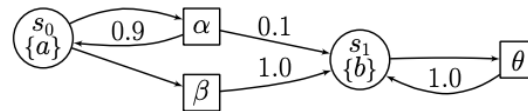
$$\pi_* = \operatorname{argmax}_{\pi} \Pr_{\pi}(\mathcal{M} \models \varphi)$$



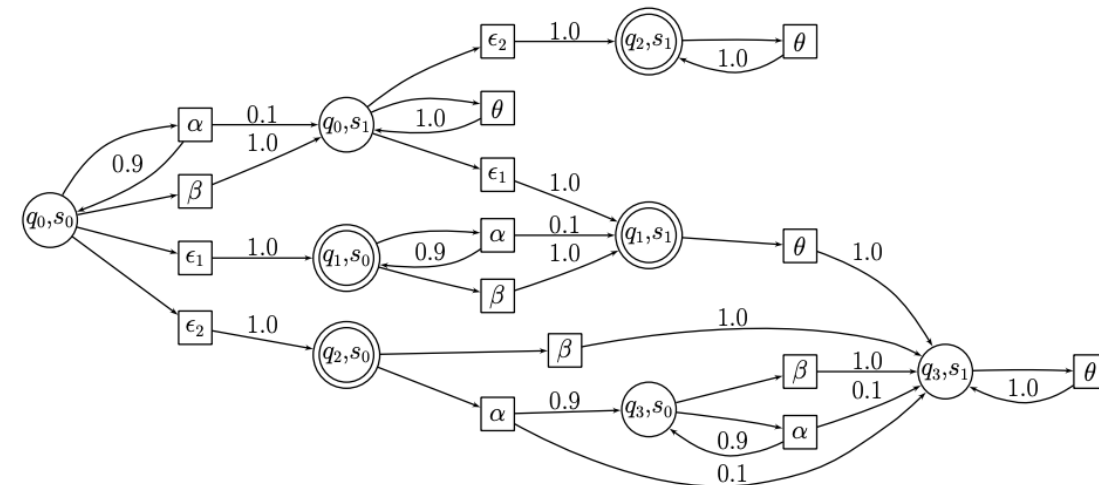
**Limit-Deterministic Büchi Automata (LDBA)** – consist of two deterministic components the *initial* and *accepting*. The only nondeterministic transitions are the  $\epsilon$ -moves from the initial component to the accepting components.



(a) A derived LDBA  $\mathcal{A}$  for the LTL formula  $\varphi = \square\square a \vee \square\square b$



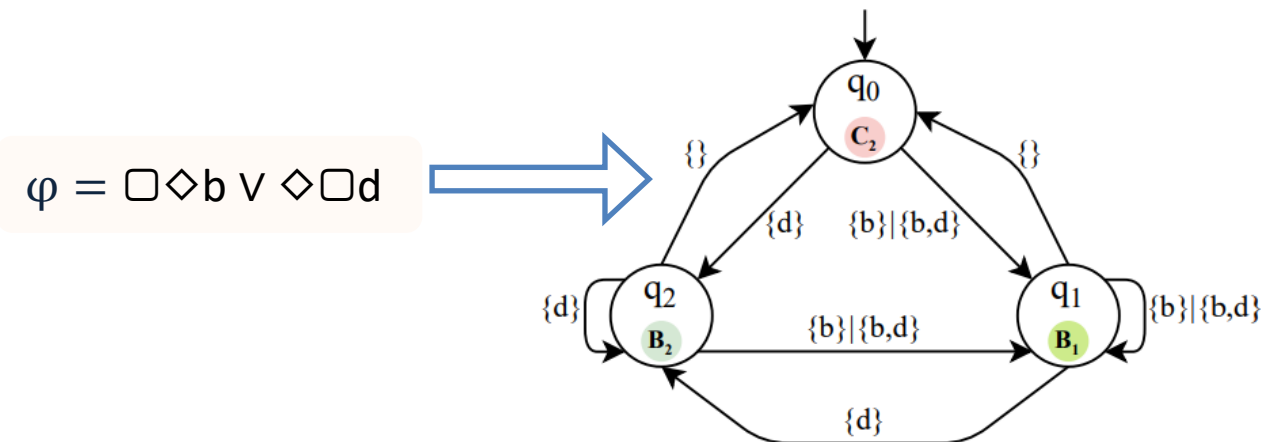
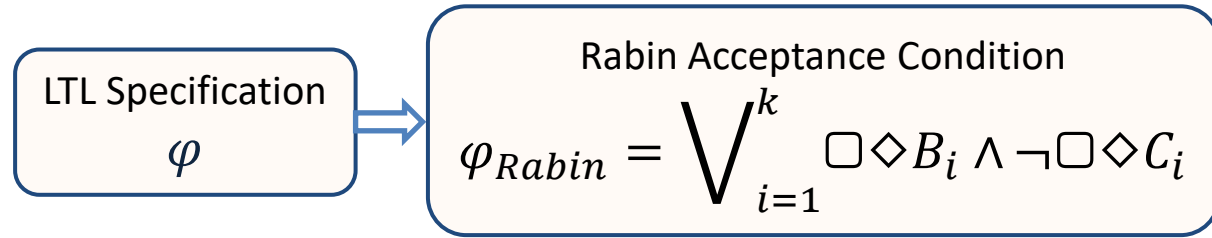
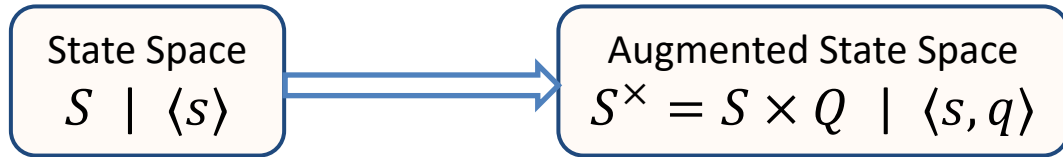
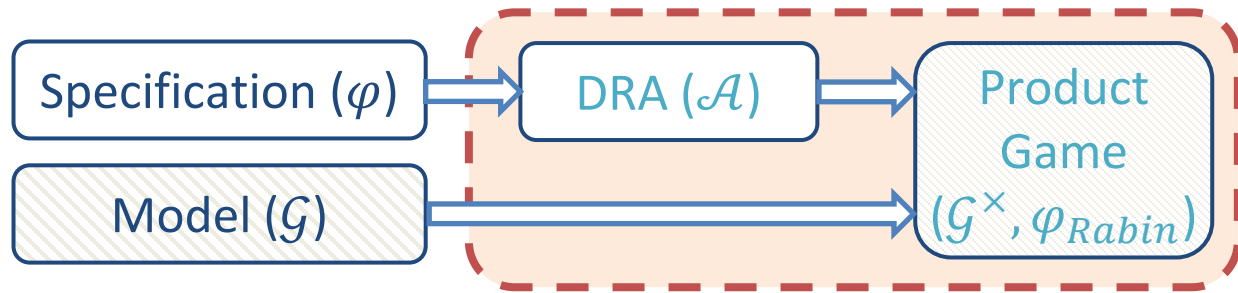
(b) An example MDP  $\mathcal{M}$ ; the circles denote MDP states, rectangles denote actions, and numbers transition probabilities



(c) The obtained product MDP

[1] A. K. Bozkurt, Y. Wang, M. M. Zavlanos, and M. Pajic, "Control Synthesis from LTL Specifications using Model-Free Reinforcement Learning", *IEEE International Conference on Robotics and Automation (ICRA)*, 2020

# Product Game Construction



$\varphi = \square \diamond b \vee \diamond \square d$

## Rabin(1) Acceptance Condition

$$\varphi_{Rabin}^{(1)} = \square \diamond B \wedge \neg \square \diamond C$$

**Pure** and **memoryless** strategies suffice for both Player 1 (**Controller**) and Player 2 (**Environment**) for  $\varphi_{Rabin}^{(1)}$ .

(Chatterjee et al., 2012)



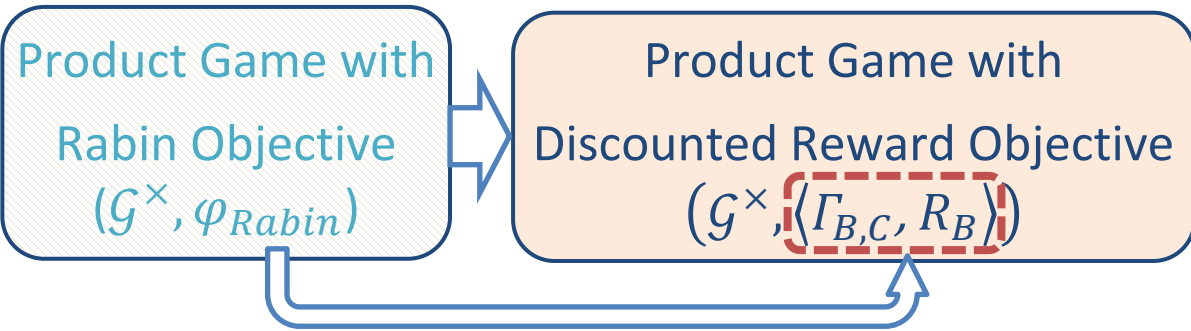
For pure and memoryless  $\mu^x$  and  $\nu^x$ :

$$\mu_*^x = \operatorname{argmax}_{\mu^x} \min_{\nu^x} Pr_{\mu^x, \nu^x} (\mathcal{G}^x \models \varphi_{Rabin}^{(1)})$$

Optimal Finite-Memory Controller Strategy

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

# Rabin(1) Acceptance Condition as Sum of Discounted Rewards

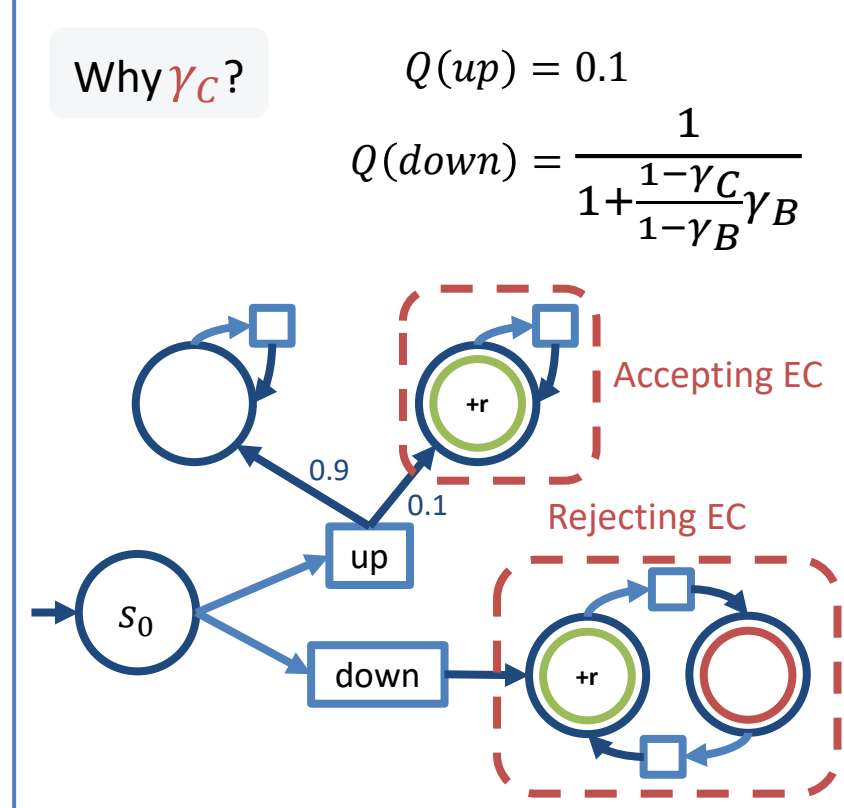
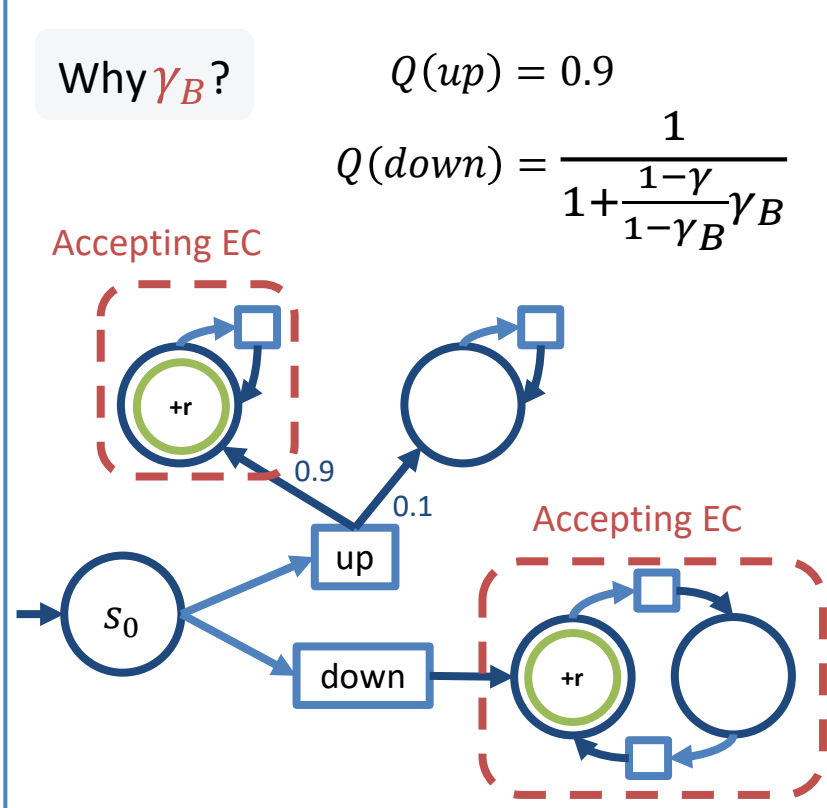
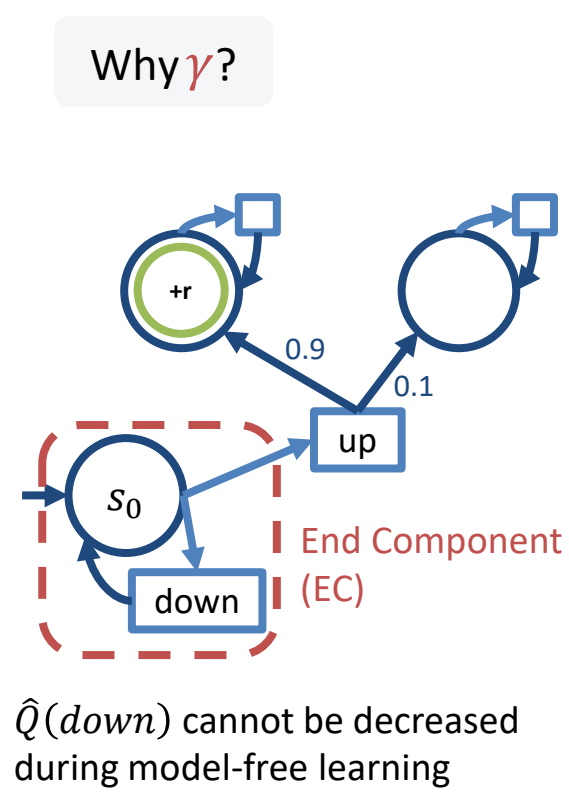
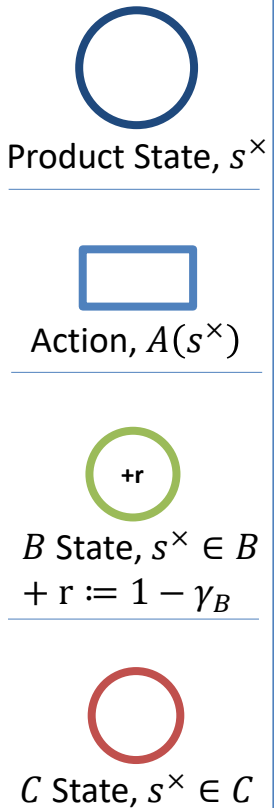


$$R_B(s) := \begin{cases} 1 - \gamma_B, & s^\times \in B \\ 0, & s^\times \notin B \end{cases}$$

Reward Function

$$\Gamma_{B,C}(s^\times) := \begin{cases} \gamma_B, & s^\times \in B \\ \gamma_C, & s^\times \in C \\ \gamma, & otherwise \end{cases}$$

Discount Function



# Main Theoretical Results

Reward Function

$$R_B^\gamma(s) := \begin{cases} 1 - \gamma_B, & s^\times \in B \\ 0, & s^\times \notin B \end{cases}$$

$$\Gamma_{B,C}^\gamma(s^\times) := \begin{cases} \gamma_B(\gamma), & s^\times \in B \\ \gamma_C(\gamma), & s^\times \in C \\ \gamma, & \text{otherwise} \end{cases}$$

Discount Function

Discounted Rewards

$$G_{B,C}^\gamma(\sigma) = \sum_{i=0}^{\infty} \left( \prod_{j=0}^{i-1} \Gamma_{B,C}^\gamma(\sigma[j]) \right) R_B^\gamma(\sigma[i])$$

Discount Constraints

$$\lim_{\gamma \rightarrow 1^-} \frac{1 - \gamma}{1 - \gamma_B(\gamma)} = \lim_{\gamma \rightarrow 1^-} \frac{1 - \gamma_B(\gamma)}{1 - \gamma_C(\gamma)} = 0$$

Theorem

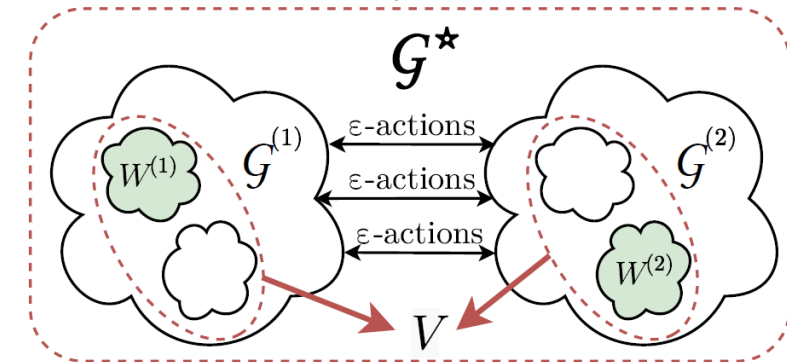
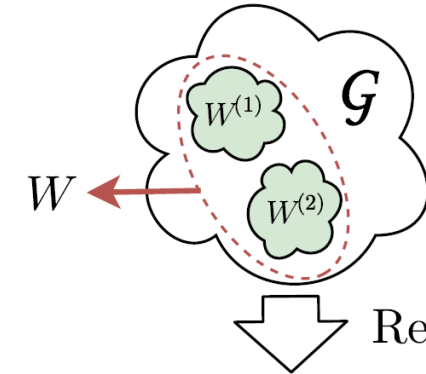
$$\lim_{\gamma \rightarrow 1^-} \mathbb{E}_{\sigma \sim \mathcal{G}_{\mu^\times, \nu^\times}^\times} [G_{B,C}^\gamma(\sigma)] = Pr_{\mu^\times, \nu^\times}(\mathcal{G}^\times \models \Box \Diamond B \wedge \neg \Box \Diamond C)$$

Corollary

$$\begin{aligned} \text{There exists } \gamma' < 1 \text{ such that for all } \gamma \geq \gamma', \\ \mu_*^\times &= \operatorname{argmax}_{\mu^\times} \min_{\nu^\times} \mathbb{E}_{\sigma \sim \mathcal{G}_{\mu^\times, \nu^\times}^\times} [G_{B,C}^\gamma(\sigma)] \\ &= \operatorname{argmax}_{\mu^\times} \min_{\nu^\times} Pr_{\mu^\times, \nu^\times}(\mathcal{G}^\times \models \Box \Diamond B \wedge \neg \Box \Diamond C) \end{aligned}$$

$\mu_*$ : Optimal Finite-Memory Controller Strategy

Multiple Rabin Pairs

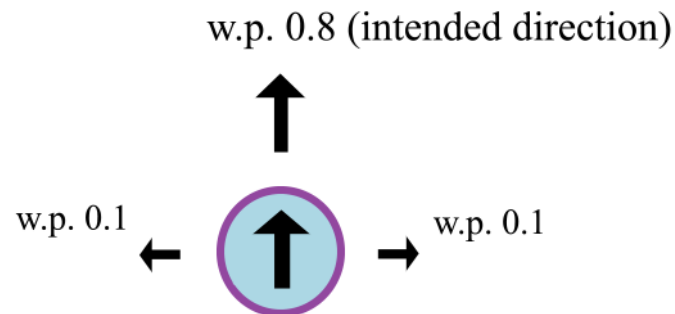


[1] A. K. Bozkurt, Y. Wang, M. M. Zavlanos, and M. Pajic, "Model-Free Reinforcement Learning for Stochastic Games with Linear Temporal Logic Objectives", arXiv:2010.01050, 2020

# Case Study: Avoiding Adversary

## Grid World:

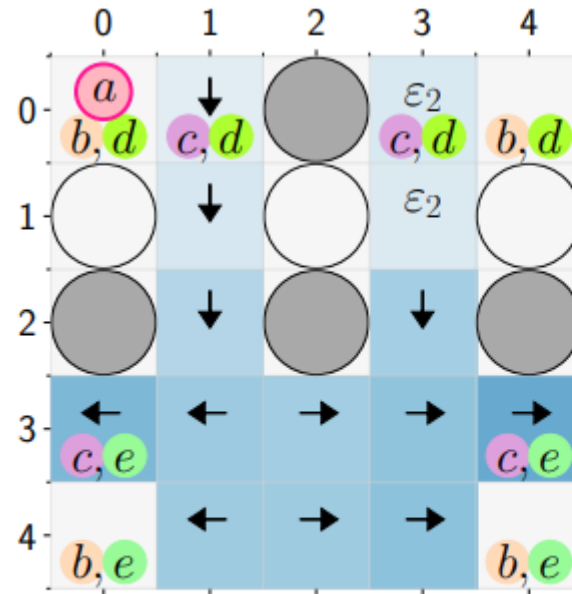
- The agent can take four actions:  
**North, South, East, West**
- The transition model :
  - The probability that the robot moves in the *intended* direction: **0.8**
  - The probability that the robot moves in a direction *orthogonal* to the intended direction: **0.2**
- Action: **North**



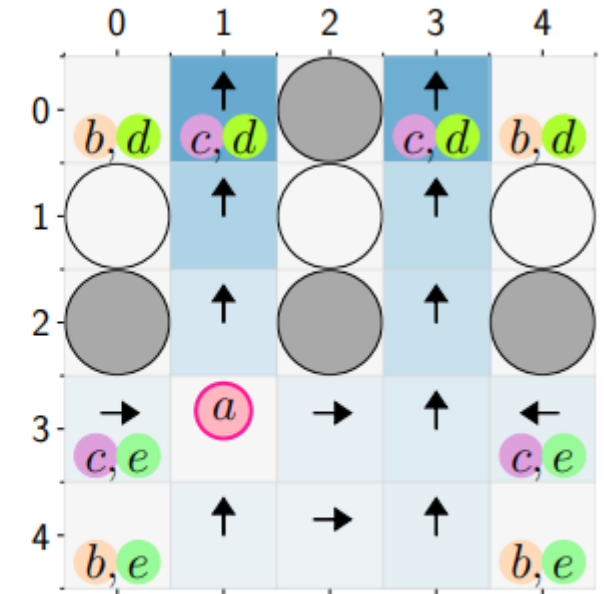
## Objective:

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

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



(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

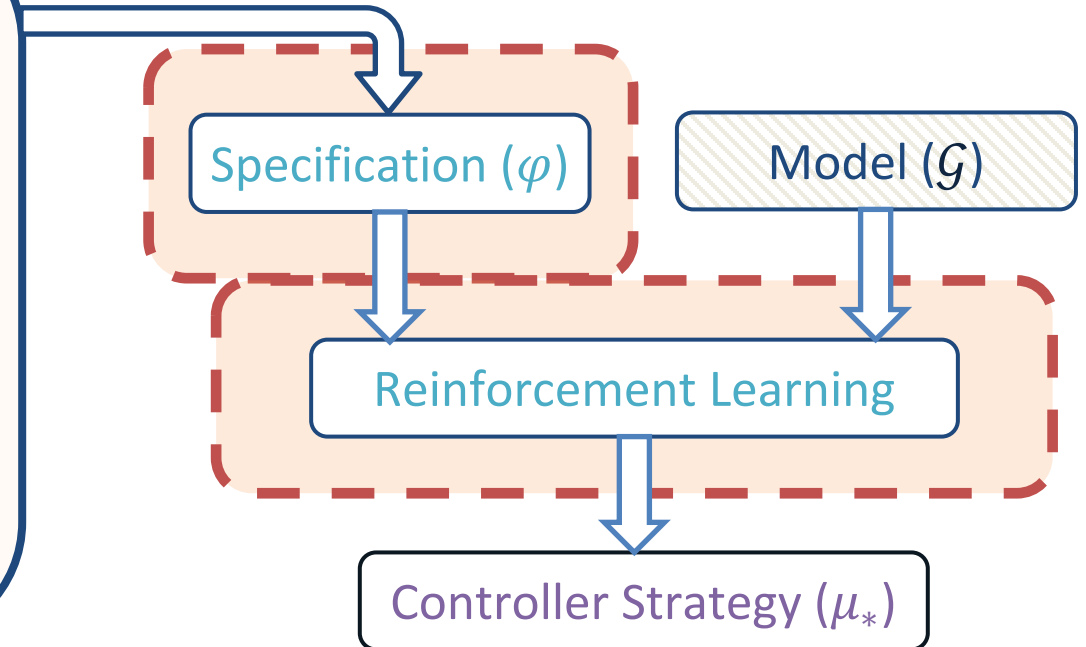
# Secure Planning Against Stealthy Attacks

## Controller:

- aims to perform a given **task**
- does **not have a model** of the environment
- has a perfect knowledge of the current state
- has an intrusion-detection system (**IDS**) that monitors anomalies
- can **detect** attacks only when the IDS raises an **alarm**

## Attacker:

- aims to prevent the controller from performing the given task
- has a **perfect knowledge** of the current state, the controller strategy and the IDS mechanism
- can attack on **actuators** unless detected
- tends to stay **stealthy**



## LTL Formulation of Controller Objective $\varphi$

- captures the controller task and the IDS mechanism
- reflects the behavior of stealthy attackers
- translates into a small DRA

$\varphi = \varphi_{IDS} \vee \varphi_{TASK}$ , where  $\varphi_{IDS}$  is a **reachability objective**

**Example: Counting-Based IDS**

$$\varphi_{IDS} = \diamond \left( \text{anomaly} \wedge \bigcirc \left( \text{anomaly} \wedge \bigcirc \diamond^{\leq 1} (\text{attack} \wedge \bigcirc \diamond \text{attack}) \right) \right)$$



# Secure Planning Case Studies

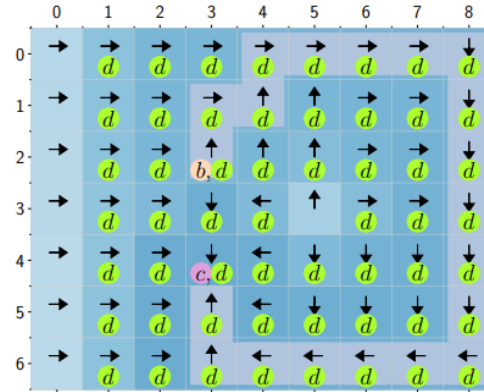
$$\varphi = \varphi_{IDS} \vee \varphi_{TASK}$$

$$\varphi_{IDS} = \diamond \left( \text{anomaly} \wedge \bigcirc \left( \text{anomaly} \wedge \bigcirc \diamond^{\leq 1} (\text{attack} \wedge \bigcirc \diamond \text{attack}) \right) \right)$$

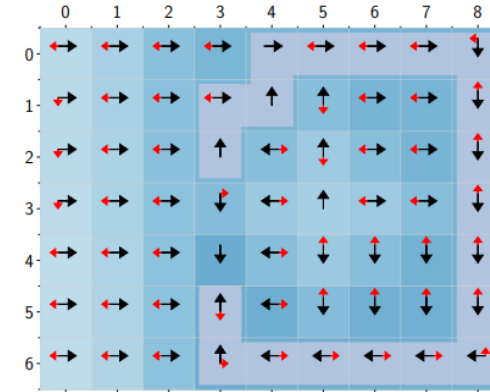
## Repeated Coverage:

- Repeatedly visit a **b** and a **c** cell
- Eventually reach a safe region labeled with **d** and do not leave

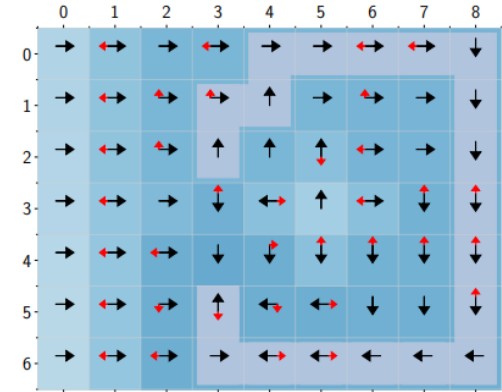
$$\varphi_{TASK} = \square \diamond \mathbf{b} \wedge \square \diamond \mathbf{c} \wedge \diamond \square \mathbf{d}$$



(a) The controller strategy from *b* to *c* and the labels of the cells



(b) The controller and the attacker strategies from *b* to *c* before any anomaly

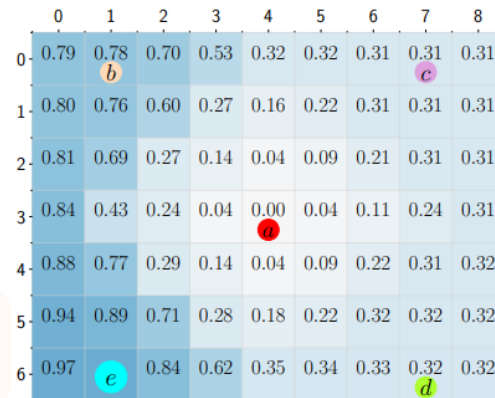


(c) The controller and the attacker strategies from *b* to *c* after one anomaly

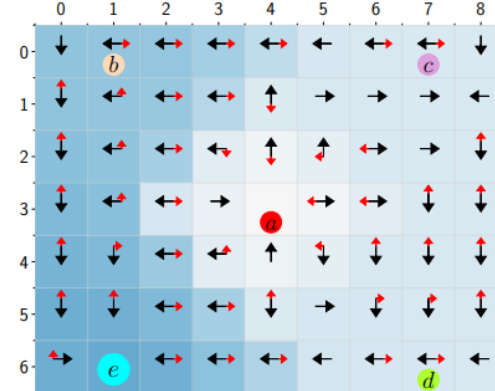
## Sequence of Tasks:

- Visit **b, c, d, e** in order
- Avoid the danger zone **a** at all costs

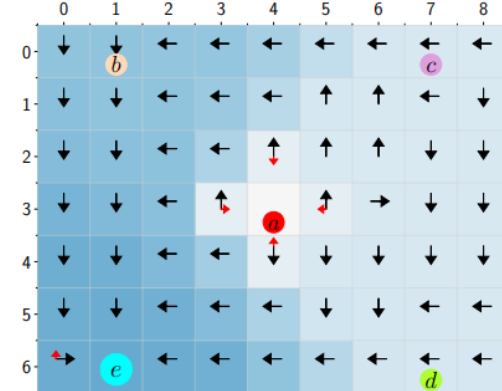
$$\varphi_{TASK} = \diamond \left( \mathbf{b} \wedge \diamond \left( \mathbf{c} \wedge \diamond \left( \mathbf{d} \wedge \diamond \mathbf{e} \right) \right) \right) \wedge \square \neg \mathbf{a}$$



(a) The controller strategy from *d* to *e* and the labels of the cells



(b) The controller and the attacker strategies from *d* to *e* right after an anomaly happens



(c) The controller and the attacker strategies from *d* to *e* right after an alarm is raised

## Summary:

- We convert a control synthesis problem in stochastic games to a reinforcement learning problem
- A controller strategy maximizing the return maximizes the satisfaction probability
- Our method does not require (or learn) the transition probabilities or the topology
- Convergence of reinforcement learning is ensured

## Future Work:

- More practical algorithms that converge to the desired strategy faster
- The use of approximate reinforcement learning to handle large state spaces

# Thank you

---



**Duke**  
UNIVERSITY

PRATT SCHOOL *of*  
**ENGINEERING**