# Learning Reward Machines for Partially Observable Reinforcement Learning

**Rodrigo Toro Icarte**    Ethan Waldie    Toryn Q. Klassen    Richard Valenzano
Margarita P. Castro    Sheila A. McIlraith

UNIVERSITY OF
TORONTO

VECTOR
INSTITUTE

E L E M E N T ^AI

**NeurIPS 2019**
December 11

**What is a**
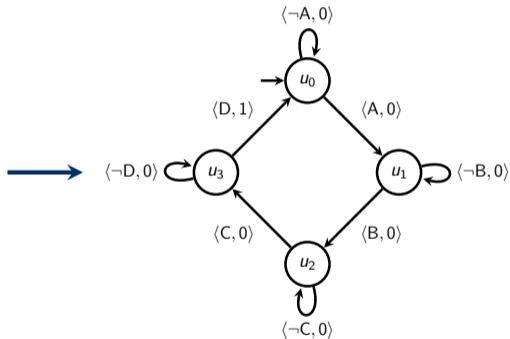# Reward Machine (RM)?

👀

# Reward Machines (RMs)

RMs are automata-based reward functions:

# Reward Machines (RMs)

RMs are automata-based reward functions:
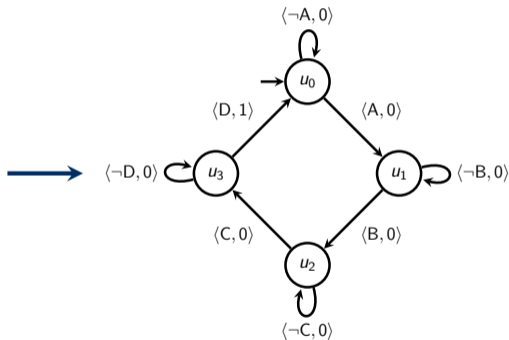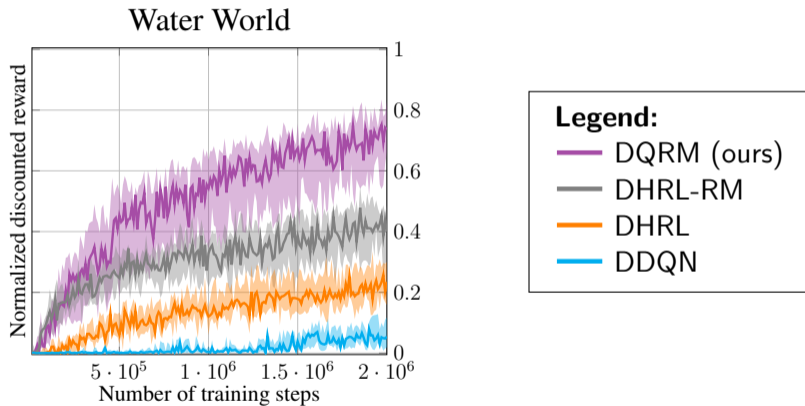
```
 1  m = 0 # global variable
 2  def get_reward(s):
 3    if m == 0 and s.at("A"):
 4      m = 1
 5    if m == 1 and s.at("B"):
 6      m = 2
 7    if m == 2 and s.at("C"):
 8      m = 3
 9    if m == 3 and s.at("D"):
10      m = 0
11      return 1
12    return 0
```
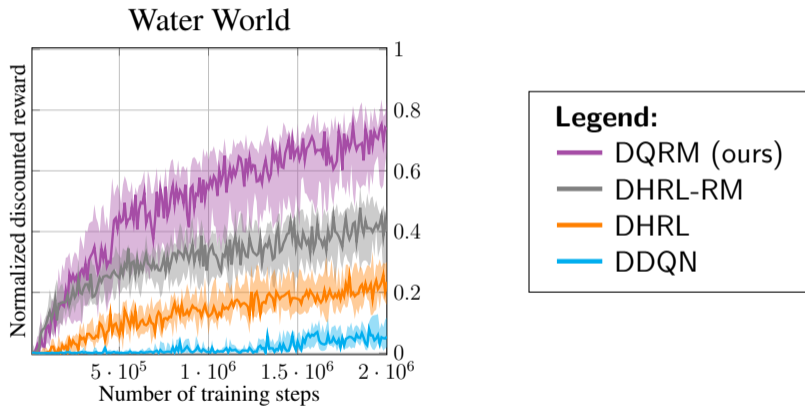
# Reward Machines (RMs)

RMs are automata-based reward functions:

```
 1  m = 0 # global variable
 2  def get_reward(s):
 3      if m == 0 and s.at("A"):
 4          m = 1
 5      if m == 1 and s.at("B"):
 6          m = 2
 7      if m == 2 and s.at("C"):
 8          m = 3
 9      if m == 3 and s.at("D"):
10          m = 0
11          return 1
12      return 0
```



... that allow for **learning policies faster**.

Water World

Legend:
- DQRM (ours)
- DHRL-RM
- DHRL
- DDQN

… but the RMs were **handcrafted**.

**This work**:

1. Shows how to learn RMs from experiences (LRM).

# Learning Reward Machines for Partially Observable RL

**This work**:

1. Shows how to learn RMs from experiences (LRM).
2. Uses RMs as memory for partially observable RL.

**This work**:

1. Shows how to learn RMs from experiences (LRM).
2. Uses RMs as memory for partially observable RL.
3. Extends QRM to work under partial observability.

**This work**:

1. Shows how to learn RMs from experiences (LRM).
2. Uses RMs as memory for partially observable RL.
3. Extends QRM to work under partial observability.
4. Provides a theoretical and empirical analysis of LRM.

# The Cookie Domain

**Agent**

# The cookie domain



Button

(Cookie)

**(+1 Reward)**

Solving the cookie domain requires **memory**!

Solving the cookie domain requires **memory**!

$$\pi^*(a|o_t) \neq \pi^*(a|o_0, \cdots, o_t)$$

## Partially Observable RL

**The most popular approach**:

Training LSTMs policies using a policy gradient method.

**The most popular approach**:

Training LSTMs policies using a policy gradient method.

**... starves in the cookie domain**.



**Legend**:
- Optimal
- ACER
- A3C
- PPO
- DDQN

# RMs as memory

If the agent can detect the color of the rooms (■, □, ■, □),

## Reward Machines as memory

If the agent can detect the color of the rooms (■, □, ■, □), and when it presses the button (○),

If the agent can detect the color of the rooms (🟥, ⬜, 🟦, 🟨), and when it presses the button (🔵), eats a cookie (☺),

If the agent can detect the color of the rooms (🟥, ⬜, 🟦, 🟨), and when it presses the button (🔵), eats a cookie (☉), and sees a cookie (🍪),

If the agent can detect the color of the rooms (■, □, ■, □), and when it presses the button (◯), eats a cookie (⊙), and sees a cookie (●), then:



... becomes a **"perfect" memory** for the cookie domain.

# Reward Machines as memory

| conditions at state $u_0$ | | |
|---|---|---|
| **if** ($\square \bigcirc$) | $\rightarrow$ | **goto** $u_1$ |
| **else** | $\rightarrow$ | **goto** $u_0$ |

| conditions at state $u_0$ | | |
|---|---|---|
| **if** ($\square$ $\bigcirc$) | $\rightarrow$ | **goto** $u_1$ |
| **else** | $\rightarrow$ | **goto** $u_0$ |

| conditions at state $u_0$ | | |
|---|---|---|
| **if** ($\square$ $\bigcirc$) | $\rightarrow$ | **goto** $u_1$ |
| **else** | $\rightarrow$ | **goto** $u_0$ |

# Reward Machines as memory

| conditions at state $u_0$ | | |
|---|---|---|
| **if** ($\square$ $\bullet$) | $\rightarrow$ | **goto** $u_1$ |
| **else** | $\rightarrow$ | **goto** $u_0$ |

# Reward Machines as memory



| conditions at state $u_0$ | | |
|---|---|---|
| **if** ($\square$ $\bigcirc$) | $\rightarrow$ | **goto** $u_1$ |
| **else** | $\rightarrow$ | **goto** $u_0$ |

| conditions at state $u_1$ | | |
| --- | --- | --- |
| **if** (🟦 **or** 🟥 🍪) | $\rightarrow$ | **goto** $u_2$ |
| **if** (🟥 **or** 🟦 🍪) | $\rightarrow$ | **goto** $u_3$ |
| **else** | $\rightarrow$ | **goto** $u_1$ |

| conditions at state $u_1$ | | |
| --- | --- | --- |
| if (🟦 or 🟥🍪) | $\rightarrow$ | goto $u_2$ |
| if (🟥 or 🟦🍪) | $\rightarrow$ | goto $u_3$ |
| else | $\rightarrow$ | goto $u_1$ |

# Reward Machines as memory



|   | conditions at state $u_1$ |   |   |
|---|---|---|---|
| if (🟦 or 🟥🍪) | → | goto $u_2$ |
| if (🟥 or 🟦🍪) | → | goto $u_3$ |
| else | → | goto $u_1$ |

| conditions at state $u_1$ | | |
|---|---|---|
| **if** (🟦 or 🟥 🍪) | $\rightarrow$ | **goto** $u_2$ |
| **if** (🟥 or 🟦 🍪) | $\rightarrow$ | **goto** $u_3$ |
| **else** | $\rightarrow$ | **goto** $u_1$ |

| conditions at state $u_1$ | | |
|---|---|---|
| if ($\blacksquare$ or $\blacksquare$ 🍪) | $\rightarrow$ | **goto** $u_2$ |
| if ($\blacksquare$ or $\blacksquare$ 🍪) | $\rightarrow$ | **goto** $u_3$ |
| **else** | $\rightarrow$ | **goto** $u_1$ |

| conditions at state $u_1$ | | |
|---|---|---|
| if (🟦 or 🟥🍪) | → | goto $u_2$ |
| if (🟥 or 🟦🍪) | → | goto $u_3$ |
| else | → | goto $u_1$ |

| conditions at state $u_1$ | | |
|---|---|---|
| if (🟦 or 🟥🍪) | $\rightarrow$ | goto $u_2$ |
| if (🟥 or 🟦🍪) | $\rightarrow$ | goto $u_3$ |
| else | $\rightarrow$ | goto $u_1$ |

# Reward Machines as memory



| conditions at state $u_1$ | | |
|---|---|---|
| **if** (🟦 or 🟥🍪) | $\rightarrow$ | **goto** $u_2$ |
| **if** (🟥 or 🟦🍪) | $\rightarrow$ | **goto** $u_3$ |
| **else** | $\rightarrow$ | **goto** $u_1$ |

| conditions at state $u_3$ | | |
|---|---|---|
| **if** ($\square \odot$) | $\rightarrow$ | **goto** $u_0$ |
| **else** | $\rightarrow$ | **goto** $u_3$ |

| conditions at state $u_3$ | | |
|---|---|---|
| **if** ($\blacksquare$ $\odot$) | $\rightarrow$ | **goto** $u_0$ |
| **else** | $\rightarrow$ | **goto** $u_3$ |

# Reward Machines as memory

# Reward Machines as memory



|          | conditions at state $u_3$ |          |
|----------|:--:|----------|
| **if** ($\blacksquare\,\odot$) | $\rightarrow$ | **goto** $u_0$ |
| **else** | $\rightarrow$ | **goto** $u_3$ |

| conditions at state $u_3$ | | |
|---|---|---|
| if ($\blacksquare$ $\odot$) | $\rightarrow$ | goto $u_0$ |
| else | $\rightarrow$ | goto $u_3$ |

| conditions at state $u_3$ | | |
|---|---|---|
| if ($\blacksquare \odot$) | $\rightarrow$ | goto $u_0$ |
| else | $\rightarrow$ | goto $u_3$ |

| conditions at state $u_3$ | | |
|---|---|---|
| **if** ($\blacksquare \odot$) | $\rightarrow$ | **goto** $u_0$ |
| **else** | $\rightarrow$ | **goto** $u_3$ |

| conditions at state $u_3$ | | |
|---|---|---|
| if ($\blacksquare$ $\odot$) | $\rightarrow$ | goto $u_0$ |
| else | $\rightarrow$ | goto $u_3$ |

# Reward Machines as memory



| conditions at state $u_3$ | | |
|---|---|---|
| if ($\square$ ☉) | $\rightarrow$ | goto $u_0$ |
| else | $\rightarrow$ | goto $u_3$ |

| conditions at state $u_3$ | | |
|---|---|---|
| if ($\square$ ☺) | $\rightarrow$ | goto $u_0$ |
| else | $\rightarrow$ | goto $u_3$ |

| conditions at state $u_3$ | | |
|---|---|---|
| **if** ($\square \odot$) | $\rightarrow$ | **goto** $u_0$ |
| **else** | $\rightarrow$ | **goto** $u_3$ |

| conditions at state $u_0$ | | |
|---|---|---|
| **if** ($\square \bigcirc$) | $\rightarrow$ | **goto** $u_1$ |
| **else** | $\rightarrow$ | **goto** $u_0$ |

# Reward Machines as memory

| conditions at state $u_0$ | | |
|---|---|---|
| if ($\square$ $\bullet$) | $\rightarrow$ | goto $u_1$ |
| else | $\rightarrow$ | goto $u_0$ |

| conditions at state $u_0$ | | |
|---|---|---|
| **if** (🟨 🔵) | $\rightarrow$ | **goto** $u_1$ |
| **else** | $\rightarrow$ | **goto** $u_0$ |

|  conditions at state $u_0$ | | |
| --- | --- | --- |
| **if** ($\square$ $\bullet$) | $\rightarrow$ | **goto** $u_1$ |
| **else** | $\rightarrow$ | **goto** $u_0$ |

|  | conditions at state $u_0$ | | |
|---|---|---|---|
| **if** ($\square$ $\bullet$) | $\rightarrow$ | **goto** $u_1$ |
| **else** | $\rightarrow$ | **goto** $u_0$ |

| conditions at state $u_0$ | | |
| --- | --- | --- |
| **if** ($\square$ $\bullet$) | $\rightarrow$ | **goto** $u_1$ |
| **else** | $\rightarrow$ | **goto** $u_0$ |

|  | conditions at state $u_0$ | | |
| --- | --- | --- | --- |
| **if** ($\square\,\bigcirc$) | $\rightarrow$ | **goto** $u_1$ | |
| **else** | $\rightarrow$ | **goto** $u_0$ | |

**Why is this a perfect memory?**

**Why is this a perfect memory?**

$$\pi^*(a|o_0, \cdots, o_t) = \pi^*(a|o_t, u_t)$$

**Why is this a perfect memory?**

$$\pi^*(a|o_0, \cdots, o_t) = \pi^*(a|o_t, u_t)$$

Hard problem $\xrightarrow{\text{RM}}$ Easy problem

# How to learn such RMs?

Given a set of detectors (e.g., $\{\blacksquare, \square, \blacksquare, \square, \bigcirc, \bullet, \odot\}$) and traces $\mathcal{T}$,

## Learning Reward Machines

Given a set of detectors (e.g., $\{■, □, ■, □, ●, ●, ☺\}$) and traces $\mathcal{T}$, learning RMs is a **discrete optimization** problem:

## Learning Reward Machines

Given a set of detectors (e.g., $\{\blacksquare, \square, \blacksquare, \square, \bigcirc, \odot, \odot\}$) and traces $\mathcal{T}$, learning RMs is a **discrete optimization** problem:

$$\underset{\langle U, u_0, \delta_u, \delta_r \rangle}{\text{minimize}} \sum_{i \in I} \sum_{t \in T_i} \log(|N_{x_{i,t}, L(e_{i,t})}|) \tag{LRM}$$

$$s.t. \ \langle U, u_0, \delta_u, \delta_r \rangle \in \mathcal{R}_\mathcal{P} \tag{1}$$

$$|U| \leq u_{\max} \tag{2}$$

$$x_{i,t} \in U \qquad \forall i \in I, t \in T_i \cup \{t_i\} \tag{3}$$

$$x_{i,0} = u_0 \qquad \forall i \in I \tag{4}$$

$$x_{i,t+1} = \delta_u(x_{i,t}, L(e_{i,t+1})) \qquad \forall i \in I, t \in T_i \tag{5}$$

$$N_{u,l} \subseteq 2^{2^\mathcal{P}} \qquad \forall u \in U, l \in 2^\mathcal{P} \tag{6}$$

$$L(e_{i,t+1}) \in N_{x_{i,t}, L(e_{i,t})} \qquad \forall i \in I, t \in T_i \tag{7}$$

## Learning Reward Machines

Given a set of detectors (e.g., $\{\blacksquare, \square, \blacksquare, \square, \bigcirc, \bullet, \odot\}$) and traces $\mathcal{T}$, learning RMs is a **discrete optimization** problem:

$$\underset{\langle U, u_0, \delta_u, \delta_r \rangle}{\text{minimize}} \sum_{i \in I} \sum_{t \in T_i} \log(|N_{x_{i,t}, L(e_{i,t})}|) \tag{LRM}$$

$$s.t. \ \langle U, u_0, \delta_u, \delta_r \rangle \in \mathcal{R}_\mathcal{P} \tag{1}$$

$$|U| \leq u_{\max} \tag{2}$$

$$x_{i,t} \in U \qquad \forall i \in I, t \in T_i \cup \{t_i\} \tag{3}$$

$$x_{i,0} = u_0 \qquad \forall i \in I \tag{4}$$

$$x_{i,t+1} = \delta_u(x_{i,t}, L(e_{i,t+1})) \qquad \forall i \in I, t \in T_i \tag{5}$$

$$N_{u,l} \subseteq 2^{2^\mathcal{P}} \qquad \forall u \in U, l \in 2^\mathcal{P} \tag{6}$$

$$L(e_{i,t+1}) \in N_{x_{i,t}, L(e_{i,t})} \qquad \forall i \in I, t \in T_i \tag{7}$$

Learn a
"*simple*"
**causal model**

## Learning Reward Machines

Given a set of detectors (e.g., $\{\blacksquare, \square, \blacksquare, \square, \bullet, \bullet, \odot\}$) and traces $\mathcal{T}$, learning RMs is a **discrete optimization** problem:

$$\underset{\langle U, u_0, \delta_u, \delta_r \rangle}{\text{minimize}} \sum_{i \in I} \sum_{t \in T_i} \log(|N_{x_{i,t}, L(e_{i,t})}|) \qquad \text{(LRM)}$$

$$\begin{aligned}
s.t. \ &\langle U, u_0, \delta_u, \delta_r \rangle \in \mathcal{R}_{\mathcal{P}} && (1) \\
&|U| \leq u_{\max} && (2) \\
&x_{i,t} \in U && \forall i \in I, t \in T_i \cup \{t_i\} \quad (3) \\
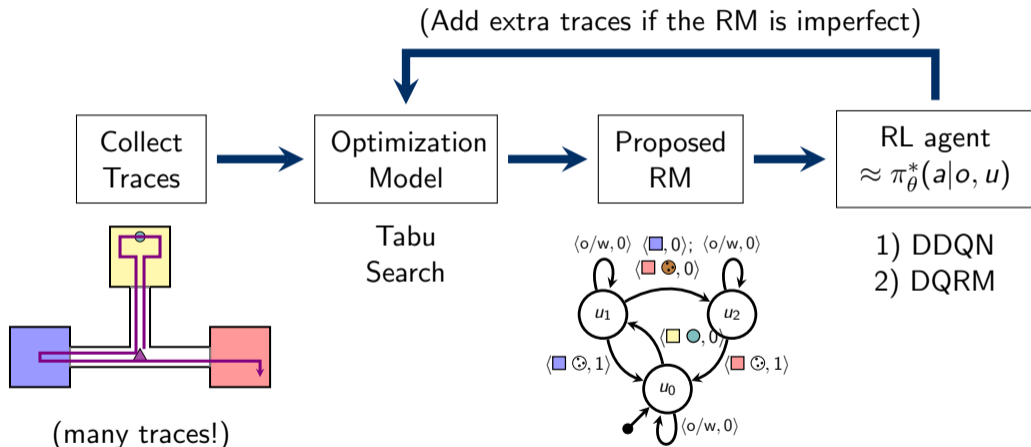&x_{i,0} = u_0 && \forall i \in I \quad (4) \\
&x_{i,t+1} = \delta_u(x_{i,t}, L(e_{i,t+1})) && \forall i \in I, t \in T_i \quad (5) \\
&N_{u,l} \subseteq 2^{2^{\mathcal{P}}} && \forall u \in U, l \in 2^{\mathcal{P}} \quad (6) \\
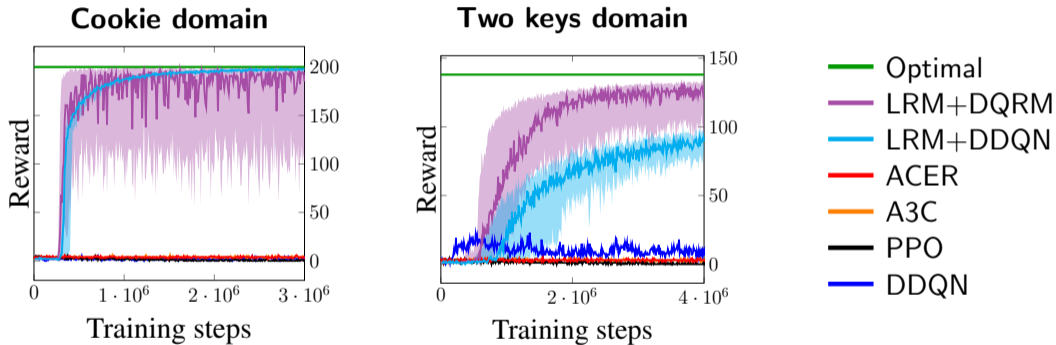&L(e_{i,t+1}) \in N_{x_{i,t}, L(e_{i,t})} && \forall i \in I, t \in T_i \quad (7)
\end{aligned}$$

… that we solved using **Tabu Search**.

# Results

# Results

## Cookie domain



## Two keys domain



- Optimal
- LRM+DQRM
- LRM+DDQN
- ACER
- A3C
- PPO
- DDQN

*Note: The detectors were also given to the baselines.

# Discussion at poster #210

`https://bitbucket.org/RToroIcarte/lrm`

**Thanks! :)**



Rodrigo  Ethan  Toryn  Rick  Margarita  Sheila