# Reinforcement Learning with Reward Machines
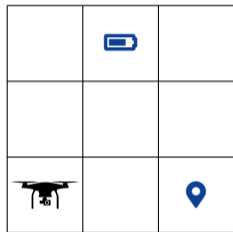
**Daniel Neider**
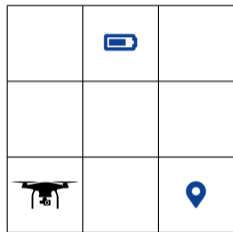
tu dortmund university

CENTER FOR TRUSTWORTHY DATA SCIENCE AND SECURITY

UA RUHR | RESEARCH ALLIANCE

agent

action

state +
reward

environment
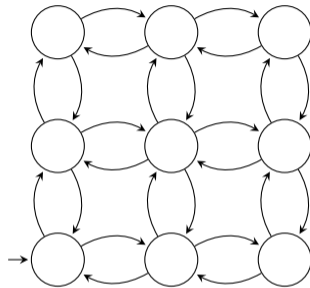
- Actions $A = \{\uparrow, \downarrow, \leftarrow, \rightarrow\}$
- Labels $\mathcal{P} = \{\text{📍}, \text{🔋}\}$

- Actions $A = \{\uparrow, \downarrow, \leftarrow, \rightarrow\}$
- Labels $\mathcal{P} = \{\text{📍}, \text{🔋}\}$

$$\mathcal{M} = (S, s_I, A, \mathcal{P}, p, L, R, \gamma)$$

- Actions $A = \{\uparrow, \downarrow, \leftarrow, \rightarrow\}$
- Labels $\mathcal{P} = \{\textcolor{blue}{\bullet}, \textcolor{blue}{\blacksquare}\}$

$\mathcal{M} = (\textcolor{orange}{S}, s_I, A, \mathcal{P}, p, L, R, \gamma)$

<span style="color:orange">set $S$ of states</span>
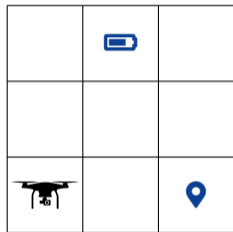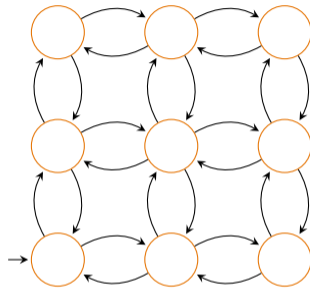
- Actions $A = \{\uparrow, \downarrow, \leftarrow, \rightarrow\}$
- Labels $\mathcal{P} = \{\text{📍}, \text{🔋}\}$
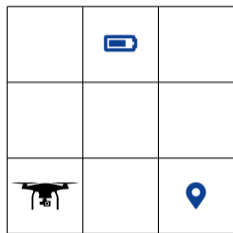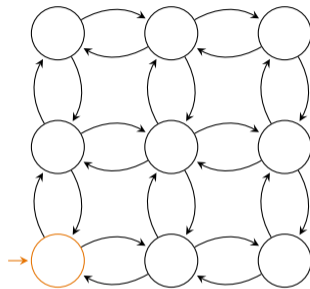
$\mathcal{M} = (S, s_I, A, \mathcal{P}, p, L, R, \gamma)$

initial state $s_I \in S$

- Actions $A = \{\uparrow, \downarrow, \leftarrow, \rightarrow\}$
- Labels $\mathcal{P} = \{$ 📍, 🔋 $\}$

$$\mathcal{M} = (S, s_I, A, \mathcal{P}, p, L, R, \gamma)$$
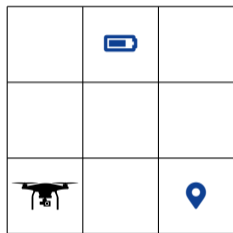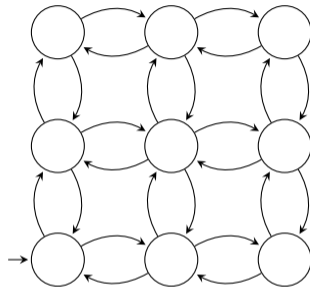
actions $A$ and labels $\mathcal{P}$
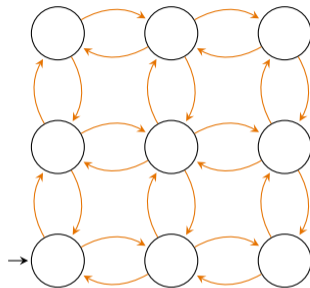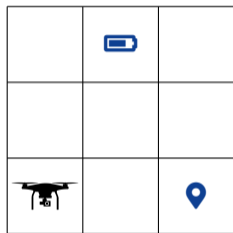
- Actions $A = \{\uparrow, \downarrow, \leftarrow, \rightarrow\}$
- Labels $\mathcal{P} = \{\text{\textpin}, \text{\textbattery}\}$

$$\mathcal{M} = (S, s_I, A, \mathcal{P}, p, L, R, \gamma)$$

transition function
$p: S \times A \times S \to [0, 1]$

- Actions $A = \{\uparrow, \downarrow, \leftarrow, \rightarrow\}$
- Labels $\mathcal{P} = \{ \text{📍}, \text{🔋} \}$

$\mathcal{M} = (S, s_I, A, \mathcal{P}, p, L, R, \gamma)$

labeling function
$L: S \times A \times S \rightarrow \mathcal{P}$

Actions $A = \{\uparrow, \downarrow, \leftarrow, \rightarrow\}$

Labels $\mathcal{P} = \{\text{📍}, \text{🔋}\}$

$\mathcal{M} = (S, s_I, A, \mathcal{P}, p, L, R, \gamma)$

reward function
$R \colon S \times A \times S \to \mathbb{R}$

- Actions $A = \{\uparrow, \downarrow, \leftarrow, \rightarrow\}$
- Labels $\mathcal{P} = \{\textcolor{blue}{\pmb{\varphi}}, \textcolor{blue}{\pmb{\boxminus}}\}$

$\mathcal{M} = (S, s_I, A, \mathcal{P}, p, L, R, \textcolor{orange}{\gamma})$

discount factor $\gamma \in (0, 1)$

Find a (probabilistic) policy $\pi \colon S \times A \to [0,1]$ maximizing the expected discounted reward

$$\mathbb{E}_\pi \left[ \sum_{i=0}^{k} \gamma^i \cdot R(s_i, a_{i+1}, s_{i+1}) \right]$$

of every trajectory $s_0 a_0 s_1 \ldots s_{k+1}$, $k \in \mathbb{N}$, through the MDP

Find a (probabilistic) policy $\pi\colon S \times A \to [0,1]$ maximizing the expected discounted reward

$$\mathbb{E}_\pi\left[\sum_{i=0}^{k} \gamma^i \cdot R(s_i, a_{i+1}, s_{i+1})\right]$$

of every trajectory $s_0 a_0 s_1 \ldots s_{k+1}$, $k \in \mathbb{N}$, through the MDP

## Q-Learning

1. Maintain a table $Q\colon S \times A \to \mathbb{R}$ (initialized to, e.g., 0)
2. Explore the environment according to $\pi$, resulting in a trajectory $s_0 a_1 s_1 a_2 s_2 \ldots$
3. In step $t$, update $Q$ by
   $Q(s_t, a_t) \leftarrow (1-\alpha) \cdot Q(s_t, a_t) + \alpha\big[R(s_t, a_{t+1}, s_{t+1}) + \gamma \max_a Q(s_{t+1}, a)\big]$
4. After each episode, update $\pi$ by $\pi(s, a) \leftarrow \arg\max_{a \in A} Q(s, a)$
5. Repeat until this process converges; $\pi$ is then the optimal policy

A; B; C; D:                    1

A; B; C; D:        1

A; B; B; C; D:        1

| A; B; C; D: | 1 |
| A; B; B; C; D: | 1 |
| A; D: | -1 |
| A; C; D: | -1 |

| | | | | |
|---|---|---|---|---|
| A | | D | | |
| | | | | B |
| | | C | | |
| | A | | | B |
| | | | | |

| | |
|---|---:|
| A; B; C; D: | 1 |
| A; B; B; C; D: | 1 |
| A; D: | -1 |
| A; C; D: | -1 |

How to handle such situations?

**Joint Inference of Reward Machines and Policies for Reinforcement Learning**

Zhe Xu,[1][*] Ivan Gavran,[2][*] Yousef Ahmad,[1] Rupak Majumdar,[2] Daniel Neider,[2] Ufuk Topcu,[1]
Bo Wu[1]

**Reinforcement Learning with Stochastic Reward Machines**

Jan Corazza[1,2], Ivan Gavran[2], Daniel Neider[2]

-sws.org

**Advice-Guided Reinforcement Learning in a non-Markovian Environment**

Daniel Neider[1], Jean-Raphael Gaglione[2], Ivan Gavran[1], Ufuk Topcu[3], Bo Wu[3], Zhe Xu[4]

[1] Max Planck Institute for Software Systems, Kaiserslautern, Germany
[2] Ecole Polytechnique, France
[3] University of Texas at Austin, Texas, USA
[4] Arizona State University, Arizona, USA

# 1. Joint Inference of Policies and Reward Machines

**(joint work with Yousef Ahmad, Ivan Gavran, Rupak Majumdar, Ufuk Topcu, Bo Wu, and Zhe Xu)**

A; B; C; D:          1

A; B; B; C; D:        1

A; D:               -1

A; C; D:            -1

Bacchus et al. (1996)
Jothimurgan et al. (2019)
Icarte et al. (2018)
Brafman et al. (2018)

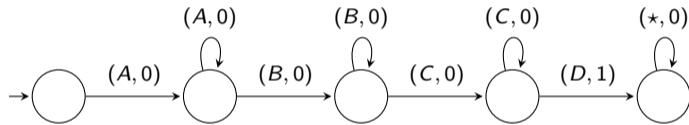"Use automata/temporal logic to capture non-Markovian rewards"

A; B; C; D:          1

A; B; B; C; D:        1

A; D:                -1

A; C; D:             -1

Bacchus et al. (1996)
Jothimurgan et al. (2019)
Icarte et al. (2018)
Brafman et al. (2018)

| A; B; C; D: | 1 |

| A; B; B; C; D: | 1 |

| A; D: | -1 |

| A; C; D: | -1 |

Bacchus et al. (1996)
Jothimurgan et al. (2019)
Icarte et al. (2018)
Brafman et al. (2018)

A; B; C; D:           1

A; B; B; C; D:        1

A; D:                 -1

A; C; D:              -1
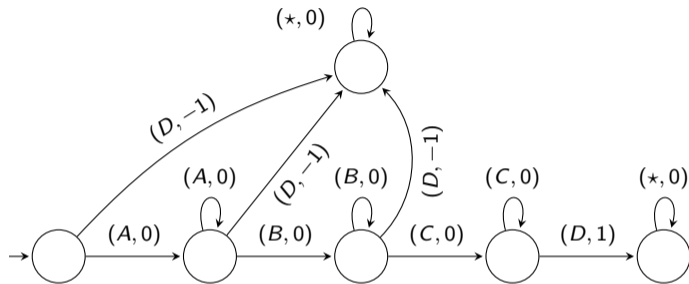
Bacchus et al. (1996)
Jothimurgan et al. (2019)
Icarte et al. (2018)
Brafman et al. (2018)

A; B; C; D:          1

A; B; B; C; D:       1

A; D:               -1

A; C; D:            -1

Bacchus et al. (1996)
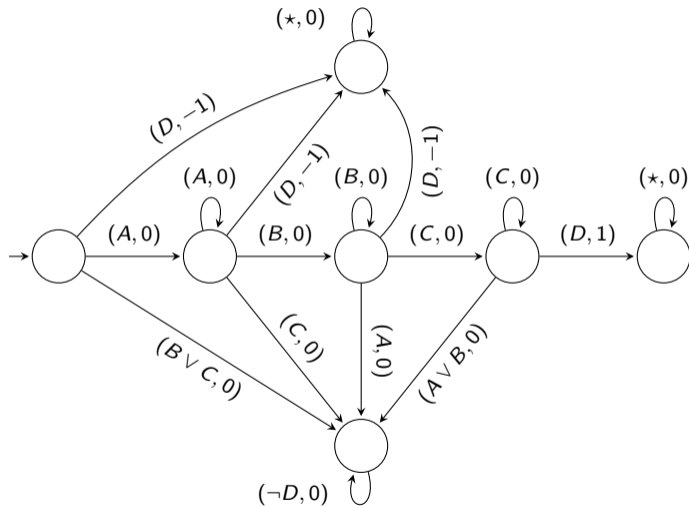Jothimurgan et al. (2019)
Icarte et al. (2018)
Brafman et al. (2018)

| A | | D | | |
|---|---|---|---|---|
| | | | | B |
| | | C | | |
| | A | | | B |
| | | | | |

A; B; C; D:        1

A; B; B; C; D:        1

A; D:        -1

A; C; D:        -1

Bacchus et al. (1996)
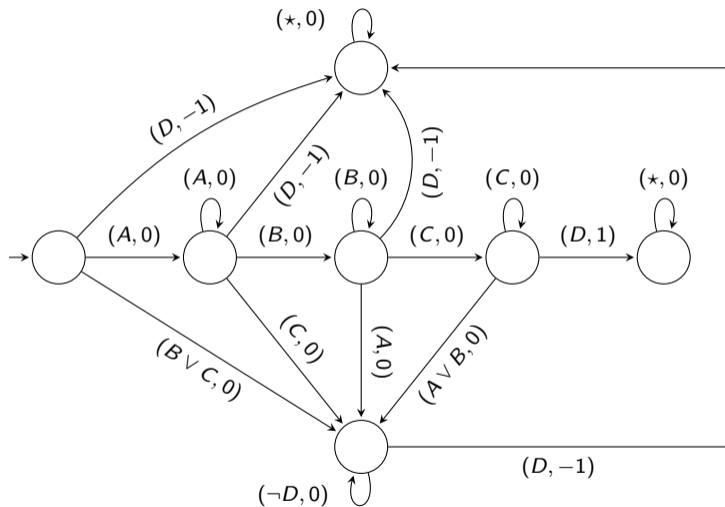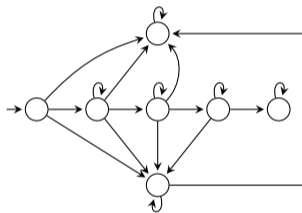Jothimurgan et al. (2019)
Icarte et al. (2018)
Brafman et al. (2018)

Icarte et al. (2018) have proposed an extension of the Q-learning algorithm, named QRM, that can handle reward machines

- ▶ avoids building the cross-product explicitly
- ▶ exploits the structure of the reward machine during exploration



ICML 2018

Icarte et al. (2018) have proposed an extension of the Q-learning algorithm, named QRM, that can handle reward machines

- ▶ avoids building the cross-product explicitly
- ▶ exploits the structure of the reward machine during exploration

## Problem

How does one construct reward machines?

- ▶ direct construction, from temporal logics, learning, . . .



ICML 2018

## Key idea

- ▶ Given the current hypothesis reward machine $H$, perform QRM and record the resulting label sequence $\lambda = \ell_1 \ldots \ell_n$ and reward sequence $\rho = r_1 \ldots r_n$
- ▶ If the pair $(\lambda, \rho)$ contradicts $H$, learn a new reward machine $H'$
- ▶ Repeat until this process converges to the "true" reward machine and an optimal policy

labeled MDP

Initialize reward machine $H$;

$H$: $\quad \rightarrow \bigcirc\!\!\!p_1 \circlearrowright (\star, 0)$

Initialize reward machine $H$;

Initialize a set $Q$ of $q$-functions;

$H$: $\rightarrow \left( p_1 \right) \circlearrowleft (\star, 0)$

$Q$: $\left\{ q^{p_1} \right\}$

Initialize reward machine $H$;

Initialize a set $Q$ of $q$-functions;

Initialize a sample $X$ of traces;

$H$:    $\rightarrow \left(p_1\right) \circlearrowleft (\star, 0)$

$Q$:    $\left\{ q^{p_1} \right\}$

$X$:    $\emptyset$

Initialize $H$, $Q$, $X$;

**repeat**

    $(\lambda, \rho, Q) \leftarrow \mathrm{QRM}(H, Q);$

$H$:    $\rightarrow \boxed{p_1} \circlearrowleft (\star, 0)$

$Q$:    $\left\{ q^{p_1} \right\}$

$X$:    $\emptyset$

Initialize $H$, $Q$, $X$;

**repeat**

$A; B; C; D$

$(\lambda, \rho, Q) \leftarrow \text{QRM}(H, Q);$

$H$:     $\rightarrow \widehat{p_1} \circlearrowleft (\star, 0)$

$Q$:     $\{ q^{p_1} \}$

$X$:     $\emptyset$

Initialize $H$, $Q$, $X$;

**repeat**

$(\lambda, \rho, Q) \leftarrow \text{QRM}(H, Q);$

$A; B; C; D$

$0; 0; 0; 1$

$H:$ → $p_1$ ⟲ $(\star, 0)$

$Q:$ $\left\{ q^{p_1} \right\}$

$X:$ $\emptyset$

Initialize $H$, $Q$, $X$;

**repeat**

$(\lambda, \rho, Q) \leftarrow \mathrm{QRM}(H, Q)$;

$\boxed{A; B; C; D}$

$\boxed{0; 0; 0; 1}$

**if** $H(\lambda) \neq \rho$ **then**

add $(\lambda, \rho)$ to $X$;

$H:$   →  $p_1$   ↺  $(\star, 0)$

$Q:$   $\{q^{p_1}\}$

$X:$   $\emptyset$

Initialize $H$, $Q$, $X$;

**repeat**

$(\lambda, \rho, Q) \leftarrow \text{QRM}(H, Q);$

$A; B; C; D$

$0; 0; 0; 1$

**if** $H(\lambda) \neq \rho$ **then**

add $(\lambda, \rho)$ to $X$;

$H: \quad \rightarrow \boxed{p_1} \; \circlearrowleft \; (\star, 0)$

$Q: \quad \{q^{p_1}\}$

$X: \quad \{(A; B; C; D / 0; 0; 0; 1)\}$

Initialize $H$, $Q$, $X$;

**repeat**

    $(\lambda, \rho, Q) \leftarrow \mathrm{QRM}(H, Q)$;

    **if** $H(\lambda) \neq \rho$ **then**

        add $(\lambda, \rho)$ to $X$;

    **if** $X$ was modified **then**

        $H \leftarrow \mathrm{infer}(X)$;

$H$:    $\rightarrow \left(\!\! p_1 \!\!\right) \circlearrowleft (\star, 0)$

$Q$:    $\left\{ q^{p_1} \right\}$

$X$:    $\left\{ (A; B; C; D/0; 0; 0; 1) \right\}$

Initialize $H$, $Q$, $X$;

**repeat**

$(\lambda, \rho, Q) \leftarrow \mathrm{QRM}(H, Q)$;

**if** $H(\lambda) \neq \rho$ **then**
add $(\lambda, \rho)$ to $X$;

**if** $X$ was modified **then**
$H \leftarrow \mathrm{infer}(X)$;

$H$:



$Q$: $\{q^{p_1}\}$

$X$: $\{(A; B; C; D/0; 0; 0; 1)\}$

Initialize $H$, $Q$, $X$;

**repeat**

    $(\lambda, \rho, Q) \leftarrow \text{QRM}(H, Q)$;

    **if** $H(\lambda) \neq \rho$ **then**

        add $(\lambda, \rho)$ to $X$;

    **if** $X$ was modified **then**

        $H \leftarrow \text{infer}(X)$;

        re-initialize $Q$ if necessary;

$H$:



$Q$:     $\{q^{p_1}\}$

$X$:     $\{(A; B; C; D/0; 0; 0; 1)\}$

Initialize $H$, $Q$, $X$;

**repeat**

$\quad (\lambda, \rho, Q) \leftarrow \text{QRM}(H, Q);$

$\quad$ **if** $H(\lambda) \neq \rho$ **then**

$\quad\quad$ add $(\lambda, \rho)$ to $X$;

$\quad$ **if** $X$ was modified **then**

$\quad\quad H \leftarrow \text{infer}(X);$

$\quad\quad$ re-initialize $Q$ if necessary;

$H$:



$Q$: $\quad \left\{ q^{p_1}, q^{p_2} \right\}$

$X$: $\quad \left\{ (A; B; C; D/0; 0; 0; 1) \right\}$

Initialize $H$, $Q$, $X$;

**repeat**

$(\lambda, \rho, Q) \leftarrow \mathsf{QRM}(H, Q)$;

**if** $H(\lambda) \neq \rho$ **then**
add $(\lambda, \rho)$ to $X$;

**if** $X$ was modified **then**
$H \leftarrow \mathsf{infer}(X)$;
re-initialize $Q$ if necessary;

It is crucial to infer **minimal** reward machines

We use two sets of propositional variables to encode reward machines:

$d_{p,\ell,q}$      encodes the transition function of the reward machine
(i.e., the machine transitions from state $p$ to state $q$ on reading symbol $\ell$)

$o_{p,\ell,r}$      encodes the output function of the reward machine
(i.e., the machine outputs reward $r$ in state $p$ on reading symbol $\ell$)

We use two sets of propositional variables to encode reward machines:

$d_{p,\ell,q}$      encodes the transition function of the reward machine
(i.e., the machine transitions from state $p$ to state $q$ on reading symbol $\ell$)

$o_{p,\ell,r}$      encodes the output function of the reward machine
(i.e., the machine outputs reward $r$ in state $p$ on reading symbol $\ell$)

## Enforcing deterministic functions

We impose pseudo-Boolean constraints to enforce for each pair of state $p$ and input $a$ that

- ▶ exactly one variable $d_{p,\ell,q}$ is set to `true`
- ▶ exactly one variable $o_{p,\ell,r}$ is set to `true`

We introduce a set of auxiliary variables:

$x_{\lambda,p}$     encodes the run of the reward machine on all prefixes of examples
           (i.e., the machine reaches state $p$ after reading the prefix $\lambda$)

We introduce a set of auxiliary variables:

$x_{\lambda,p}$      encodes the run of the reward machine on all prefixes of examples
(i.e., the machine reaches state $p$ after reading the prefix $\lambda$)

## Enforcing consistency with the examples

$$\left[ \bigwedge_{u \in Pref(X)} \text{one}(x_{u,q_1}, \ldots, x_{u,q_n}) \right] \wedge x_{\varepsilon, q_I}$$

reward machine:

$$(x_{\lambda,p} \wedge d_{p,\ell,q}) \rightarrow x_{\lambda\ell,q}$$



$$x_{\lambda,p} \rightarrow o_{p,\ell,r}$$

$x_{\varepsilon,q_I}$

We introduce a set of auxiliary variables:

$x_{\lambda,p}$      encodes the <span style="color:red">run</span> of the reward machine on all prefixes of examples
(i.e., the machine reaches state $p$ after reading the prefix $\lambda$)

## Enforcing consistency with the examples

$$\left[ \bigwedge_{u \in Pref(X)} \text{one}(x_{u,q_1}, \ldots, x_{u,q_n}) \right] \wedge x_{\varepsilon,q_I}$$

reward machine:



$$(x_{\lambda,p} \wedge d_{p,\ell,q}) \to x_{\lambda\ell,q}$$

$$x_{\lambda,p} \to o_{p,\ell,r}$$

We introduce a set of auxiliary variables:

$x_{\lambda,p}$   encodes the run of the reward machine on all prefixes of examples
(i.e., the machine reaches state $p$ after reading the prefix $\lambda$)

## Enforcing consistency with the examples

$$\left[ \bigwedge_{u \in Pref(X)} \text{one}(x_{u,q_1}, \ldots, x_{u,q_n}) \right] \wedge x_{\varepsilon,q_I}$$

$$(x_{\lambda,p} \wedge d_{p,\ell,q}) \rightarrow x_{\lambda\ell,q}$$

$$x_{\lambda,p} \rightarrow o_{p,\ell,r}$$

reward machine:



$$\rightarrow \boxed{q_I} \dashrightarrow^{(\lambda, \rho)} \boxed{p} \xrightarrow{(\ell, r)} \boxed{q}$$

$x_{\varepsilon,q_I}$          $x_{\lambda,p}$        $x_{\lambda\ell,q}$

### Theorem (Ahmad, Gavran, Majumdar, N., Topcu, Wu, and Xu)

*Given*

- *a sufficient episode length*
- *an $\varepsilon$-greedy exploration strategy*

*we have the following:*

1. *JIRP almost surely learns the "true" reward machine*
2. *JIRP almost surely converges to an optimal policy*

## Office World Scenario (Icarte et al., 2018)



## Conclusion

▶ JIRP is the only method that converges to an optimal policy

▶ JIRP converges faster than any of the competing methods

**2.** **Reinforcement Learning with Stochastic Reward Machines**

(joint work with Jan Corazza and Ivan Gavran)

T: tools 🔧    M: market 💲

S: silver mine 🔺    G: gold mine 🔺

T; G; M:  1.9

T: tools 🛠
M: market 💲
S: silver mine 🗻
G: gold mine 🔺

T: tools ⚒

M: market 💲

S: silver mine ⛰

G: gold mine ⛰

| T; G; M: | 1.9 |
| T; G; M: | 2.2 |
| T; S; M: | 1.2 |
| T; S; M: | 0.9 |

T: tools 🔧     M: market 💲

S: silver mine 🔺     G: gold mine 🔺

| T; G; M: | 1.9 |
| T; G; M: | 2.2 |
| T; S; M: | 1.2 |
| T; S; M: | 0.9 |

T: tools 🛠    M: market 💲

S: silver mine 🗻    G: gold mine 🗻

- If the label sequences are identical, no reward machine matches both traces
- If the label sequences are different, the resulting reward machine explodes in size

| T; G; M: | 1.9 |
| T; G; M: | 2.2 |
| T; S; M: | 1.2 |
| T; S; M: | 0.9 |

Outputs are bounded continuous distributions

Two SRMs $A$ and $B$ are equivalent in expectation ($A \sim_E B$) if they output sequences of distributions with equal expected values for each label sequence

Two SRMs $A$ and $B$ are equivalent in expectation ($A \sim_E B$) if they output sequences of distributions with equal expected values for each label sequence

## Corollary

If two SRMs are equivalent in expectation, then they induce the same optimal policy in an environment

## A naive algorithm

1. Collect many samples
2. Take the average reward in every position of the same trajectory
3. Construct an ordinary reward machine based on the average rewards

| | |
|---|---|
| T; G; M: | 1.9 |
| T; G; M: | 2.2 |
| T; G; M: | 2.1 |
| T; G; M: | 2.0 |

average: 2.05

## A naive algorithm

1. Collect many samples
2. Take the average reward in every position of the same trajectory
3. Construct an ordinary reward machine based on the average rewards

## Problem

Collecting samples is too slow!

| T; G; M: | 1.9 |
| T; G; M: | 2.2 |
| T; G; M: | 2.1 |
| T; G; M: | 2.0 |

average: 2.05

1. Probability distributions are continuous
   and have bounded support with
   "width" $\delta$

1. Probability distributions are continuous and have bounded support with "width" $\delta$

2. The noise from one distribution does not fully conceal the signal from another one (except in symmetric circumstances)

we will eventually observe
enough rewards to differentiate

1. Probability distributions are continuous and have bounded support with "width" $\delta$

2. The noise from one distribution does not fully conceal the signal from another one (except in symmetric circumstances)

no differentiation necessary:
equal expectations

1. Probability distributions are continuous and have bounded support with "width" $\delta$

2. The noise from one distribution does not fully conceal the signal from another one (except in symmetric circumstances)

Initialize $H$, $Q$, $X$, $A$;

**repeat**

    $(\lambda, \rho, Q) \leftarrow \mathrm{QRM}(H, Q)$;

    add $(\lambda, \rho)$ to $A$;

one would keep a
moving average in practice

Initialize $H$, $Q$, $X$, $A$;

**repeat**

    $(\lambda, \rho, Q) \leftarrow \mathsf{QRM}(H, Q)$;

    add $(\lambda, \rho)$ to $A$;

    **if** $H$ is not $\delta$-consistent with $(\lambda, \rho)$ **then**

        add $(\lambda, \rho)$ to $X$;

        $H' \leftarrow \mathsf{infer}(X)$;

infers a minimal
$\delta$-consistent "proto"-SRM

(only cares for $\delta$-consistency,
not estimating distribution)

Initialize $H$, $Q$, $X$, $A$;

**repeat**

    $(\lambda, \rho, Q) \leftarrow$ QRM$(H, Q)$;

    add $(\lambda, \rho)$ to $A$;

    **if** $H$ is not $\delta$-consistent with $(\lambda, \rho)$ **then**

        add $(\lambda, \rho)$ to $X$;

        $H' \leftarrow$ infer$(X)$;

        $H \leftarrow$ estimate$(H', A)$;

        re-initialize $Q$ if necessary;

corrects outputs of $H'$
by estimating distribution
parameters from samples in $A$

sample $X$ $\longrightarrow$

$n \leftarrow 0$

$n \leftarrow n + 1$

create SMT
formula $\Psi_n^X$

is $\Psi_n^X$ satisfiable?

no

yes $\longrightarrow$ derive "proto"-SRM
of size $n$ from model

We use propositional and real-valued variables to encode a "proto"-SRM:

$d_{p,\ell,q} \in \mathbb{B}$      encodes the transition function of the reward machine

$x_{\lambda,p} \in \mathbb{B}$      encodes the run of the reward machine on prefixes from $X$

$o_{p,\ell} \in \mathbb{R}$      encodes a "conjectured mean" of an output distribution
                (i.e., the distr. returned in state $p$ on reading symbol $\ell$ has mean $o_{p,\ell}$)

## Enforcing consistency with the examples

$$x_{\lambda,p} \to |o_{p,\ell} - r| \le \frac{\delta}{2}$$

SRM:

## Theorem (Corazza, Gavran, N.)

*Given*

- ▶ *a sufficient episode length*
- ▶ *an $\varepsilon$-greedy exploration strategy*
- ▶ *Assumptions 1 and 2 hold for the "true" (environment) SRM*

*we have the following:*

1. *SRMI almost surely learns a SRM that is equivalent in expectation to the "true" SRM*
2. *SRMI almost surely converges to an optimal policy*

## Mining Environment



stochastic rewards

deterministic rewards

## Conclusion

- ▶ SRMI converges faster than the baseline method
- ▶ SRMI's performance does not degrade in the case of deterministic rewards

# Inferring Probabilistic Reward Machines from Non-Markovian Reward Signals for Reinforcement Learning

**Taylor Dohmen[1*], Noah Topper[2*], George Atia[2], Andre Beckus[3],**
**Ashutosh Trivedi[1], Alvaro Velasquez[3]**

[1] University of Colorado Boulder
[2] University of Central Florida
[3] Air Force Research Laboratory

**Abstract**

The success of reinforcement learning in typical settings is

2021). They also serve as a memory mechanism for reasoning over partially observable environments (Icarte et al. 2019), are useful for reward shaping to mitigate sparse re-

**3.** **Advice-Guided Reinforcement Learning**

(joint work with Jean-Raphaël Gaglione, Ivan Gavran, Ufuk Topcu, Bo Wu, and Zhe Xu)

reward
machine
is given

Icarte et al., 2018

reward
machine
is inferred

Icarte et a., 2019
Furelos-Blanco et al., 2020
Gaon & Brafman, 2020
Xu et al., 2020

We formalize advice by means of regular languages:

- ► Deterministic Finite Automata (DFA)
- ► Regular expressions
- ► Linear Temporal Logic
- ► . . .

We formalize advice by means of regular languages:

▶ Deterministic Finite Automata (DFA)
▶ Regular expressions
▶ Linear Temporal Logic
▶ . . .



"every $A$ is
followed by $B$"

We formalize advice by means of regular languages:

- ▶ Deterministic Finite Automata (DFA)
- ▶ Regular expressions
- ▶ Linear Temporal Logic
- ▶ . . .

## Compatibility of advice DFAs (i.e., semantics)

*A reward can only be positive (negative/non-zero) if the advice DFA accepts the label sequence*

- ▶ A reward machine satisfying this property is called compatible



"every $A$ is followed by $B$"

$H$:　$\rightarrow$( $p_1$ ) $\circlearrowright$ $(\star, 0)$

Initialize reward machine $H$;

Initialize reward machine $H$;

Initialize a set $Q$ of $q$-functions;

$H$:    $\rightarrow \widehat{p_1} \circlearrowleft (\star, 0)$

$Q$:    $\{q^{p_1}\}$

Initialize reward machine $H$;

Initialize a set $Q$ of $q$-functions;

Initialize a sample $X$ of traces;

$H$:  $p_1$ $(\star, 0)$

$Q$: $\left\{ q^{p_1} \right\}$

$X$: $\emptyset$

Initialize reward machine $H$;

Initialize a set $Q$ of $q$-functions;

Initialize a sample $X$ of traces;

Initialize a set $D$ of advice DFAs;

$H$:    $(\star, 0)$

$Q$:   $\left\{ q^{p_1} \right\}$

$X$:   $\emptyset$

$D$:

Initialize $H$, $Q$, $X$, $D$;

**repeat**

$\quad (\lambda, \rho, Q) \leftarrow \text{QRM}(H, Q);$

$H:$ 

$Q: \quad \left\{ q^{p_1} \right\}$

$X: \quad \emptyset$

$D: \quad$

Initialize $H$, $Q$, $X$, $D$;

**repeat**

$(\lambda, \rho, Q) \leftarrow \text{QRM}(H, Q);$

$A; B; C; D$

$H$:  $\rightarrow p_1 \circlearrowleft (\star, 0)$

$Q$: $\left\{ q^{p_1} \right\}$

$X$: $\emptyset$

$D$:

Initialize $H$, $Q$, $X$, $D$;

**repeat**

$(\lambda, \rho, Q) \leftarrow \text{QRM}(H, Q);$

$A; B; C; D$

$0; 0; 0; 1$

$H$:  $\rightarrow \left(p_1\right) \circlearrowright (\star, 0)$

$Q$: $\left\{ q^{p_1} \right\}$

$X$: $\emptyset$

$D$:

Initialize $H$, $Q$, $X$, $D$;

**repeat**

$A; B; C; D$

$(\lambda, \rho, Q) \leftarrow \text{QRM}(H, Q);$

$0; 0; 0; 1$

    **if** $H(\lambda) \neq \rho$ **then**

        add $(\lambda, \rho)$ to $X$;

$H$:  $\rightarrow p_1 \;\; (\star, 0)$

$Q$: $\{q^{p_1}\}$

$X$: $\emptyset$

$D$:

Initialize $H$, $Q$, $X$, $D$;

**repeat**

$A; B; C; D$

$(\lambda, \rho, Q) \leftarrow \text{QRM}(H, Q);$

$0; 0; 0; 1$

    **if** $H(\lambda) \neq \rho$ **then**

        add $(\lambda, \rho)$ to $X$;

$H$:  $p_1$ $(\star, 0)$

$Q$: $\left\{ q^{p_1} \right\}$

$X$: $\left\{ (A; B; C; D / 0; 0; 0; 1) \right\}$

$D$:

Initialize $H$, $Q$, $X$, $D$;

**repeat**

$A; B; C; D$

$(\lambda, \rho, Q) \leftarrow \text{QRM}(H, Q);$

$0; 0; 0; 1$

**if** $H(\lambda) \neq \rho$ **then**

add $(\lambda, \rho)$ to $X$;

**if** $(\lambda, \rho)$ is not compatible with
some $\mathcal{D} \in D$ **then**

remove $\mathcal{D}$ from $D$;

$H$:  $\rightarrow (p_1) \circlearrowright (\star, 0)$

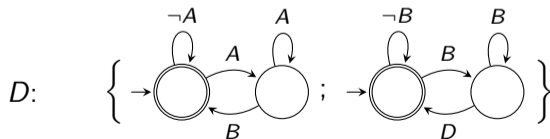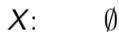$Q$: $\{q^{p_1}\}$

$X$: $\{(A; B; C; D/0; 0; 0; 1)\}$

$D$:

Initialize $H$, $Q$, $X$, $D$;

**repeat**

    $(\lambda, \rho, Q) \leftarrow \text{QRM}(H, Q)$;

    $A; B; C; D$

    $0; 0; 0; 1$

    **if** $H(\lambda) \neq \rho$ **then**

        add $(\lambda, \rho)$ to $X$;

    **if** $(\lambda, \rho)$ is not compatible with
        some $\mathcal{D} \in D$ **then**

        remove $\mathcal{D}$ from $D$;

$H$:  $(\star, 0)$ with state $p_1$

$Q$: $\left\{ q^{p_1} \right\}$

$X$: $\left\{ (A; B; C; D / 0; 0; 0; 1) \right\}$

$D$:

Initialize $H$, $Q$, $X$, $D$;

**repeat**

    $(\lambda, \rho, Q) \leftarrow \text{QRM}(H, Q)$;

    **if** $H(\lambda) \neq \rho$ **then**
        add $(\lambda, \rho)$ to $X$;

    **if** $(\lambda, \rho)$ is not compatible with
        some $\mathcal{D} \in D$ **then**
        remove $\mathcal{D}$ from $D$;

    **if** $X$ or $D$ were modified **then**
        $H \leftarrow \text{infer}(X, D)$;

$H$:  $(\star, 0)$

$Q$: $\left\{ q^{p_1} \right\}$

$X$: $\left\{ (A; B; C; D/0; 0; 0; 1) \right\}$

$D$:

Initialize $H$, $Q$, $X$, $D$;

**repeat**

    $(\lambda, \rho, Q) \leftarrow \text{QRM}(H, Q)$;

    **if** $H(\lambda) \neq \rho$ **then**
        add $(\lambda, \rho)$ to $X$;

    **if** $(\lambda, \rho)$ is not compatible with
        some $\mathcal{D} \in D$ **then**
        remove $\mathcal{D}$ from $D$;

    **if** $X$ or $D$ were modified **then**
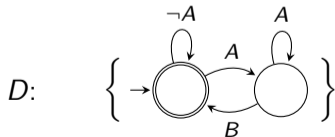        $H \leftarrow \text{infer}(X, D)$;

$H$:



$Q$:     $\{q^{p_1}\}$

$X$:     $\{(A; B; C; D/0; 0; 0; 1)\}$

$D$:     

Initialize $H$, $Q$, $X$, $D$;

**repeat**

$(\lambda, \rho, Q) \leftarrow \text{QRM}(H, Q)$;

**if** $H(\lambda) \neq \rho$ **then**
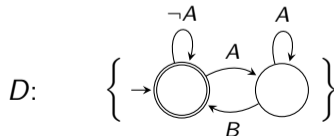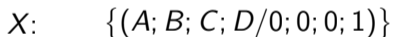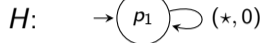add $(\lambda, \rho)$ to $X$;

**if** $(\lambda, \rho)$ is not compatible with
some $\mathcal{D} \in D$ **then**
remove $\mathcal{D}$ from $D$;

**if** $X$ or $D$ were modified **then**
$H \leftarrow \text{infer}(X, D)$;
re-initialize $Q$ if necessary;

$H$:



$Q$:  $\left\{ q^{p_1} \right\}$

$X$:  $\{(A; B; C; D/0; 0; 0; 1)\}$

$D$:

Initialize $H$, $Q$, $X$, $D$;

**repeat**

$(\lambda, \rho, Q) \leftarrow \text{QRM}(H, Q)$;

**if** $H(\lambda) \neq \rho$ **then**

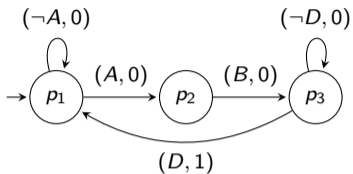add $(\lambda, \rho)$ to $X$;

**if** $(\lambda, \rho)$ is not compatible with
some $\mathcal{D} \in D$ **then**

remove $\mathcal{D}$ from $D$;

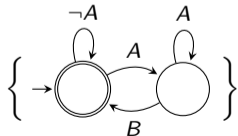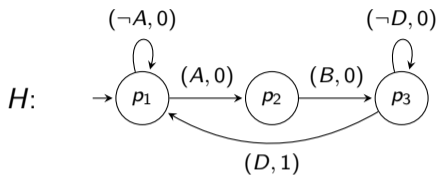**if** $X$ or $D$ were modified **then**

$H \leftarrow \text{infer}(X, D)$;

re-initialize $Q$ if necessary;

$H$:



$Q$: $\left\{ q^{p_1}, q^{p_2}, q^{p_3} \right\}$

$X$: $\{(A; B; C; D/0; 0; 0; 1)\}$

$D$:

sample $X$
advice DFAs $D$

$n \leftarrow 0$

$n \leftarrow n + 1$

create propositional formula $\Phi_n^{X,D}$

is $\Phi_n^{X,D}$ satisfiable?

no

yes

derive reward machine of size $n$ from model

## Theorem (N., Gaglione, Gavran, Topcu, Wu, Xu)

*Given*

- ▶ *a sufficient episode length*
- ▶ *an $\varepsilon$-greedy exploration strategy*

*we have the following:*

1. *AdvisoRL almost surely learns the "true" reward machine*
2. *AdvisoRL almost surely converges to an optimal policy*

## Conclusion

▶ AdvisoRL's performance improves with the "quality" of the given advice

▶ AdvisoRL is robust to incorrect advice

Office World Scenario (Icarte et al., 2018)



AdvisoRL

# Conclusion

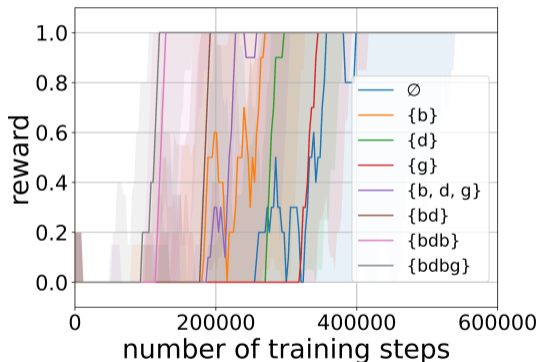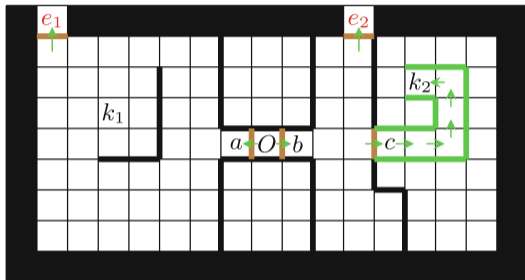# Reinforcement Learning with Temporal-Logic-Based Causal Diagrams

Yash Paliwal[1], Rajarshi Roy[2(✉)], Jean-Raphaël Gaglione[3],
Nasim Baharisangari[1], Daniel Neider[4,5], Xiaoming Duan[6], Ufuk Topcu[3],
and Zhe Xu[1]

[1] Arizona State University, Tempe, AZ, USA
xzhe1@asu.edu
[2] Max Planck Institute for Software Systems, Kaiserslautern, Germany
rajarshi008@gmail.com
[3] University of Texas at Austin, Austin, TX, USA

$$b \blacktriangleright \mathsf{G} \neg e_1$$
$$c \blacktriangleright \mathsf{X X X X X} k_2$$
$$k_2 \blacktriangleright \mathsf{G} \neg e_2$$

$$b \blacktriangleright \mathsf{G} \neg e_1$$

$$c \blacktriangleright \mathsf{X}\,\mathsf{X}\,\mathsf{X}\,\mathsf{X}\,\mathsf{X}\, k_2$$

$$k_2 \blacktriangleright \mathsf{G} \neg e_2$$



Performance Comparison

With TL-CD
Without TL-CD

### Summary

- ▶ We have been on a journey through reinforcement learning with reward machines
- ▶ There are several extension (often by other research groups)
    - ▶ partial observability, active automata learning, etc.

### Future work

- ▶ Incorporating (temporal) causal information
- ▶ Automatically synthesizing high level propositions
- ▶ More expressive classes of finite-state machines (e.g., counter)

**CENTER FOR TRUSTWORTHY DATA SCIENCE AND SECURITY**

UA RUHR | **RESEARCH ALLIANCE**

▶ **Three universities:** University of Duisburg-Essen, University of Bochum, TU Dortmund University

▶ **Four disciplines:** Computer science, IT Security, Statistics, Psychology

**We offer opportunities . . .**

▶ Collaborations with academia and industry

▶ Open positions for research group leaders, postdocs, Ph.D.s, students

▶ Internship program

▶ . . .