

# Learning Verified Monitors for Hidden Markov Models

Luko van der Maas, Sebastian Junges

October 30, 2025

Radboud Universiteit









When should the light turn on?

# Learning a monitor

- Monitor is a classifier



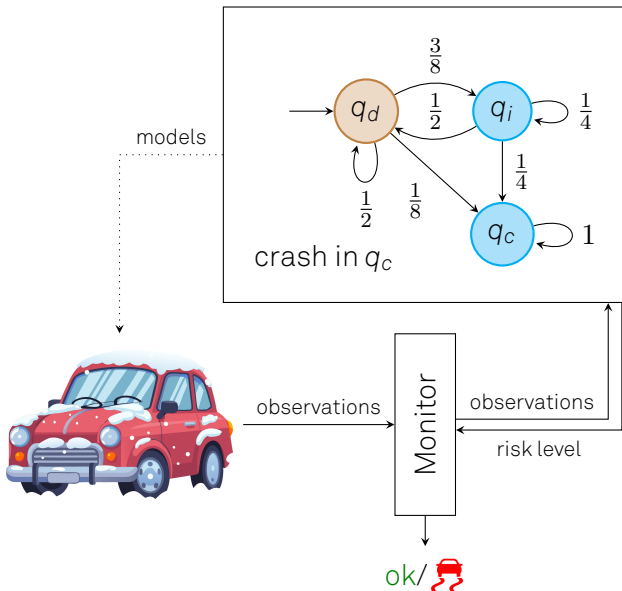
# Learning a monitor

- Monitor is a classifier
- Classic learning approaches do not **guarantee safety**



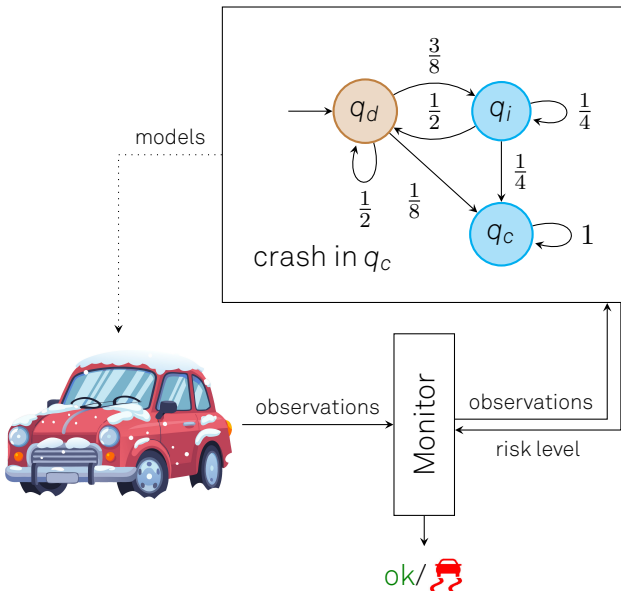
# Learning a monitor

- Monitor is a classifier
- Classic learning approaches do not **guarantee safety**
- We introduce a model



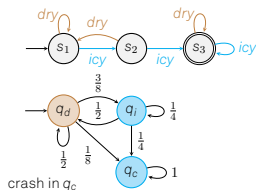
# Learning a monitor

- Monitor is a classifier
- Classic learning approaches do not **guarantee safety**
- We introduce a model
- **Verify** the monitor is safe on the model





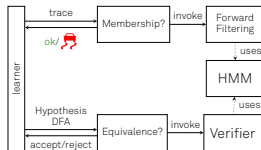
# Overview



Problem statement



Verification approach



Learning approach

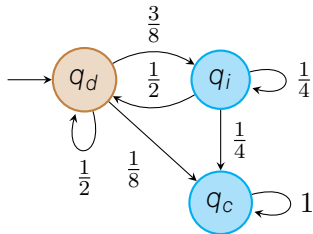


Correctness results

# What is our System Model?

## Definition (Hidden Markov Model)

- States:  $S$
- Transition function:  $\mathbf{P} : S \rightarrow \Delta S$
- Observations:  $Z$
- Observation function:  $obs : S \rightarrow Z$



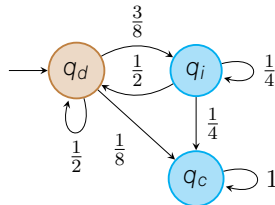
$$S = \{q_d, q_i, q_c\}$$

$$Z = \{\text{dry: } \text{brown circle}, \text{icy: } \text{blue circle}\}$$

# Monitoring Observations

## Question

Having observed the observation sequence  $\tau$ , what is the probability of being in  $q_c$ ?

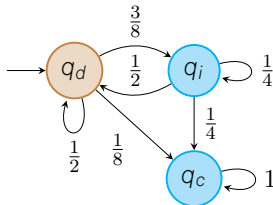


# Monitoring Observations

## Question

Having observed the observation sequence  $\tau$ , what is the probability of being in  $q_c$ ?


$$\tau_1 = \text{brown circle}$$



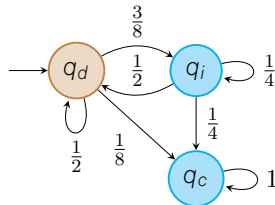
# Monitoring Observations

## Question

Having observed the observation sequence  $\tau$ , what is the probability of being in  $q_c$ ?

$\tau_1 =$  


0



# Monitoring Observations

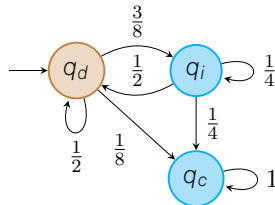
## Question

Having observed the observation sequence  $\tau$ , what is the probability of being in  $q_c$ ?

$\tau_1 =$  

0

$\tau_2 =$   



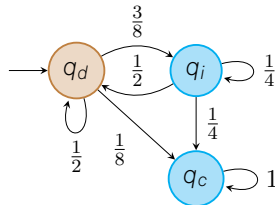
# Monitoring Observations

## Question

Having observed the observation sequence  $\tau$ , what is the probability of being in  $q_c$ ?

$$\tau_1 = \text{brown circle} \quad 0$$

$$\tau_2 = \text{brown circle} \text{ blue circle} \quad 1/4$$



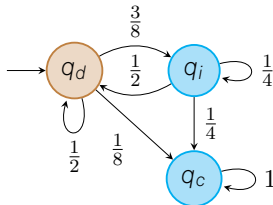
# Monitoring Observations

## Question

Having observed the observation sequence  $\tau$ , what is the probability of being in  $q_c$ ?

$$\begin{aligned}\tau_1 &= \text{brown circle} & 0 \\ \tau_2 &= \text{brown circle, blue circle} & 1/4\end{aligned}$$

$$\tau_3 = \text{brown circle, blue circle, blue circle}$$





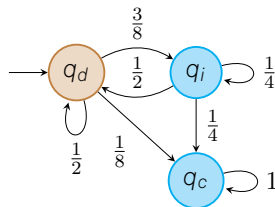
# Monitoring Observations

## Question

Having observed the observation sequence  $\tau$ , what is the probability of being in  $q_c$ ?

$$\begin{aligned}\tau_1 &= \text{brown circle} && 0 \\ \tau_2 &= \text{brown circle, blue circle} && 1/4\end{aligned}$$

$$\tau_3 = \text{brown circle, blue circle, blue circle} \quad 5/8$$



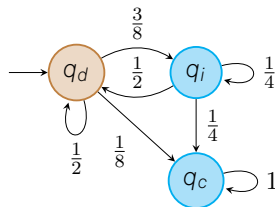
# Monitoring Observations

## Question

Having observed the observation sequence  $\tau$ , what is the probability of being in  $q_c$ ?

$$\begin{aligned}\tau_1 &= \text{brown circle} & 0 \\ \tau_2 &= \text{brown circle, blue circle} & 1/4\end{aligned}$$

$$\begin{aligned}\tau_3 &= \text{brown circle, blue circle, blue circle} & 5/8 \\ \tau_4 &= \text{blue circle, brown circle} & \end{aligned}$$



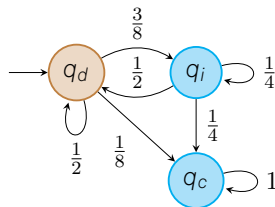
# Monitoring Observations

## Question

Having observed the observation sequence  $\tau$ , what is the probability of being in  $q_c$ ?

$$\begin{aligned}\tau_1 &= \text{brown circle} && 0 \\ \tau_2 &= \text{brown circle, blue circle} && 1/4\end{aligned}$$

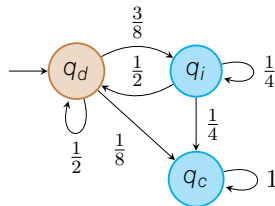
$$\begin{aligned}\tau_3 &= \text{brown circle, blue circle, blue circle} && 5/8 \\ \tau_4 &= \text{blue circle, brown circle} && ?\end{aligned}$$



# Monitoring Observations

## Question

Probability above  $\lambda = 0.3$  is **unsafe**.  
Should the warning light go on?



$$\tau_1 = \text{orange circle}$$

ok

$$\tau_2 = \text{orange circle, blue circle}$$

ok

$$\tau_3 = \text{orange circle, blue circle, blue circle}$$



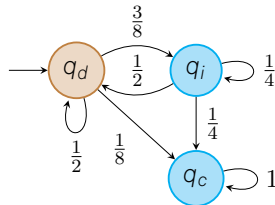
$$\tau_4 = \text{blue circle, orange circle}$$

?

# Monitoring Observations

## Question

Probability above  $\lambda = 0.3$  is **unsafe**.  
Should the warning light go on?



$$\tau_1 = \text{orange circle}$$

ok

$$\tau_2 = \text{orange circle, blue circle}$$

ok

$$\tau_3 = \text{orange circle, blue circle, blue circle}$$



$$\tau_4 = \text{blue circle, orange circle}$$

?

$$\mathbb{U}_\lambda = \{\tau_3, \dots\}$$

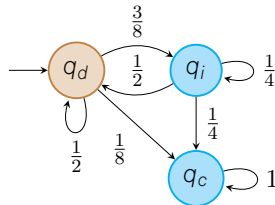
# Monitoring Observations

## Question

Probability below  $\lambda_s = 0.1$  is **safe**.

Probability above  $\lambda_u = 0.3$  is **unsafe**.

Should the warning light go on?



$$\tau_1 = \text{brown circle}$$

ok

$$\tau_2 = \text{brown circle, blue circle}$$

?

$$\tau_3 = \text{brown circle, blue circle, blue circle}$$



$$\tau_4 = \text{blue circle, brown circle}$$

?

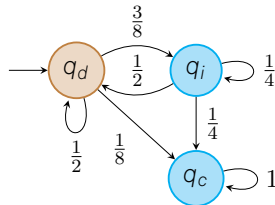
# Monitoring Observations

## Question

Probability below  $\lambda_s = 0.1$  is **safe**.

Probability above  $\lambda_u = 0.3$  is **unsafe**.

Should the warning light go on?



$$\tau_1 = \text{orange circle}$$

ok

$$\tau_2 = \text{orange circle, blue circle}$$

?

$$\tau_3 = \text{orange circle, blue circle, blue circle}$$



$$\tau_4 = \text{blue circle, orange circle}$$

?

$$\mathbb{U}_{\lambda_u} = \{\tau_3, \dots\}$$

$$\mathbb{S}_{\lambda_s} = \{\tau_1, \dots\}$$

# Monitoring Observations

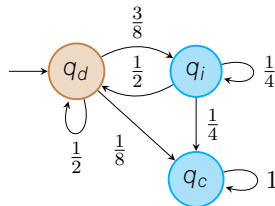
## Question

Probability below  $\lambda_s = 0.1$  is **safe**.

Probability above  $\lambda_u = 0.3$  is **unsafe**.

Horizon of 3 observations.

Should the warning light go on?



$\tau_1 =$   ok

$\tau_2 =$    ?

$\tau_3 =$     

$\tau_4 =$    ?

$$\mathbb{U}_{\lambda_u}^{\leq 3} = \{\tau_3\}$$

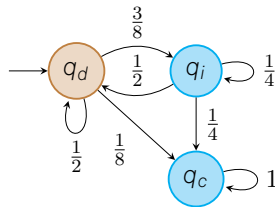
$$\mathbb{S}_{\lambda_s}^{\leq 3} = \{\tau_1, \dots\}$$



# What is a Monitor?

$$\mathbb{U}_{\lambda_u}^{\leq 3} = \{\text{brown circle}, \text{blue circle}, \text{blue circle}\}$$

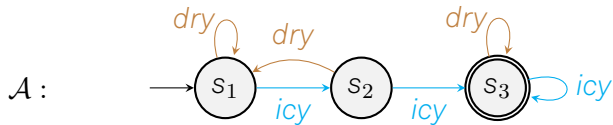
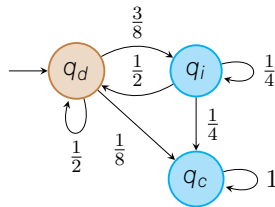
$$\mathbb{S}_{\lambda_s}^{\leq 3} = \{\text{brown circle}, \dots\}$$



# What is a Monitor?

$$\mathbb{U}_{\lambda_u}^{\leq 3} = \{\text{brown circle}, \text{blue circle}, \text{blue circle}\}$$

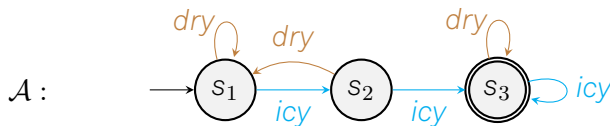
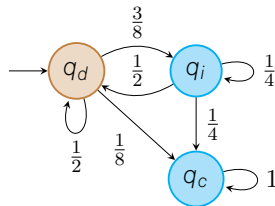
$$\mathbb{S}_{\lambda_s}^{\leq 3} = \{\text{brown circle}, \dots\}$$



# What is a Monitor?

$$\mathbb{U}_{\lambda_u}^{\leq 3} = \{\text{brown circle}, \text{blue circle}, \text{blue circle}\}$$

$$\mathbb{S}_{\lambda_s}^{\leq 3} = \{\text{brown circle}, \dots\}$$



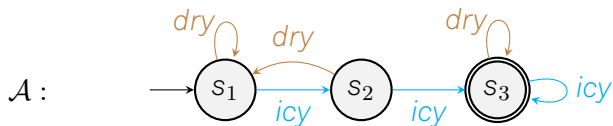
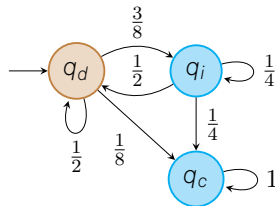
## Monitor Correctness

$$\mathbb{U}_{\lambda_u}^{\leq 3} \subseteq \mathcal{L}(\mathcal{A}) \subseteq \Sigma^* \setminus \mathbb{S}_{\lambda_s}^{\leq 3}$$

# What is a Monitor?

$$\mathbb{U}_{\lambda_u}^{\leq 3} = \{\text{brown circle}, \text{blue circle}, \text{blue circle}\}$$

$$\mathbb{S}_{\lambda_s}^{\leq 3} = \{\text{brown circle}, \dots\}$$



## Monitor Correctness

$$\mathbb{U}_{\lambda_u}^{\leq 3} \subseteq \mathcal{L}(\mathcal{A}) \subseteq \Sigma^* \setminus \mathbb{S}_{\lambda_s}^{\leq 3}$$

# Verifying Monitors (this paper)

## No Missed Alarms Problem

Given a HMM generating a set of traces  $\mathbb{U}_{\lambda_u}^{\leq h}$ , and a monitor  $\mathcal{A}$ , verify that

$$\forall \tau \in \mathbb{U}_{\lambda_u}^{\leq h}. \tau \in \mathcal{L}(\mathcal{A})$$

# Verifying Monitors (this paper)

## No Missed Alarms Problem

Given a HMM generating a set of traces  $\mathbb{U}_{\lambda_u}^{\leq h}$ , and a monitor  $\mathcal{A}$ , verify that

$$\forall \tau \in \mathbb{U}_{\lambda_u}^{\leq h}. \tau \in \mathcal{L}(\mathcal{A})$$

↓ Find a counter example

## Find Missed Alarm Problem

Given a HMM generating a set of traces  $\mathbb{U}_{\lambda_u}^{\leq h}$ , and a monitor  $\mathcal{A}$ ,

$$\exists \tau \in \mathbb{U}_{\lambda_u}^{\leq h}. \tau \notin \mathcal{L}(\mathcal{A})$$

# Verifying Monitors (this paper)

## No Missed Alarms Problem

Given a HMM generating a set of traces  $\mathbb{U}_{\lambda_u}^{\leq h}$ , and a monitor  $\mathcal{A}$ , verify that

$$\forall \tau \in \mathbb{U}_{\lambda_u}^{\leq h}. \tau \in \mathcal{L}(\mathcal{A})$$

↓ Find a counter example

## Find Missed Alarm Problem

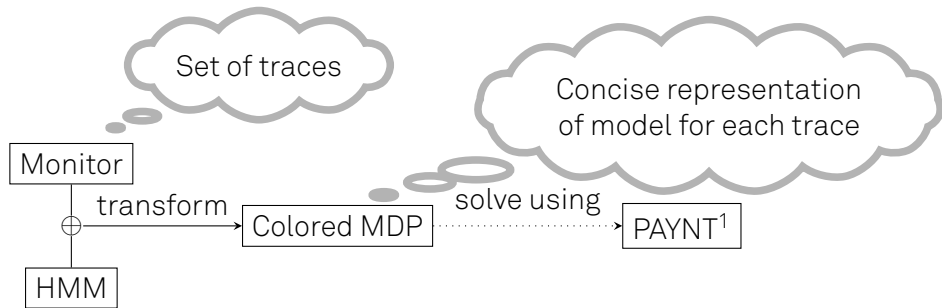
Given a HMM generating a set of traces  $\mathbb{U}_{\lambda_u}^{\leq h}$ , and a monitor  $\mathcal{A}$ ,

$$\exists \tau \in \mathbb{U}_{\lambda_u}^{\leq h}. \tau \notin \mathcal{L}(\mathcal{A})$$

## Complexity

Finding a missed alarm is NP-complete (proof in the paper).

# Searching for Missed Alarms



- Writing **conditional probability properties** using reachability, by Baier et al.<sup>2</sup>.
- Equate **traces** in the HMM to **policies** in the colored MDP, by Badings et al.<sup>3</sup>.

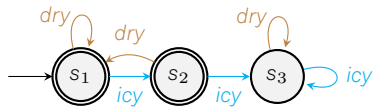
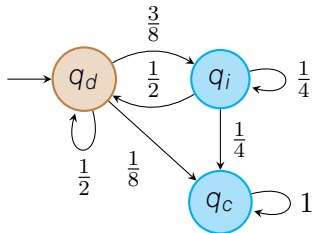
<sup>1</sup>R. Andriushchenko et al., “PAYNT: A tool for inductive synthesis of probabilistic programs,” 2021.

<sup>2</sup>C. Baier et al., “Computing conditional probabilities in markovian models efficiently,” 2014

<sup>3</sup>T. S. Badings et al., “Ctmcs with imprecisely timed observations,” 2024

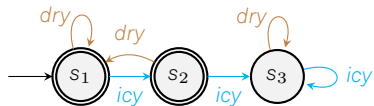
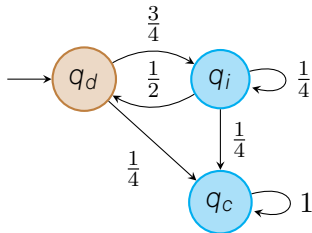


## Transformation 1/4



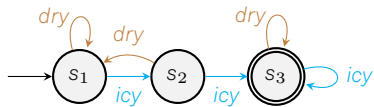
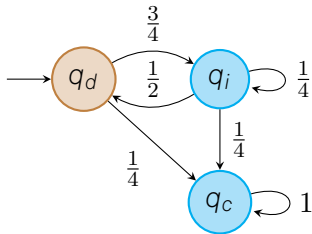
# Transformation 1/4

Find an unsafe trace which is not in the monitor



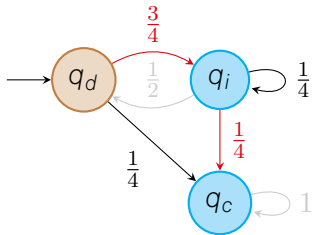
## Transformation 1/4

Find an unsafe trace which is not in the monitor




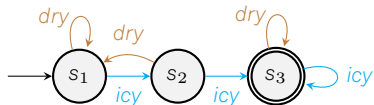
## Transformation 1/4

Find an unsafe trace which is not in the monitor



→ path:  $q_d \rightarrow q_i \rightarrow q_c$


→ trace: 

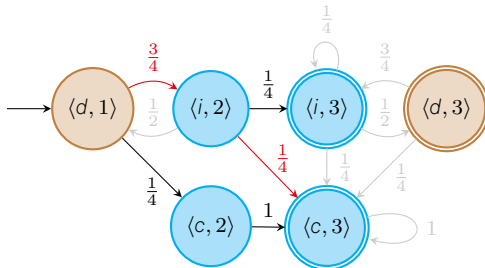
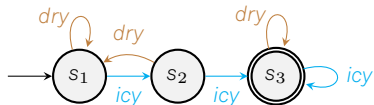
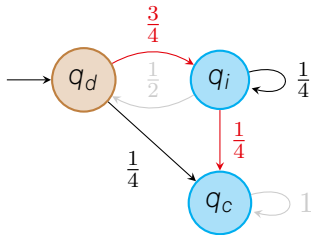


## Transformation 2/4

Find an unsafe trace which is not in the monitor


→ path:  $q_d \rightarrow q_i \rightarrow q_c$

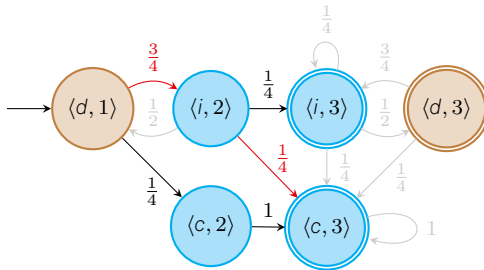
→ trace: 



## Transformation 3/4

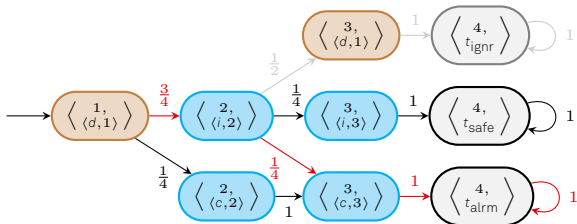
→ path:  $q_d \rightarrow q_i \rightarrow q_c$

→ trace: 



Find a trace which

- does not end in  $t_{\text{ignr}}$ ,
- has probability  $> \lambda_u$  to reach  $\langle 4, t_{\text{alarm}} \rangle$ .



# Transformation 4/4

Find a trace which


- does not end in  $t_{\text{ignr}}$ ,
- has probability  $> \lambda_u$  to reach  $\langle 4, t_{\text{alarm}} \rangle$ .

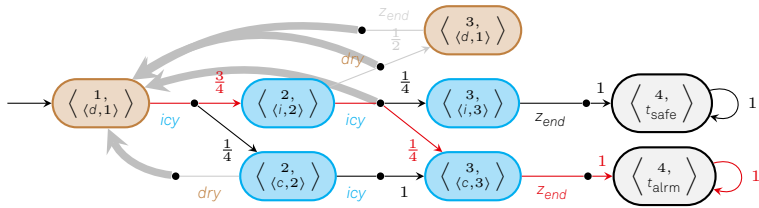
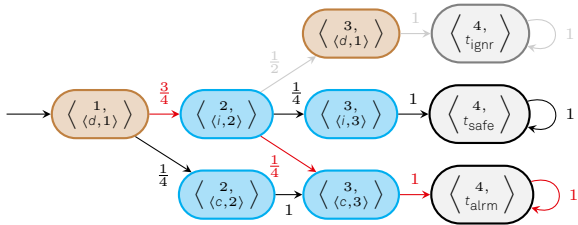
Find a policy

$\sigma: \mathbb{N}^{\leq 4} \rightarrow Z$  s.t.

- we reach an end state,
- reach  $t_{\text{alarm}}$  with prob.  $\geq \lambda_u$ .


→ path:  $q_d \rightarrow q_i \rightarrow q_c$

→ trace: 



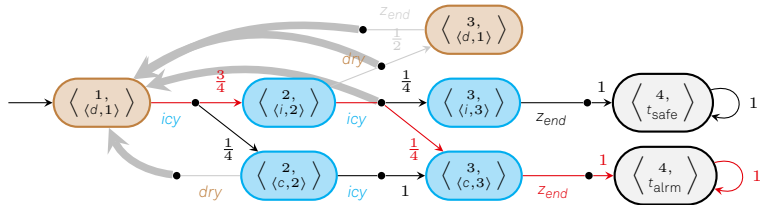
## Transformation 4/4

→ path:  $q_d \rightarrow q_i \rightarrow q_c$

→ trace: 

Find a policy  
 $\sigma: \mathbb{N}^{\leq 4} \rightarrow Z$  s.t.

- we reach an end state,
- reach  $t_{\text{alarm}}$  with prob.  $\geq \lambda_u$ .




Solvable by PAYNT



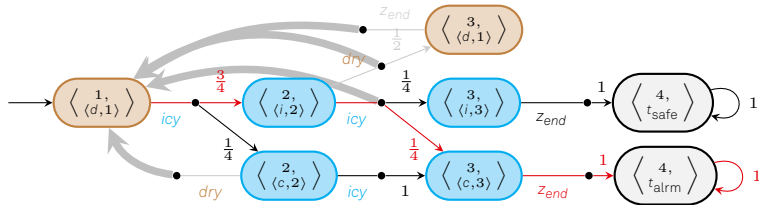
## Transformation 4/4

→ path:  $q_d \rightarrow q_i \rightarrow q_c$

→ trace: 

Find a policy  
 $\sigma: \mathbb{N}^{\leq 4} \rightarrow Z$  s.t.

- we reach an end state,
- reach  $t_{\text{alarm}}$  with prob.  $\geq \lambda_U$ .



Solvable by PAYNT

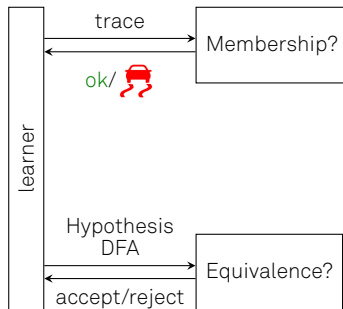
### Theorem 1

“The transformation is correct”

# **Learning** Verified Monitors for Hidden Markov Models

# Learning a monitor

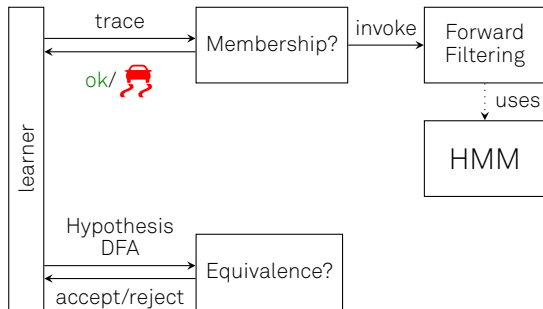
- Active automata learning using  $L^*$ .



# Learning a monitor

- Active automata learning using  $L^*$ .
- MQ: **Forward Filtering** implemented by Premise<sup>4</sup> on the HMM with threshold  $\lambda_l$ .

$$\lambda_s \leq \lambda_l \leq \lambda_u$$



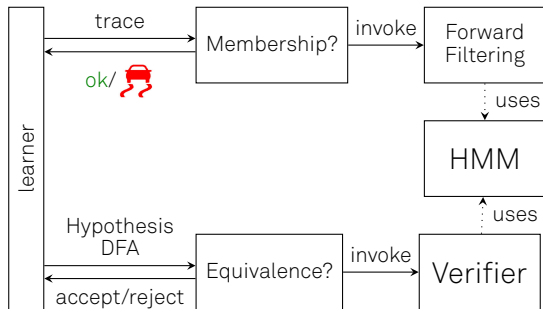
<sup>4</sup>S. Junges et al., “Runtime monitors for markov decision processes,” 2021

# Learning a monitor

- Active automata learning using  $L^*$ .
- MQ: **Forward Filtering** implemented by Premise<sup>4</sup> on the HMM with threshold  $\lambda_l$ .

$$\lambda_s \leq \lambda_l \leq \lambda_u$$

- EQ: is a candidate monitor correct.



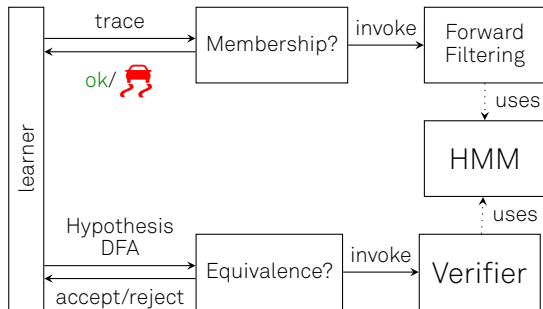
<sup>4</sup>S. Junges et al., “Runtime monitors for markov decision processes,” 2021

# Learning a monitor

- Active automata learning using  $L^*$ .
- MQ: **Forward Filtering** implemented by Premise<sup>4</sup> on the HMM with threshold  $\lambda_l$ .

$$\lambda_s \leq \lambda_l \leq \lambda_u$$

- EQ: is a candidate monitor correct.



## Theorem 2

“Monitors learned using our verification algorithm are correct.”

<sup>4</sup>S. Junges et al., “Runtime monitors for markov decision processes,” 2021

# Correctness Experiments

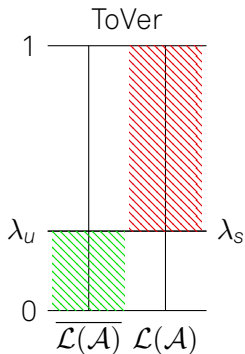
## Benchmark A-63/64

290 states, 1258 transitions,  
50 observations, horizon of 10  
 $\lambda_s = \lambda_l = \lambda_u = 0.3$

# Correctness Experiments

## Benchmark A-63/64

290 states, 1258 transitions,  
50 observations, horizon of 10  
 $\lambda_s = \lambda_l = \lambda_u = 0.3$



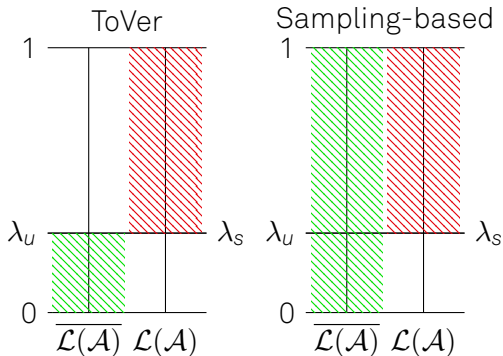


# Correctness Experiments

## Benchmark A-63/64

290 states, 1258 transitions,  
50 observations, horizon of 10

$$\lambda_s = \lambda_l = \lambda_u = 0.3$$

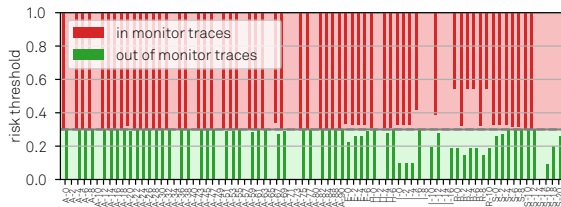
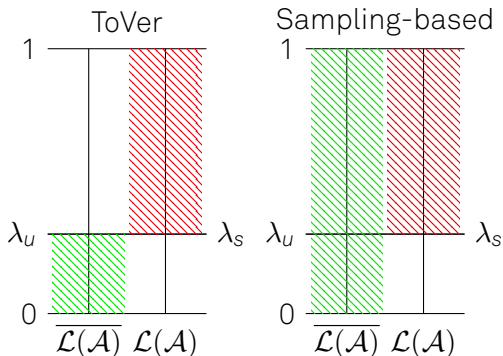


# Correctness Experiments

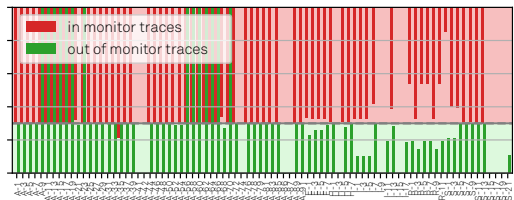
## Benchmark A-63/64

290 states, 1258 transitions,  
50 observations, horizon of 10

$$\lambda_s = 0.1 \quad \lambda_l = 0.3 \quad \lambda_u = 0.35$$



(a) Learning with verification,  $\lambda_s = \lambda_l = \lambda_u$

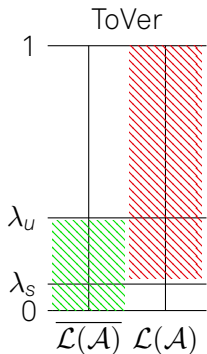


(b) Learning with sampling-based verification

# Correctness Experiments

## Benchmark A-63/64

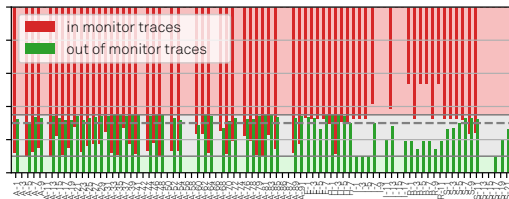
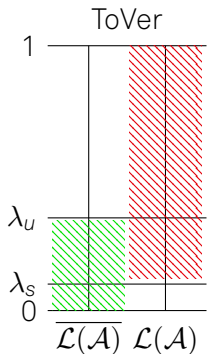
290 states, 1258 transitions,  
50 observations, horizon of 10  
 $\lambda_s = 0.1$     $\lambda_l = 0.3$     $\lambda_u = 0.35$



# Correctness Experiments

## Benchmark A-63/64

290 states, 1258 transitions,  
50 observations, horizon of 10  
 $\lambda_s = 0.1$     $\lambda_l = 0.3$     $\lambda_u = 0.35$



(c) Learning with verification,  $\lambda_s < \lambda_l < \lambda_u$

# Conclusion



## Summary

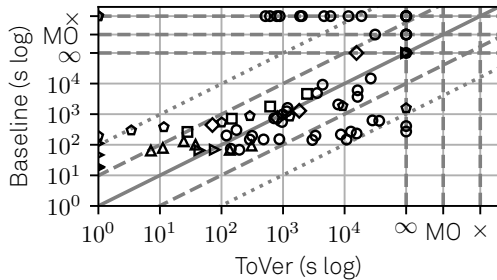
- We present a **verification algorithm** for HMM monitors.
- We prove the verification problem is **coNP-complete**.
- We integrate it with active automata learning to **learn correct monitors**.
- We learn monitors with up to **1500 states** in **11 hours** on models with **100s of states**.

## Future interests

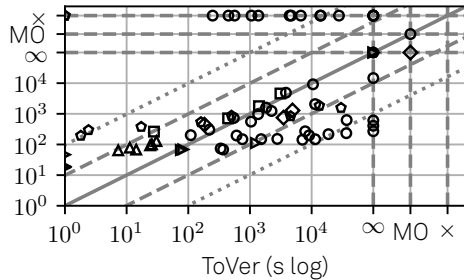
- Ideas to adapt AAL more to our specific problem.
- Adapt colored MDP model checking more to our specific problem of conditional probabilities.
- Learn models from data such that they are useful for monitoring.

Email: `luko.vandermaas@ru.nl`

# Results: Runtime

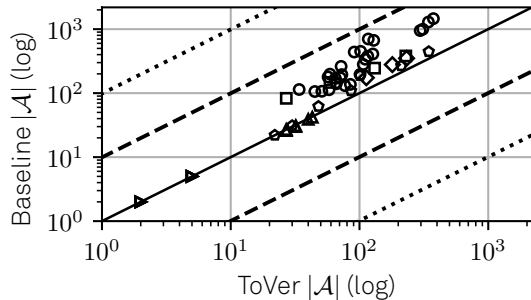


(a)  $\lambda_s < \lambda_l < \lambda_u$ , Runtime

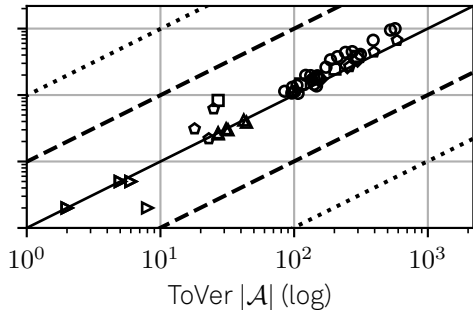


(b)  $\lambda_s = \lambda_l = \lambda_u$ , Runtime

## Results: Monitor Size



(a)  $\lambda_s < \lambda_l < \lambda_u$ , Size of  $\mathcal{A}$



(b)  $\lambda_s = \lambda_l = \lambda_u$ , Size of  $\mathcal{A}$