

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

ELEMENT^{AI}

KR 2020
September 16

Hi, I'm Rodrigo :)

Hi, I'm an AI researcher

Hi, I'm an AI researcher

“The ultimate goal of AI is to create computer programs that can solve problems in the world as well as humans.”

— John McCarthy

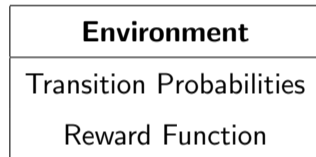
Hi, I'm an AI researcher

“The ultimate goal of AI is to create computer programs that can solve problems in the world as well as humans.”

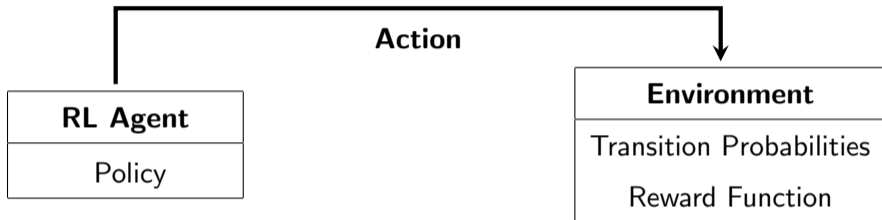
— John McCarthy

Our research incorporates insights from **knowledge**, **reasoning**, and **learning**,
in service of building general-purpose agents.

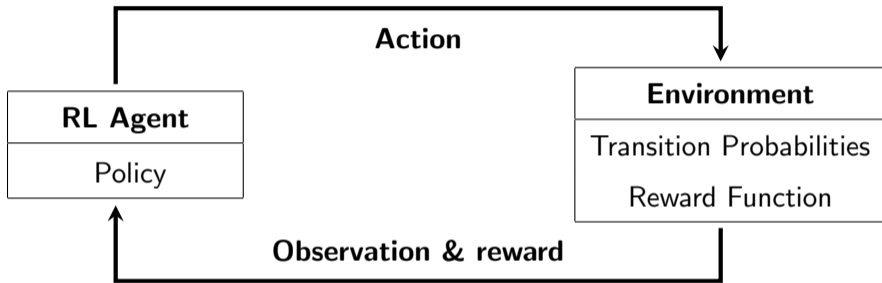
Reinforcement Learning (RL)



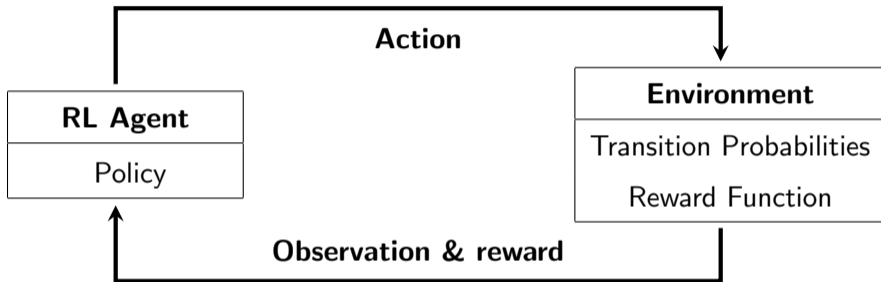
Reinforcement Learning (RL)



Reinforcement Learning (RL)

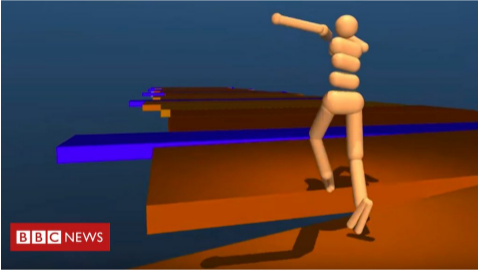


Reinforcement Learning (RL)

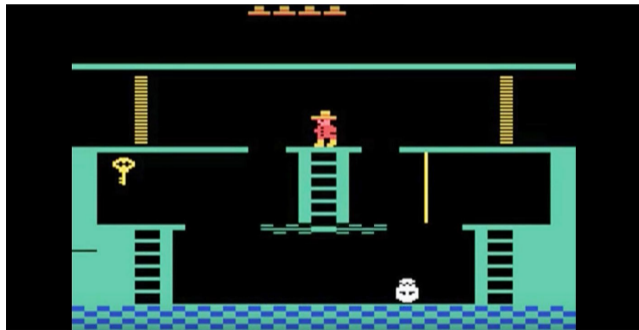


This learning process captures some aspects of human intelligence.

Reinforcement Learning (RL)



Reinforcement Learning (RL)



How to enhance RL with KR

Reinforcement Learning (RL)

Long-standing RL problems that we tackled using KR:

- Reward specification.
- Sample efficiency.
- Memory.
- ...

Reward specification

Make a bridge: get wood, iron, and use the factory

Make a bridge: get wood, iron, and use the factory

LTL specifications¹:

$\exists(\text{got_wood} \wedge \exists \text{used_factory}) \wedge \exists(\text{got_iron} \wedge \exists \text{used_factory})$

¹ Teaching Multiple Tasks to an RL Agent using LTL (AAMAS-18).

Make a bridge: get wood, iron, and use the factory

LTL specifications¹:

$\exists(\text{got_wood} \wedge \exists \text{used_factory}) \wedge \exists(\text{got_iron} \wedge \exists \text{used_factory})$

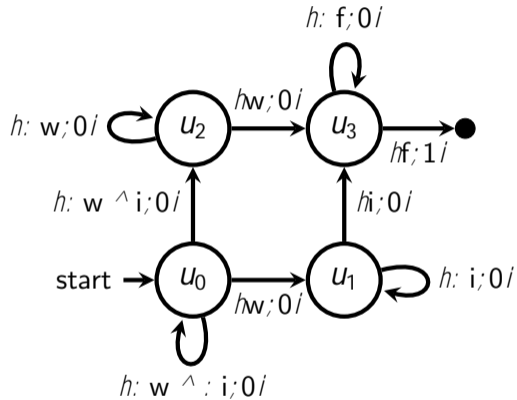
Reward machines²:

Automata-based reward functions

¹ Teaching Multiple Tasks to an RL Agent using LTL (AAMAS-18).

² Using Reward Machines for High-Level Task Specification and Decomposition in RL (ICML-18).

Reward machine



Make a bridge: get wood, iron, and use the factory

Make a bridge: get wood, iron, and use the factory

LTL specifications¹:

$\exists(\text{got_wood} \wedge \exists \text{used_factory}) \wedge \exists(\text{got_iron} \wedge \exists \text{used_factory})$

Reward machines²:

Automata-based reward functions

Formal languages³:

Many formal languages / Reward machines.

¹ Teaching Multiple Tasks to an RL Agent using LTL (AAMAS-18).

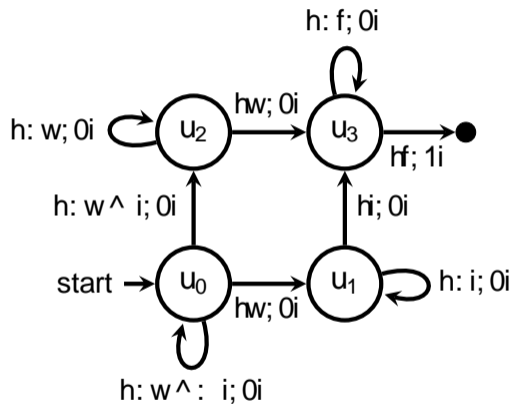
² Using Reward Machines for High-Level Task Specification and Decomposition in RL (ICML-18).

³ LTL and Beyond: Formal Languages for Reward Function Specification in RL (IJCAI-19).

Sample efficiency

Sample efficiency

Reward machine



How to exploit the reward machine's structure:

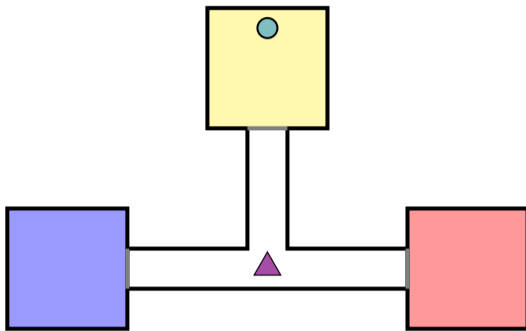
- CRM: Counterfactual reasoning.
- HRM: Task decomposition.
- RS: Reward shaping.

Sample efficiency

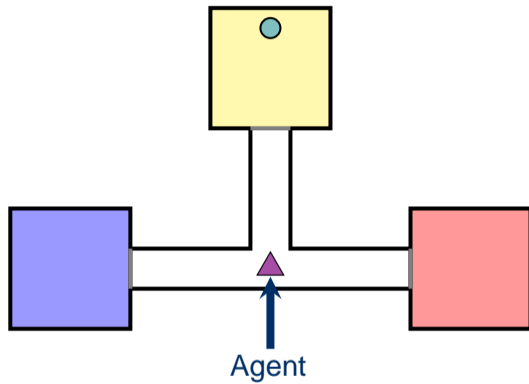
Sample efficiency

Memory

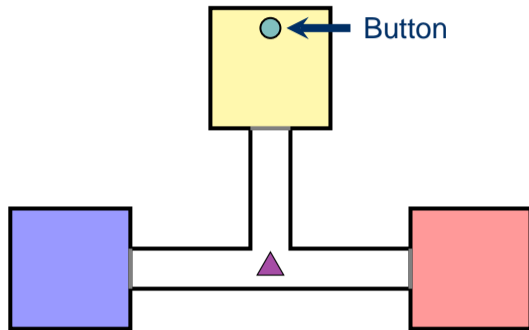
Memory



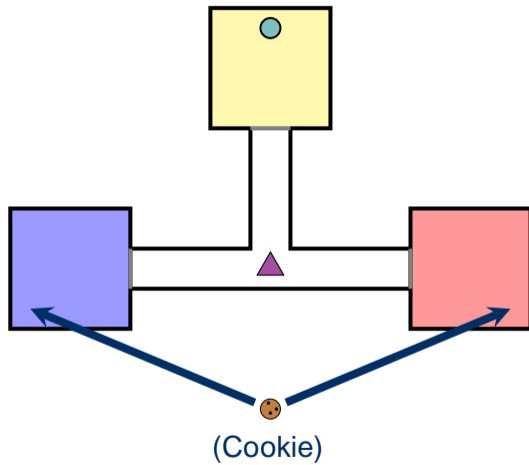
Memory



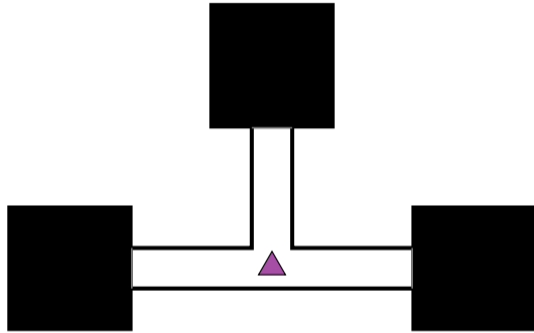
Memory



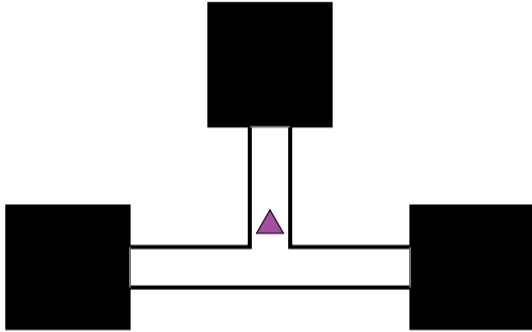
Memory



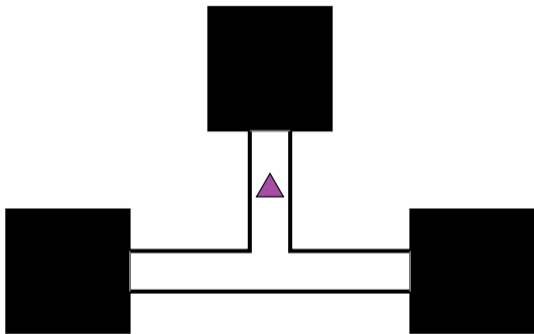
Memory



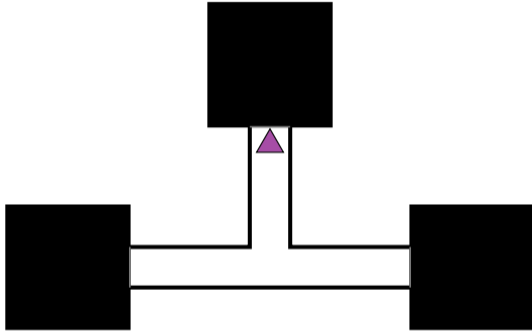
Memory



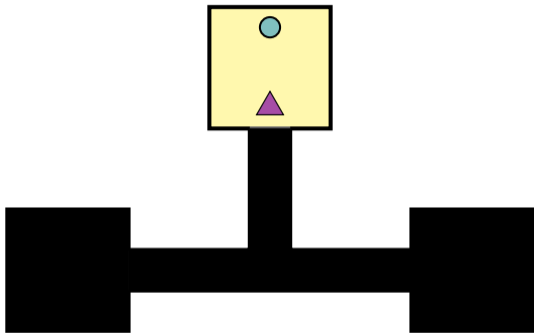
Memory



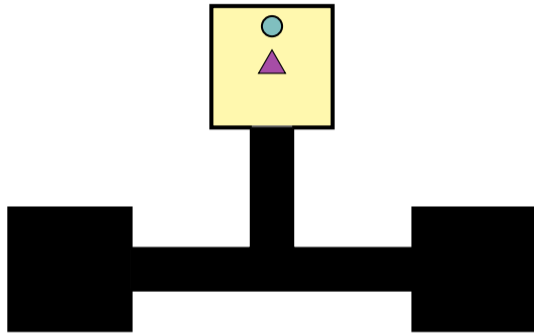
Memory



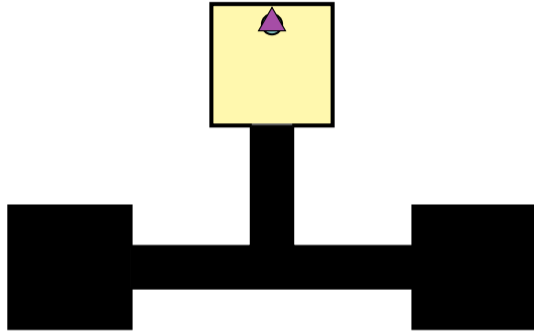
Memory



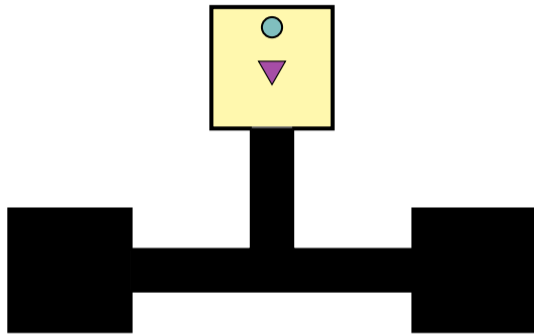
Memory



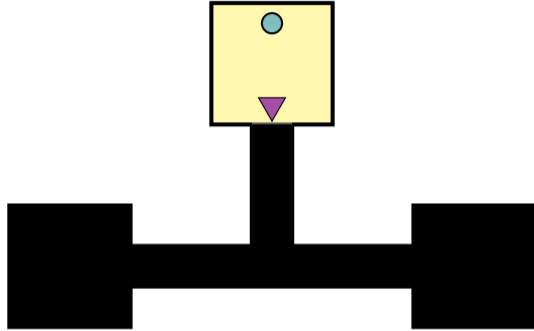
Memory



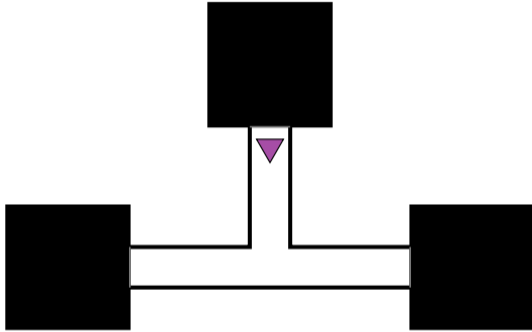
Memory



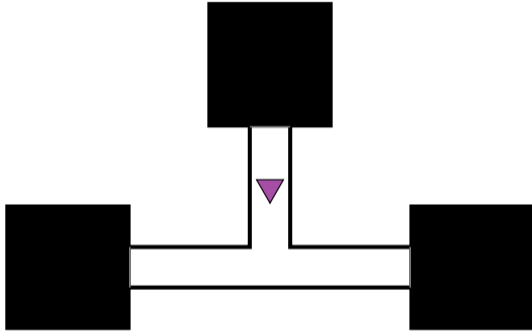
Memory



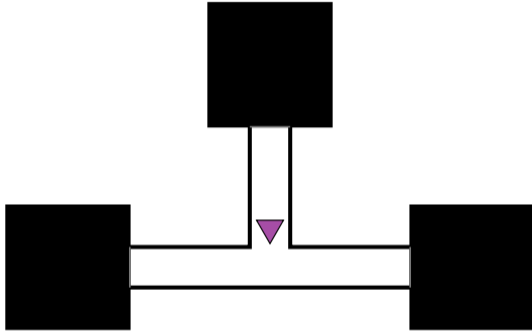
Memory



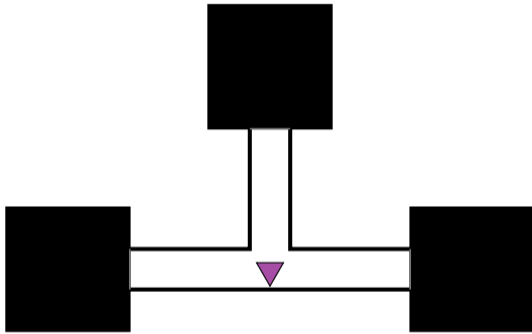
Memory



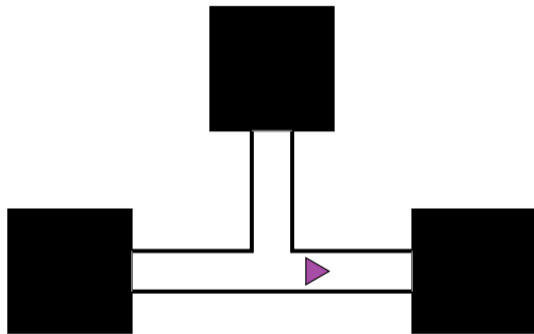
Memory



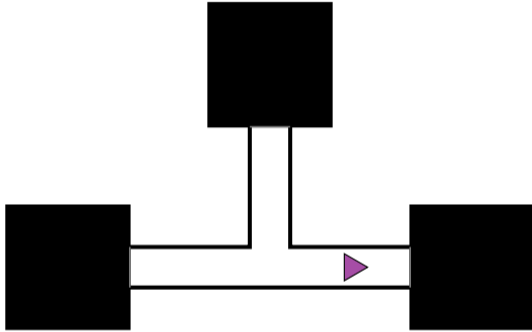
Memory



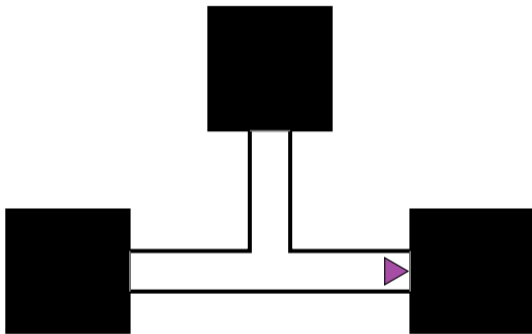
Memory



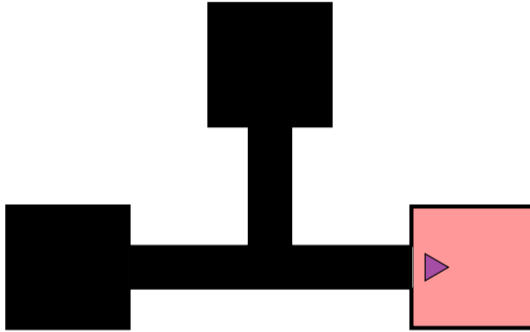
Memory



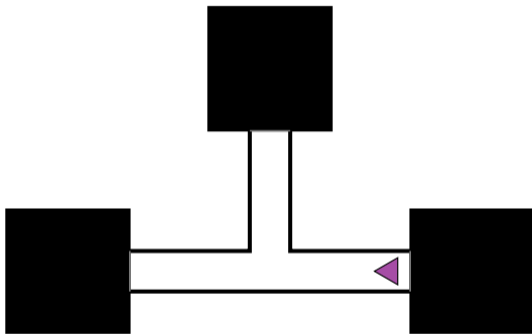
Memory



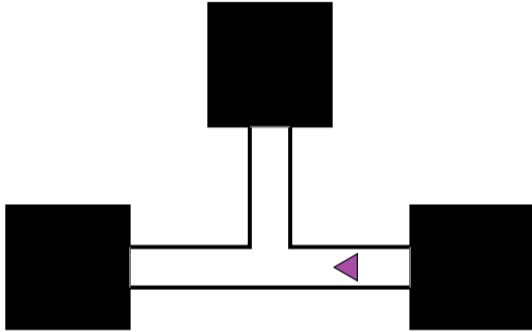
Memory



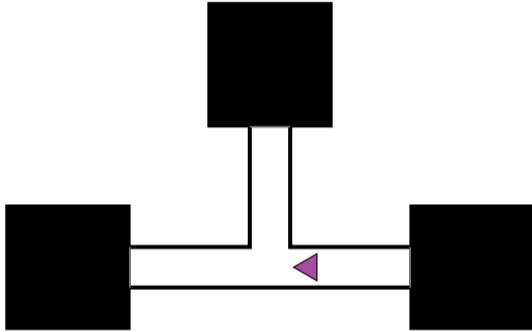
Memory



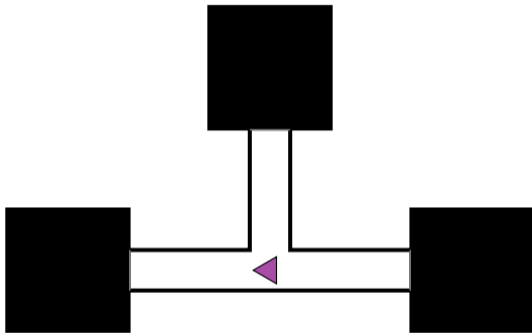
Memory



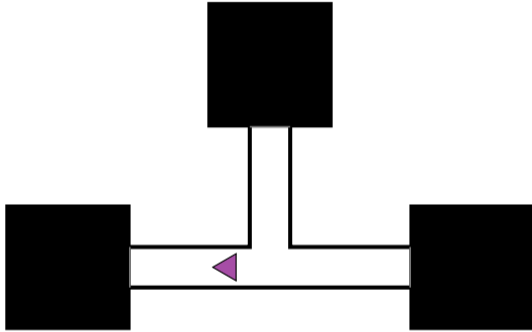
Memory



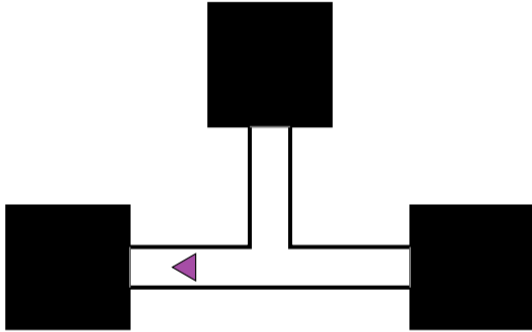
Memory



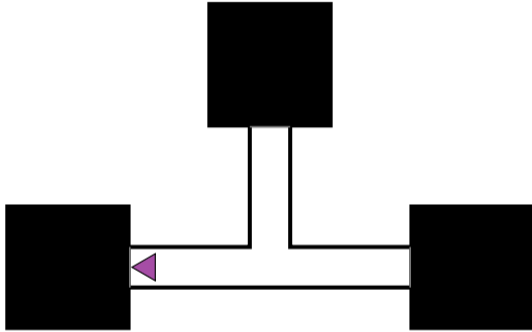
Memory



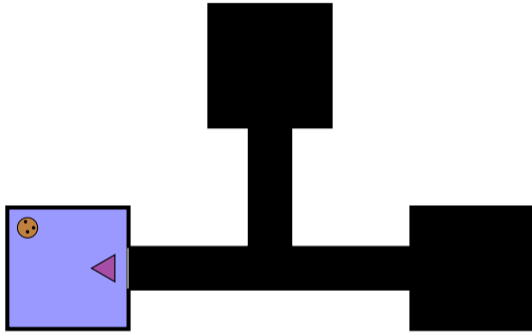
Memory



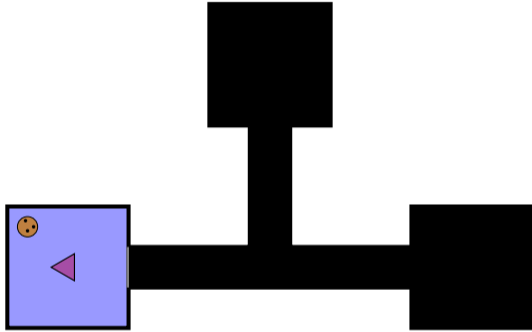
Memory



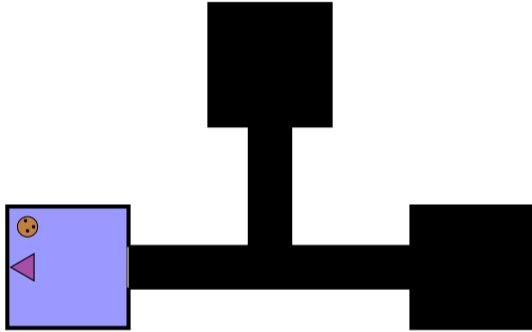
Memory



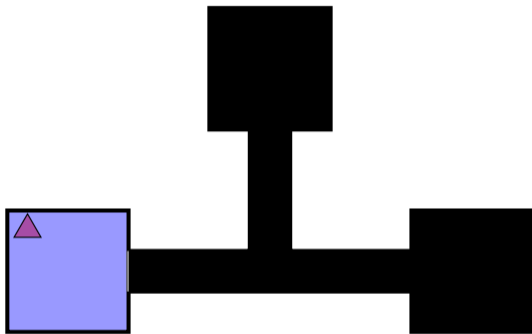
Memory



Memory

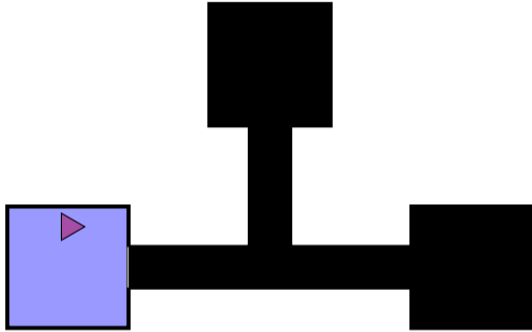


Memory

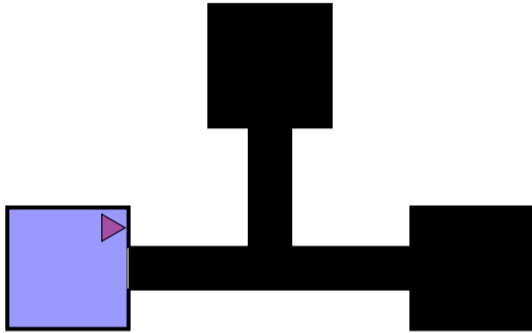


(+1 Reward)

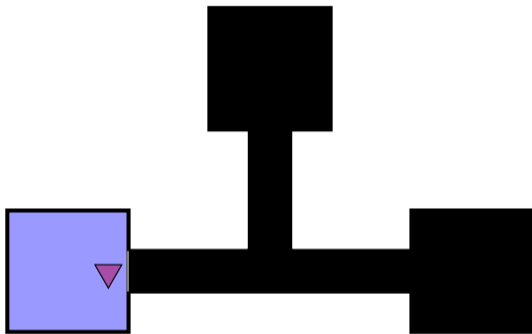
Memory



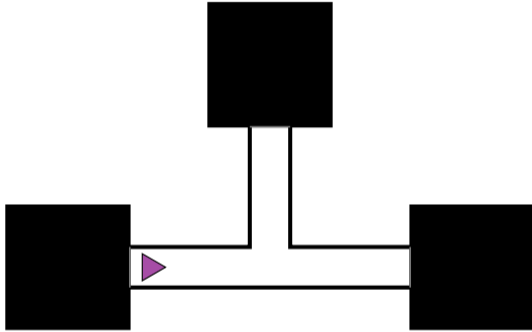
Memory



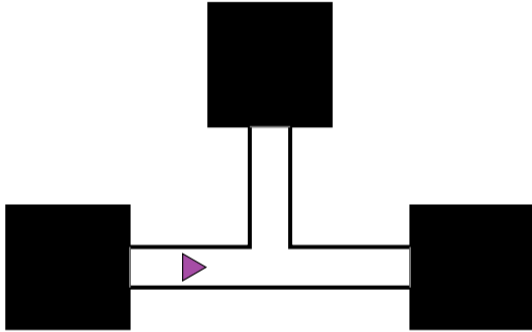
Memory



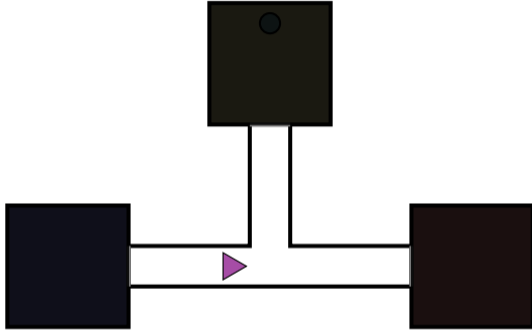
Memory



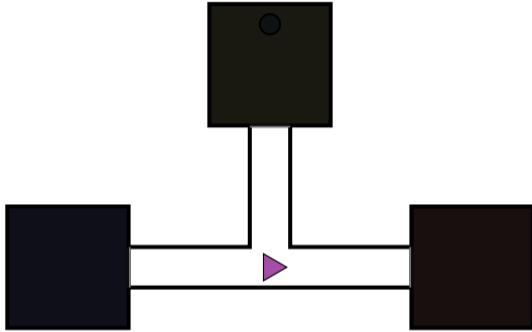
Memory



Memory



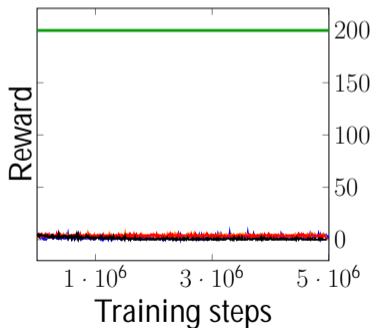
Memory



The most popular approach:

Training LSTMs policies using a policy gradient method.

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

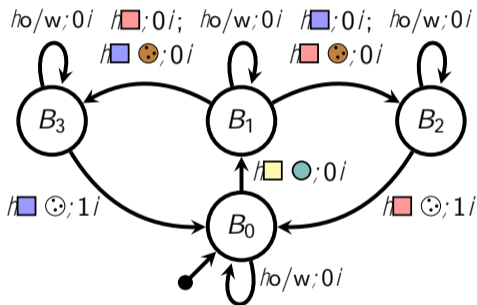


Legend:

- Optimal
- ACER
- A3C
- PPO
- DDQN

Reward Machines as memory

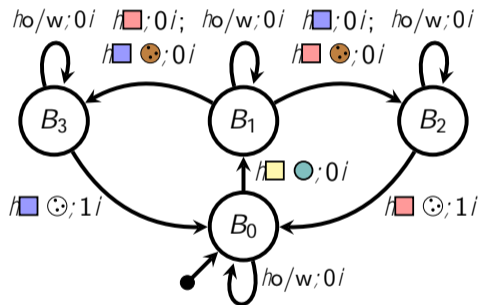
If the agent can detect the color of the rooms (■;□;■;■), and when it presses the button (●), eats a cookie (☺), and sees a cookie (⦿), then:



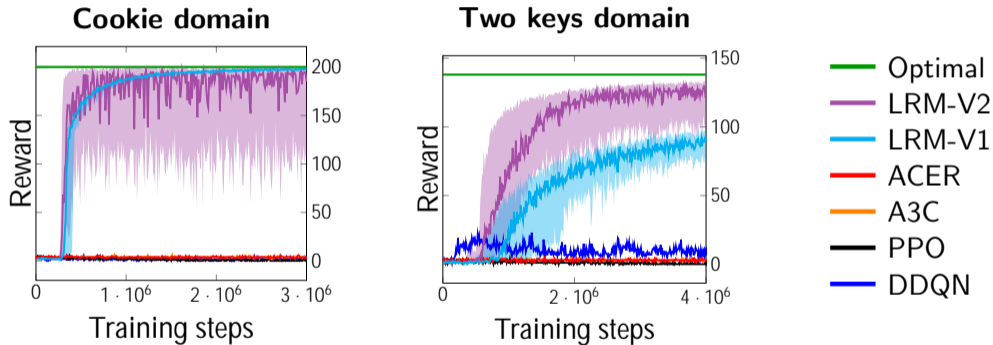
... becomes a **“perfect” memory** for the cookie domain.

Reward Machines as memory

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.



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

Summary

Summary

If you are interested in $KR \setminus RL$, consider reading our papers:

Summary

If you are interested in KR \ RL, consider reading our papers:

Advice-Based Exploration in Model-Based Reinforcement Learning (Canadian AI-18)

Teaching Multiple Tasks to an RL Agent using LTL (AAMAS-18)

Using Reward Machines for High-Level Task Specification and Decomposition in RL (ICML-18)

LTL and Beyond: Formal Languages for Reward Function Specification in RL (IJCAI-19)

Learning Reward Machines for Partially Observable RL (NeurIPS-19)

Symbolic Plans as High-Level Instructions for Reinforcement Learning (ICAPS-20)

Summary

If you are interested in KR \ RL, consider reading our papers:

Advice-Based Exploration in Model-Based Reinforcement Learning (Canadian AI-18)

Teaching Multiple Tasks to an RL Agent using LTL (AAMAS-18)

Using Reward Machines for High-Level Task Specification and Decomposition in RL (ICML-18)

LTL and Beyond: Formal Languages for Reward Function Specification in RL (IJCAI-19)

Learning Reward Machines for Partially Observable RL (NeurIPS-19)

Symbolic Plans as High-Level Instructions for Reinforcement Learning (ICAPS-20)

Code: <https://bitbucket.org/RTorolcarte/>

Summary

If you are interested in KR \ RL, consider reading our papers:

Advice-Based Exploration in Model-Based Reinforcement Learning (Canadian AI-18)

Teaching Multiple Tasks to an RL Agent using LTL (AAMAS-18)

Using Reward Machines for High-Level Task Specification and Decomposition in RL (ICML-18)

LTL and Beyond: Formal Languages for Reward Function Specification in RL (IJCAI-19)

Learning Reward Machines for Partially Observable RL (NeurIPS-19)

Symbolic Plans as High-Level Instructions for Reinforcement Learning (ICAPS-20)

Code: <https://bitbucket.org/RTorolcarte/>

Thanks! :)