

REWARD MACHINES FOR COOPERATIVE MULTI-AGENT REINFORCEMENT LEARNING

Cyrus Neary, Zhe Xu, Bo Wu, and Ufuk Topcu

The University of Texas at Austin



The 20th International Conference on Autonomous Agents and Multiagent Systems 3-7 May 2021, Online

Cooperative multi-agent RL: A team of agents learn interact in a shared environment to achieve a common objective.

Cooperative multi-agent RL: A team of agents learn interact in a shared environment to achieve a common objective.



Source: Blizzard Entertainment



Source: OpenAl



Source: Starship Technologies

Cooperative multi-agent RL: A team of agents learn interact in a shared environment to achieve a common objective.



Source: Blizzard Entertainment

Source: OpenAl



Source: Starship Technologies

Cooperative multi-agent RL: A team of agents learn interact in a shared environment to achieve a common objective.



Cooperative multi-agent RL: A team of agents learn interact in a shared environment to achieve a common objective.



Cooperative multi-agent RL: A team of agents learn interact in a shared environment to achieve a common objective.

| Objectives: | |
|---|--|
| 1) Present a framework for specifying structured | |
| representations of team tasks. | |
| 2) Use this specification to decompose problem into necessary | |
| individual behaviors. | |
| 3) Present a reinforcement learning algorithm that uses | |
| decomposition to simplify multi-agent learning. | |



Team objective: Have A_1 safely reach the goal location *Goal*.



Team objective: Have A_1 safely reach the goal location Goal.

- Colored walls block the agents' paths.



Team objective: Have A_1 safely reach the goal location Goal.

- Colored walls block the agents' paths.
- Buttons remove the wall of the corresponding color.







Team objective:

Have A_1 safely reach the goal location *Goal*.

Colored walls block the agents' paths.

Buttons remove the wall of the corresponding color.



Team objective:

Have A_1 safely reach the goal location *Goal*.

Colored walls block the agents' paths.

Buttons remove the wall of the corresponding color.



Team objective:

Have A_1 safely reach the goal location *Goal*.

Colored walls block the agents' paths.

Buttons remove the wall of the corresponding color.



Team objective:

Have A_1 safely reach the goal location *Goal*.

Colored walls block the agents' paths.

Buttons remove the wall of the corresponding color.





 $U = \{u_I, u_1, u_2, ..., u_7\}$ – Set of states





$$U = \{u_I, u_1, u_2, \dots, u_7\} - \text{Set of states}$$

$$\Sigma = \{Y_B, G_B, R_B, A_2^{R_B}, A_2^{\neg R_B}, A_3^{R_B}, A_3^{\neg R_B}, Goal\} - \text{Set of events}$$





$$\begin{aligned} U &= \{u_I, u_1, u_2, \dots, u_7\} - \text{Set of states} \\ \Sigma &= \{Y_B, G_B, R_B, A_2^{R_B}, A_2^{\neg R_B}, A_3^{R_B}, A_3^{\neg R_B}, Goal\} - \text{Set of events} \\ \delta &: U \times \Sigma \to U - \text{Transition function} \end{aligned}$$





 $U = \{u_{I}, u_{1}, u_{2}, ..., u_{7}\} - \text{Set of states}$ $\Sigma = \{Y_{B}, G_{B}, R_{B}, A_{2}^{R_{B}}, A_{2}^{\neg R_{B}}, A_{3}^{R_{B}}, A_{3}^{\neg R_{B}}, Goal\} - \text{Set of events}$ $\delta : U \times \Sigma \rightarrow U - \text{Transition function}$ $\sigma : U \times U \rightarrow \mathbb{R} - \text{Output function}$ F - Set of final states





$$\Sigma = \{Y_B, G_B, R_B, A_2^{R_B}, A_2^{\neg R_B}, A_3^{\neg R_B}, A_3^{\neg R_B}, Goal\}$$





$$\Sigma = \{ Y_B, G_B, R_B, A_2^{R_B}, A_2^{\neg R_B}, A_3^{\neg R_B}, A_3^{\neg R_B}, Goal \}$$





$$\Sigma = \{Y_B, G_B, R_B, A_2^{R_B}, A_2^{\neg R_B}, A_3^{\neg R_B}, A_3^{\neg R_B}, Goal\}$$





$$\Sigma = \{Y_B, G_B, R_B, A_2^{R_B}, A_2^{\neg R_B}, A_3^{\neg R_B}, A_3^{\neg R_B}, Goal\}$$





$$\Sigma = \{Y_B, G_B, R_B, A_2^{R_B}, A_2^{\neg R_B}, A_3^{\neg R_B}, A_3^{\neg R_B}, \textbf{Goal}\}$$



Connecting Environment States to a Reward Machine

The labeling function: Relate the environment state to collections of high-level events.



High-Level Task (Model: Reward Machine)

- Reward machine states U
- High-level Events Σ

Connecting Environment States to a Reward Machine

The labeling function: Relate the environment state to collections of high-level events.



High-Level Task (Model: Reward Machine)

- Reward machine states U
- High-level Events Σ



Environment Dynamics Model: Transition distribution $p(\cdot | s, a)$

- Environment states $S = S_1 \times S_2 \times S_3$
- Environment actions $\mathbf{A} = A_1 \times A_2 \times A_3$

Connecting Environment States to a Reward Machine

The labeling function: Relate the environment state to collections of high-level events.





















How does the reward machine change if one only has access to events from $\Sigma_i \subseteq \Sigma$?



Example: Event subset for A_1 is $\Sigma_1 = \{Y_B, R_B, Goal\}$

How does the reward machine change if one only has access to events from $\Sigma_i \subseteq \Sigma$?



Example: Event subset for A_1 is $\Sigma_1 = \{Y_B, R_B, Goal\}$

• Transitions triggered by missing events are not observed.

How does the reward machine change if one only has access to events from $\Sigma_i \subseteq \Sigma$?



Example: Event subset for A_1 is $\Sigma_1 = \{Y_B, R_B, Goal\}$

- Transitions triggered by missing events are not observed.
- Merge states that cannot be differentiated by events from Σ_1 .

How does the reward machine change if one only has access to events from $\Sigma_i \subseteq \Sigma$?



Projected Reward Machine R_1 start u_1^1 u_2^1

 R_B

Goa

- Transitions triggered by missing events are not observed.
- Merge states that cannot be differentiated by events from Σ_1 .

How does the reward machine change if one only has access to events from $\Sigma_i \subseteq \Sigma$?



How does the reward machine change if one only has access to events from $\Sigma_i \subseteq \Sigma$?

Projected reward machines encode the sub-tasks of individual agents who only observe events in Σ_i . $\rightarrow (u_1^3) \xrightarrow{*} (u_2^3)$

Problem Equivalence

When is the task described by the team reward machine equivalent to the composition of its projections?



Problem Equivalence

When is the task described by the team reward machine equivalent to the composition of its projections?



Local Labeling Functions

Connecting environment dynamics with projected reward machines?



Local Labeling Functions

Connecting environment dynamics with projected reward machines?



Local Labeling Functions

Connecting environment dynamics with projected reward machines?



Problem Equivalence



Problem Equivalence









While ensuring the resulting composite behavior accomplishes the team task.



While ensuring the resulting composite behavior accomplishes the team task.





While ensuring the resulting composite behavior accomplishes the team task.





While ensuring the resulting composite behavior accomplishes the team task.





Rendezvous Experiment

| A_1 | | A_2 | | | | A_4 | G_5 |
|----------|--|-------|---|-------|-------|-------|----------|
| | | | | | | | |
| A_3 | | | | | | | G_3 |
| | | | R | | | | |
| A_6 | | | | | | | A_8 |
| G_8 | | | | | | | A_{10} |
| | | | | | | | G_9 |
| A_7 | | | | | | | G_2 |
| G_{10} | | | | | | | |
| A_5 | | | | A_9 | G_1 | | G_4 |

Team Task: Each agent must *simultaneously* occupy the rendezvous location \mathbb{R} , before proceeding to their respective goals \mathbb{G}_i .

Rendezvous Experiment



Team Task: Each agent must *simultaneously* occupy the rendezvous location \mathbb{R} , before proceeding to their respective goals \mathbb{G}_i .

- Two agent rendezvous.
- Ten agent rendezvous.

Rendezvous Experiment



Team Task: Each agent must *simultaneously* occupy the rendezvous location \mathbb{R} , before proceeding to their respective goals G_i .

- Two agent rendezvous.
- Ten agent rendezvous.

Rendezvous Experiment



Team Task: Each agent must *simultaneously* occupy the rendezvous location \mathbb{R} , before proceeding to their respective goals \mathbb{G}_i .

• Two agent rendezvous.

• Ten agent rendezvous.

Rendezvous Experiment



Team Task: Each agent must *simultaneously* occupy the rendezvous location \mathbb{R} , before proceeding to their respective goals \mathbb{G}_i .

• Two agent rendezvous.

• Ten agent rendezvous.

DQPRM scales well with the number of agents.

Each agent learns in the absence of its teammates.

| A1 | | |
|----|--|--|
| | | |
| | | |
| | | |
| | | |
| | | |
| | | |
| | | |
| | | |
| | | |





| | A4 |
|--|----|
| | |
| | |
| | |
| | |
| | |
| | |
| | |
| | |
| | |

| | _ | _ | |
|----------------|----------|-------|--|
| | | _ | |
| | + + | | |
| | + + | | |
| | + + | | |
| | | | |
| A ₅ | | | |

| | A ₇ | |
|--|----------------|--|
| | | |



| Δ | |
|----|--|
| Ag | |

DQPRM scales well with the number of agents.

Each agent learns in the absence of its teammates.

| 41 | | |
|----|--|--|
| | | |
| | | |
| | | |
| | | |
| | | |
| | | |
| | | |
| | | |
| | | |

| A ₂ | | |
|----------------|--|---|
| | | |
| | | _ |
| | | |
| | | |
| | | |
| | | |







| | _ | | |
|----|---|--|---|
| | | | |
| ۸. | | | |
| 76 | _ | | |
| | _ | | |
| | | | |
| | | | |
| | | | 1 |

| A ₇ |
|----------------|
| |

| | - | | | |
|---|---|---|---|----------------|
| | | | | A ₈ |
| | | | | |
| | | | | |
| | _ | _ | _ | |
| _ | _ | | | |

| | _ | |
|-----|-------|--|
| | Г | |
| | | |
| | | |
| | | |
| | | |
| | | |
| | | |
| | | |
| | | |
| | | |
| Δ. | | |
| 719 | | |

Composite behavior **accomplishes the team task.**



DQPRM scales well with the number of agents.



1.5

Two agent rendezvous

Ten agent rendezvous



Reward machines to **specify** and **decompose** cooperative RL tasks.

Cyrus Neary cneary@utexas.edu

Reward machines to **specify** and **decompose** cooperative RL tasks.

Conditions guaranteeing equivalence between original and decomposed tasks.



Reward machines to **specify** and **decompose** cooperative RL tasks.

Conditions guaranteeing equivalence between original and decomposed tasks.



Decentralized Q-learning algorithm that trains each agent separately.

Reward machines to **specify** and **decompose** cooperative RL tasks.

Conditions guaranteeing equivalence between original and decomposed tasks.



Decentralized Q-learning algorithm that trains each agent separately.





Team

Cyrus Neary cneary@utexas.edu