Teaching Multiple Tasks to a Reinforcement Learning Agent using Linear Temporal Logic (LTL)

Rodrigo Toro IcarteToryn Q. KlassenRichard ValenzanoSheila A. McIlraith



ELEMENT<sup>AI</sup> V<sup>VECTOR</sup> INSTITUT

AAMAS 2018 July 11

#### 1 Motivation

- 2 Related work
- 3 Approach:
  - Why (and how) we use LTL for specifying tasks in RL.
  - Why (and how) LPOPL speeds up learning of multiple tasks.
- 4 Results
- 5 Concluding remarks



## 1 Motivation

- 2 Related work
- 3 Approach:
  - Why (and how) we use LTL for specifying tasks in RL.
  - Why (and how) LPOPL speeds up learning of multiple tasks.
- 4 Results
- 5 Concluding remarks











#### Luigi can collect raw materials:







#### Luigi can collect raw materials:



... and make new objects in:



factory



toolshed



workbench





#### Luigi can collect raw materials:



factory

toolshed

workbench

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



Task type	Example
Single goal	get wood
Sequence of goals	get wood and then use the factory
Disjunctive goals	get wood or iron
Conjunctive goals	get grass and iron
Safety constraints	do not leave the shelter at night



en use the factory
า
shelter at night



 $\mathsf{Description} \neq \mathsf{Reward} \ \mathsf{function}$ 



Description  $\neq$  Reward function

 $\mathsf{Language} \to \mathsf{Reward} \ \mathsf{function}$ 



Description  $\neq$  Reward function

 $\mathsf{Language} \to \mathsf{Reward} \ \mathsf{function}$ 

Why would we want such a language?



Description  $\neq$  Reward function

 $\mathsf{Language} \to \mathsf{Reward} \ \mathsf{function}$ 

Why would we want such a language?

To define new task faster.

To transfer learning between tasks.



### 1 Motivation

- 2 Related work
- 3 Approach:
  - Why (and how) we use LTL for specifying tasks in RL.
  - Why (and how) LPOPL speeds up learning of multiple tasks.
- 4 Results
- 5 Concluding remarks



### 1 Motivation

## 2 Related work

- 3 Approach:
  - Why (and how) we use LTL for specifying tasks in RL.
  - Why (and how) LPOPL speeds up learning of multiple tasks.
- 4 Results
- 5 Concluding remarks



### Hindsight Experience Replay by Andrychowicz et al. (NIPS-17)

Language: Single goal condition Advantage: Learn to achieve goals in parallel (off-policy RL)



### Hindsight Experience Replay by Andrychowicz et al. (NIPS-17)

Language: Single goal condition Advantage: Learn to achieve goals in parallel (off-policy RL)

Task	HER
get wood	1
get wood and then use the factory	
get wood or iron	
get grass and iron	
do not leave the shelter at night	



Modular Multitask RL with Policy Sketches by Andreas et al. (ICML-17)

**Language**: Sequence of sub-goals (called sketch) **Advantage**: Decompose the problem using the sketch.



Modular Multitask RL with Policy Sketches by Andreas et al. (ICML-17)

**Language**: Sequence of sub-goals (called sketch) **Advantage**: Decompose the problem using the sketch.

Task	HER	Sketches
get wood	✓	1
get wood and then use the factory		$\checkmark$
get wood or iron		
get grass and iron		
do not leave the shelter at night		



### Teaching Multiple Tasks to an RL Agent using LTL

Language: Linear Temporal Logic (LTL) Advantage: Decompose and use off-policy RL to learn subtasks.



### Teaching Multiple Tasks to an RL Agent using LTL

#### Language: Linear Temporal Logic (LTL) Advantage: Decompose and use off-policy RL to learn subtasks.

Task	HER	Sketches	LTL
get wood	1	1	✓
get wood and then use the factory		1	$\checkmark$
get wood or iron			$\checkmark$
get grass and iron			$\checkmark$
do not leave the shelter at night			$\checkmark$



### Teaching Multiple Tasks to an RL Agent using LTL

Language: Linear Temporal Logic (LTL) Advantage: Decompose and use off-policy RL to learn subtasks.

Task	HER	Sketches	LTL
get wood	1	1	1
get wood and then use the factory		$\checkmark$	$\checkmark$
get wood or iron			$\checkmark$
get grass and iron			$\checkmark$
do not leave the shelter at night			$\checkmark$
Off-policy learning	1		1
Task decomposition		$\checkmark$	$\checkmark$



Other works using variants of LTL in RL:

- Reinforcement learning with temporal logic rewards by Li et al. (IROS, 2017)
- Environment-independent task specifications via GLTL by Littman et al. (arXiv, 2017)
- Logically constrained reinforcement learning by Hasanbeig et al. (arXiv, 2018)



Other works using variants of LTL in RL:

- Reinforcement learning with temporal logic rewards by Li et al. (IROS, 2017)
- Environment-independent task specifications via GLTL by Littman et al. (arXiv, 2017)
- Logically constrained reinforcement learning by Hasanbeig et al. (arXiv, 2018)

We exploit task decomposition and off-policy RL.



### 1 Motivation

## 2 Related work

- 3 Approach:
  - Why (and how) we use LTL for specifying tasks in RL.
  - Why (and how) LPOPL speeds up learning of multiple tasks.
- 4 Results
- 5 Concluding remarks



#### 1 Motivation

- 2 Related work
- 3 Approach:
  - Why (and how) we use LTL for specifying tasks in RL.
  - Why (and how) LPOPL speeds up learning of multiple tasks.
- 4 Results
- 5 Concluding remarks



Idea: Let's give the RL agent a set of high-level event detectors.



Idea: Let's give the RL agent a set of high-level event detectors.

#### Example

 $\mathcal{P} = \{ \texttt{got\_wood, got\_iron, got\_grass, used\_workbench, used\_factory, is\_night, at\_shelter, ...} \}$ 



Idea: Let's give the RL agent a set of high-level event detectors.

#### Example

Use LTL to define tasks by composing occurrences of events in  ${\mathcal P}$ 



# Defining tasks using LTL

### Linear Temporal Logic (syntax)

LTL augments propositional logic with temporal operators  $\bigcirc$  (*next*),  $\Diamond$  (*eventually*), and U (*until*):

 $\varphi ::= p \mid \neg \varphi \mid \varphi_1 \land \varphi_2 \mid \bigcirc \varphi \mid \Diamond \varphi \mid \varphi_1 \, \mathsf{U} \, \varphi_2 \text{ with } p \in \mathcal{P}$ 



# Defining tasks using LTL

#### Linear Temporal Logic (syntax)

LTL augments propositional logic with temporal operators  $\bigcirc$  (*next*),  $\Diamond$  (*eventually*), and U (*until*):

 $\varphi ::= p \mid \neg \varphi \mid \varphi_1 \land \varphi_2 \mid \bigcirc \varphi \mid \Diamond \varphi \mid \varphi_1 \, \mathsf{U} \, \varphi_2 \text{ with } p \in \mathcal{P}$ 

#### Examples

```
◊got_wood
◊(got_grass ∧ ◊used_factory)
◊got_wood ∨ ◊got_iron
◊got_grass ∧ ◊got_iron
(is_night → at_shelter) Ugot_wood
```



# Defining tasks using LTL

### Linear Temporal Logic (syntax)

LTL augments propositional logic with temporal operators  $\bigcirc$  (*next*),  $\Diamond$  (*eventually*), and U (*until*):

 $\varphi ::= p \mid \neg \varphi \mid \varphi_1 \land \varphi_2 \mid \bigcirc \varphi \mid \Diamond \varphi \mid \varphi_1 \, \mathsf{U} \, \varphi_2 \text{ with } p \in \mathcal{P}$ 

#### Examples

eventually got\_wood eventually (got\_grass and eventually used\_factory) eventually got\_wood or eventually got\_iron eventually got\_grass and eventually got\_iron (is\_night→at\_shelter) until got\_wood





**Detected events** 

 $at_shelter$ 

$$\label{eq:phi} \begin{split} \varphi = \textbf{eventually}(\texttt{got\_iron} \text{ and } \textbf{eventually} \text{ used\_factory}) \\ & \text{and } \textbf{eventually} \text{ got\_gold} \end{split}$$





**Detected events** 

none

#### $\varphi = eventually(got_iron and eventually used_factory)$ and eventually got\_gold





**Detected events** 

none

#### $\varphi = eventually(got_iron and eventually used_factory)$ and eventually got\_gold




none





got\_wood

$$\label{eq:phi} \begin{split} \varphi = \textbf{eventually}(\texttt{got\_iron} \text{ and } \textbf{eventually} \text{ used\_factory}) \\ & \text{and } \textbf{eventually} \text{ got\_gold} \end{split}$$





none





got\_grass





none





none





used\_factory





none





got\_iron

$$\label{eq:phi} \begin{split} \varphi = \textbf{eventually}(\texttt{got\_iron} \text{ and } \textbf{eventually} \text{ used\_factory}) \\ & \text{and } \textbf{eventually} \text{ got\_gold} \end{split}$$





got\_iron



# LTL progression

# LTL progression

Given an LTL formula  $\varphi$  and state *s*, we can *progress*  $\varphi$  using *s*:

- $prog(s, p) = true \text{ if } p \in L(s)$ , where  $p \in \mathcal{P}$
- $prog(s, p) = false if p \notin L(s)$ , where  $p \in \mathcal{P}$

• 
$$\operatorname{prog}(s, \neg \varphi) = \neg \operatorname{prog}(s, \varphi)$$

$$\mathsf{prog}(s,\varphi_1 \land \varphi_2) = \mathsf{prog}(s,\varphi_1) \land \mathsf{prog}(s,\varphi_2)$$

• 
$$\operatorname{prog}(s, \bigcirc \varphi) = \varphi$$

• 
$$\mathsf{prog}(s,\Diamond \varphi) = \mathsf{prog}(s,\varphi) \lor \Diamond \varphi$$

 $\bullet \operatorname{prog}(s, \varphi_1 \cup \varphi_2) = \operatorname{prog}(s, \varphi_2) \vee (\operatorname{prog}(s, \varphi_1) \land \varphi_1 \cup \varphi_2)$ 



# LTL progression

# LTL progression

Given an LTL formula  $\varphi$  and state *s*, we can *progress*  $\varphi$  using *s*:

- $prog(s, p) = true \text{ if } p \in L(s)$ , where  $p \in \mathcal{P}$
- $prog(s, p) = false if p \notin L(s)$ , where  $p \in \mathcal{P}$

• 
$$\operatorname{prog}(s, \neg \varphi) = \neg \operatorname{prog}(s, \varphi)$$

• 
$$\operatorname{prog}(s, \varphi_1 \land \varphi_2) = \operatorname{prog}(s, \varphi_1) \land \operatorname{prog}(s, \varphi_2)$$

• 
$$\operatorname{prog}(s, \bigcirc \varphi) = \varphi$$

• 
$$\mathsf{prog}(s,\Diamond \varphi) = \mathsf{prog}(s,\varphi) \lor \Diamond \varphi$$

 $\bullet \operatorname{prog}(s, \varphi_1 \cup \varphi_2) = \operatorname{prog}(s, \varphi_2) \vee (\operatorname{prog}(s, \varphi_1) \land \varphi_1 \cup \varphi_2)$ 

### This is a correct and well-defined procedure!





got\_iron





none





none





none





none





none





none





none





got\_gold

 $\varphi = \mathbf{eventually} \text{ used_factory and } \mathbf{eventually} \text{ got_gold}$ 





got\_gold





none





none





got\_grass





none





none





used\_factory





used\_factory

$$\varphi = \mathsf{true} (+1 \mathsf{ reward})$$





used\_factory

### LTL formulas $\rightarrow$ Rewards



Toro Icarte et al: Teaching Multiple Tasks to an RL Agent using LTL



used\_factory

( $\varphi$  can be learned using standard RL)





used\_factory

### We can do better than this!



# 1 Motivation

- 2 Related work
- 3 Approach:
  - Why (and how) we use LTL for specifying tasks in RL.
  - Why (and how) LPOPL speeds up learning of multiple tasks.
- 4 Results
- 5 Concluding remarks



- 1 Motivation
- 2 Related work
- 3 Approach:
  - Why (and how) we use LTL for specifying tasks in RL.
  - Why (and how) LPOPL speeds up learning of multiple tasks.
- 4 Results
- 5 Concluding remarks



# LPOPL: An example

Suppose Luigi has to learn two tasks:



Suppose Luigi has to learn two tasks:

 $\varphi_1 = eventually(got_iron and eventually used_factory) and eventually got_gold$ 

 $arphi_2 = {f eventually}([{f got_grass or got_wood}] \ {f and eventually} \ {f used_factory})$ 


Suppose Luigi has to learn two tasks:

 $\varphi_1 = eventually(got_iron and eventually used_factory) and eventually got_gold$ 

 $arphi_2 = extbf{eventually}([ extbf{got_grass} extbf{or got_wood}] extbf{and eventually} used_factory)$ 

 $\ldots$  and begins by trying to solve  $\varphi_1$ 





**Detected events** 

at\_shelter





**Detected events** 

none





**Detected events** 

none





**Detected events** 

none





**Detected events** 

got\_wood

# $$\begin{split} \varphi_1 &= \text{eventually}(\texttt{got\_iron and eventually used\_factory}) \\ & \text{and eventually got\_gold} \\ \varphi_2 &= \texttt{eventually}([\texttt{got\_grass or got\_wood}] \text{ and eventually} \\ & \text{used\_factory}) \end{split}$$





**Detected events** 

got\_wood







**Detected events** 

got\_wood



LPOPL learns all the tasks in parallel!



 $\label{eq:step1} Step \ 1: \ {\tt Decompose the tasks into subtasks using LTL progression}.$ 



**Step 1**: Decompose the tasks into subtasks using LTL progression.



**Step 1**: Decompose the tasks into subtasks using LTL progression.

- $\varphi_1 = eventually(got_iron and eventually used_factory)$ and eventually got\_gold
- $\varphi_2 = eventually([got_grass or got_wood])$ and eventually used\_factory)
- $\varphi_3 = eventually(\texttt{got_iron} and eventually used_factory)$
- $arphi_4 = {f eventually} \; {\tt used\_factory} \; {\tt and} \; {f eventually} \; {\tt got\_gold}$
- $\varphi_5 = {f eventually} \ {\tt used\_factory}$
- $\varphi_6 = {\it eventually got_gold}$
- $\varphi_7 = \mathsf{true}$



 $\label{eq:step1} Step \ 1: \ {\tt Decompose the tasks into subtasks using LTL progression}.$ 

$$\begin{array}{l} \varphi_1 = \Diamond (\texttt{got\_iron} \land \Diamond \texttt{used\_factory}) \\ \land \Diamond \texttt{got\_gold} \\ \varphi_2 = \Diamond ([\texttt{got\_grass} \lor \texttt{got\_wood}] \\ \land \Diamond \texttt{used\_factory}) \\ \varphi_3 = \Diamond (\texttt{got\_iron} \land \Diamond \texttt{used\_factory}) \\ \varphi_4 = \Diamond \texttt{used\_factory} \land \Diamond \texttt{got\_gold} \\ \varphi_5 = \Diamond \texttt{used\_factory} \\ \varphi_6 = \Diamond \texttt{got\_gold} \\ \varphi_7 = \texttt{true} \end{array}$$



Step 1: Decompose the tasks into subtasks using LTL progression.

$$\begin{array}{l} \varphi_1 = \Diamond (\texttt{got\_iron} \land \Diamond \texttt{used\_factory}) \\ \land \Diamond \texttt{got\_gold} \\ \varphi_2 = \Diamond ([\texttt{got\_grass} \lor \texttt{got\_wood}] \\ \land \Diamond \texttt{used\_factory}) \\ \varphi_3 = \Diamond (\texttt{got\_iron} \land \Diamond \texttt{used\_factory}) \\ \varphi_4 = \Diamond \texttt{used\_factory} \land \Diamond \texttt{got\_gold} \\ \varphi_5 = \Diamond \texttt{used\_factory} \\ \varphi_6 = \Diamond \texttt{got\_gold} \\ \varphi_7 = \texttt{true} \end{array}$$

Step 2: Learn one policy per subtask using off-policy learning.



# LPOPL using tabular q-learning

Standard q-learning update given experience (s, a, r, s'):

$$Q(s, a) \stackrel{lpha}{\longleftarrow} r + \gamma \max_{a'} Q(s', a')$$

... where  $x \xleftarrow{\alpha} y$  is a shorthand notation for  $x \leftarrow x + \alpha \cdot (y - x)$ .



# LPOPL using tabular q-learning

Standard q-learning update given experience (s, a, r, s'):

$$Q(s, a) \stackrel{lpha}{\longleftarrow} r + \gamma \max_{a'} Q(s', a')$$

... where  $x \xleftarrow{\alpha} y$  is a shorthand notation for  $x \leftarrow x + \alpha \cdot (y - x)$ .

LPOPL update using q-learning given experience (s, a, s'):

$$\begin{split} Q_{arphi}(s,a) & \stackrel{lpha}{\longleftarrow} r_{arphi} + \gamma \max_{a'} Q_{arphi'}(s',a') \ ... \ \text{where} \ arphi' = \mathrm{prog}(s',arphi) \ \text{and} \ r_{arphi} = 1 \ \mathrm{iff} \ arphi 
eq arphi' = \mathrm{true}. \end{split}$$



# LPOPL using tabular q-learning

Standard q-learning update given experience (s, a, r, s'):

$$Q(s, a) \stackrel{lpha}{\longleftarrow} r + \gamma \max_{a'} Q(s', a')$$

... where  $x \xleftarrow{\alpha} y$  is a shorthand notation for  $x \leftarrow x + \alpha \cdot (y - x)$ .

LPOPL update using q-learning given experience (s, a, s'):

$$\begin{split} Q_{\varphi}(s,a) & \xleftarrow{\alpha} r_{\varphi} + \gamma \max_{a'} Q_{\varphi'}(s',a') \\ \dots \text{ where } \varphi' = \operatorname{prog}(s',\varphi) \text{ and } r_{\varphi} = 1 \text{ iff } \varphi \neq \varphi' = \operatorname{true.} \end{split}$$

## Theorem

LPOPL using tabular q-learning converges to an optimal policy.





$$\begin{split} &\varphi_1 = \diamondsuit(\texttt{got\_iron} \land \diamondsuit\texttt{used\_factory}) \land \diamondsuit\texttt{got\_gold} \\ &\varphi_2 = \diamondsuit(\texttt{[got\_grass} \lor \texttt{got\_wood}] \land \diamondsuit\texttt{used\_factory}) \\ &\varphi_3 = \diamondsuit(\texttt{got\_iron} \land \diamondsuit\texttt{used\_factory}) \\ &\varphi_4 = \diamondsuit\texttt{used\_factory} \land \diamondsuit\texttt{got\_gold} \\ &\varphi_5 = \diamondsuit\texttt{used\_factory} \\ &\varphi_6 = \diamondsuit\texttt{got\_gold} \end{split}$$

Detected events: at\_shelter





$$\begin{array}{l} \varphi_1 = \diamondsuit(\texttt{got\_iron} \land \diamondsuit\texttt{used\_factory}) \land \diamondsuit\texttt{got\_gold} \\ \varphi_2 = \diamondsuit(\texttt{[got\_grass} \lor \texttt{got\_wood}] \land \diamondsuit\texttt{used\_factory}) \\ \varphi_3 = \diamondsuit(\texttt{got\_iron} \land \diamondsuit\texttt{used\_factory}) \\ \varphi_4 = \diamondsuit\texttt{used\_factory} \land \diamondsuit\texttt{got\_gold} \\ \varphi_5 = \diamondsuit\texttt{used\_factory} \\ \varphi_6 = \diamondsuit\texttt{got\_gold} \end{array}$$

#### Detected events: none

$$\begin{array}{ll} Q_{\varphi_1}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_1}(s',a') & Q_{\varphi_4}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_4}(s',a') \\ Q_{\varphi_2}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_2}(s',a') & Q_{\varphi_5}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_5}(s',a') \\ Q_{\varphi_3}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_3}(s',a') & Q_{\varphi_6}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_6}(s',a') \end{array}$$





$$\begin{array}{l} \varphi_1 = \diamondsuit(\texttt{got\_iron} \land \diamondsuit\texttt{used\_factory}) \land \diamondsuit\texttt{got\_gold} \\ \varphi_2 = \diamondsuit(\texttt{[got\_grass} \lor \texttt{got\_wood}] \land \diamondsuit\texttt{used\_factory}) \\ \varphi_3 = \diamondsuit(\texttt{got\_iron} \land \diamondsuit\texttt{used\_factory}) \\ \varphi_4 = \diamondsuit\texttt{used\_factory} \land \diamondsuit\texttt{got\_gold} \\ \varphi_5 = \diamondsuit\texttt{used\_factory} \\ \varphi_6 = \diamondsuit\texttt{got\_gold} \end{array}$$

#### Detected events: none

$$\begin{array}{ll} Q_{\varphi_1}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_1}(s',a') & Q_{\varphi_4}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_4}(s',a') \\ Q_{\varphi_2}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_2}(s',a') & Q_{\varphi_5}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_5}(s',a') \\ Q_{\varphi_3}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_3}(s',a') & Q_{\varphi_6}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_6}(s',a') \end{array}$$





$$\begin{array}{l} \varphi_1 = \diamondsuit(\texttt{got\_iron} \land \diamondsuit\texttt{used\_factory}) \land \diamondsuit\texttt{got\_gold} \\ \varphi_2 = \diamondsuit(\texttt{[got\_grass} \lor \texttt{got\_wood}] \land \diamondsuit\texttt{used\_factory}) \\ \varphi_3 = \diamondsuit(\texttt{got\_iron} \land \diamondsuit\texttt{used\_factory}) \\ \varphi_4 = \diamondsuit\texttt{used\_factory} \land \diamondsuit\texttt{got\_gold} \\ \varphi_5 = \diamondsuit\texttt{used\_factory} \\ \varphi_6 = \diamondsuit\texttt{got\_gold} \end{array}$$

#### Detected events: none

$$\begin{array}{ll} Q_{\varphi_1}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_1}(s',a') & Q_{\varphi_4}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_4}(s',a') \\ Q_{\varphi_2}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_2}(s',a') & Q_{\varphi_5}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_5}(s',a') \\ Q_{\varphi_3}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_3}(s',a') & Q_{\varphi_6}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_6}(s',a') \end{array}$$





$$\begin{array}{l} \varphi_1 = \diamondsuit(\texttt{got\_iron} \land \diamondsuit\texttt{used\_factory}) \land \diamondsuit\texttt{got\_gold} \\ \varphi_2 = \diamondsuit(\texttt{[got\_grass} \lor \texttt{got\_wood}] \land \diamondsuit\texttt{used\_factory}) \\ \varphi_3 = \diamondsuit(\texttt{got\_iron} \land \diamondsuit\texttt{used\_factory}) \\ \varphi_4 = \diamondsuit\texttt{used\_factory} \land \diamondsuit\texttt{got\_gold} \\ \varphi_5 = \diamondsuit\texttt{used\_factory} \\ \varphi_6 = \diamondsuit\texttt{got\_gold} \end{array}$$

Detected events: got\_wood

$$\begin{array}{ll} Q_{\varphi_1}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_1}(s',a') & Q_{\varphi_4}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_4}(s',a') \\ Q_{\varphi_2}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_5}(s',a') & Q_{\varphi_5}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_5}(s',a') \\ Q_{\varphi_3}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_3}(s',a') & Q_{\varphi_6}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_6}(s',a') \end{array}$$





$$\begin{array}{l} \varphi_1 = \diamondsuit(\texttt{got\_iron} \land \diamondsuit\texttt{used\_factory}) \land \diamondsuit\texttt{got\_gold} \\ \varphi_2 = \diamondsuit(\texttt{[got\_grass} \lor \texttt{got\_wood}] \land \diamondsuit\texttt{used\_factory}) \\ \varphi_3 = \diamondsuit(\texttt{got\_iron} \land \diamondsuit\texttt{used\_factory}) \\ \varphi_4 = \diamondsuit\texttt{used\_factory} \land \diamondsuit\texttt{got\_gold} \\ \varphi_5 = \diamondsuit\texttt{used\_factory} \\ \varphi_6 = \diamondsuit\texttt{got\_gold} \end{array}$$

Detected events: got\_wood

$$\begin{array}{l} Q_{\varphi_1}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_1}(s',a') \\ Q_{\varphi_2}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_5}(s',a') \\ Q_{\varphi_3}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_3}(s',a') \end{array}$$

$$\begin{array}{l} Q_{\varphi_4}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_4}(s',a') \\ Q_{\varphi_5}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_5}(s',a') \\ Q_{\varphi_6}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_6}(s',a') \end{array}$$





$$\begin{array}{l} \varphi_1 = & \left( \texttt{got\_iron} \land \left\texttt{oused\_factory} \right) \land \left\texttt{ogt\_gold} \right. \\ \varphi_2 = & \left( [\texttt{got\_grass} \lor \texttt{got\_wood}] \land \left\texttt{oused\_factory} \right) \\ \varphi_3 = & \left( \texttt{got\_iron} \land \left\texttt{oused\_factory} \right) \\ \varphi_4 = & \left\texttt{oused\_factory} \land \left\texttt{ogt\_gold} \right. \\ \varphi_5 = & \left\texttt{oused\_factory} \right. \\ \varphi_6 = & \left\texttt{ogt\_gold} \right. \end{array}$$

#### Detected events: none

$$\begin{array}{ll} Q_{\varphi_1}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_1}(s',a') & Q_{\varphi_4}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_4}(s',a') \\ Q_{\varphi_2}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_2}(s',a') & Q_{\varphi_5}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_5}(s',a') \\ Q_{\varphi_3}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_3}(s',a') & Q_{\varphi_6}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_6}(s',a') \end{array}$$





$$\begin{array}{l} \varphi_1 = \diamondsuit(\texttt{got\_iron} \land \diamondsuit\texttt{used\_factory}) \land \diamondsuit\texttt{got\_gold} \\ \varphi_2 = \diamondsuit(\texttt{[got\_grass} \lor \texttt{got\_wood}] \land \diamondsuit\texttt{used\_factory}) \\ \varphi_3 = \diamondsuit(\texttt{got\_iron} \land \diamondsuit\texttt{used\_factory}) \\ \varphi_4 = \diamondsuit\texttt{used\_factory} \land \diamondsuit\texttt{got\_gold} \\ \varphi_5 = \diamondsuit\texttt{used\_factory} \\ \varphi_6 = \diamondsuit\texttt{got\_gold} \end{array}$$

Detected events: got\_grass

$$\begin{array}{ll} Q_{\varphi_1}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_1}(s',a') & Q_{\varphi_4}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_4}(s',a') \\ Q_{\varphi_2}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_5}(s',a') & Q_{\varphi_5}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_5}(s',a') \\ Q_{\varphi_3}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_3}(s',a') & Q_{\varphi_6}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_6}(s',a') \end{array}$$





$$\begin{array}{l} \varphi_1 = \diamondsuit(\texttt{got\_iron} \land \diamondsuit\texttt{used\_factory}) \land \diamondsuit\texttt{got\_gold} \\ \varphi_2 = \diamondsuit(\texttt{[got\_grass} \lor \texttt{got\_wood}] \land \diamondsuit\texttt{used\_factory}) \\ \varphi_3 = \diamondsuit(\texttt{got\_iron} \land \diamondsuit\texttt{used\_factory}) \\ \varphi_4 = \diamondsuit\texttt{used\_factory} \land \diamondsuit\texttt{got\_gold} \\ \varphi_5 = \diamondsuit\texttt{used\_factory} \\ \varphi_6 = \diamondsuit\texttt{got\_gold} \end{array}$$

Detected events: got\_grass

$$\begin{array}{l} Q_{\varphi_1}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_1}(s',a') \\ Q_{\varphi_2}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_5}(s',a') \\ Q_{\varphi_3}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_3}(s',a') \end{array}$$

$$\begin{array}{c} Q_{\varphi_4}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_4}(s',a') \\ Q_{\varphi_5}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_5}(s',a') \\ Q_{\varphi_6}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_6}(s',a') \end{array}$$





$$\begin{array}{l} \varphi_1 = & \left( \texttt{got\_iron} \land \left\texttt{oused\_factory} \right) \land \left\texttt{ogt\_gold} \right. \\ \varphi_2 = & \left( [\texttt{got\_grass} \lor \texttt{got\_wood}] \land \left\texttt{oused\_factory} \right) \\ \varphi_3 = & \left( \texttt{got\_iron} \land \left\texttt{oused\_factory} \right) \\ \varphi_4 = & \left\texttt{oused\_factory} \land \left\texttt{ogt\_gold} \right. \\ \varphi_5 = & \left\texttt{oused\_factory} \right. \\ \varphi_6 = & \left\texttt{ogt\_gold} \right. \end{array}$$

#### Detected events: none

$$\begin{array}{ll} Q_{\varphi_1}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_1}(s',a') & Q_{\varphi_4}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_4}(s',a') \\ Q_{\varphi_2}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_2}(s',a') & Q_{\varphi_5}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_5}(s',a') \\ Q_{\varphi_3}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_3}(s',a') & Q_{\varphi_6}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_6}(s',a') \end{array}$$





$$\begin{array}{l} \varphi_1 = & \left( \texttt{got\_iron} \land \left\texttt{oused\_factory} \right) \land \left\texttt{ogt\_gold} \right. \\ \varphi_2 = & \left( [\texttt{got\_grass} \lor \texttt{got\_wood}] \land \left\texttt{oused\_factory} \right) \\ \varphi_3 = & \left( \texttt{got\_iron} \land \left\texttt{oused\_factory} \right) \\ \varphi_4 = & \left\texttt{oused\_factory} \land \left\texttt{ogt\_gold} \right. \\ \varphi_5 = & \left\texttt{oused\_factory} \right. \\ \varphi_6 = & \left\texttt{ogt\_gold} \right. \end{array}$$

#### Detected events: none

$$\begin{array}{ll} Q_{\varphi_1}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_1}(s',a') & Q_{\varphi_4}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_4}(s',a') \\ Q_{\varphi_2}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_2}(s',a') & Q_{\varphi_5}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_5}(s',a') \\ Q_{\varphi_3}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_3}(s',a') & Q_{\varphi_6}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_6}(s',a') \end{array}$$





$$\begin{array}{l} \varphi_1 = \diamondsuit(\texttt{got\_iron} \land \diamondsuit\texttt{used\_factory}) \land \diamondsuit\texttt{got\_gold} \\ \varphi_2 = \diamondsuit(\texttt{[got\_grass} \lor \texttt{got\_wood}] \land \diamondsuit\texttt{used\_factory}) \\ \varphi_3 = \diamondsuit(\texttt{got\_iron} \land \diamondsuit\texttt{used\_factory}) \\ \varphi_4 = \diamondsuit\texttt{used\_factory} \land \diamondsuit\texttt{got\_gold} \\ \varphi_5 = \diamondsuit\texttt{used\_factory} \\ \varphi_6 = \diamondsuit\texttt{got\_gold} \end{array}$$

Detected events: used\_factory

$$\begin{array}{ll} Q_{\varphi_1}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_1}(s',a') & Q_{\varphi_4}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_6}(s',a') \\ Q_{\varphi_2}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_2}(s',a') & Q_{\varphi_5}(s,a) \xleftarrow{\alpha} 1 \\ Q_{\varphi_3}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_3}(s',a') & Q_{\varphi_6}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_6}(s',a') \end{array}$$





$$\begin{array}{l} \varphi_1 = \diamondsuit(\texttt{got\_iron} \land \diamondsuit\texttt{used\_factory}) \land \diamondsuit\texttt{got\_gold} \\ \varphi_2 = \diamondsuit(\texttt{[got\_grass} \lor \texttt{got\_wood}] \land \diamondsuit\texttt{used\_factory}) \\ \varphi_3 = \diamondsuit(\texttt{got\_iron} \land \diamondsuit\texttt{used\_factory}) \\ \varphi_4 = \diamondsuit\texttt{used\_factory} \land \diamondsuit\texttt{got\_gold} \\ \varphi_5 = \diamondsuit\texttt{used\_factory} \\ \varphi_6 = \diamondsuit\texttt{got\_gold} \end{array}$$

Detected events: used\_factory

$$\begin{array}{c} Q_{\varphi_1}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_1}(s',a') \\ Q_{\varphi_2}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_2}(s',a') \\ Q_{\varphi_3}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_3}(s',a') \end{array}$$

$$\begin{array}{l} Q_{\varphi_4}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_6}(s',a') \\ Q_{\varphi_5}(s,a) \xleftarrow{\alpha} 1 \\ Q_{\varphi_6}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_6}(s',a') \end{array}$$





$$\begin{array}{l} \varphi_1 = \diamondsuit(\texttt{got\_iron} \land \diamondsuit\texttt{used\_factory}) \land \diamondsuit\texttt{got\_gold} \\ \varphi_2 = \diamondsuit(\texttt{[got\_grass} \lor \texttt{got\_wood}] \land \diamondsuit\texttt{used\_factory}) \\ \varphi_3 = \diamondsuit(\texttt{got\_iron} \land \diamondsuit\texttt{used\_factory}) \\ \varphi_4 = \diamondsuit\texttt{used\_factory} \land \diamondsuit\texttt{got\_gold} \\ \varphi_5 = \diamondsuit\texttt{used\_factory} \\ \varphi_6 = \diamondsuit\texttt{got\_gold} \end{array}$$

Detected events: used\_factory

$$\begin{array}{ll} Q_{\varphi_1}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_1}(s',a') & Q_{\varphi_4}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_6}(s',a') \\ Q_{\varphi_2}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_2}(s',a') & Q_{\varphi_5}(s,a) \xleftarrow{\alpha} 1 \\ Q_{\varphi_3}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_3}(s',a') & Q_{\varphi_6}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_6}(s',a') \end{array}$$





$$\begin{array}{l} \varphi_1 = & \left( \texttt{got\_iron} \land \left\texttt{oused\_factory} \right) \land \left\texttt{ogt\_gold} \right. \\ \varphi_2 = & \left( [\texttt{got\_grass} \lor \texttt{got\_wood}] \land \left\texttt{oused\_factory} \right) \\ \varphi_3 = & \left( \texttt{got\_iron} \land \left\texttt{oused\_factory} \right) \\ \varphi_4 = & \left\texttt{oused\_factory} \land \left\texttt{ogt\_gold} \right. \\ \varphi_5 = & \left\texttt{oused\_factory} \right. \\ \varphi_6 = & \left\texttt{ogt\_gold} \right. \end{array}$$

#### Detected events: none

$$\begin{array}{ll} Q_{\varphi_1}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_1}(s',a') & Q_{\varphi_4}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_4}(s',a') \\ Q_{\varphi_2}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_2}(s',a') & Q_{\varphi_5}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_5}(s',a') \\ Q_{\varphi_3}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_3}(s',a') & Q_{\varphi_6}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_6}(s',a') \end{array}$$





$$\begin{array}{l} \varphi_1 = \diamondsuit(\texttt{got\_iron} \land \diamondsuit\texttt{used\_factory}) \land \diamondsuit\texttt{got\_gold} \\ \varphi_2 = \diamondsuit(\texttt{[got\_grass} \lor \texttt{got\_wood}] \land \diamondsuit\texttt{used\_factory}) \\ \varphi_3 = \diamondsuit(\texttt{got\_iron} \land \diamondsuit\texttt{used\_factory}) \\ \varphi_4 = \diamondsuit\texttt{used\_factory} \land \diamondsuit\texttt{got\_gold} \\ \varphi_5 = \diamondsuit\texttt{used\_factory} \\ \varphi_6 = \diamondsuit\texttt{got\_gold} \end{array}$$

Detected events: got\_iron

$$\begin{array}{ll} Q_{\varphi_1}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_4}(s',a') & Q_{\varphi_4}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_4}(s',a') \\ Q_{\varphi_2}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_2}(s',a') & Q_{\varphi_5}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_5}(s',a') \\ Q_{\varphi_3}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_5}(s',a') & Q_{\varphi_6}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_6}(s',a') \end{array}$$





$$\begin{array}{l} \varphi_1 = \Diamond (\texttt{got\_iron} \land \Diamond \texttt{used\_factory}) \land \Diamond \texttt{got\_gold} \\ \varphi_2 = \Diamond ([\texttt{got\_grass} \lor \texttt{got\_wood}] \land \Diamond \texttt{used\_factory}) \\ \varphi_3 = \Diamond (\texttt{got\_iron} \land \Diamond \texttt{used\_factory}) \\ \varphi_4 = \Diamond \texttt{used\_factory} \land \Diamond \texttt{got\_gold} \\ \varphi_5 = \Diamond \texttt{used\_factory} \\ \varphi_6 = \Diamond \texttt{got\_gold} \end{array}$$

Detected events: got\_iron

$$\begin{array}{ll} Q_{\varphi_1}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_4}(s',a') & Q_{\varphi_4}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_4}(s',a') \\ Q_{\varphi_2}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_2}(s',a') & Q_{\varphi_5}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_5}(s',a') \\ Q_{\varphi_3}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_5}(s',a') & Q_{\varphi_6}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_6}(s',a') \end{array}$$





$$\begin{array}{l} \varphi_1 = \Diamond (\texttt{got\_iron} \land \Diamond \texttt{used\_factory}) \land \Diamond \texttt{got\_gold} \\ \varphi_2 = \Diamond ([\texttt{got\_grass} \lor \texttt{got\_wood}] \land \Diamond \texttt{used\_factory}) \\ \varphi_3 = \Diamond (\texttt{got\_iron} \land \Diamond \texttt{used\_factory}) \\ \varphi_4 = \Diamond \texttt{used\_factory} \land \Diamond \texttt{got\_gold} \\ \varphi_5 = \Diamond \texttt{used\_factory} \\ \varphi_6 = \Diamond \texttt{got\_gold} \end{array}$$

Detected events: got\_iron

$$\begin{array}{ll} Q_{\varphi_1}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_4}(s',a') & Q_{\varphi_4}(s,a') \\ Q_{\varphi_2}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_2}(s',a') & Q_{\varphi_5}(s,a') \\ Q_{\varphi_3}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_5}(s',a') & Q_{\varphi_6}(s,a) \end{array}$$

$$\begin{array}{l} Q_{\varphi_4}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_4}(s',a') \\ Q_{\varphi_5}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_5}(s',a') \\ Q_{\varphi_6}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_6}(s',a') \end{array}$$





 $\begin{array}{l} \varphi_1 = \Diamond (\texttt{got\_iron} \land \Diamond \texttt{used\_factory}) \land \Diamond \texttt{got\_gold} \\ \varphi_2 = \Diamond ([\texttt{got\_grass} \lor \texttt{got\_wood}] \land \Diamond \texttt{used\_factory}) \\ \varphi_3 = \Diamond (\texttt{got\_iron} \land \Diamond \texttt{used\_factory}) \\ \varphi_4 = \Diamond \texttt{used\_factory} \land \Diamond \texttt{got\_gold} \\ \varphi_5 = \Diamond \texttt{used\_factory} \\ \varphi_6 = \Diamond \texttt{got\_gold} \end{array}$ 

Detected events: got\_iron

$$\begin{array}{ll} Q_{\varphi_1}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_4}(s',a') & Q_{\varphi_4}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_4}(s',a') \\ Q_{\varphi_2}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_2}(s',a') & Q_{\varphi_5}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_5}(s',a') \\ Q_{\varphi_3}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_5}(s',a') & Q_{\varphi_6}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_6}(s',a') \end{array}$$





$$\begin{array}{l} \varphi_1 = \Diamond (\texttt{got\_iron} \land \Diamond \texttt{used\_factory}) \land \Diamond \texttt{got\_gold} \\ \varphi_2 = \Diamond ([\texttt{got\_grass} \lor \texttt{got\_wood}] \land \Diamond \texttt{used\_factory}) \\ \varphi_3 = \Diamond (\texttt{got\_iron} \land \Diamond \texttt{used\_factory}) \\ \varphi_4 = \Diamond \texttt{used\_factory} \land \Diamond \texttt{got\_gold} \\ \varphi_5 = \Diamond \texttt{used\_factory} \\ \varphi_6 = \Diamond \texttt{got\_gold} \end{array}$$

Detected events: none

$$\begin{array}{ll} Q_{\varphi_1}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_1}(s',a') & Q_{\varphi_4}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_4}(s',a') \\ Q_{\varphi_2}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_2}(s',a') & Q_{\varphi_5}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_5}(s',a') \\ Q_{\varphi_3}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_3}(s',a') & Q_{\varphi_6}(s,a) \xleftarrow{\alpha} \gamma \max_{a'} Q_{\varphi_6}(s',a') \end{array}$$


### 1 Motivation

- 2 Related work
- 3 Approach:
  - Why (and how) we use LTL for specifying tasks in RL.
  - Why (and how) LPOPL speeds up learning of multiple tasks.
- 4 Results
- 5 Concluding remarks



### 1 Motivation

- 2 Related work
- 3 Approach:
  - Why (and how) we use LTL for specifying tasks in RL.
  - Why (and how) LPOPL speeds up learning of multiple tasks.

## 4 Results

5 Concluding remarks



## Experiments

## Goal

- Study LPOPL + DQN
- Compare with standard RL
- Compare with alternative decomposition methods for RL



## Experiments

## Goal

- Study LPOPL + DQN
- Compare with standard RL
- Compare with alternative decomposition methods for RL

### Baselines

- DQN-L: Standard DQN<sup>1</sup> (no decomposition).
- HRL-E: Hierarchical Deep RL proposed by Kulkarni et al.<sup>2</sup>
- HRL-L: HRL-E but exploiting LTL to prune *useless* options.

<sup>1</sup>Human-level control through deep reinforcement learning (Nature-15)
<sup>2</sup>HDRL: Integrating Temporal Abstraction and Intrinsic Motivation (NIPS-16)



# Baselines: HRL-E



# Baselines: HRL-L



# Baselines: HRL-L



# Hierarchical RL might converge to suboptimal policies





# Hierarchical RL might converge to suboptimal policies





# Hierarchical RL might converge to suboptimal policies



We tested over 5 random grids and 5 adversarial grids.



### Description

10 tasks defined as sequence of subgoals (Andreas et al., 2017) e.g. get iron, then get wood, then use factory.



# Experiment 1: subgoal sequences

### Description

10 tasks defined as sequence of subgoals (Andreas et al., 2017) e.g. get iron, then get wood, then use factory.



### Description

Same set of 10 tasks but removing unnecessary order constraints. e.g. (get iron and get wood), then use factory.



# Experiment 2: Interleaving subtasks

#### Description

Same set of 10 tasks but removing unnecessary order constraints. e.g. (get iron and get wood), then use factory.



### Description

Same set of 10 tasks but including the safety constraint of being at the shelter during the night.



# Experiment 3: Safety constraints

### Description

Same set of 10 tasks but including the safety constraint of being at the shelter during the night.





### 1 Motivation

- 2 Related work
- 3 Approach:
  - Why (and how) we use LTL for specifying tasks in RL.
  - Why (and how) LPOPL speeds up learning of multiple tasks.
- 4 Results
- 5 Concluding remarks



### 1 Motivation

- 2 Related work
- 3 Approach:
  - Why (and how) we use LTL for specifying tasks in RL.
  - Why (and how) LPOPL speeds up learning of multiple tasks.
- 4 Results
- 5 Concluding remarks



# Concluding remarks



## Problem

How to tell an RL agent what to do.

### Our approach

Define tasks using LTL. Decompose tasks using LTL progression. Learn subtasks using off-policy RL.



# Concluding remarks



#### Problem

How to tell an RL agent what to do.

#### Our approach

Define tasks using LTL. Decompose tasks using LTL progression. Learn subtasks using off-policy RL.

#### Experiments:

- LPOPL outperformed HRL and DQN over a variety of tasks.
- https://bitbucket.org/RToroIcarte/lpopl



# Thanks!

