

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

Rodrigo Toro Icarte Toryn Q. Klassen
Richard Valenzano Sheila A. McIlraith



Computer Science
UNIVERSITY OF TORONTO

ELEMENT^{AI}



Learning by Instruction Workshop,
NeurIPS 2018

- 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

How can we instruct RL agents?

Running example



Running example



Luigi can collect raw materials:



wood



grass



iron



gold



gems

... and make new objects in:



factory



toolshed



workbench

Running example



Luigi can collect raw materials:



wood



grass



iron



gold



gems

... and make new objects in:



factory



toolshed



workbench

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

Running example

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

Running example

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

How can we instruct RL agents?

Motivation

Question: How can we instruct RL agents?

Proposal: Let's use language to describe tasks!

Question: How can we instruct RL agents?

Proposal: Let's use language to describe tasks!

Desirable properties:

- It is expressive.
- RL agents understand it:
 - Task description / Reward function.
 - Task description / Learn faster.
- Humans understand it.

- 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)

Task	HER
get wood	3
get wood and then use the factory	
get wood or iron	
get grass and iron	
do not leave the shelter at night	
Off-policy learning	3

Related work

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	3	3
get wood and then use the factory		3
get wood or iron		
get grass and iron		
do not leave the shelter at night		
Off-policy learning	3	
Task decomposition		3

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	3	3	3
get wood and then use the factory		3	3
get wood or iron			3
get grass and iron			3
do not leave the shelter at night			3
Off-policy learning	3		3
Task decomposition		3	3

- 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.

Defining tasks using LTL

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

Example

$$P = f_{\text{got_wood}, \text{got_iron}, \text{got_grass}, \text{used_workbench}, \text{used_factory}, \text{is_night}, \text{at_shelter}, \dots} g$$

Defining tasks using LTL

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

Example

$$P = \{ \text{got_wood}, \text{got_iron}, \text{got_grass}, \text{used_workbench}, \\ \text{used_factory}, \text{is_night}, \text{at_shelter}, \dots \}$$

Use LTL to define tasks by composing occurrences of events in P

Defining tasks using LTL

Linear Temporal Logic (syntax)

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

$$\phi ::= p \mid \neg \phi \mid \phi_1 \wedge \phi_2 \mid \phi_1 \vee \phi_2 \mid \phi_1 \text{U} \phi_2 \text{ with } p \in P$$

Defining tasks using LTL

Linear Temporal Logic (syntax)

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

$$\phi ::= p \mid j : \phi_1 \wedge \phi_2 \mid \phi_1 \text{ U } \phi_2 \text{ with } p \in P$$

Examples

```
got_wood
(got_grass ^ used_factory)
got_wood _ got_iron
got_grass ^ got_iron
(is_night / at_shelter) U got_wood
```


Defining tasks using LTL

Linear Temporal Logic (syntax)

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

$$\phi ::= p \mid j \mid \neg \phi \mid \phi_1 \wedge \phi_2 \mid \phi_1 \vee \phi_2 \mid \phi_1 U \phi_2 \text{ with } p \in P$$

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

Example



Detected events
at_shel ter

' = **eventually**(got_i ron **and eventually** used_factory)
and eventually got_gold

Example



Detected events

none

' = **eventually**(got_i ron **and eventually** used_factory)
and eventually got_gold

Example



Detected events

none

' = **eventually**(got_i ron **and eventually** used_factory)
and eventually got_gold

Example



Detected events

none

' = **eventually**(got_i ron **and eventually** used_factory)
and eventually got_gold

Example



Detected events

got_wood

' = **eventually**(got_iron **and eventually** used_factory)
and eventually got_gold

Example



Detected events

none

' = **eventually**(got_i ron **and eventually** used_factory)
and eventually got_gol d

Example



Detected events

none

' = **eventually**(got_i ron **and eventually** used_factory)
and eventually got_gol d

Example



Detected events

got_i ron

' = **eventually**(got_i ron **and eventually** used_factory)
and eventually got_gol d

Example



Detected events

got_i ron

' = **eventually**(~~got_i ron~~ and **eventually** used_factory)
and eventually got_gold

Example



Detected events

got_i ron

' = **eventually** used_factory **and eventually** got_gol d

LTL progression

Given an LTL formula ϕ and state s , we can *progress* ϕ using s :

- $\text{prog}(s; p) = \text{true}$ if $p \in L(s)$, where $p \in P$
- $\text{prog}(s; p) = \text{false}$ if $p \notin L(s)$, where $p \in P$
- $\text{prog}(s; \neg \phi) = \neg \text{prog}(s; \phi)$
- $\text{prog}(s; \phi_1 \wedge \phi_2) = \text{prog}(s; \phi_1) \wedge \text{prog}(s; \phi_2)$
- $\text{prog}(s; \phi) = \phi$
- $\text{prog}(s; \neg \phi) = \neg \text{prog}(s; \phi)$
- $\text{prog}(s; \phi_1 \cup \phi_2) = \text{prog}(s; \phi_2) \cup (\text{prog}(s; \phi_1) \wedge \neg \phi_1 \cup \phi_2)$

LTL progression

Given an LTL formula ϕ and state s , we can *progress* ϕ using s :

- $\text{prog}(s; p) = \text{true}$ if $p \in L(s)$, where $p \in P$
- $\text{prog}(s; p) = \text{false}$ if $p \notin L(s)$, where $p \in P$
- $\text{prog}(s; \neg \phi) = \neg \text{prog}(s; \phi)$
- $\text{prog}(s; \phi_1 \wedge \phi_2) = \text{prog}(s; \phi_1) \wedge \text{prog}(s; \phi_2)$
- $\text{prog}(s; \phi) = \phi$
- $\text{prog}(s; \neg \phi) = \neg \text{prog}(s; \phi)$
- $\text{prog}(s; \phi_1 \cup \phi_2) = \text{prog}(s; \phi_2) \cup (\text{prog}(s; \phi_1) \wedge \neg \phi_2)$

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

Example



Detected events

got_i ron

' = **eventually** used_factory **and eventually** got_gol d

Example



Detected events

none

' = **eventually** used_factory **and eventually** got_gold

Example



Detected events

none

' = **eventually** used_factory **and eventually** got_gold

Example



Detected events

none

' = **eventually** used_factory **and eventually** got_gold

Example



Detected events

none

' = **eventually** used_factory **and eventually** got_gold

Example



Detected events

none

' = **eventually** used_factory **and eventually** got_gold

Example



Detected events
got_gol d

' = **eventually** used_factory **and eventually** got_gol d

Example



Detected events

got_gold

' = **eventually** used_factory **and eventually** ~~got_gold~~

Example



Detected events

got_gold

' = **eventually** used_factory

Example



Detected events

none

' = **eventually** used_factory

Example



Detected events

none

' = **eventually** used_factory

Example



Detected events

none

' = **eventually** used_factory

Example



Detected events
used_factory

' = **eventually** used_factory

Example



Detected events
used_factory

' = **eventually** ~~used_factory~~

Example



Detected events
used_factory

' = true (+1 reward)

Example



Detected events
used_factory

LTL formulas ! Rewards

Example



Detected events
used_factory

(' can be learned using standard RL)

Example



Detected events
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:

LPOPL: An example

Suppose Luigi has to learn two tasks:

'₁ = **eventually**(got_iron **and** **eventually** used_factory) **and**
eventually got_gold

'₂ = **eventually**([got_grass **or** got_wood] **and** **eventually**
used_factory)

LPOPL: An example

Suppose Luigi has to learn two tasks:

'₁ = **eventually**(got_iron **and** **eventually** used_factory) **and**
eventually got_gold

'₂ = **eventually**([got_grass **or** got_wood] **and** **eventually**
used_factory)

... and begins by trying to solve '₁

LPOPL: An example



Detected events
at_shel ter

- '₁ = **eventually**(got_i ron **and eventually** used_factory)
and eventually got_gol d
- '₂ = **eventually**([got_grass **or** got_wood] **and eventually**
used_factory)

LPOPL: An example



Detected events

none

- '₁ = **eventually**(got_iron **and eventually** used_factory)
and eventually got_gold
- '₂ = **eventually**([got_grass **or** got_wood] **and eventually**
used_factory)

LPOPL: An example



Detected events

none

- '₁ = **eventually**(got_iron **and eventually** used_factory)
and eventually got_gold
- '₂ = **eventually**([got_grass **or** got_wood] **and eventually**
used_factory)

LPOPL: An example



Detected events

none

- '₁ = **eventually**(got_iron **and eventually** used_factory)
and eventually got_gold
- '₂ = **eventually**([got_grass **or** got_wood] **and eventually**
used_factory)

LPOPL: An example



Detected events
got_wood

- '₁ = **eventually**(got_iron **and eventually** used_factory)
and eventually got_gold
- '₂ = **eventually**([got_grass **or** got_wood] **and eventually**
used_factory)

LPOPL: An example



Detected events
got_wood

- '₁ = **eventually**(got_iron **and eventually** used_factory)
and eventually got_gold
- '₂ = **eventually**([got_grass **or** got_wood] **and eventually**
used_factory)

LPOPL: An example



Detected events
got_wood

- '₁ = **eventually**(got_iron **and eventually** used_factory)
and eventually got_gold
- '₂ = **eventually**([got_grass **or** got_wood] **and eventually**
used_factory)

LPOPL learns all the tasks in parallel!

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

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

'₁ = **eventually**(got_iron **and eventually** used_factory)
and eventually got_gold

'₂ = **eventually**([got_grass **or** got_wood]
and eventually used_factory)

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

- '₁ = **eventually**(got_i ron **and eventually** used_factory)
and eventually got_gol d
- '₂ = **eventually**([got_grass **or** got_wood]
and eventually used_factory)
- '₃ = **eventually**(got_i ron **and eventually** used_factory)
- '₄ = **eventually** used_factory **and eventually** got_gol d
- '₅ = **eventually** used_factory
- '₆ = **eventually** got_gol d
- '₇ = true

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

$$\tau_1 = (\text{got_iron} \wedge \text{used_factory}) \\ \wedge \text{got_gold}$$

$$\tau_2 = ([\text{got_grass} _ \text{got_wood}] \\ \wedge \text{used_factory})$$

$$\tau_3 = (\text{got_iron} \wedge \text{used_factory})$$

$$\tau_4 = \text{used_factory} \wedge \text{got_gold}$$

$$\tau_5 = \text{used_factory}$$

$$\tau_6 = \text{got_gold}$$

$$\tau_7 = \text{true}$$

LPOPL overview

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

$$\tau_1 = (\text{got_iron} \wedge \text{used_factory}) \\ \wedge \text{got_gold}$$

$$\tau_2 = ([\text{got_grass} _ \text{got_wood}] \\ \wedge \text{used_factory})$$

$$\tau_3 = (\text{got_iron} \wedge \text{used_factory})$$

$$\tau_4 = \text{used_factory} \wedge \text{got_gold}$$

$$\tau_5 = \text{used_factory}$$

$$\tau_6 = \text{got_gold}$$

$$\tau_7 = \text{true}$$

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

Example



Subtasks

-
- '₁ = (got_iron ^ used_factory) ^ got_gold
 - '₂ = ([got_grass _ got_wood] ^ used_factory)
 - '₃ = (got_iron ^ used_factory)
 - '₄ = **used_factory** ^ **got_gold**
 - '₅ = used_factory
 - '₆ = got_gold

Detected events: none

Example



Subtasks

-
- '₁ = (got_iron ^ used_factory) ^ got_gold
 - '₂ = ([got_grass _ got_wood] ^ used_factory)
 - '₃ = (got_iron ^ used_factory)
 - '₄ = **used_factory ^ got_gold**
 - '₅ = used_factory
 - '₆ = got_gold

Detected events: none

Q-updates

$$Q'_1(s; a)$$

$$Q'_2(s; a)$$

$$Q'_3(s; a)$$

$$Q'_4(s; a)$$

$$Q'_5(s; a)$$

$$Q'_6(s; a)$$

Example



Subtasks

-
- '₁ = (got_iron ^ used_factory) ^ got_gold
 - '₂ = ([got_grass _ got_wood] ^ used_factory)
 - '₃ = (got_iron ^ used_factory)
 - '₄ = **used_factory ^ got_gold**
 - '₅ = used_factory
 - '₆ = got_gold

Detected events: none

Q-updates

$$Q'_1(s; a)$$

$$Q'_2(s; a)$$

$$Q'_3(s; a)$$

$$Q'_4(s; a)$$

$$Q'_5(s; a)$$

$$Q'_6(s; a)$$

Example

Subtasks

- '₁ = (got_iron ^ used_factory) ^ got_gold
- '₂ = ([got_grass _ got_wood] ^ used_factory)
- '₃ = (got_iron ^ used_factory)
- '₄ = used_factory ^ got_gold
- '₅ = used_factory
- '₆ = got_gold

Detected events: got_gold

Q-updates

$$Q'_1(s; a)$$

$$Q'_2(s; a)$$

$$Q'_3(s; a)$$

$$Q'_4(s; a)$$

$$Q'_5(s; a)$$

$$Q'_6(s; a)$$

Example

Subtasks

- '₁ = (got_iron ^ used_factory) ^ got_gold
- '₂ = ([got_grass _ got_wood] ^ used_factory)
- '₃ = (got_iron ^ used_factory)
- '₄ = used_factory ^ got_gold
- '₅ = **used_factory**
- '₆ = got_gold

Detected events: got_gold

Q-updates

$$Q'_{1}(s; a)$$

$$Q'_{2}(s; a)$$

$$Q'_{3}(s; a)$$

$$Q'_{4}(s; a)$$

$$Q'_{5}(s; a)$$

$$Q'_{6}(s; a)$$

Example

Subtasks

- '₁ = (got_iron ^ used_factory) ^ got_gold
- '₂ = ([got_grass _ got_wood] ^ used_factory)
- '₃ = (got_iron ^ used_factory)
- '₄ = used_factory ^ got_gold
- '₅ = **used_factory**
- '₆ = got_gold

Detected events: got_gold

Q-updates

$Q'_1(s; a)$

$Q'_2(s; a)$

$Q'_3(s; a)$

$Q'_4(s; a)$

$Q'_5(s; a)$

$Q'_6(s; a)$

Example

Subtasks

- '₁ = (got_iron ^ used_factory) ^ got_gold
- '₂ = ([got_grass _ got_wood] ^ used_factory)
- '₃ = (got_iron ^ used_factory)
- '₄ = used_factory ^ got_gold
- '₅ = used_factory
- '₆ = got_gold

Detected events: got_gold

Q-updates

$Q'_{1}(s; a)$	$Q'_{4}(s; a)$	$\max_{a^{\theta}} Q'_{5}(s^{\theta}; a^{\theta})$
$Q'_{2}(s; a)$	$Q'_{5}(s; a)$	
$Q'_{3}(s; a)$	$Q'_{6}(s; a)$	

Example

Subtasks

- '₁ = (got_iron ^ used_factory) ^ got_gold
- '₂ = ([got_grass _ got_wood] ^ used_factory)
- '₃ = (got_iron ^ used_factory)
- '₄ = used_factory ^ got_gold
- '₅ = used_factory
- '₆ = got_gold

Detected events: got_gold

Q-updates

$Q'_{1}(s; a)$	$Q'_{4}(s; a)$	$\max_{a^{\theta}} Q'_{5}(s^{\theta}; a^{\theta})$
$Q'_{2}(s; a)$	$Q'_{5}(s; a)$	
$Q'_{3}(s; a)$	$Q'_{6}(s; a)$	

Example

Subtasks

-
- '₁ = (got_iron ^ used_factory) ^ got_gold
 - '₂ = ([got_grass _ got_wood] ^ used_factory)
 - '₃ = (got_iron ^ used_factory)
 - '₄ = used_factory ^ got_gold
 - '₅ = used_factory
 - '₆ = got_gold

Detected events: got_gold

Q-updates

$Q'_{1}(s; a)$	$Q'_{4}(s; a)$	$\max_{a^{\theta}} Q'_{5}(s^{\theta}; a^{\theta})$
$Q'_{2}(s; a)$	$Q'_{5}(s; a)$	
$Q'_{3}(s; a)$	$Q'_{6}(s; a)$	1

Example

Subtasks

-
- '₁ = (got_iron ^ used_factory) ^ got_gold
 - '₂ = ([got_grass _ got_wood] ^ used_factory)
 - '₃ = (got_iron ^ used_factory)
 - '₄ = used_factory ^ got_gold
 - '₅ = used_factory
 - '₆ = got_gold

Detected events: got_gold

Q-updates

$Q'_{1}(s; a)$	$\max_{a^0} Q'_{3}(s^0; a^0)$	$Q'_{4}(s; a)$	$\max_{a^0} Q'_{5}(s^0; a^0)$
$Q'_{2}(s; a)$	$\max_{a^0} Q'_{2}(s^0; a^0)$	$Q'_{5}(s; a)$	$\max_{a^0} Q'_{5}(s^0; a^0)$
$Q'_{3}(s; a)$	$\max_{a^0} Q'_{3}(s^0; a^0)$	$Q'_{6}(s; a)$	1

Example

Subtasks

- '₁ = (got_iron ^ used_factory) ^ got_gold
- '₂ = ([got_grass _ got_wood] ^ used_factory)
- '₃ = (got_iron ^ used_factory)
- '₄ = used_factory ^ got_gold
- '₅ = used_factory
- '₆ = got_gold

Detected events: got_gold

Theorem

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

- 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

Goal

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

Baselines

- DQN-L: Standard DQN¹ (no decomposition).
- HRL-E: Hierarchical Deep RL proposed by Kulkarni et al.²
- HRL-L: HRL-E but exploiting LTL to prune *useless* options.

¹Human-level control through deep reinforcement learning (Nature-15)

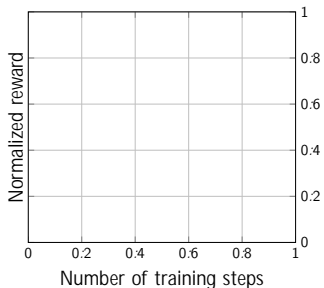
²HDRL: Integrating Temporal Abstraction and Intrinsic Motivation (NIPS-16)

Experiment 1: Interleaving subtasks

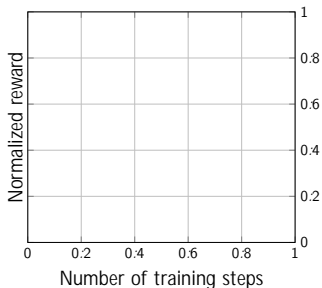
Description

Tasks from (Andreas et al., 2017) w/o unneeded order constraints.
e.g. (get iron and get wood), then use factory.

5 random maps



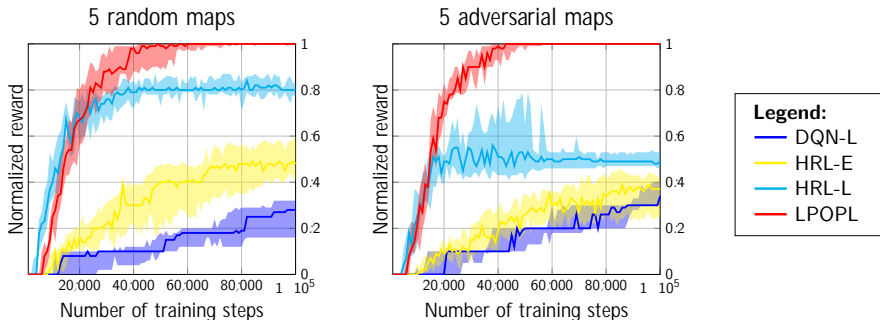
5 adversarial maps



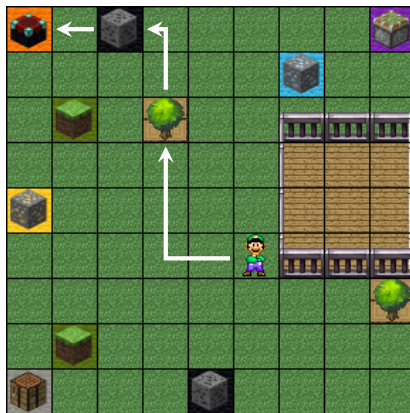
Experiment 1: Interleaving subtasks

Description

Tasks from (Andreas et al., 2017) w/o unneeded order constraints.
e.g. (get iron and get wood), then use factory.



Hierarchical RL might converge to suboptimal policies

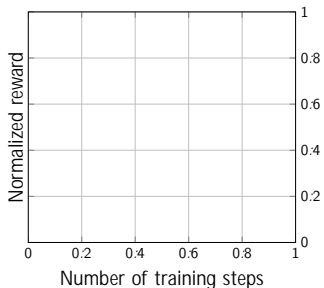


Experiment 2: Safety constraints

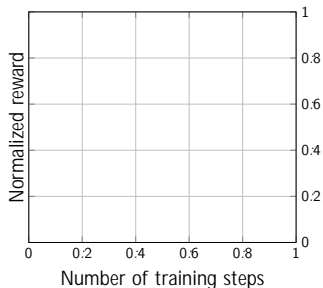
Description

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

5 random maps



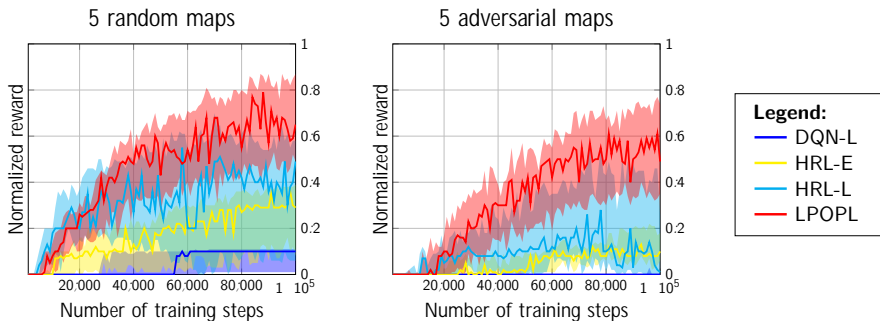
5 adversarial maps



Experiment 2: 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

Question: How can we instruct RL agents?

Concluding remarks

Question: How can we instruct RL agents?

Proposal: Define tasks using Linear Temporal Logic.

Concluding remarks

Question: How can we instruct RL agents?

Proposal: Define tasks using Linear Temporal Logic.

Desirable properties:

- It is expressive.

Concluding remarks

Question: How can we instruct RL agents?

Proposal: Define tasks using Linear Temporal Logic.

Desirable properties:

- It is expressive.
- RL agents understand it:
 - Task description / Reward function.
 - Task description / Learn faster.

Concluding remarks

Question: How can we instruct RL agents?

Proposal: Define tasks using Linear Temporal Logic.

Desirable properties:

- It is expressive.
- RL agents understand it:
 - Task description / Reward function.
 - Task description / Learn faster.
- Humans understand it.

Concluding remarks

Question: How can we instruct RL agents?

Proposal: Define tasks using Linear Temporal Logic.

Desirable properties:

- It is expressive.
- RL agents understand it:
 - Task description / Reward function.
 - Task description / Learn faster.
- Humans understand it.
 - NL / LTL by Dzifcak et al. (ICRA-09).

Concluding remarks

Question: How can we instruct RL agents?

Proposal: Define tasks using Linear Temporal Logic.

Desirable properties:

- It is expressive.
- RL agents understand it:
 - Task description / Reward function.
 - Task description / Learn faster.
- Humans understand it.
 - NL / LTL by Dzifcak et al. (ICRA-09).

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