

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}



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

How do you describe to an RL agent what to do?

Running example



Running example



Luigi can collect raw materials:



wood



grass



iron



gold



gems

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 do you describe to an RL agent what to do?

How do you describe to an RL agent what to do?

Description \neq Reward function

Motivation

How do you describe to an RL agent what to do?

Description \neq Reward function

Language \neq Reward function

How do you describe to an RL agent what to do?

Description \neq Reward function

Language \neq Reward function

Why would we want such a language?

How do you describe to an RL agent what to do?

Description \neq Reward function

Language \neq Reward 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	3
get wood and then use the factory	
get wood or iron	
get grass and iron	
do not leave the shelter at night	

Related works

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.

Related works

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		

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	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

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

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

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

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*):

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

Examples

```
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*):

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

Examples

eventually got_wood

eventually (got_grass and **eventually** used_factory)

eventually got_wood or **eventually** got_i ron

eventually got_grass and **eventually** got_i ron

(i_s_night! at_shel ter) **until** got_wood

Example



Detected events
at_shel ter

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

Example



Detected events

none

' = **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_gold

Example



Detected events

none

' = **eventually**(got_iron 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_gold

Example



Detected events
got_grass

' = **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_gold

Example



Detected events

none

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

Example



Detected events
used_factory

' = **eventually**(got_i_iron 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_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_gold

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 $p \in P$
- $\text{prog}(s; \neg \phi) = \neg \text{prog}(s; \phi)$
- $\text{prog}(s; \phi_1 \wedge \phi_2) = \text{prog}(s; \phi_1) \wedge \text{prog}(s; \phi_2)$
- $\text{prog}(s; \phi) = \phi$
- $\text{prog}(s; \neg \phi) = \neg \text{prog}(s; \phi)$
- $\text{prog}(s; \phi_1 \cup \phi_2) = \text{prog}(s; \phi_2) \cup (\text{prog}(s; \phi_1) \wedge \neg \phi_2)$

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_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
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 _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
got _grass

' = 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

' = true (+1 reward)

Example

Detected events
used_factory

LTL formulas! Rewards

Example

Detected events
used_factory

(' can be learned using standard RL)

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 solve₁

Detected events
at _shelter

- '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
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.

- ' ϕ_1 = eventually (got _iron and eventually used_factory)
and eventually got _gold
- ' ϕ_2 = eventually ([got _grass or got _wood]
and eventually used_factory)

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

- ' $\tau_1 = \text{eventually}(\text{got_iron} \ \text{and} \ \text{eventually} \ \text{used_factory} \)$
and eventually got _gold
- ' $\tau_2 = \text{eventually}([\text{got_grass} \ \text{or} \ \text{got_wood}]$
and eventually used_factory)
- ' $\tau_3 = \text{eventually}(\text{got_iron} \ \text{and} \ \text{eventually} \ \text{used_factory} \)$
- ' $\tau_4 = \text{eventually} \ \text{used_factory} \ \ \text{and} \ \text{eventually} \ \text{got_gold}$
- ' $\tau_5 = \text{eventually} \ \text{used_factory}$
- ' $\tau_6 = \text{eventually} \ \text{got_gold}$
- ' $\tau_7 = \text{true}$

LPOPL overview

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

'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

'₇ = true

LPOPL overview

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

```
' 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  
' 7 = true
```

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

LPOPL using tabular q-learning

Standard q-learning update given experience $(s, a, r; s^0)$:

$$Q(s; a) \leftarrow r + \max_{a^0} Q(s^0; a^0)$$

... where $x \leftarrow y$ is a shorthand notation for $x = x + (y - x)$.

LPOPL using tabular q-learning

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

$$Q(s; a) \leftarrow r + \max_{a'} Q(s', a')$$

... where $x \leftarrow y$ is a shorthand notation for $x = y + (y - x)$.

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

$$Q_t(s; a) \leftarrow r_t + \max_{a'} Q_{t-1}(s', a')$$

... where $s' = \text{prog}(s, a)$ and $r_t = 1$ if $s' \in \mathcal{G}$ else $r_t = 0$.

LPOPL using tabular q-learning

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

$$Q(s; a) \leftarrow r + \max_{a'} Q(s', a')$$

... where $x \leftarrow y$ is a shorthand notation for $x = y + (x - y)$.

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

$$Q_t(s; a) \leftarrow r_t + \max_{a'} Q_{t-1}(s', a')$$

... where $s' = \text{prog}(s, a)$ and $r_t = 1$ if $s' \in \mathcal{G}$.

Theorem

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

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 : at _shelter

Example

Subtasks

-
- ' $1 = (\text{got_iron} \wedge \text{used_factory}) \wedge \text{got_gold}$
 - ' $2 = ([\text{got_grass} \ _ \ \text{got_wood}] \wedge \text{used_factory})$
 - ' $3 = (\text{got_iron} \wedge \text{used_factory})$
 - ' $4 = \text{used_factory} \wedge \text{got_gold}$
 - ' $5 = \text{used_factory}$
 - ' $6 = \text{got_gold}$

Detected events: none

Q-updates

$Q_1(s; a)$	$\max_{a^0} Q_1(s^0; a^0)$	$Q_4(s; a)$	$\max_{a^0} Q_4(s^0; a^0)$
$Q_2(s; a)$	$\max_{a^0} Q_2(s^0; a^0)$	$Q_5(s; a)$	$\max_{a^0} Q_5(s^0; a^0)$
$Q_3(s; a)$	$\max_{a^0} Q_3(s^0; a^0)$	$Q_6(s; a)$	$\max_{a^0} Q_6(s^0; a^0)$

Example

Subtasks

-
- ' $1 = (\text{got_iron} \wedge \text{used_factory}) \wedge \text{got_gold}$
 - ' $2 = ([\text{got_grass} _ \text{got_wood}] \wedge \text{used_factory})$
 - ' $3 = (\text{got_iron} \wedge \text{used_factory})$
 - ' $4 = \text{used_factory} \wedge \text{got_gold}$
 - ' $5 = \text{used_factory}$
 - ' $6 = \text{got_gold}$

Detected events: none

Q-updates

$Q_1(s; a)$	$\max_{a^0} Q_1(s^0; a^0)$	$Q_4(s; a)$	$\max_{a^0} Q_4(s^0; a^0)$
$Q_2(s; a)$	$\max_{a^0} Q_2(s^0; a^0)$	$Q_5(s; a)$	$\max_{a^0} Q_5(s^0; a^0)$
$Q_3(s; a)$	$\max_{a^0} Q_3(s^0; a^0)$	$Q_6(s; a)$	$\max_{a^0} Q_6(s^0; a^0)$

Example

Subtasks

-
- ' $1 = (\text{got_iron} \wedge \text{used_factory}) \wedge \text{got_gold}$
 - ' $2 = ([\text{got_grass} _ \text{got_wood}] \wedge \text{used_factory})$
 - ' $3 = (\text{got_iron} \wedge \text{used_factory})$
 - ' $4 = \text{used_factory} \wedge \text{got_gold}$
 - ' $5 = \text{used_factory}$
 - ' $6 = \text{got_gold}$

Detected events: none

Q-updates

$Q_1(s; a)$	$\max_{a^0} Q_1(s^0; a^0)$	$Q_4(s; a)$	$\max_{a^0} Q_4(s^0; a^0)$
$Q_2(s; a)$	$\max_{a^0} Q_2(s^0; a^0)$	$Q_5(s; a)$	$\max_{a^0} Q_5(s^0; a^0)$
$Q_3(s; a)$	$\max_{a^0} Q_3(s^0; a^0)$	$Q_6(s; a)$	$\max_{a^0} Q_6(s^0; a^0)$

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_wood

Q-updates

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

Example

Subtasks

-
- ' $1 = (\text{got_iron} \wedge \text{used_factory}) \wedge \text{got_gold}$
 - ' $2 = ([\text{got_grass} _ \text{got_wood}] \wedge \text{used_factory})$
 - ' $3 = (\text{got_iron} \wedge \text{used_factory})$
 - ' $4 = \text{used_factory} \wedge \text{got_gold}$
 - ' $5 = \text{used_factory}$
 - ' $6 = \text{got_gold}$

Detected events: got_wood

Q-updates

$Q_1(s; a)$	$\max_{a^0} Q_1(s^0; a^0)$	$Q_4(s; a)$	$\max_{a^0} Q_4(s^0; a^0)$
$Q_2(s; a)$	$\max_{a^0} Q_5(s^0; a^0)$	$Q_5(s; a)$	$\max_{a^0} Q_5(s^0; a^0)$
$Q_3(s; a)$	$\max_{a^0} Q_3(s^0; a^0)$	$Q_6(s; a)$	$\max_{a^0} Q_6(s^0; a^0)$

Example

Subtasks

-
- ' $1 = (\text{got_iron} \wedge \text{used_factory}) \wedge \text{got_gold}$
 - ' $2 = ([\text{got_grass} _ \text{got_wood}] \wedge \text{used_factory})$
 - ' $3 = (\text{got_iron} \wedge \text{used_factory})$
 - ' $4 = \text{used_factory} \wedge \text{got_gold}$
 - ' $5 = \text{used_factory}$
 - ' $6 = \text{got_gold}$

Detected events: none

Q-updates

$Q_1(s; a)$	$\max_{a^0} Q_1(s^0; a^0)$	$Q_4(s; a)$	$\max_{a^0} Q_4(s^0; a^0)$
$Q_2(s; a)$	$\max_{a^0} Q_2(s^0; a^0)$	$Q_5(s; a)$	$\max_{a^0} Q_5(s^0; a^0)$
$Q_3(s; a)$	$\max_{a^0} Q_3(s^0; a^0)$	$Q_6(s; a)$	$\max_{a^0} Q_6(s^0; a^0)$

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_grass

Q-updates

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

Example

Subtasks

- ' $1 = (\text{got_iron} \wedge \text{used_factory}) \wedge \text{got_gold}$
- ' $2 = ([\text{got_grass} _ \text{got_wood}] \wedge \text{used_factory})$
- ' $3 = (\text{got_iron} \wedge \text{used_factory})$
- ' $4 = \text{used_factory} \wedge \text{got_gold}$
- ' $5 = \text{used_factory}$
- ' $6 = \text{got_gold}$

Detected events: got_grass

Q-updates

$Q_1(s; a)$	$\max_{a^0} Q_1(s^0; a^0)$	$Q_4(s; a)$	$\max_{a^0} Q_4(s^0; a^0)$
$Q_2(s; a)$	$\max_{a^0} Q_5(s^0; a^0)$	$Q_5(s; a)$	$\max_{a^0} Q_5(s^0; a^0)$
$Q_3(s; a)$	$\max_{a^0} Q_3(s^0; a^0)$	$Q_6(s; a)$	$\max_{a^0} Q_6(s^0; a^0)$

Example

Subtasks

- ' $1 = (\text{got_iron} \wedge \text{used_factory}) \wedge \text{got_gold}$
- ' $2 = ([\text{got_grass} _ \text{got_wood}] \wedge \text{used_factory})$
- ' $3 = (\text{got_iron} \wedge \text{used_factory})$
- ' $4 = \text{used_factory} \wedge \text{got_gold}$
- ' $5 = \text{used_factory}$
- ' $6 = \text{got_gold}$

Detected events: none

Q-updates

$Q_1(s; a)$	$\max_{a^0} Q_1(s^0; a^0)$	$Q_4(s; a)$	$\max_{a^0} Q_4(s^0; a^0)$
$Q_2(s; a)$	$\max_{a^0} Q_2(s^0; a^0)$	$Q_5(s; a)$	$\max_{a^0} Q_5(s^0; a^0)$
$Q_3(s; a)$	$\max_{a^0} Q_3(s^0; a^0)$	$Q_6(s; a)$	$\max_{a^0} Q_6(s^0; a^0)$

Example

Subtasks

-
- ' $1 = (\text{got_iron} \wedge \text{used_factory}) \wedge \text{got_gold}$
 - ' $2 = ([\text{got_grass} _ \text{got_wood}] \wedge \text{used_factory})$
 - ' $3 = (\text{got_iron} \wedge \text{used_factory})$
 - ' $4 = \text{used_factory} \wedge \text{got_gold}$
 - ' $5 = \text{used_factory}$
 - ' $6 = \text{got_gold}$

Detected events: none

Q-updates

$Q_1(s; a)$	$\max_{a^0} Q_1(s^0; a^0)$	$Q_4(s; a)$	$\max_{a^0} Q_4(s^0; a^0)$
$Q_2(s; a)$	$\max_{a^0} Q_2(s^0; a^0)$	$Q_5(s; a)$	$\max_{a^0} Q_5(s^0; a^0)$
$Q_3(s; a)$	$\max_{a^0} Q_3(s^0; a^0)$	$Q_6(s; a)$	$\max_{a^0} Q_6(s^0; a^0)$

Example

Subtasks

-
- ' $1 = (\text{got_iron} \wedge \text{used_factory}) \wedge \text{got_gold}$
 - ' $2 = ([\text{got_grass} \ _ \ \text{got_wood}] \wedge \text{used_factory})$
 - ' $3 = (\text{got_iron} \wedge \text{used_factory})$
 - ' $4 = \text{used_factory} \wedge \text{got_gold}$
 - ' $5 = \text{used_factory}$
 - ' $6 = \text{got_gold}$

Detected events: used_factory

Q-updates

$Q_1(s; a)$	$\max_{a^0} Q_1(s^0; a^0)$	$Q_4(s; a)$	$\max_{a^0} Q_6(s^0; a^0)$
$Q_2(s; a)$	$\max_{a^0} Q_2(s^0; a^0)$	$Q_5(s; a)$	1
$Q_3(s; a)$	$\max_{a^0} Q_3(s^0; a^0)$	$Q_6(s; a)$	$\max_{a^0} Q_6(s^0; a^0)$

Example

Subtasks

-
- ' $1 = (\text{got_iron} \wedge \text{used_factory}) \wedge \text{got_gold}$
 - ' $2 = ([\text{got_grass} \ _ \ \text{got_wood}] \wedge \text{used_factory})$
 - ' $3 = (\text{got_iron} \wedge \text{used_factory})$
 - ' $4 = \text{used_factory} \wedge \text{got_gold}$
 - ' $5 = \text{used_factory}$
 - ' $6 = \text{got_gold}$

Detected events: used_factory

Q-updates

$Q_1(s; a)$	$\max_{a^0} Q_1(s^0; a^0)$	$Q_4(s; a)$	$\max_{a^0} Q_6(s^0; a^0)$
$Q_2(s; a)$	$\max_{a^0} Q_2(s^0; a^0)$	$Q_5(s; a)$	1
$Q_3(s; a)$	$\max_{a^0} Q_3(s^0; a^0)$	$Q_6(s; a)$	$\max_{a^0} Q_6(s^0; a^0)$

Example

Subtasks

-
- ' $1 = (\text{got_iron} \wedge \text{used_factory}) \wedge \text{got_gold}$
 - ' $2 = ([\text{got_grass} \ _ \ \text{got_wood}] \wedge \text{used_factory})$
 - ' $3 = (\text{got_iron} \wedge \text{used_factory})$
 - ' $4 = \text{used_factory} \wedge \text{got_gold}$
 - ' $5 = \text{used_factory}$
 - ' $6 = \text{got_gold}$

Detected events: used_factory

Q-updates

$Q_1(s; a)$	$\max_{a^0} Q_1(s^0; a^0)$	$Q_4(s; a)$	$\max_{a^0} Q_6(s^0; a^0)$
$Q_2(s; a)$	$\max_{a^0} Q_2(s^0; a^0)$	$Q_5(s; a)$	1
$Q_3(s; a)$	$\max_{a^0} Q_3(s^0; a^0)$	$Q_6(s; a)$	$\max_{a^0} Q_6(s^0; a^0)$

Example

Subtasks

- ' $1 = (\text{got_iron} \wedge \text{used_factory}) \wedge \text{got_gold}$
- ' $2 = ([\text{got_grass} \ _ \ \text{got_wood}] \wedge \text{used_factory})$
- ' $3 = (\text{got_iron} \wedge \text{used_factory})$
- ' $4 = \text{used_factory} \wedge \text{got_gold}$
- ' $5 = \text{used_factory}$
- ' $6 = \text{got_gold}$

Detected events: none

Q-updates

$Q_1(s; a)$	$\max_{a^0} Q_1(s^0; a^0)$	$Q_4(s; a)$	$\max_{a^0} Q_4(s^0; a^0)$
$Q_2(s; a)$	$\max_{a^0} Q_2(s^0; a^0)$	$Q_5(s; a)$	$\max_{a^0} Q_5(s^0; a^0)$
$Q_3(s; a)$	$\max_{a^0} Q_3(s^0; a^0)$	$Q_6(s; a)$	$\max_{a^0} Q_6(s^0; a^0)$

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_iron

Q-updates

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

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_iron

Q-updates

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

Example

Subtasks

-
- ' $1 = (\text{got_iron} \wedge \text{used_factory}) \wedge \text{got_gold}$
 - ' $2 = ([\text{got_grass} \ _ \ \text{got_wood}] \wedge \text{used_factory})$
 - ' $3 = (\text{got_iron} \wedge \text{used_factory})$
 - ' $4 = \text{used_factory} \wedge \text{got_gold}$
 - ' $5 = \text{used_factory}$
 - ' $6 = \text{got_gold}$

Detected events: got_iron

Q-updates

$Q_1(s; a)$	$\max_{a^0} Q_4(s^0; a^0)$	$Q_4(s; a)$	$\max_{a^0} Q_4(s^0; a^0)$
$Q_2(s; a)$	$\max_{a^0} Q_2(s^0; a^0)$	$Q_5(s; a)$	$\max_{a^0} Q_5(s^0; a^0)$
$Q_3(s; a)$	$\max_{a^0} Q_5(s^0; a^0)$	$Q_6(s; a)$	$\max_{a^0} Q_6(s^0; a^0)$

Example

Subtasks

- ' $1 = (\text{got_iron} \wedge \text{used_factory}) \wedge \text{got_gold}$
- ' $2 = ([\text{got_grass} \ _ \ \text{got_wood}] \wedge \text{used_factory})$
- ' $3 = (\text{got_iron} \wedge \text{used_factory})$
- ' $4 = \text{used_factory} \wedge \text{got_gold}$
- ' $5 = \text{used_factory}$
- ' $6 = \text{got_gold}$

Detected events: got_iron

Q-updates

$Q_1(s; a)$	$\max_{a^0} Q_4(s^0; a^0)$	$Q_4(s; a)$	$\max_{a^0} Q_4(s^0; a^0)$
$Q_2(s; a)$	$\max_{a^0} Q_2(s^0; a^0)$	$Q_5(s; a)$	$\max_{a^0} Q_5(s^0; a^0)$
$Q_3(s; a)$	$\max_{a^0} Q_5(s^0; a^0)$	$Q_6(s; a)$	$\max_{a^0} Q_6(s^0; a^0)$

Example

Subtasks

-
- ' $1 = (\text{got_iron} \wedge \text{used_factory}) \wedge \text{got_gold}$
 - ' $2 = ([\text{got_grass} _ \text{got_wood}] \wedge \text{used_factory})$
 - ' $3 = (\text{got_iron} \wedge \text{used_factory})$
 - ' $4 = \text{used_factory} \wedge \text{got_gold}$
 - ' $5 = \text{used_factory}$
 - ' $6 = \text{got_gold}$

Detected events: none

Q-updates

$Q_1(s; a)$	$\max_{a^0} Q_1(s^0; a^0)$	$Q_4(s; a)$	$\max_{a^0} Q_4(s^0; a^0)$
$Q_2(s; a)$	$\max_{a^0} Q_2(s^0; a^0)$	$Q_5(s; a)$	$\max_{a^0} Q_5(s^0; a^0)$
$Q_3(s; a)$	$\max_{a^0} Q_3(s^0; a^0)$	$Q_6(s; a)$	$\max_{a^0} Q_6(s^0; a^0)$

- 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

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

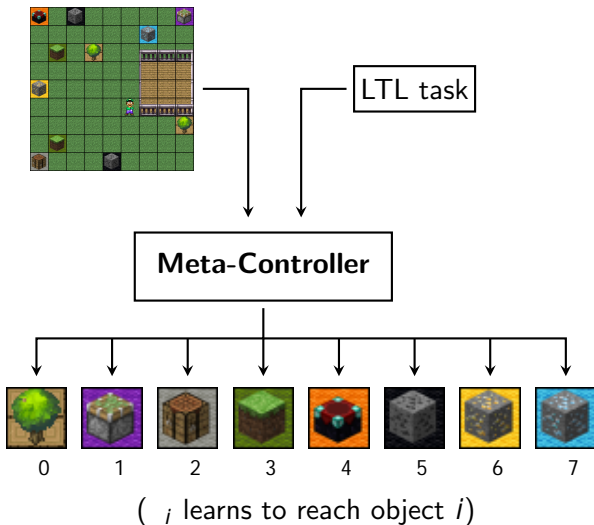
Baselines

- 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.

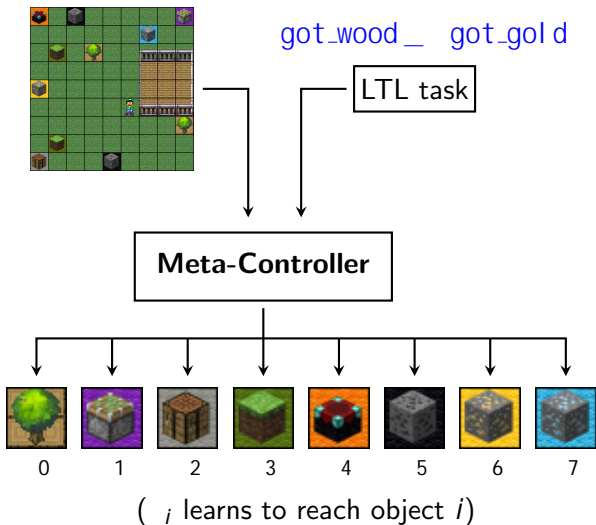
¹Human-level control through deep reinforcement learning (Nature-15)

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

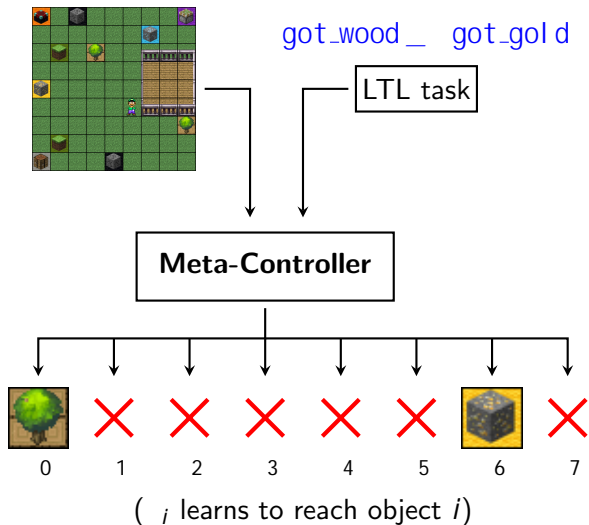
Baselines: HRL-E



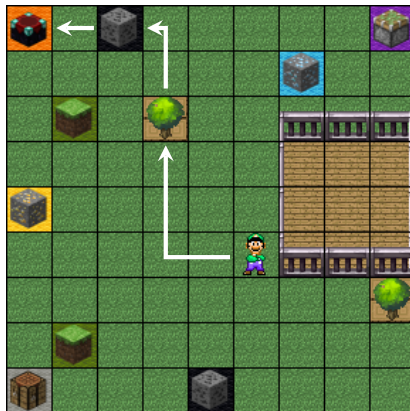
Baselines: HRL-L



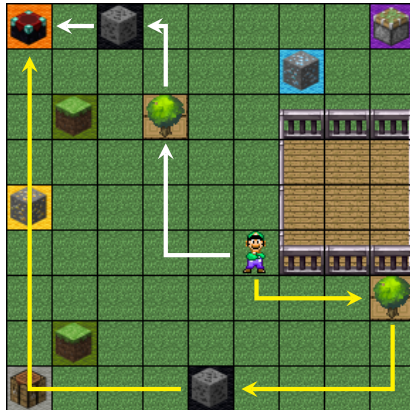
Baselines: HRL-L



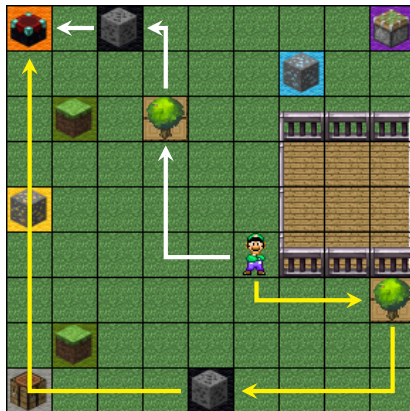
Hierarchical RL might converge to suboptimal policies



Hierarchical RL might converge to suboptimal policies



Hierarchical RL might converge to suboptimal policies



We tested over 5 random grids and 5 adversarial grids.

Experiment 1: subgoal sequences

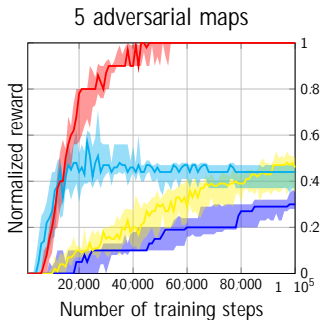
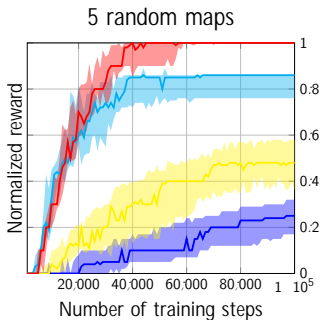
Description

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

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.



Legend:

- DQN-L
- HRL-E
- HRL-L
- LPOPL

Experiment 2: Interleaving subtasks

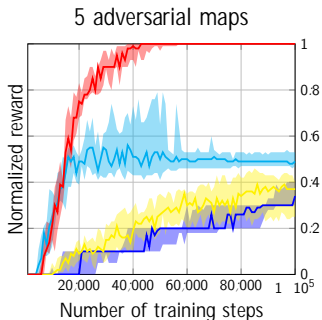
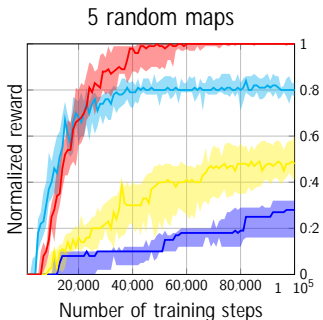
Description

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

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.



Legend:

- DQN-L
- HRL-E
- HRL-L
- LPOPL

Experiment 3: Safety constraints

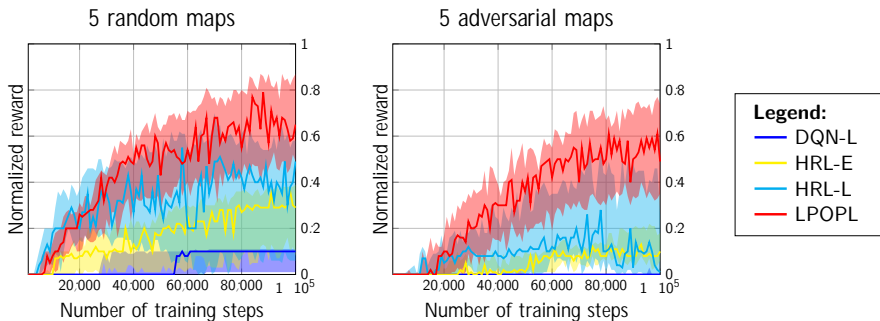
Description

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!