Leveraging Formal Task Structure in Reinforcement Learning Ameesh Shah, UC Berkeley

Simons Institute, 5.01

Reinforcement learning for autonomous behavior



Agent

Actions

Environment

Consequences: Observations **Rewards**

Reward design for reinforcement learning is hard

- The success of Reinforcement Learning (RL) relies on designing good *reward functions*
 - Rewards should incentivize desirable behavior from our agents and discourage undesirable behavior
 - However, "reward engineering" is...
 - Tedious at best
 - Difficult to impossible at worst
- The key question: how do we represent desirable behavior?



A classic example of a reward function leading to undesirable behavior

How do we represent desirable behavior?

- of contexts
- To verify computational systems, FM has developed representations for φ for short)



Formal methods has been asking (and answering) this question in a number

desirable properties in the form of formal task specifications (formal specs or



Can we use formal specifications in RL?

- A natural question emerges:
 - Can we use formal specifications as objectives in RL?
 - Yes! (Lots of prior work)
 - The more salient question: When should we use formal specifications as objectives in RL?

Pros and Cons of using φ **in RL** Pros



★ Provides compositionality (can break down and build up separate specs)

★ Easy to capture non-Markovian semantics e.g. temporally ordered tasks

Cons

- Difficult to optimize and provides sparse feedback
- Solution space for specs is non-continuous and difficult to search
- Difficult to construct by hand

Pros and Cons of using φ **in RL**

Pros

★ Precise notion of satisfaction

★ Provides compositionality (can break down and build up separate specs)

★ Easy to capture non-Markovian semantics e.g. temporally ordered tasks

Cons

- Difficult to optimize and provides sparse feedback
- Solution space for specs is non-continuous and difficult to search
- Difficult to construct by hand

Using ϕ in RL: Differentiating "correct" and "optimal"

- In correctness-critical settings, we want to learn policies that certainly achieve the task at hand
- However, we may have **optimality** conditions alongside our task
- Example: A patrolling quadruped



Correctness: Patrol the three regions consistently!

Optimality: Navigate the terrain safely (avoid collisions!)



Defining "Correct and Optimal" using φ and r

- It's difficult to express both the correctness condition *and* the optimality condition in a single expression.
- However, the precision and Boolean satisfaction of formal specifications allows us to define "correct" behavior
- And the expressivity of Markovian reward functions allow us to define "optimality" within correct behaviors.
- Combining the two allows us to ask: which "correct" policy is the most "optimal"?



Demonstrating Example: Correctness and Optimality

Task φ : Patrol red, green, and yellow, all while avoiding blue.

Optimality condition *r***: +1** for each visit to a purple circle.

I *must* accomplish my task. Amongst all ways of accomplishing my task, choose the optimal one.



"FlatWorld"

How do we specify φ and r in our setting?

Task specification φ

Buchi Automata (Linear Temporal Logic)



Satisfaction condition: visit accepting states (green) infinitely often.

Optimality Condition *r*

Use the standard definition of MDP reward:

$$r(s, a, s') = \begin{cases} +1 & \text{if s' in purple} \\ 0 & \text{otherwise} \end{cases}$$

At each transition, receive a scalar reward.



Approach Outline

1. How do I combine φ with r in a form that is readily optimized by (deep) RL?

2. Does the objective from (1) work in practice? If not, what can we do?



1. Compile ϕ down to a proxy reward function that, when optimized, closely approximates satisfaction of the original property.



Step 1: Create a proxy reward for φ

- acceptance condition.
- Instead, we use a proxy reward:

 ${}^{r}\mathsf{LTL}(b_t) = \begin{cases} 1 & \text{if } (b_t \in \mathcal{B}^*) \\ 0 & \text{otherwise} \end{cases}$

get a reward every time I visit the accepting state!

There is no direct way to enforce satisfaction of LTL in RL due to the infinite





- approximates satisfaction of the original property.
- 2. Add the φ -proxy reward and the optimality reward (MDP reward)

$$r_{\text{DUAL}}(s, b, a, s', b') = r_{\text{MDP}}(s, a, s') + \lambda r_{\text{LTL}}(b')$$

For large enough λ , the policy that optimizes this reward will satisfy ϕ with maximum probability AND achieve optimal reward.

1. Compile φ down to a proxy reward function that, when optimized, closely

together, and multiply the φ -proxy reward by a Lagrange multiplier λ :



- 1. Compile φ down to a proxy reward function that, when optimized, closely approximates satisfaction of the original property.
- 2. Add the φ -proxy reward and the optimality reward (MDP reward) together, and multiply the φ -proxy reward by a Lagrange multiplier λ : ∇
 - $r_{\text{DUAL}}(s, b, a, s', b') = r_{\text{MDP}}(s, a, s') + \lambda r_{\text{LT}}(b')$
- 3. Learn a policy that optimizes the combined reward above.

We can just use our favorite RL algorithm for this!



- approximates satisfaction of the original property.
- and multiply the φ -proxy reward by a Lagrange multiplier λ :

$$r_{\text{DUAL}}(s, b, a, s', b') =$$

- 3. Learn a policy that optimizes the combined reward above. ∇
- 4. Profit?

1. Compile φ down to a proxy reward function that, when optimized, closely

2. Add the φ -proxy reward and the optimality reward (MDP reward) together,

 $r_{\text{MDP}}(s, a, s') + \lambda r_{\text{ITI}}(b')$



Approach Outline

1. How do I combine φ with r in a form that is readily optimized by (deep) RL?

2. Does the objective from (1) work in practice? If not, what can we do?



Getting LTL-Constrained Deep RL to work in practice

Task φ : Patrol red, green, and yellow, all while avoiding blue.

Optimality condition *r*: +1 for each visit to a purple circle.

$$r_{\mathsf{DUAL}}(s, b, a, s', b') = r_{\mathsf{MDP}}(s, a, s')$$

Let's train a policy on this reward and see what happens in this example.

 $+ \lambda r_{ITI}(b')$







Getting LTL-Constrained Deep RL to work in practice

Why does this happen?

$r_{|\mathsf{T}|}$ is sparse.

if *r*_{MDP} is **dense**, a policy can get 'distracted' during learning, and ignore r_{ITI} .







Addressing the sparsity issue of reward proxies



Can we shape $r_{|T|}$ to make it more dense?



Addressing the sparsity issue of reward proxies Key idea: every visit to an accepting state must traverse an "accepting cycle" within our BA.





An accepting trajectory for our example.



Addressing the sparsity issue of reward proxies So, why not give a reward every time we progress along an accepting cycle?





Reward Shaping with Cycle Experience Replay (CyclER)

- Counterfactually reason over all possible accepting cycles at the end of each trajectory to find the accepting cycle that provides the most dense reward
- Retroactively update the rewards from this trajectory based on the most dense accepting cycle
- Notably: we can incorporate Quantitative Semantics (QS) into CyclER and densify reward even further!

LTL-Constrained Policy Optimization With Cycle Experience Replay. Shah, Voloshin, Yang, Verma, Chaudhuri, Seshia. TMLR, 2025.

hah

CyclER: Example





t=0

CyclER: Experimental Evaluation

- Evaluated in environments with both discrete and continuous state/action spaces
- Tasks: Navigation tasks of varying complexity
- Baselines: Existing reward shaping methods for LTL-guided RL



Example MuJoCo Safety-Gymnasium Environments



CyclER: Experimental Evaluation

	FlatWorld		ZonesEnv		ButtonsEnv	
	$r_{ m LTL}$	$r_{ m MDP}$	$r_{ m LTL}$	$r_{ m MDP}$	$r_{ m LTL}$	$r_{ m MDP}$
CyclER	2.0 ± 0.5	45.3 ± 8.5	1.8 ± 0.4	-27.8 ± 4.55	2.6 ± 0.3	30.4 ± 5.6
LCER	0.0 ± 0.0	103.4 ± 76.6	0.0 ± 0.0	-3.8 ± 1.9	0.0 ± 0.0	118.8 ± 143.2
LCER, no r_{MDP}	0.8 ± 0.4	30.8 ± 10.1	0.0 ± 0.0	-2.7 ± 0.9	0.6 ± 0.4	13.2 ± 8.53
TLTL	-	_	0.0 ± 0.0	-4.0 ± 2.2	0.0 ± 0.0	9.6 ± 6.7
BHNR	-	_	0.0 ± 0.0	-0.8 ± 0.6	0.0 ± 0.0	35.8 ± 7.9

CyclER succeeds where other baselines cannot accomplish the LTL task!

Even a baseline (blue) that completely ignores MDP reward cannot accomplish the LTL task, which means that sparsity is an issue even in the absence of $r_{\rm MDP}$.



Pros and Cons of using φ **in RL** Pros



Provides compositionality (can break down and build up separate specs)

★ Easy to capture non-Markovian semantics e.g. temporally ordered tasks

Cons

- Difficult to optimize and provides sparse feedback
- Solution space for specs is non-continuous and difficult to search
- Difficult to construct by hand

The "Group Project" Problem in Multi-Agent RL





When a team of agents share reward in a cooperative setting, issues may arise:

Teams. Shah, Lauffer, Chen, Pitta, Seshia. AAMAS, 2025.





We can verify that a decomposition is valid: completing each sub-task will complete the overall task.

eceive separate rewards for accomplishing their own tasks

Learning Symbolic Task Decompositions for Multi-Agent Teams. Shah, Lauffer, Chen, Pitta, Seshia. AAMAS, 2025.







Problem: many possible valid decompositions exist





Summary of our problem

- Decomposing symbolic task representations enables assigning sub-tasks to individual agents with the guarantee that completing all sub-tasks will complete the overall task
 - This allows for learning more efficient strategies for completing cooperative tasks.
- The problem: many valid decompositions may exist. We don't know which one solves the task most efficiently (optimally)
 - because we don't have a-priori knowledge of the environment.



Our work: Learning Task Decompositions On-The-Fly

As we collect experience, use it to inform which decomposition is the best!





Overview of our procedure Learning Optimal Task Decompositions (LOTaD)

- 1. Generate a set of *candidate* decompositions.
- 2. Use a task-conditioned policy to accelerate learning across different candidates.
- 3. Use the Upper Confidence Bound (UCB) algorithm to balance exploring and exploiting different candidates.





Experiments: Evaluating LOTaD

- **Baselines**:
 - ATAD: Prior work on generating task decompositions that uses heuristics to select and fix a single decomposition for use during MARL
 - Monolithic: All agents are given the original symbolic task as their task
- **Example environments:**





Learning Symbolic Task Decompositions for Multi-Agent Teams. Shah, Lauffer, Chen, Pitta, Seshia. AAMAS, 2025.



Experiments: Evaluating LOTaD





LOTaD strongly outperforms existing baselines!







Conclusions and Next Steps

- Formal task specifications, such as automata, can act as a precise and composable means of identifying objectives in RL
- When should you use formal specs?
 - When you have a precise satisfaction condition in correctness-critical settings
 - When your task is a *composition* of separable sub-tasks
- Next steps: exploring how formal structure can augment foundation model capabilities



Thanks!

- Papers discussed today:
 - LTL-Constrained Policy Optimization With Cycle Experience Replay. Shah, Voloshin, Yang, Verma, Chaudhuri, Seshia. TMLR, 2025.
 - Learning Symbolic Task Decompositions for Multi-Agent Teams. Shah, Lauffer, Chen, Pitta, Seshia. AAMAS, 2025.
- reach out to me: <u>ameesh@berkeley.edu</u>
- Thanks to my collaborators!
 - Berkeley: Niklas Lauffer, Thomas Chen, Nikhil Pitta, Sanjit A. Seshia
 - Non-Berkeley: Cameron Voloshin, Chenxi Yang, Abhinav Verma, Swarat Chaudhuri

r,