



McGill
UNIVERSITY



Blog on HRL

Hierarchical Reinforcement Learning

Ray Luo

WITH SPECIAL THANKS TO

*Doina Precup, Xujie Si, Martin Klissarov, Akhil Bagaria, George Konidakis, Marlos C. Machado, Tianwei Ni,
Xutong Zhao, Nishanth Anand Vermgal, Zijing Wu, Zihan Wang, Yivan Zhang, Isabeau Premont-schwarz,
David Abel, Padideh Nouri, Jonathan Concalo Carr*

What is “Hierarchical” RL?

- An arrangement of **items** that are represented as being “above”, “below”, or “at the same **level** as” one another.
(Wikipedia)
- What are **items** in RL?
 - **Actions/policy** (control; [temporal](#)):
 - managers propose high-level goals, workers take concrete low-level acts
 - States (perception; spacial):
 - manager knows less details than workers
- How to decide “above” / “below”?
 - Determined by **abstraction levels**
 - The amount of details ignored



Blog on HRL

HIGH-LEVEL GOAL: Gather Nuts for Winter



SUB-GOAL 1:

Climb Tree



Grip Bark



Grip Bark

Move Branch

Ascend

SUB-GOAL 2:

Collect Acorns



Identify Nut



Move Paws

Pick Up Nut

Put in Pouch

SUB-GOAL 3:

Store in Drey



Find Drey



Find Drey

Enter Nut

Deposit Nut

Temporal abstractions (option/skill/subgoal)



Desiderata of HRL

- HRL exploits **temporal structure** that organize **knowledge** across time
 - **Procedural** knowledge (policies, skills / **goal-driven behavior**)
 - **Predictive** knowledge (value functions / models)

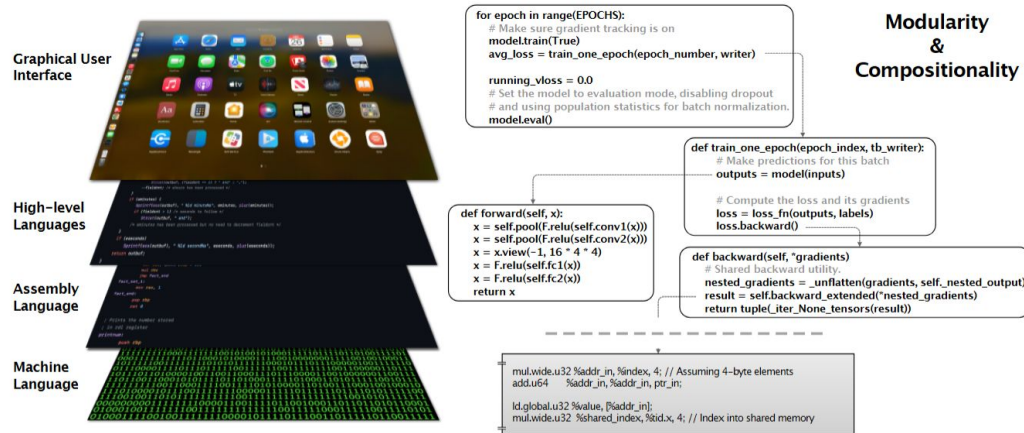


Building knowledge within RL

HRL exploits **temporal structure** that organize **knowledge** across time

Knowledge should be:

- **Expressive**: can represent many things (objects, goals, high-level strategies)
- **Learnable**: from streaming experience, ideally without supervision
- **Composable**: fast planning with assembling existing pieces
- **Abstract**: managing complexity and omit useless details



Summary: **Why** HRL?

- Intelligence (animals) reasons and acts across **extended timescales**.
- Growing knowledge in **open-ended environments**.
- May manage those knowledge with **temporal abstractions!**
 - **Standard RL**: states -> primitive actions
 - Making a decision at **every timestep**
 - **Hierarchical RL**: formalizing the idea of flexibly reasoning over different timescales at multiple levels of **temporal abstraction**



How and When does temporal abstraction help in RL?

And...

What behaviors/skills/options/structures/abstractions do we desire?

What can good temporal abstractions achieve?

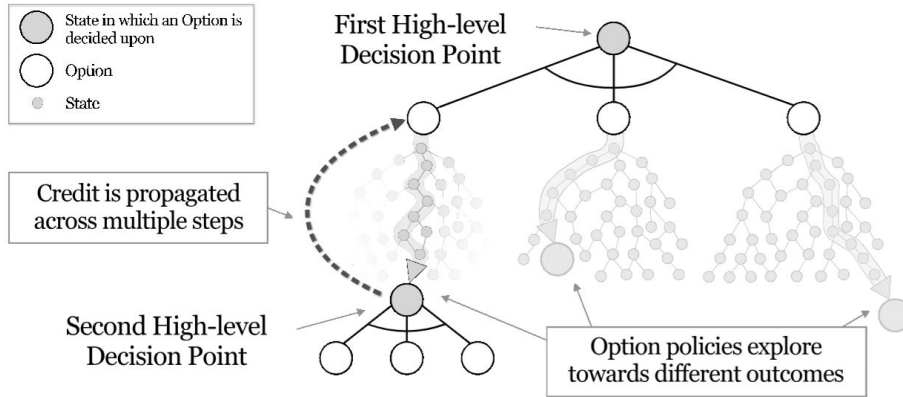
The dual question: what are the **benefits** of useful temporal abstractions?

With good/useful abstraction, we can achieve better:

- Exploration
 - Credit Assignment
 - Transferability (Generalizability)
 - Interpretability
- } Faster, sample efficient learning and planning

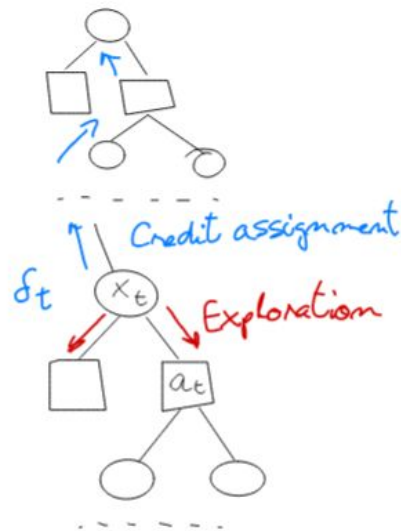
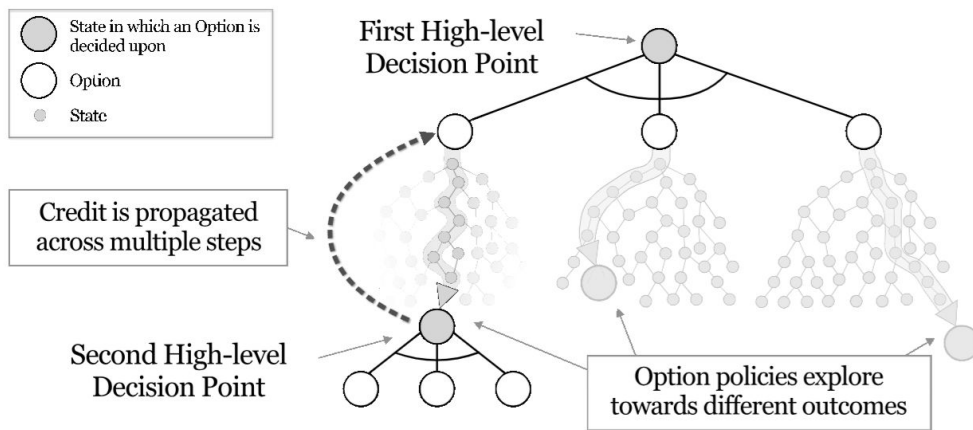
Structured Exploration

- Explore for pursuing subgoals, in different meaningful directions
- Explore within more abstract state/action spaces
- Explore over extended time, not just single action



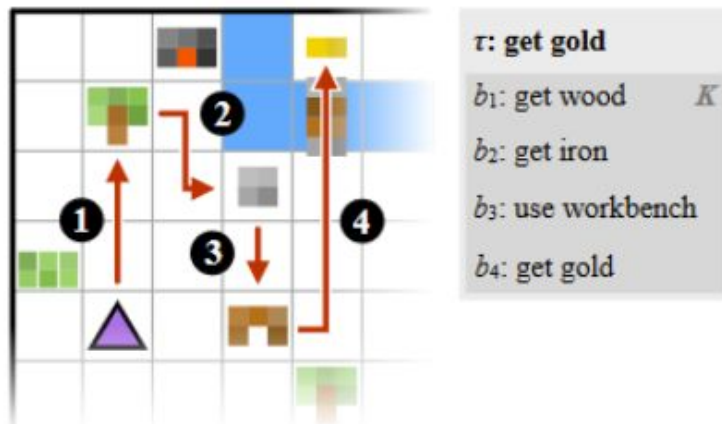
Effective Credit Assignment

- *CA problem*: determining which actions (or decisions) in a trajectory are responsible for later rewards
 - Allow learning signals (TD-errors/gradients) to be assigned at higher levels of abstraction
 - Identify “critical” decisions (bottlenecks, subgoals)



Transferability and Generalizability

- Options can be **reused** across different tasks
- Options can **generalize** in similar but novel conditions
 - Similar tasks (goals) require similar/same options to complete



Interpretability

- HRL provides an **interface** via a more abstract formulation
- More abstract representations are (sometimes) easier for humans to understand
 - As the system is smaller
 - Mind over-abstraction!

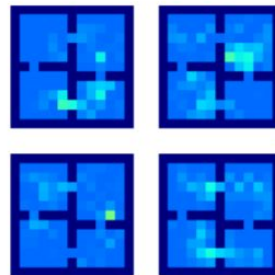
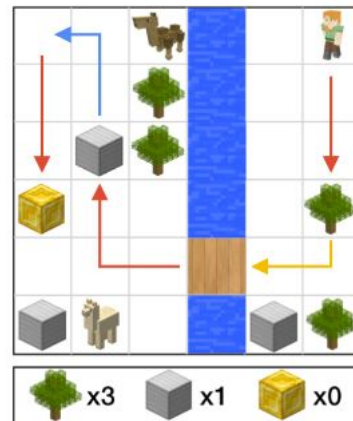


Figure 3: Termination probabilities for the option-critic agent learning with 4 options. The darkest color represents the walls in the environment while lighter colors encode higher termination probabilities.

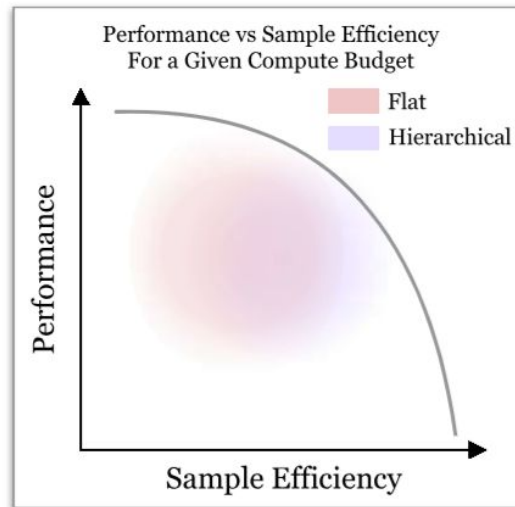
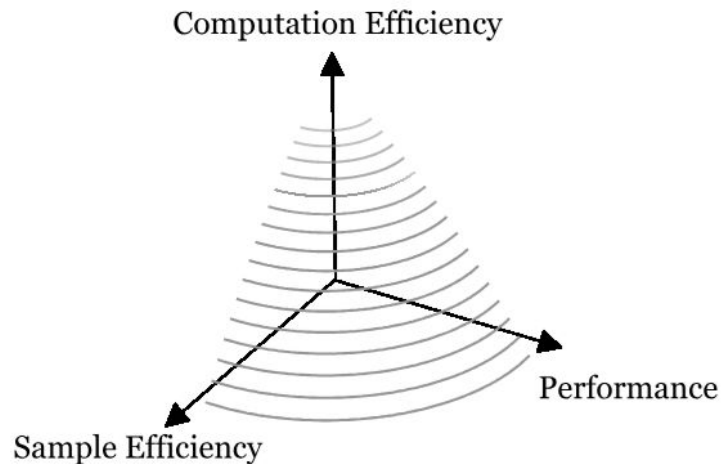
Program

```
def run():  
    if is_there[River]:  
        mine(Wood)  
        build_bridge()  
        if agent[Iron]<3:  
            mine(Iron)  
            place(Iron, 1, 1)  
        else:  
            goto(4, 2)  
        while env[Gold]>0:  
            mine(Gold)
```



HRL aims for desirable trade-offs

- No free lunch! (Wolpert and Macready, 1997)
- Flat RL can also achieve optimal policies
 - But it may take forever!
- HRL agents face a trade-off between **performance, sample efficiency, and computation efficiency**
- Another hidden axis: Sometimes, some **existing knowledge** is easy to get (offline dataset, foundation models, ...)





Formalism

- ❑ Option, SMDP
- ❑ Alternative Terminologies
- ❑ Option Discovery Problem
- ❑ Scope of HRL

Formalize the term “knowledge” and “temporal abstraction/structure”

Option framework,

and friends

$$\pi : \mathcal{S} \times \mathcal{O} \rightarrow \Delta(\mathcal{A}) \quad \text{Intra-option policy}$$

defines how the agent acts given state s under option o

$$\beta : \mathcal{S} \times \mathcal{O} \rightarrow [0, 1] \quad \text{Termination function}$$

probability of ending option o upon reaching state s

$$\mathcal{I} : \mathcal{S} \times \mathcal{O} \rightarrow [0, 1] \quad \text{Initiation function (set)}$$

indicates where an option o can start, “affordances”

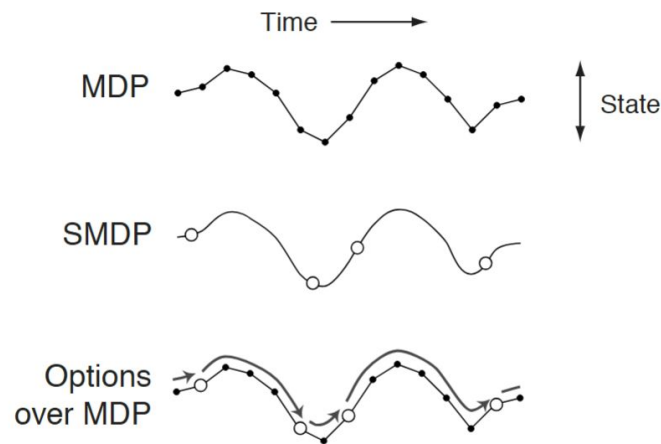
$$r^o : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$$

Option reward function

$$\mu : \mathcal{S} \rightarrow \Delta(\mathcal{O} \cup \mathcal{A})$$

High-level policy

Decision-making with Options



MDP vs. Semi-MDP (Sutton, et al., 1998)

Per-timestep vs. Multiple-timestep decision

$$P(s' | s, a) \quad \text{vs.} \quad P(s', \tau | s, o)$$

$$R(s, a) \quad \text{vs.} \quad R(s, a, \tau)$$

r.v. of # timestep
induced by an option