

HIERARCHIES OF REWARD MACHINES

Daniel Furelos-Blanco¹, Mark Law², Anders Jonsson³, Krysia Broda¹, and Alessandra Russo¹

¹Imperial College London, UK ²ILASP Limited, UK ³Universitat Pompeu Fabra, Barcelona, Spain

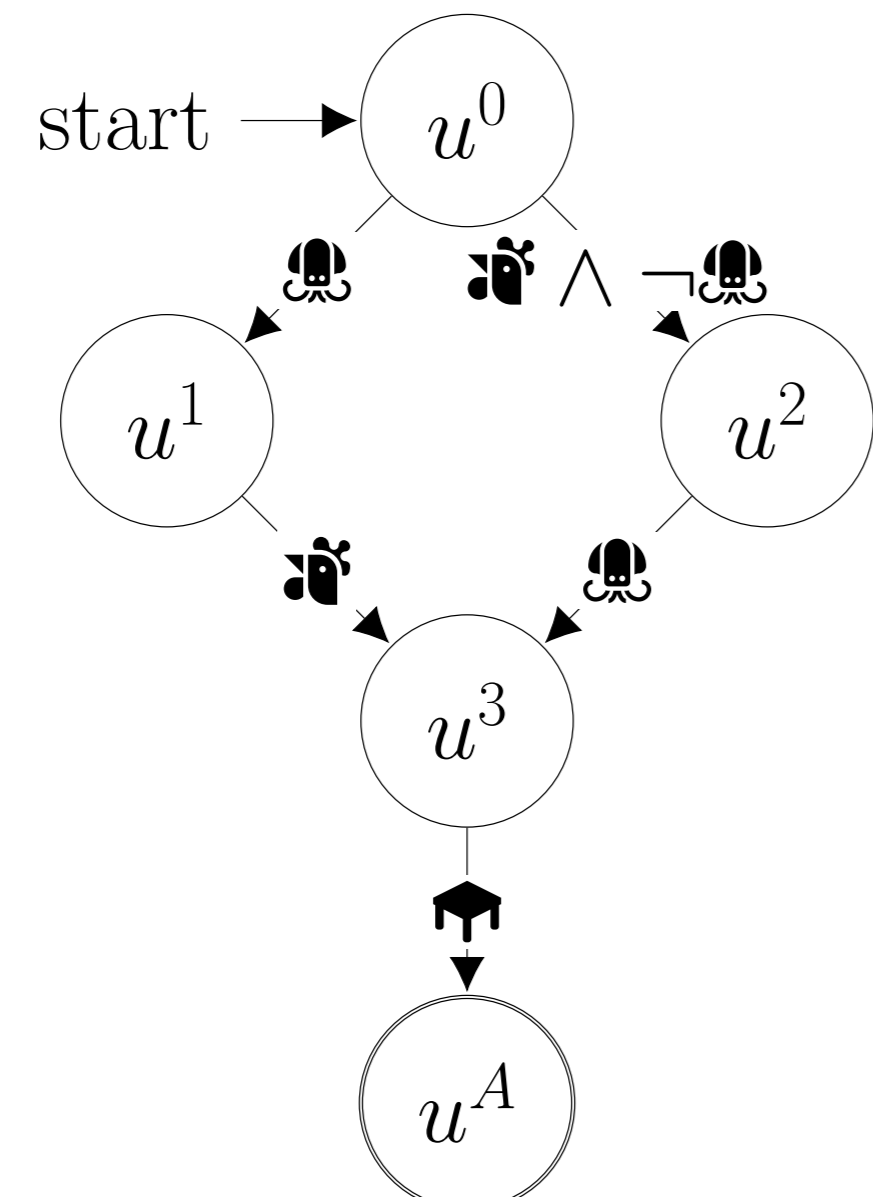
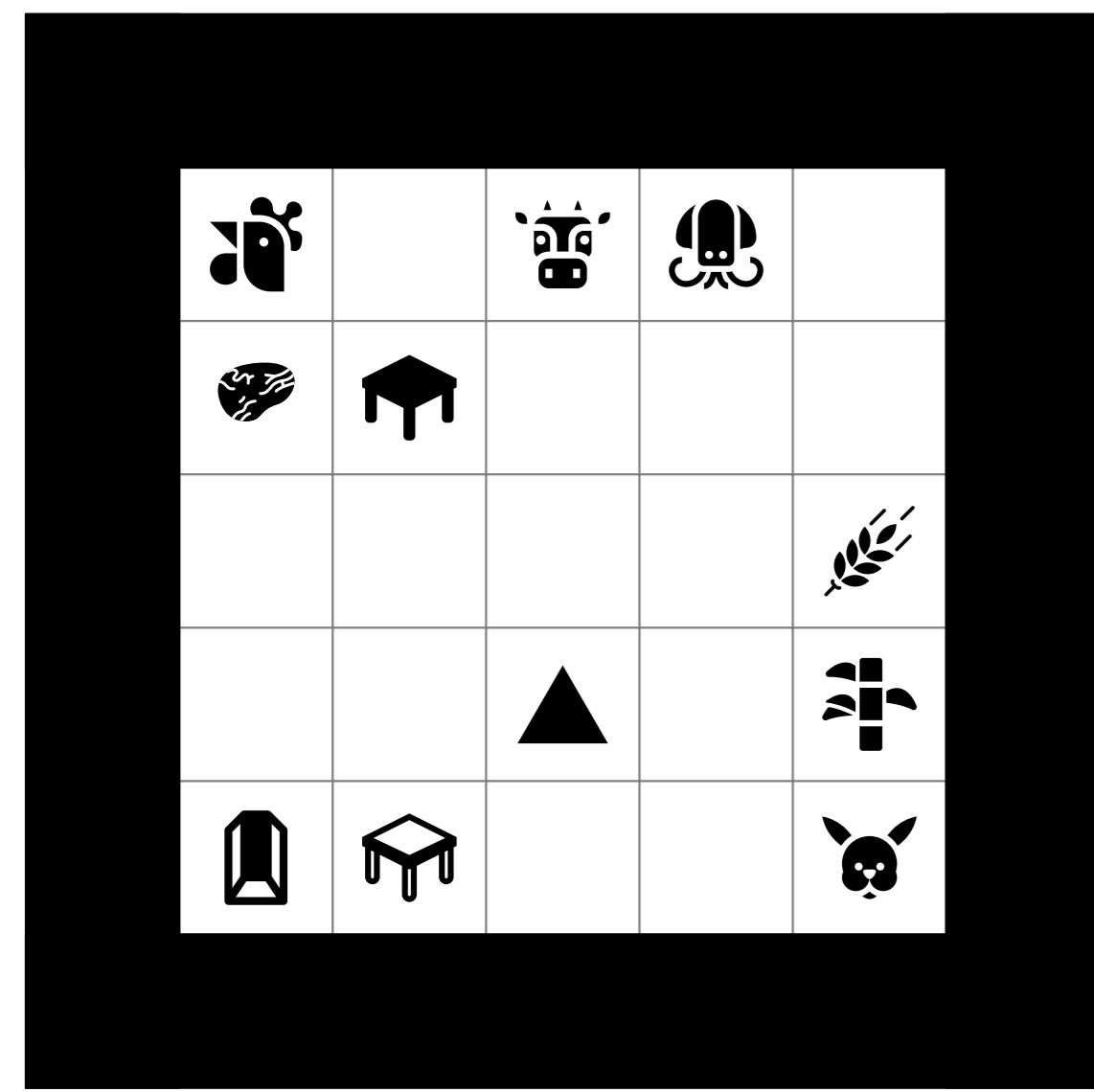
Motivation

- Learn at *multiple time scales* simultaneously.
- Learn with a rich *structure* of events and durations.

— How? —

Using a type of finite-state machines called Reward Machines [3].

What is a Reward Machine (RM)?

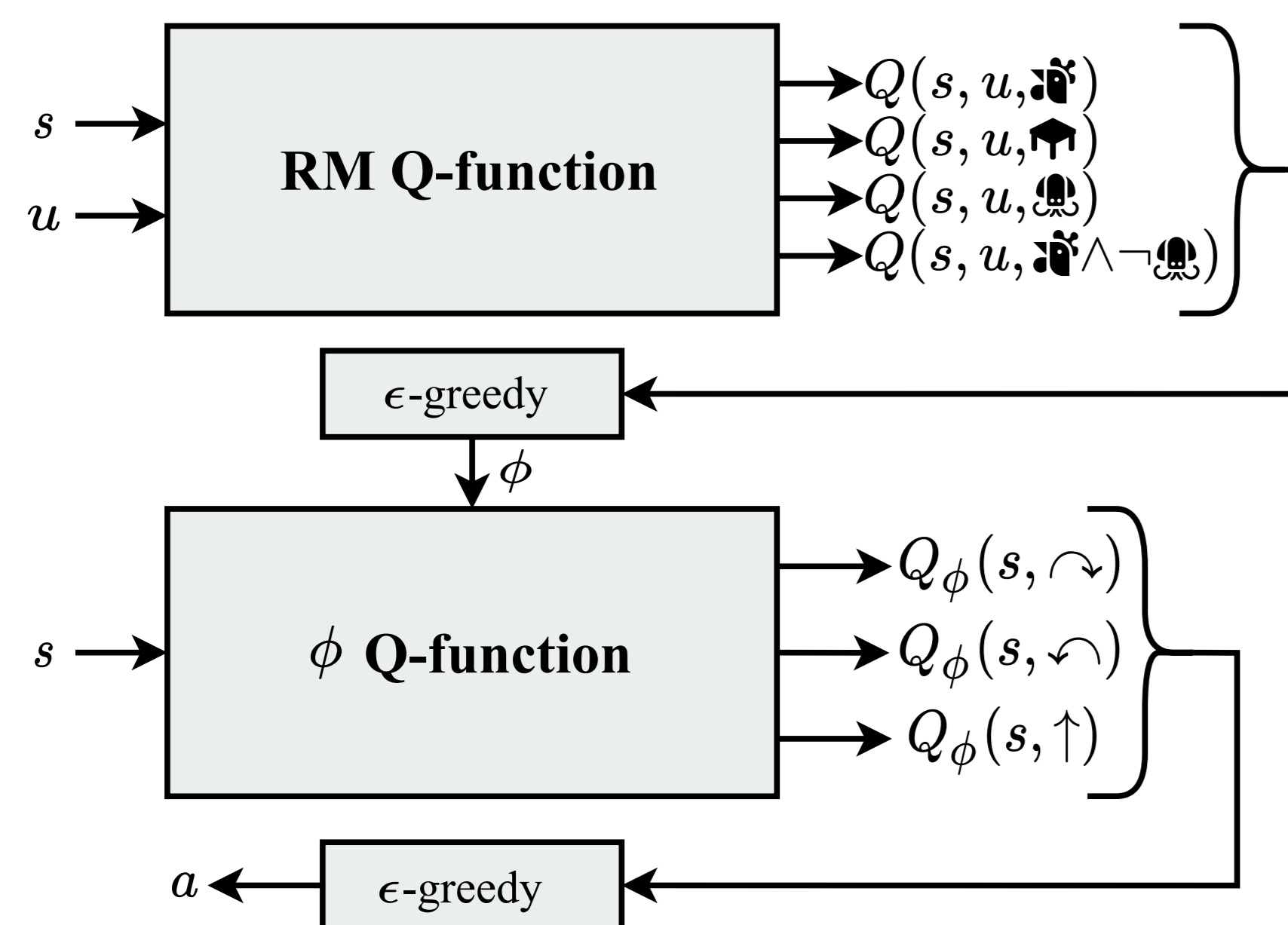


$\mathcal{P} = \{\text{BATTER}, \text{BUCKET}, \text{COMPASS}, \text{LEATHER}, \text{PAPER}, \text{QUILL}, \text{SUGAR}, \text{MILKBUCKET}, \text{BOOK}, \text{MILKB.SUGAR}, \text{CAKE}\}$
Task: Observe BATTER and BUCKET in any order, then CAKE .

$$r(u, u') = \mathbb{1}[u \neq u^A \wedge u' = u^A]$$

Policy Learning in RMs

Using the *options* framework [2] for hierarchical reinforcement learning.



Limitations of RMs

- Lack of *modularity*.
- *Hard to learn* when they contain more than a few states.

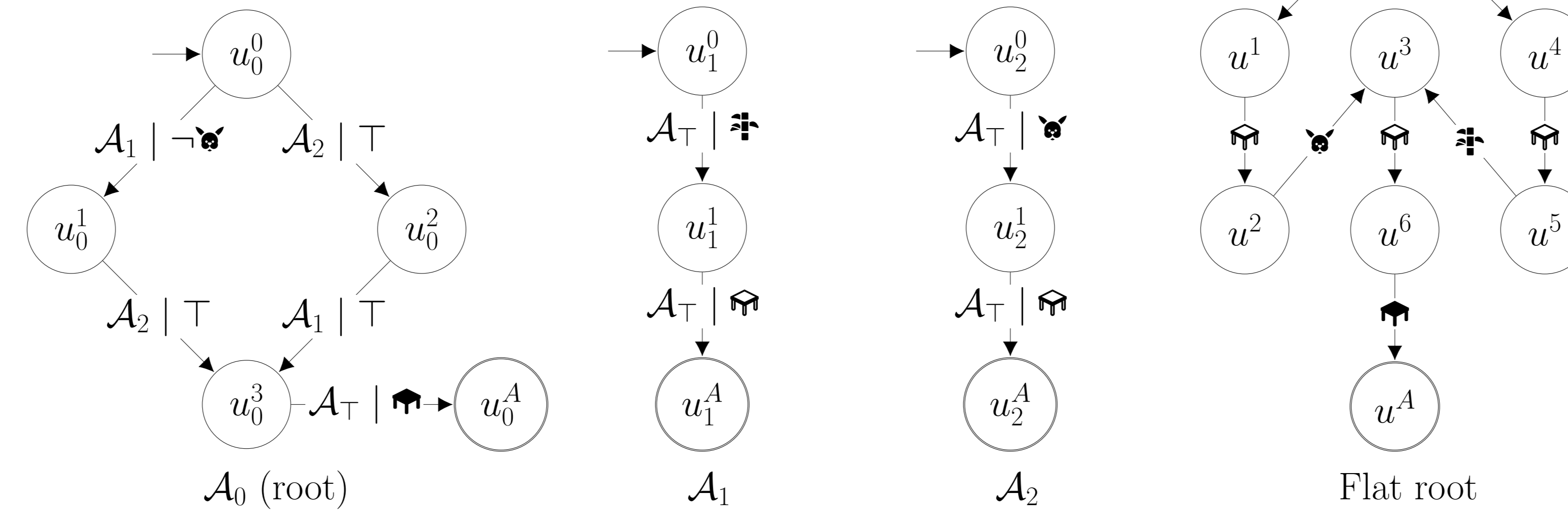
How to Address These?

Compose RMs into *hierarchies*.

Contributions

Hierarchies of Reward Machines (HRMs)

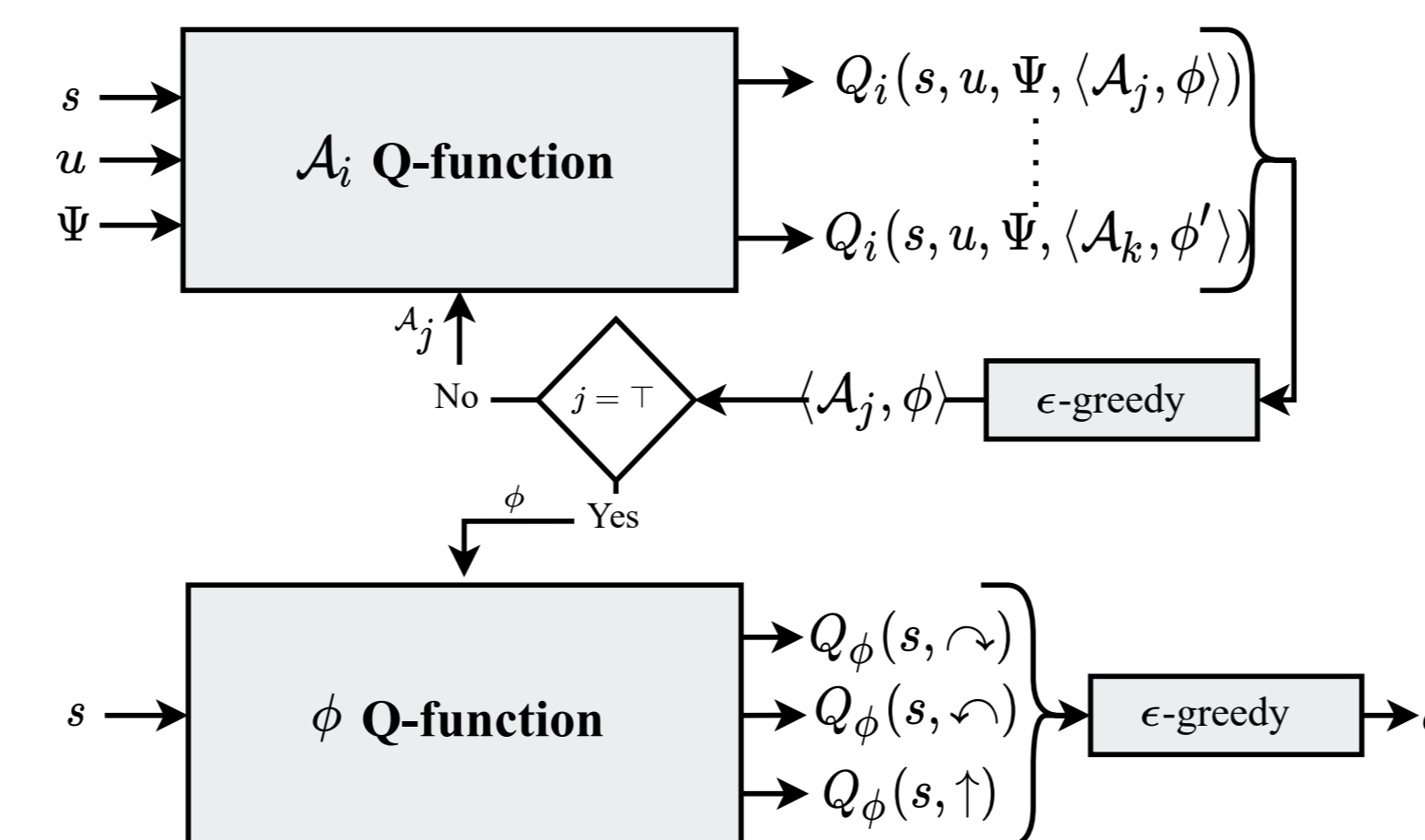
- Endue RMs with the ability to *call* each other.
- Theory:
 1. Given an HRM, there exists an *equivalent flat* HRM.
 2. The number of states and edges in an equivalent flat HRM may be *exponential* in the number of levels of the HRM.



Task	h	Description	Task	h	Description	Task	h	Description
BATTER	1	(BATTER & BUCKET); CAKE	QUILL	1	(QUILL & SUGAR); CAKE	BOOK&QUILL	3	BOOK & QUILL
BUCKET	1	BUCKET ; CAKE	SUGAR	1	SUGAR ; CAKE	MILKB.SUGAR	3	MILKBUCKET & SUGAR
COMPASS	1	(COMPASS & LEATHER); CAKE	BOOK	2	(PAPER & LEATHER); CAKE	CAKE	4	BATTER; MILKB.SUGAR; CAKE
LEATHER	1	LEATHER ; CAKE	MAP	2	(PAPER & COMPASS); CAKE			
PAPER	1	PAPER ; CAKE	MILKBUCKET	2	BUCKET; CAKE			

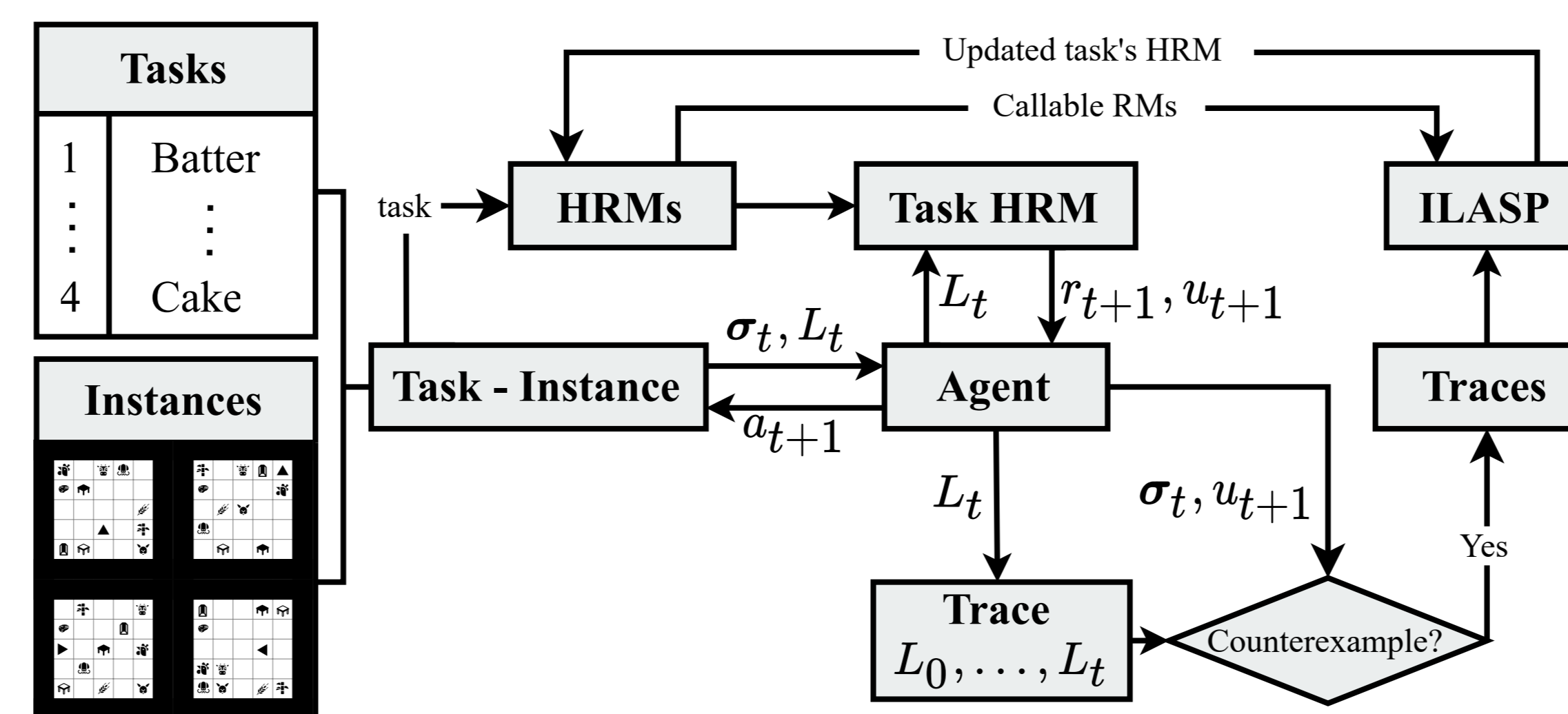
Learning Policies in HRMs

- Option types: *formula* options (associated with calls to \mathcal{A}_T) and *call* options (other calls).



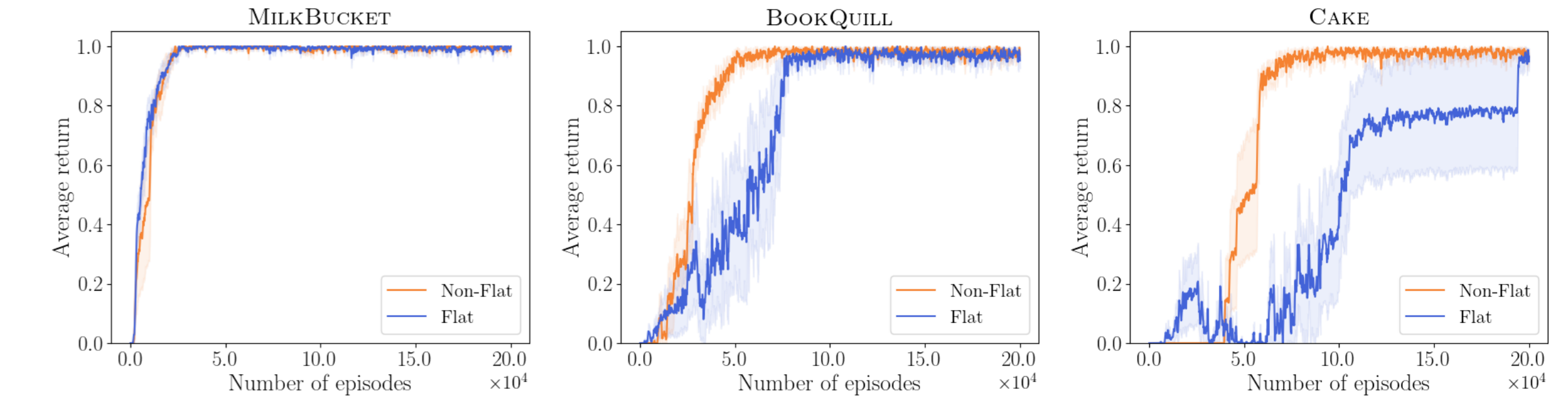
Interleaving Policy and HRM Learning

- Task-instance pairs are selected following a *curriculum learning* method.
- HRMs are learned using the *ILASP* inductive logic programming system [1].

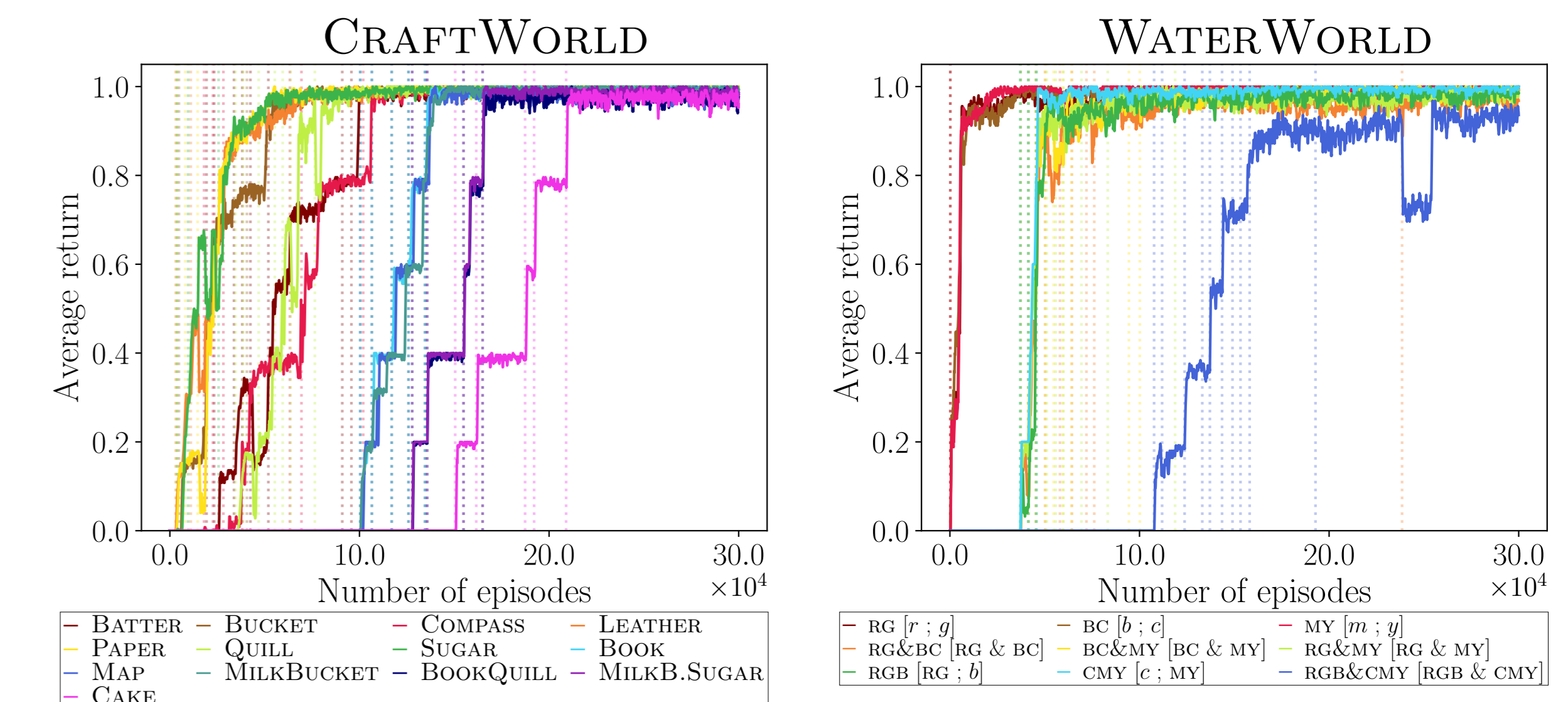


Results

Policy Learning can be Faster in Non-Flat HRMs



Learning Non-Flat HRMs



- A *restricted set of callable RMs* speeds up HRM learning by 5-7 \times .
- Using *options to explore* helps observing counterexamples faster ($\approx 128\times$ less episodes in some CW settings for the easiest tasks).

Learning Flat HRMs

- We compare our method for learning a non-flat HRM against:
 1. Our method for learning a flat HRM.
 2. Existing RM learning methods that label edges with proposition sets instead of formulas.
- Learning a non-flat HRM is more scalable than learning a flat HRM since it reuses previously learned RMs.
 - \Rightarrow The root may consist of less states and edges.
 - \Rightarrow Easier to learn!
- Abstraction through *formulas* is key in WATERWORLD.

Future Work

Relax some *assumptions* (e.g., handcrafted propositions, the level of each task) and increase *generalization* across instances.

References

- [1] M. Law, A. Russo, and K. Broda. The ILASP System for Learning Answer Set Programs, 2015.
- [2] R. S. Sutton, D. Precup, and S. P. Singh. Between MDPs and Semi-MDPs: A Framework for Temporal Abstraction in Reinforcement Learning. *Artif. Intell.*, 112(1-2):181–211, 1999.
- [3] R. Toro Icarte, T. Q. Klassen, R. A. Valenzano, and S. A. McIlraith. Using Reward Machines for High-Level Task Specification and Decomposition in Reinforcement Learning. In *ICML*, 2018.