
Hierarchies of Reward Machines

Daniel Furelos-Blanco¹ Mark Law^{1,2} Anders Jonsson³ Kryssia Broda¹ Alessandra Russo¹

Abstract

Reward machines (RMs) are a recent formalism for representing the reward function of a reinforcement learning task through a finite-state machine whose edges encode subgoals of the task using high-level events. The structure of RMs enables the decomposition of a task into simpler and independently solvable subtasks that help tackle long-horizon and/or sparse reward tasks. We propose a *formalism* for further abstracting the subtask structure by endowing an RM with the ability to call other RMs, thus composing a hierarchy of RMs (HRM). We *exploit* HRMs by treating each call to an RM as an independently solvable subtask using the options framework, and describe a curriculum-based method to *learn* HRMs from traces observed by the agent. Our experiments reveal that exploiting a handcrafted HRM leads to faster convergence than with a flat HRM, and that learning an HRM is feasible in cases where its equivalent flat representation is not.

1. Introduction

More than a decade ago, Dietterich et al. (2008) argued for the need to “learn at *multiple time scales* simultaneously, and with a rich *structure* of events and durations”. Finite-state machines (FSMs) are a simple yet powerful formalism for abstractly representing temporal tasks in a structured manner. One of the most prominent recent types of FSMs used in reinforcement learning (RL; Sutton & Barto, 2018) are reward machines (RMs; Toro Icarte et al., 2018; 2022), which compactly represent state-action histories in terms of high-level events; specifically, each edge is labeled with (i) a formula over a set of high-level events that capture a task’s subgoal, and (ii) a reward for satisfying the formula. Hence, RMs fulfill the need for structuring events and durations, and keep track of the achieved and pending subgoals.

¹Imperial College London, UK ²ILASP Limited, UK ³Universitat Pompeu Fabra, Spain. Correspondence to: Daniel Furelos-Blanco <d.furelos-blanco18@imperial.ac.uk>.

Hierarchical reinforcement learning (HRL; Barto & Mahadevan, 2003) frameworks, such as options (Sutton et al., 1999), have been used to *exploit* RMs by learning policies at two levels of abstraction: (i) select a formula (i.e., subgoal) from a given RM state, and (ii) select an action to (eventually) satisfy the chosen formula (Toro Icarte et al., 2018; Furelos-Blanco et al., 2021). The subtask decomposition powered by HRL enables learning at multiple scales simultaneously, and eases the handling of long-horizon and sparse reward tasks. In addition, several works have considered the problem of *learning* the RMs themselves from interaction (e.g., Toro Icarte et al., 2019; Xu et al., 2020; Furelos-Blanco et al., 2021; Hasanbeig et al., 2021). A common problem among methods learning minimal RMs is that they scale poorly as the number of states grows.

In this work, we make the following *contributions*:

1. Enhance the abstraction power of RMs by defining *hierarchies of RMs (HRMs)*, where constituent RMs can call other RMs (Section 3). We prove that any HRM can be transformed into an *equivalent* flat HRM that behaves exactly like the original RMs. We show that under certain conditions, the equivalent flat HRM can have exponentially more states and edges.
2. Propose an HRL algorithm to *exploit* HRMs by treating each call as a subtask (Section 4). Learning policies in HRMs further fulfills the desiderata posed by Dietterich et al. since (i) there is an arbitrary number of time scales to learn across (not only two), and (ii) there is a richer range of increasingly abstract events and durations. Besides, hierarchies enable *modularity* and, hence, the *reusability* of the RMs and policies. Empirically, we show that leveraging a handcrafted HRM enables faster convergence than an equivalent flat HRM.
3. Introduce a curriculum-based method for *learning* HRMs from traces given a set of composable tasks (Section 5). In line with the theory (Contribution 1), our experiments reveal that decomposing an RM into several is *crucial* to make its learning feasible (i.e., the flat HRM cannot be efficiently learned from scratch) since (i) the constituent RMs are simpler (i.e., they have fewer states and edges), and (ii) previously learned RMs can be used to efficiently *explore* the environment in the search for traces in more complex tasks.

2. Background

Given a finite set \mathcal{X} , we use $\Delta(\mathcal{X})$ to denote the probability simplex over \mathcal{X} , \mathcal{X}^* to denote (possibly empty) sequences of elements from \mathcal{X} , and \mathcal{X}^+ to denote non-empty sequences. We use \perp and \top to denote the truth values false and true, respectively. $\mathbb{1}[A]$ is the indicator function of event A .

Reinforcement Learning. We represent RL tasks as *episodic* labeled Markov decision processes (MDPs; Xu et al., 2020), each consisting of a set of states \mathcal{S} , a set of actions \mathcal{A} , a transition function $p : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$, a reward function $r : (\mathcal{S} \times \mathcal{A})^+ \times \mathcal{S} \rightarrow \mathbb{R}$, a discount factor $\gamma \in [0, 1)$, a finite set of *propositions* \mathcal{P} representing high-level events, a *labeling function* $l : \mathcal{S} \rightarrow 2^{\mathcal{P}}$ mapping states to proposition subsets called *labels*, and a *termination function* $\tau : (\mathcal{S} \times \mathcal{A})^* \times \mathcal{S} \rightarrow \{\perp, \top\} \times \{\perp, \top\}$. Hence the transition function p is Markovian, but the reward function r and termination function τ are not. Given a *history* $h_t = \langle s_0, a_0, \dots, s_t \rangle \in (\mathcal{S} \times \mathcal{A})^* \times \mathcal{S}$, a *label trace* (or trace, for short) $\lambda_t = \langle l(s_0), \dots, l(s_t) \rangle \in (2^{\mathcal{P}})^+$ assigns labels to all states in h_t . We assume (λ_t, s_t) captures all relevant information about h_t ; thus, the reward and transition information can be written $r(h_t, a_t, s_{t+1}) = r(h_{t+1}) = r(\lambda_{t+1}, s_{t+1})$ and $\tau(h_t) = \tau(\lambda_t, s_t)$, respectively. We aim to find a *policy* $\pi : (2^{\mathcal{P}})^+ \times \mathcal{S} \rightarrow \mathcal{A}$, a mapping from traces-states to actions, that maximizes the expected cumulative discounted reward (or *return*) $R_t = \mathbb{E}_\pi[\sum_{k=t}^n \gamma^{k-t} r(\lambda_{k+1}, s_{k+1})]$, where n is the last episode’s step.

At time t , the trace is $\lambda_t \in (2^{\mathcal{P}})^+$, and the agent observes a tuple $\mathbf{s}_t = \langle s_t, s_t^T, s_t^G \rangle$, where $s_t \in \mathcal{S}$ is the state and $(s_t^T, s_t^G) = \tau(\lambda_t, s_t)$ is the termination information, with s_t^T and s_t^G indicating whether or not the history (λ_t, s_t) is terminal or a goal, respectively. The agent also observes a label $\mathcal{L}_t = l(s_t)$. If the history is non-terminal, the agent runs action $a_t \in \mathcal{A}$, and the environment transitions to state $s_{t+1} \sim p(\cdot | s_t, a_t)$. The agent then observes tuple \mathbf{s}_{t+1} and label \mathcal{L}_{t+1} , extends the trace as $\lambda_{t+1} = \lambda_t \oplus \mathcal{L}_{t+1}$, and receives reward $r_{t+1} = r(\lambda_{t+1}, s_{t+1})$. A trace λ_t is a *goal trace* if $(s_t^T, s_t^G) = (\top, \top)$, a *dead-end trace* if $(s_t^T, s_t^G) = (\top, \perp)$, and an *incomplete trace* if $s_t^T = \perp$. We assume that the reward is $r(\lambda_{t+1}, s_{t+1}) = \mathbb{1}[\tau(\lambda_{t+1}, s_{t+1}) = (\top, \top)]$, i.e. 1 for goal traces and 0 otherwise.

Example 1. The CRAFTWORLD domain (cf. Figure 1a) is used as a running example. In this domain, the agent (\blacktriangle) can move forward or rotate 90°, staying put if it moves towards a wall. Locations are labeled with propositions from $\mathcal{P} = \{\text{L}, \text{R}, \text{F}, \text{B}, \text{D}, \text{S}, \text{W}, \text{E}, \text{N}, \text{O}\}$. The agent observes propositions that it steps on, e.g. $\{\text{D}\}$ in the top-left corner. Table 1 lists the tasks we consider, which consist of observing a sequence of propositions.¹ For the task BOOK, two goal traces are $\langle \{\text{F}\}, \{\text{R}\}, \{\text{W}\}, \{\text{N}\}, \{\text{O}\} \rangle$ and

¹The tasks are based on those by Andreas et al. (2017) and Toro Icarte et al. (2018), but definable in terms of each other.

$\langle \{\text{W}\}, \{\text{R}\}, \{\text{F}\}, \{\text{N}\}, \{\text{O}\} \rangle$ (they could contain irrelevant labels in between). The traces $\langle \{\text{W}\}, \{\text{R}\} \rangle$ and $\langle \{\text{D}\} \rangle$ are incomplete. There are no dead-end traces in this scenario.

Options (Sutton et al., 1999) address temporal abstraction in RL. Given an episodic labeled MDP, an option is a tuple $\omega = \langle \mathcal{I}_\omega, \pi_\omega, \beta_\omega \rangle$, where $\mathcal{I}_\omega \subseteq \mathcal{S}$ is the option’s initiation set, $\pi_\omega : \mathcal{S} \rightarrow \mathcal{A}$ is the option’s policy, and $\beta_\omega : \mathcal{S} \rightarrow [0, 1]$ is the option’s termination condition. An option is available in $s \in \mathcal{S}$ if $s \in \mathcal{I}_\omega$, selects actions according to π_ω , and terminates in $s \in \mathcal{S}$ with probability $\beta_\omega(s)$.

Reward Machines. A (*simple*) *reward machine* (RM; Toro Icarte et al., 2018; 2022) is a tuple $\langle \mathcal{U}, \mathcal{P}, \varphi, r, u^0, \mathcal{U}^A, \mathcal{U}^R \rangle$, where \mathcal{U} is a finite set of states; \mathcal{P} is a finite set of propositions; $\varphi : \mathcal{U} \times \mathcal{U} \rightarrow \text{DNF}_{\mathcal{P}}$ is a state transition function such that $\varphi(u, u')$ denotes the disjunctive normal form (DNF) formula over \mathcal{P} to be satisfied to transition from u to u' ; $r : \mathcal{U} \times \mathcal{U} \rightarrow \mathbb{R}$ is a reward function such that $r(u, u')$ is the reward for transitioning from u to u' ; $u^0 \in \mathcal{U}$ is an initial state; $\mathcal{U}^A \subseteq \mathcal{U}$ is a set of accepting states denoting the task’s goal achievement; and $\mathcal{U}^R \subseteq \mathcal{U}$ is a set of rejecting states denoting the unfeasibility of achieving the goal. The state transition function is *deterministic*, i.e. at most one formula from each state is satisfied. To verify if a formula is satisfied by a label $\mathcal{L} \subseteq \mathcal{P}$, \mathcal{L} is used as a truth assignment where propositions in \mathcal{L} are true, and false otherwise (e.g., $\{a\} \models a \wedge \neg b$). If no transition formula is satisfied, the state remains unchanged.

Ideally, RM states should capture traces, such that (i) pairs (u, s) of an RM state and an MDP state make termination and reward Markovian, (ii) the reward $r(u, u')$ matches the underlying MDP’s reward, and (iii) goal traces end in an accepting state, rejecting traces end in a rejecting state, and incomplete traces do not end in accepting or rejecting states. As per the previous reward assumption, the reward transition functions are $r(u, u') = \mathbb{1}[u \notin \mathcal{U}^A \wedge u' \in \mathcal{U}^A]$.

Example 2. Figure 1b shows an RM for BOOK. The state transition function φ is deterministic since no label satisfies both $\varphi(u^0, u^1) = \text{F} \wedge \neg \text{W}$ and $\varphi(u^0, u^4) = \text{W}$; in contrast, φ becomes non-deterministic if $\varphi(u^0, u^1) = \text{F}$ since $\{\text{F}, \text{W}\}$ satisfies both formulas. Note that RMs compactly represent traces in terms of key events, e.g. u^2 indicates that a label satisfying $\text{F} \wedge \neg \text{W}$ followed by another satisfying W were observed (everything else is ignored).

3. Formalization of HRMs

RMs are the building blocks of our formalism. To constitute a hierarchy of RMs, we need to endow RMs with the ability to call each other. We redefine the *state transition function* as $\varphi : \mathcal{U} \times \mathcal{U} \times \mathcal{M} \rightarrow \text{DNF}_{\mathcal{P}}$, where \mathcal{M} is a set of RMs. The expression $\varphi(u, u', M)$ denotes the DNF formula over \mathcal{P} that must be satisfied to transition from $u \in \mathcal{U}$ to $u' \in \mathcal{U}$ by

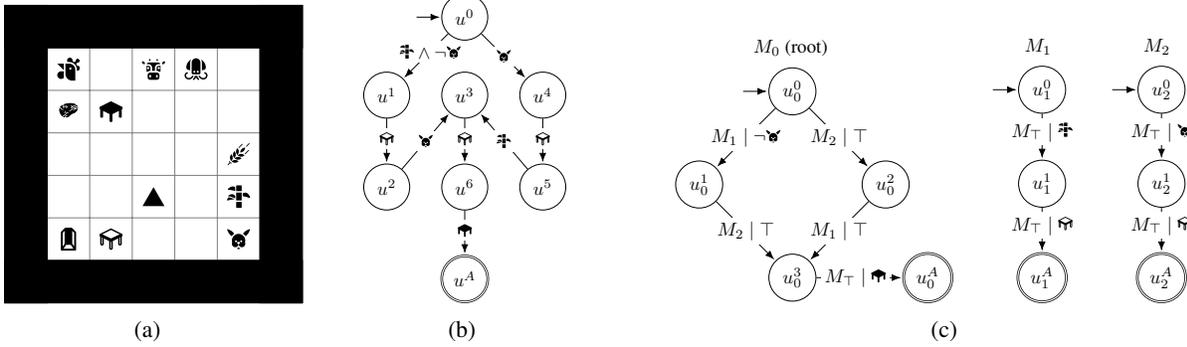


Figure 1: A CRAFTWORLD grid (a), and a flat HRM (b) and a non-flat one (c) for BOOK. In (b), an edge from u to u' is labeled $\varphi(u, u')$. In (c), an edge from u to u' in RM M_i is labeled $M_j \mid \varphi_i(u, u', M_j)$. In both cases, accepting states are double circled, and loops are omitted.

Table 1: List of CRAFTWORLD tasks. Descriptions “ $x ; y$ ” express sequential order (observe/do x then y), descriptions “ x & y ” express that x and y can be observed/done in any order, and h is the root RM’s height.

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

calling RM $M \in \mathcal{M}$. We refer to the formulas $\varphi(u, u', M)$ as *contexts* since they represent conditions under which calls are made. As we shall see later, contexts help preserve determinism and must be satisfied to start a call (a necessary but not sufficient condition). The hierarchies we consider contain an RM M_\top called the *leaf* RM, which solely consists of an accepting state (i.e., $\mathcal{U}_\top = \mathcal{U}_\top^A = \{u_\top^0\}$), and immediately returns control to the RM that calls it.

Definition 3.1. A *hierarchy of reward machines (HRM)* is a tuple $H = \langle \mathcal{M}, M_r, \mathcal{P} \rangle$, where $\mathcal{M} = \{M_0, \dots, M_{m-1}\} \cup \{M_\top\}$ is a set of m RMs and the leaf RM M_\top , $M_r \in \mathcal{M} \setminus \{M_\top\}$ is the root RM, and \mathcal{P} is a finite set of propositions used by all constituent RMs.

We make the following *assumptions*: (i) HRMs do not have circular dependencies (i.e., an RM cannot be called back from itself, including recursion), (ii) rejecting states are global (i.e., cause the root task to fail), (iii) accepting and rejecting states do not have transitions to other states, and (iv) the reward function of the root corresponds to the reward obtained in the underlying MDP. Given assumption (i), each RM M_i has a *height* h_i , which corresponds to the maximum number of nested calls needed to reach the leaf. Formally, if $i = \top$, then $h_i = 0$; otherwise, $h_i = 1 + \max_j h_j$, where j ranges over all RMs called by M_i (i.e., there exists $(u, v) \in \mathcal{U}_i \times \mathcal{U}_i$ such that $\varphi_i(u, v, M_j) \neq \perp$).

Example 3. Figure 1c shows BOOK’s HRM, whose root has height 2. The PAPER and LEATHER RMs, which have height 1 and consist of observing a two-proposition se-

quence, can be run in any order followed by observing ♣ . The context $\rightarrow \text{♠}$ in the call to M_1 preserves determinism, as detailed later.

In the following paragraphs, we describe how an HRM processes a label trace. To indicate where the agent is in an HRM, we define the notion of *hierarchy states*.

Definition 3.2. Given an HRM $H = \langle \mathcal{M}, M_r, \mathcal{P} \rangle$, a *hierarchy state* is a tuple $\langle M_i, u, \Phi, \Gamma \rangle$, where $M_i \in \mathcal{M}$ is an RM, $u \in \mathcal{U}_i$ is a state, $\Phi \in \text{DNF}_{\mathcal{P}}$ is an accumulated context, and Γ is a call stack.

Definition 3.3. Given an HRM $H = \langle \mathcal{M}, M_r, \mathcal{P} \rangle$, a *call stack* Γ contains tuples $\langle u, v, M_i, M_j, \phi, \Phi \rangle$, each denoting a call where $u \in \mathcal{U}_i$ is the state from which the call is made; $v \in \mathcal{U}_i$ is the next state in the calling RM $M_i \in \mathcal{M}$ after reaching an accepting state of the called RM $M_j \in \mathcal{M}$; $\phi \in \text{DNF}_{\mathcal{P}}$ are the disjuncts of $\varphi_i(u, v, M_j)$ satisfied by a label; and $\Phi \in \text{DNF}_{\mathcal{P}}$ is the accumulated context.

Call stacks determine where to resume the execution. Each RM appears in the stack at most once since, by assumption, HRMs have no circular dependencies. We use $\Gamma \oplus \langle u, v, M_i, M_j, \phi, \Phi \rangle$ to denote a stack recursively defined by a stack Γ and a top element $\langle u, v, M_i, M_j, \phi, \Phi \rangle$, where the *accumulated context* Φ is the condition under which a call from a state u is made. The initial hierarchy state of an HRM $H = \langle \mathcal{M}, M_r, \mathcal{P} \rangle$ is $\langle M_r, u_r^0, \top, [] \rangle$: we are in the initial state of the root, there is no accumulated context, and the stack is empty.

At the beginning of this section, we mentioned that satisfying the context of a call is a necessary but not sufficient condition to start the call. We now introduce a sufficient condition, called *exit condition*.

Definition 3.4. Given an HRM $H = \langle \mathcal{M}, M_r, \mathcal{P} \rangle$ and a hierarchy state $\langle M_i, u, \Phi, \Gamma \rangle$, the *exit condition* $\xi_{i,u,\Phi} \in \text{DNF}_{\mathcal{P}}$ is the formula that must be satisfied to leave that hierarchy state. Formally,

$$\xi_{i,u,\Phi} = \begin{cases} \Phi & \text{if } i = \top, \\ \bigvee_{\substack{\phi = \varphi_i(u,v,M_j), \\ \phi \neq \perp, v \in \mathcal{U}_i, M_j \in \mathcal{M}}} \xi_{j,u_j^0, \text{DNF}(\Phi \wedge \phi)} & \text{otherwise,} \end{cases}$$

where $\text{DNF}(\Phi \wedge \phi)$ is $\Phi \wedge \phi$ in DNF. The formula is Φ if $M_i = M_{\top}$ since it always returns control once called. Otherwise, the formula is recursively defined as the disjunction of the exit conditions from the initial state of the called RM. For instance, the exit condition for the initial hierarchy state in Figure 1c is $(\neg \heartsuit \wedge \clubsuit) \vee \heartsuit$.

We can now define the *hierarchical transition function* δ_H , which maps a hierarchy state $\langle M_i, u, \Phi, \Gamma \rangle$ into another given a label \mathcal{L} . There are three cases:

1. If u is an accepting state of M_i and the stack Γ is non-empty, pop the top element of Γ and return control to the previous RM, recursively applying δ_H in case several accepting states are reached simultaneously. Formally, the next hierarchy state is $\delta_H(\langle M_j, u', \top, \Gamma' \rangle, \perp)$ if $u \in \mathcal{U}_i^A$, $|\Gamma| > 0$, where $\Gamma = \Gamma' \oplus \langle \cdot, u', M_j, M_i, \cdot, \cdot \rangle$, \perp denotes a label that cannot satisfy any formula, and \cdot denotes something unimportant for the case.
2. If \mathcal{L} satisfies the context of a call and the exit condition from the initial state of the called RM, push the call onto the stack and recursively apply δ_H until M_{\top} is reached. Formally, the next hierarchy state is $\delta_H(\langle M_j, u_j^0, \Phi', \Gamma \oplus \langle u, u', M_i, M_j, \phi, \Phi \rangle \rangle, \mathcal{L})$ if $\mathcal{L} \models \xi_{j,u_j^0,\Phi'}$, where $\phi = \varphi_i(u, u', M_j)(\mathcal{L})$ and $\Phi' = \text{DNF}(\Phi \wedge \phi)$. Here, $\varphi(\mathcal{L})$ denotes the disjuncts of a DNF formula $\varphi \in \text{DNF}_{\mathcal{P}}$ satisfied by \mathcal{L} .
3. If none of the above holds, the hierarchy state remains unchanged.

The state transition functions φ of the RMs must be such that δ_H is *deterministic*, i.e. a label cannot simultaneously satisfy the contexts and exit conditions associated with two triplets $\langle u, v, M_i \rangle$ and $\langle u, v', M_j \rangle$ such that either (i) $v = v'$ and $i \neq j$, or (ii) $v \neq v'$. Contexts help enforce determinism by making formulas mutually exclusive. For instance, if the call to M_1 from the initial state of M_0 in Figure 1c had context \top instead of $\neg \heartsuit$, then M_1 and M_2 could be both

started if $\{\heartsuit, \spadesuit\}$ was observed, thus making the HRM non-deterministic. Finally, we introduce *hierarchy traversals*, which determine how a label trace is processed by an HRM.

Definition 3.5. Given a label trace $\lambda = \langle \mathcal{L}_0, \dots, \mathcal{L}_n \rangle$, a *hierarchy traversal* $H(\lambda) = \langle v_0, v_1, \dots, v_{n+1} \rangle$ is a unique sequence of hierarchy states such that (i) $v_0 = \langle M_r, u_r^0, \top, \perp \rangle$, and (ii) $\delta_H(v_i, \mathcal{L}_i) = v_{i+1}$ for $i = 0, \dots, n$. An HRM H *accepts* λ if $v_{n+1} = \langle M_r, u, \top, \perp \rangle$ and $u \in \mathcal{U}_r^A$ (i.e., an accepting state of the root is reached). Analogously, H *rejects* λ if $v_{n+1} = \langle M_k, u, \cdot, \cdot \rangle$ and $u \in \mathcal{U}_k^R$ for any $k \in [0, m-1]$ (i.e., a rejecting state in the HRM is reached).

Example 4. The HRM in Figure 1c accepts $\lambda = \langle \{\heartsuit\}, \{\heartsuit\}, \{\}, \{\heartsuit\}, \{\heartsuit\} \rangle$ since $H(\lambda) = \langle \langle M_0, u_0^0, \top, \perp \rangle, \langle M_1, u_1^1, \top, [\langle u_0^0, u_0^1, M_0, M_1, \neg \heartsuit, \top \rangle] \rangle, \langle M_0, u_0^1, \top, \perp \rangle, \langle M_0, u_0^1, \top, \perp \rangle, \langle M_2, u_2^2, \top, [\langle u_0^1, u_0^3, M_0, M_2, \top, \top \rangle] \rangle, \langle M_0, u_0^3, \top, \perp \rangle, \langle M_0, u_0^4, \top, \perp \rangle \rangle$. Appendix A.1 shows the step-by-step application of δ_H omitted here.

The behavior of an HRM H can be reproduced by an *equivalent flat* HRM \bar{H} ; that is, (i) the root of \bar{H} has height 1 and, (ii) \bar{H} accepts a trace iff H accepts it, rejects a trace iff H rejects it, and neither accepts nor rejects a trace iff H does not accept it nor reject it. Flat HRMs thus capture the original RM definition, e.g. Figure 1b is a flat HRM for BOOK. We formally define equivalence and prove the following equivalence theorem by construction in Appendix A.2.1.

Theorem 3.6. *Given an HRM H , there exists an equivalent flat HRM \bar{H} .*

Given the construction used in Theorem 3.6, we show that the number of states and edges of the resulting flat HRM can be *exponential* in the height of the root (see Theorem 3.7). We prove this in Appendix A.2.2 through an instance of a general HRM parametrization where the constituent RMs are *highly reused*, hence illustrating the convenience of HRMs to succinctly compose existing knowledge. In line with the theory, learning a non-flat HRM can take a few seconds, whereas learning an equivalent flat HRM is often unfeasible (see Section 6).

Theorem 3.7. *Let $H = \langle \mathcal{M}, M_r, \mathcal{P} \rangle$ be an HRM and h_r be the height of its root M_r . The number of states and edges in an equivalent flat HRM \bar{H} can be exponential in h_r .*

4. Policy Learning in HRMs

In what follows, we explain how to *exploit* the temporal structure of an HRM $H = \langle \mathcal{M}, M_r, \mathcal{P} \rangle$ using two types of *options*. We also describe (i) how to learn the policies of these options, (ii) when these options terminate, and (iii) an option selection algorithm that ensures the currently running options and the current hierarchy state are aligned.

Option Types. Given an RM $M_i \in \mathcal{M}$, a state $u \in \mathcal{U}_i$ and a context Φ , an option $\omega_{i,u,\Phi}^{j,\phi}$ is derived for each non-false

disjunct ϕ of each transition $\varphi_i(u, v, M_j)$, where $v \in \mathcal{U}_i$ and $M_j \in \mathcal{M}$. An option is either (i) a *formula option* if $j = \top$ (i.e., M_\top is called), or (ii) a *call option* otherwise. A formula option attempts to reach a label that satisfies $\phi \wedge \Phi$ through primitive actions, whereas a call option aims to reach an accepting state of the called RM M_j under context $\phi \wedge \Phi$ by invoking other options.

Policies. Policies are ϵ -greedy during training, and greedy during evaluation. A *formula option's policy* is derived from a Q-function $q_{\phi \wedge \Phi}(s, a; \theta_{\phi \wedge \Phi})$ approximated by a deep Q-network (DQN; Mnih et al., 2015) with parameters $\theta_{\phi \wedge \Phi}$, which outputs the Q-value of each action given an MDP state. We store all options' experiences $(s_t, a, s_{t+1}, \mathcal{L}_{t+1})$ in a single replay buffer \mathcal{D} , thus performing intra-option learning (Sutton et al., 1998). The Q-learning update uses the following loss function:

$$\mathbb{E}_{(s_t, a, s_{t+1}, \mathcal{L}_{t+1}) \sim \mathcal{D}} \left[(y_{\phi \wedge \Phi} - q_{\phi \wedge \Phi}(s_t, a; \theta_{\phi \wedge \Phi}))^2 \right], \quad (1)$$

where $y_{\phi \wedge \Phi} = r_{\phi \wedge \Phi}(\mathcal{L}_{t+1}) + \gamma \max_{a'} q_{\phi \wedge \Phi}(s_{t+1}, a'; \theta_{\phi \wedge \Phi}^-)$. The reward $r_{\phi \wedge \Phi}(\mathcal{L}_{t+1})$ is 1 if $\phi \wedge \Phi$ is satisfied by \mathcal{L}_{t+1} and 0 otherwise; the term $q_{\phi \wedge \Phi}(s_{t+1}, a'; \theta_{\phi \wedge \Phi}^-)$ is 0 when $\phi \wedge \Phi$ is satisfied or a dead-end is reached (i.e., $s_{t+1}^T = \top$ and $s_{t+1}^G = \perp$); and $\theta_{\phi \wedge \Phi}^-$ are the parameters of a fixed target network.

A *call option's policy* is induced by a Q-function $q_i(s, u, \Phi, \langle M_j, \phi \rangle; \theta_i)$ associated with the called RM M_i and approximated by a DQN with parameters θ_i that outputs the Q-value of each call in the RM given an MDP state, an RM state and a context. We store experiences $(s_t, \omega_{i, u, \Phi}^{j, \phi}, s_{t+k})$ in a replay buffer \mathcal{D}_i associated with M_i , and perform SMDP Q-learning using the following loss:

$$\mathbb{E}_{(s_t, \omega_{i, u, \Phi}^{j, \phi}, s_{t+k}) \sim \mathcal{D}_i} \left[(y_i - q_i(s_t, u, \Phi, \langle M_j, \phi \rangle; \theta_i))^2 \right],$$

where $y_i = r + \gamma^k \max_{j', \phi'} q_i(s_{t+k}, u', \Phi', \langle M_{j'}, \phi' \rangle; \theta_i^-)$; k is the number of steps between s_t and s_{t+k} ; r is the sum of discounted rewards during this time; u' and Φ' are the RM state and context after running the option; $M_{j'}$ and ϕ' correspond to an outgoing transition from u' , i.e. $\phi' \in \varphi_i(u', \cdot, M_{j'})$; and θ_i^- are the parameters of a fixed target network. The term $q_i(s_{t+k}, u', \dots)$ is 0 if u' is accepting or rejecting. Following the definition of δ_H , Φ' is \top if the hierarchy state changes; thus, $\Phi' = \top$ if $u' \neq u$, and $\Phi' = \Phi$ otherwise. Given our assumption on the MDP reward, we define reward transition functions as $r_i(u, u') = \mathbb{1}[u \notin \mathcal{U}_i^A \wedge u' \in \mathcal{U}_i^A]$. Learning a call option's policy and lower-level option policies at once can be unstable due to *non-stationarity* (Levy et al., 2019), e.g. a lower-level option may not achieve its goal at times. To relax the issue, experiences are added to the buffer only when options achieve their goal (i.e., call options assume lower-level options terminate successfully). The policies will be

recursively optimal (Dietterich, 2000) as each subtask is optimized individually; however, since the Q-functions are approximated, policies may only be approximately optimal. The implementation details are discussed in Appendix B.1.

Termination. An option terminates in two cases. First, if the episode ends in a goal or dead-end state. Second, if the hierarchy state changes and either successfully completes the option or interrupts the option. Concretely, a formula option $\omega_{i, u, \Phi}^{\top, \phi}$ is only applicable in a hierarchy state $\langle M_i, u, \Phi, \Gamma \rangle$, while a call option $\omega_{i, u, \Phi}^{j, \phi}$ always corresponds to a stack item $\langle u, \cdot, M_i, M_j, \phi, \Phi \rangle$. We can thus analyze the hierarchy state to see if an option is still executing or should terminate.

Algorithm. An *option stack* Ω_H stores the currently executing options. Initially, Ω_H is empty. At each step, Ω_H is filled (if needed) by repeatedly choosing options starting from the current hierarchy state using call option policies until a formula option is selected. Since HRMs have, by assumption, no circular dependencies, a formula option will eventually be chosen. An action is then selected using the formula option's policy. Once the action is applied, the DQNs associated with formula options are updated. The new hierarchy state is then used to determine which options in Ω_H have terminated. Experiences for the terminated options that achieved their goal are pushed into the corresponding buffers, and the DQNs associated with the call options are updated. Finally, Ω_H is updated to *match the call stack* of the new hierarchy state (if needed) by mapping each call stack item into an option, and adding it to Ω_H if it is not already there. By aligning the option stack with the call stack, we can update DQNs for options that ended up being run in *hindsight* and which would have been otherwise ignored. We refer the reader to Appendix B.2 for the pseudo-code and step-by-step examples.

Example 5. Given the HRM in Figure 1c, let us assume the agent had chosen to run M_1 from u_0^0 . The option running M_1 is interrupted if the agent observes $\{\blacklozenge\}$ before $\{\blackcross\}$ since M_2 is started; thus, Ω_H is updated and indicates that the agent is now acting according to M_2 . In contrast, if $\{\blackcross\}$ is observed, the agent gets into M_1 as originally decided and, hence, the corresponding option is not interrupted.

5. Learning HRMs from Traces

In the previous section, we explained how a *given* HRM can be exploited using options; however, engineering an HRM is impractical. We here describe LHRM, a method that *interleaves* policy learning with HRM learning from interaction. We consider a *multi-task* setting. Given T tasks and I instances (e.g., grids) of an environment, the agent learns (i) an HRM for each task using traces from several instances for better accuracy, and (ii) general policies to

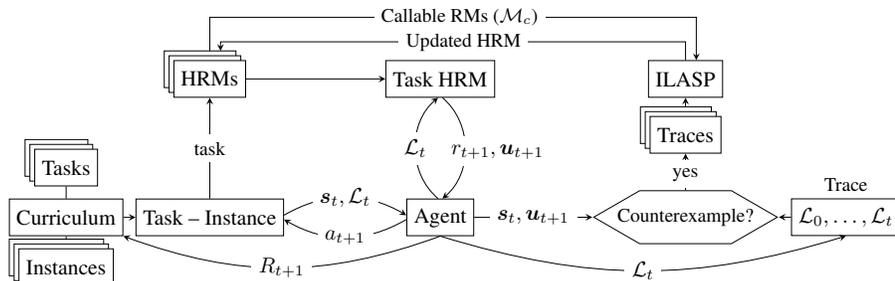


Figure 2: Overview of the interleaving algorithm. Given a set of tasks and a set of instances, the curriculum selects a task-instance pair at the start of an episode, and the HRM for the chosen task is taken from the bank of HRMs. At each step, the agent observes a tuple s_t and a label \mathcal{L}_t from the task-instance environment, and performs an action a_{t+1} . The label is used to (i) determine the next hierarchy state u_{t+1} and the reward r_{t+1} , and (ii) update the trace $\langle \mathcal{L}_0, \dots, \mathcal{L}_t \rangle$. If the trace is a counterexample, it is added to the task’s counterexample set and ILASP learns a new HRM (perhaps using previously learned RMs). The learned HRM replaces the old one in the bank of HRMs. If no counterexample is observed during the episode, the curriculum is updated using the undiscounted return R_{t+1} . Further details are described in the main text.

reach the goal in each task-instance pair. Namely, the agent interacts with $T \times I$ MDPs \mathbb{M}_{ij} , where $i \in [1, T]$ and $j \in [1, I]$. The learning proceeds from simpler to harder tasks such that HRMs for the latter build on the former.

In what follows, we detail LHRM’s components. We assume that (i) all MDPs share propositions \mathcal{P} and actions \mathcal{A} , and those defined on a given instance share states \mathcal{S} and labeling function l ; (ii) to stabilize policy learning, dead-end traces are shared across tasks;² (iii) the root’s height of a task’s HRM (or *task level*, for brevity) is known (see Table 1 for CRAFTWORLD); and (iv) without loss of generality, each RM has a single accepting state and a single rejecting state.

Curriculum Learning (Bengio et al., 2009). LHRM learns the tasks’ HRMs from lower to higher levels akin to Pierrot et al. (2019). Before starting an episode, LHRM selects an MDP \mathbb{M}_{ij} , where $i \in [1, T]$ and $j \in [1, I]$. The probability of selecting an MDP \mathbb{M}_{ij} is determined by an estimate of its average undiscounted return R_{ij} such that lower returns are mapped into higher probabilities (see details in Appendix C.1). Initially, only level 1 MDPs can be chosen. When the minimum average return across MDPs up to the current level surpasses a given threshold, the current level increases by 1, hence ensuring the learned HRMs and their associated policies are reusable in higher level tasks.

Learning an HRM. The learning of an HRM is analogous to the learning of a flat RM (Toro Icarte et al., 2019; Xu et al., 2020; Furelos-Blanco et al., 2021; Hasanbeig et al., 2021). The objective is to learn the state transition function φ_r of the root M_r with height h_r given (i) a set of states \mathcal{U}_r , (ii) a set of label traces $\Lambda = \Lambda^G \cup \Lambda^D \cup \Lambda^I$, (iii) a set of propositions \mathcal{P} , (iv) a set of RMs \mathcal{M} with lower heights

²The term $q_{\phi \wedge \Phi}(s_{t+1}, \dots)$ in Equation 1 is 0 if $(s_{t+1}^T, s_{t+1}^G) = (\top, \perp)$. Since experiences $(s_t, a, s_{t+1}, \mathcal{L}_{t+1})$ are shared through the buffer, evaluating the condition differently causes instabilities.

than h_r , (v) a set of callable RMs $\mathcal{M}_C \subseteq \mathcal{M}$ (by default, $\mathcal{M}_C = \mathcal{M}$), and (vi) the maximum number of disjuncts κ in the DNF formulas labeling the edges. The learned state transition function φ_r is such that the resulting HRM $H = \langle \mathcal{M} \cup \{M_r\}, M_r, \mathcal{P} \rangle$ accepts all goal traces Λ^G , rejects all dead-end traces Λ^D , and neither accepts or rejects incomplete traces Λ^I . The transition functions can be represented as sets of logic rules, which are learned using the ILASP (Law et al., 2015) inductive logic programming system (see Appendix C.2 for details on the ILASP encoding).

Interleaving Algorithm. LHRM *interleaves* the induction of HRMs with policy learning akin to Furelos-Blanco et al. (2021). Figure 2 illustrates the core blocks of the algorithm. Initially, the HRM’s root of each task $i \in [1, T]$ consists of 3 states (the initial, accepting, and rejecting states) and neither accepts nor rejects anything. A new HRM is learned when an episode’s label trace is not correctly recognized by the current HRM (i.e., if a goal trace is not accepted, a dead-end trace is not rejected, or an incomplete trace is accepted or rejected). The number of states in \mathcal{U}_r increases by 1 when an HRM that covers the examples cannot be learned, hence guaranteeing that the root has the smallest possible number of states (i.e., it is *minimal*) for a specific value of κ . When an HRM for task i is learned, the returns R_{ij} in the curriculum are set to 0 for all $j \in [1, I]$. Analogously to some RM learning methods (Toro Icarte et al., 2019; Xu et al., 2020; Hasanbeig et al., 2021), the first HRM for a task is learned using a set of traces; in our case, the ρ_s shortest traces from a set of ρ goal traces are used (empirically, short traces speed up learning). LHRM leverages learned options to *explore* the environment during the goal trace collection, accelerating the process when labels are sparse; specifically, options from lower height RMs are sequentially selected uniformly at random, and their greedy policy is run until termination. We describe other details in Appendix C.3.

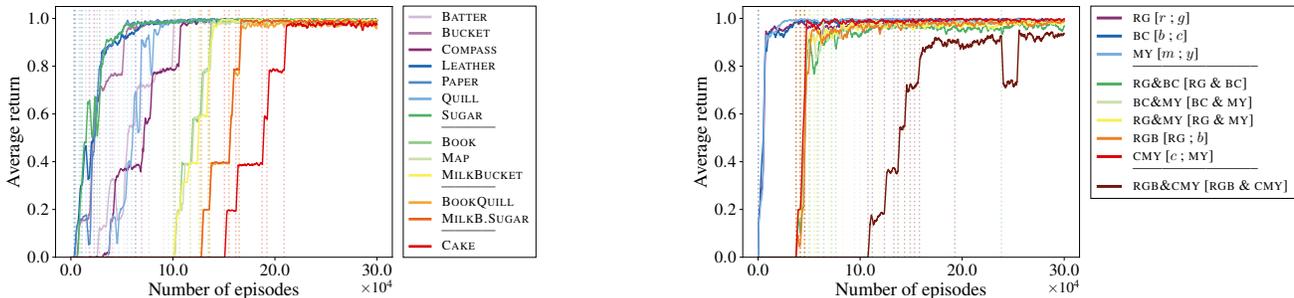


Figure 3: LHRM learning curves for CRAFTWORLD (FRL) and WATERWORLD (WD). The legends separate tasks by level. The WATERWORLD legend describes the subtask order in brackets following the specification introduced in Table 1. The dotted vertical lines correspond to episodes in which an HRM is learned.

6. Experimental Results

We evaluate the policy and HRM learning components in two *domains*: CRAFTWORLD and WATERWORLD. We consider four grid types for CRAFTWORLD (see Section 2): an open plan 7×7 grid (OP, Figure 1a), an open plan 7×7 grid with a lava location (OPL), a 13×13 four rooms grid (FR; Sutton et al., 1999), and a 13×13 four rooms grid with a lava location per room (FRL). The lava proposition must be avoided. WATERWORLD (Karpathy, 2015; Sidor, 2016; Toro Icarte et al., 2018) consists of a 2D box containing 12 balls of 6 different colors (2 per color) each moving at a constant speed in a fixed direction. The agent ball can change its velocity in any cardinal direction. The propositions $\mathcal{P} = \{r, g, b, c, m, y\}$ are the balls’ colors. Labels consist of the color of the balls the agent overlaps with and, unlike CRAFTWORLD, they may contain multiple propositions. The tasks consist in observing color sequences. We consider two settings: without dead-ends (WOD) and with dead-ends (WD). In WD, the agent must avoid 2 balls of an extra color. Further details are described in Appendix D.1.

We report the average performance across 5 runs, each using a different set of 10 random instances. The learning curves show the average undiscounted return obtained by the greedy policy every 100 episodes across instances. For other metrics (e.g., learning times), we present the average and the standard error (the latter in brackets). In HRM learning experiments, we set a 2-hour limit to learn the HRMs. The code is available at <https://github.com/ertsiger/hrm-learning>.

6.1. Learning of Non-Flat HRMs

Figure 3 shows the LHRM learning curves for CRAFTWORLD (FRL) and WATERWORLD (WD). These settings are the most challenging due to the inclusion of dead-ends since (i) they hinder the observation of goal examples in level 1 tasks using random walks, (ii) the RMs must include rejecting states, (iii) formula options must avoid dead-ends,

and (iv) call options must avoid invoking options leading to rejecting states. In line with the curriculum method, LHRM does not start learning a level h task until performance in tasks from levels $1, \dots, h - 1$ is sufficiently good. The convergence for high-level tasks is often fast due to the reuse of lower level HRMs and policies.

The average time (in seconds) exclusively spent on learning *all* HRMs is 1009.8 (122.3) for OP, 1622.6 (328.7) for OPL, 1031.6 (150.3) for FR, 1476.8 (175.3) for FRL, 35.4 (2.0) for WOD, and 67.0 (6.2) for WD (see Tables 3 and 6 in Appendix D.3.1). Dead-ends (OPL, FRL, WD) incur longer times since (i) there is one more proposition, (ii) there are edges to the rejecting state(s), and (iii) there are dead-end traces to cover. We observe that the complexity of learning an HRM does not necessarily correspond with the task complexity (e.g., the times for OP and FRL are close). Learning in WATERWORLD is faster than in CRAFTWORLD since the RMs have fewer states and there are fewer callable RMs.

Ablations. By *restricting the callable RMs* to those required by the HRM (e.g., *just* using PAPER and LEATHER RMs to learn BOOK’s), there are fewer ways to label the edges of the induced RM. Learning is 5-7 \times faster using 20% fewer calls to the learner (i.e., fewer examples) in CRAFTWORLD, and 1.5 \times faster in WATERWORLD (see Tables 4 and 7 in Appendix D.3.1); thus, HRM learning becomes less scalable as the number of tasks and levels grows. This is an instance of the *utility* problem (Minton, 1988). Refining the callable RM set prior to HRM learning is an avenue for future work.

We evaluate the performance of *exploration with options* using the number of episodes needed to collect the ρ goal traces for a given task since the activation of its level. Intuitively, the agent rarely moves far from a region of the state space using primitive actions only, thus taking longer to collect the traces; in contrast, options enable the agent to explore the state space efficiently. In CRAFTWORLD’s FRL setting, using primitive actions requires $128.1 \times$ more episodes than options in MILKBUCKET, the only level 2 task for which ρ traces are collected. Likewise, primitive

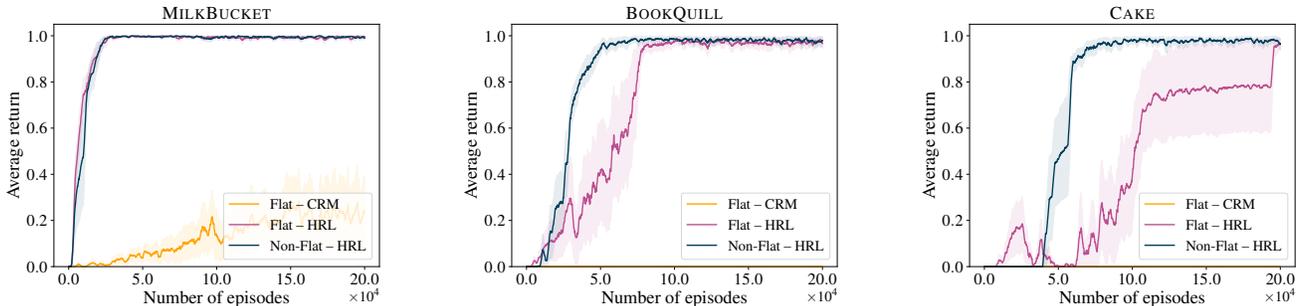


Figure 4: Learning curves for three CRAFTWORLD (FRL) tasks using handcrafted HRMs.

actions take $53.1\times$ and $10.1\times$ more episodes in OPL and WD respectively. In OP and WOD, options are not as beneficial since episodes are relatively long (1000 steps), there are no dead-ends and it is easy to observe the different propositions. See Tables 5 and 8 in Appendix D.3.1 for detailed results.

Learning the first HRMs using a *single goal trace* ($\rho = \rho_s = 1$) incurs timeouts in all CRAFTWORLD settings, thus showing the value of using many short traces instead.

6.2. Learning of Flat HRMs

Learning a flat HRM is often less scalable than learning a non-flat equivalent since (i) already learned HRMs cannot be reused, and (ii) a flat HRM usually has more states and edges (as shown in Theorem 3.7, growth can be exponential). We compare the performance of learning (from interaction) a non-flat HRM using LHRM with that of an equivalent flat HRM using LHRM, DeepSynth (Hasanbeig et al., 2021), JIRP (Xu et al., 2020) and LRM (Toro Icarte et al., 2019). LHRM and JIRP induce RMs with explicit accepting states, while DeepSynth and LRM do not. We use OP and WOD instances in CRAFTWORLD and WATERWORLD respectively.

A non-flat HRM for MILKBUCKET (level 2) is learned in 1.5 (0.2) seconds, whereas flat HRMs take longer: 3.2 (0.6) w/LHRM, 325.6 (29.7) w/DeepSynth, 17.1 (5.5) w/JIRP and 347.5 (64.5) w/LRM. LHRM and JIRP learn minimal RMs, hence producing the same RM consisting of 4 states and 3 edges. DeepSynth and LRM do not learn a minimal RM but one that is good at predicting the next possible label given the current one. In domains like ours where propositions can be observed anytime (i.e., without temporal dependencies between them), these methods tend to ‘overfit’ the input traces and output large RMs that barely reflect the task’s structure, e.g. DeepSynth learns RMs with 13.4 (0.4) states and 93.2 (1.7) edges. In contrast, methods learning minimal RMs only from observable traces may suffer from *over-generalization* (Angluin, 1980) in other domains (e.g., with temporally-dependent propositions). These observations apply to more complex tasks (i.e., involving more high-level temporal steps and multiple paths to the goal), such as

BOOK (level 2), BOOKQUILL (level 3) and CAKE (level 4). LHRM learns non-flat HRMs (e.g., see Figure 1c) for these tasks in (at most) a few minutes, while learning an informative flat HRM (e.g., see Figure 1b) is unfeasible. We refer the reader to Table 9 in Appendix D.3.2 for details.

DeepSynth, JIRP and LRM perform poorly in WATERWORLD. Unlike LHRM, these learn RMs whose edges are not labeled by formulas but proposition sets; hence, the RMs may have exponentially more edges (e.g., 64 instead of 2 for RG), and become unfeasible to learn. Indeed, flat HRM learners time out in RG&BC and RGB&CMY, while LHRM needs a few seconds (see Table 9 in Appendix D.3.2).

6.3. Policy Learning in Handcrafted HRMs

We compare the performance of policy learning in handcrafted non-flat HRMs against in flat equivalents, which are guaranteed to exist by Theorem 3.6. For fairness, the flat HRMs are minimal. To exploit the flat HRMs, we apply our HRL algorithm (Section 4) and CRM (Toro Icarte et al., 2022), which learns a Q-function over $\mathcal{S} \times \mathcal{U}$ using synthetic counterfactual experiences for each RM state. Figure 4 shows the learning curves for some CRAFTWORLD tasks in the FRL setting. The convergence rate is similar in the simplest task (MILKBUCKET), but higher for non-flat HRMs in the hardest ones. Unlike the HRL approaches, CRM does not decompose the subtask into independently solvable subtasks and, hence, deals with sparser rewards that result in a slower convergence. In the case of the HRL approaches, since both use the same set of formula option policies, differences arise from flat HRMs’ lack of modularity. Call options, which are not present in flat HRMs, form independent modules that reduce reward sparsity. MILKBUCKET involves fewer high-level steps than BOOKQUILL and CAKE, thus reward is less sparse and non-flat HRMs are not as beneficial. The efficacy of non-flat HRMs is also limited when (i) the task’s goal is reachable regardless of the chosen options (e.g., if there are no rejecting states, like in OP and FR), and (ii) the reward is not sparse, like in OPL (the grid is small) or WATERWORLD (the balls easily get near the agent). See Appendix D.3.3 for additional results.

7. Related Work

RMs and Composability. Our RMs differ from the original (Toro Icarte et al., 2018; 2022) in that (i) an RM can call other RMs, (ii) there are explicit accepting and rejecting states (Xu et al., 2020; Furelos-Blanco et al., 2021), and (iii) transitions are labeled with propositional logic formulas instead of proposition sets (Furelos-Blanco et al., 2021). Recent works *derive* RMs (or similar FSMs) from formal language specifications (Camacho et al., 2019; Araki et al., 2021) and expert demonstrations (Camacho et al., 2021), or *learn* them from experience using discrete optimization (Toro Icarte et al., 2019; Christoffersen et al., 2020), SAT solving (Xu et al., 2020; Corazza et al., 2022), active learning (Gaon & Brafman, 2020; Xu et al., 2021; Dohmen et al., 2022), state-merging (Xu et al., 2019; Gaon & Brafman, 2020), program synthesis (Hasanbeig et al., 2021) or inductive logic programming (Furelos-Blanco et al., 2021; Ardon et al., 2023). A prior way of composing RMs consists in merging the state and reward transition functions (De Giacomo et al., 2020). Other works have considered settings where the labeling function is noisy (Li et al., 2022; Verginis et al., 2022), the RM transitions and/or rewards are stochastic (Corazza et al., 2022; Dohmen et al., 2022) or defined over predicates (Zhou & Li, 2022), and multiple agents interact with the world (Neary et al., 2021; Dann et al., 2022; Ardon et al., 2023). High-probability regret bounds have been derived for RMs (Bourel et al., 2023).

Alternative methods for modeling task composability include subtask sequences (Andreas et al., 2017), context-free grammars (Chevalier-Boisvert et al., 2019), formal languages (Jothimurugan et al., 2019; Illanes et al., 2020; León et al., 2020; Wang et al., 2020) and logic-based algebras (Nangue Tasse et al., 2020).

Hierarchical RL. Our method for exploiting HRMs resembles a hierarchy of DQNs (Kulkarni et al., 2016). Akin to option discovery methods, LHRM induces a set of options from experience. While LHRM’s options are a byproduct of finding an HRM that compactly captures label traces, usual option discovery methods explicitly look for them (e.g., options that reach novel states). LHRM requires a set of propositions and tasks, which bound the number of discoverable options; similarly, some of these methods impose an explicit bound (Bacon et al., 2017; Machado et al., 2017). LHRM requires each task to be solved at least once before learning an HRM (and, hence, options), just like other methods (McGovern & Barto, 2001; Stolle & Precup, 2002). The problem of discovering options for exploration has been considered before (Bellemare et al., 2016; Machado et al., 2017; Jinnai et al., 2019; Dabney et al., 2021). While our options are not discovered for exploration, we leverage them to find goal traces in new tasks. Levy et al. (2019) learn policies from multiple hierarchical levels in parallel by training

each level as if the lower levels were optimal; likewise, we train call option policies from experiences where invoked options achieve their goal.

HRMs are close to hierarchical abstract machines (HAMs; Parr & Russell, 1997) since both are hierarchies of FSMs, but there are two core differences. First, HAMs do not have reward transition functions. Second, (H)HRMs decouple the traversal from the policies, i.e. independently of the agent’s choices, the (H)RM is followed; thus, an agent using an (H)RM must be able to interrupt its choices (see Section 4). While HAMs do not support interruption, Programmable HAMs (Andre & Russell, 2000) extend them to support it along with other program-like features. Despite the similarity, there are few works on learning HAMs (Leonetti et al., 2012) and many on learning RMs, as outlined before.

Curriculum Learning. Pierrot et al. (2019) learn hierarchies of neural programs given the level of each program, akin to our RMs’ height; likewise, Andreas et al. (2017) prioritize tasks consisting of fewer high-level steps. The ‘online’ method by Matiisen et al. (2020) also keeps an estimate of each task’s average return, but it is not applied in an HRL scenario. Wang et al. (2020) learn increasingly complex temporal logic formulas leveraging previously learned formulas using a set of templates.

8. Conclusions and Future Work

We have here proposed (1) HRMs, a *formalism* that composes RMs in a hierarchy by enabling them to call each other, (2) an HRL method that *exploits* the structure of an HRM, and (3) a curriculum-based method for *learning* a set of HRMs from traces. Non-flat HRMs have significant advantages over their flat equivalents. Theoretically, a flat equivalent of a given HRM can have exponentially more states and edges. Empirically, (i) our HRL method converges faster given a non-flat HRM instead of a flat equivalent one, and (ii) in line with the theory, learning an HRM is feasible in cases where a flat equivalent is not.

LHRM *assumes* the proposition set is known, shared dead-end indicators across tasks, and a fixed set of tasks. Relaxing these assumptions by forming the propositions from raw data, conditioning policies to dead-ends, and letting the agent propose its own composable tasks are promising directions for future work. Other directions include *non-episodic* settings and learning *globally optimal* policies over HRMs.

Acknowledgements

We thank the reviewers, as well as Hadeel Al-Negheimish, Nuri Cingillioglu, and Alex F. Spies for their comments. Anders Jonsson is partially funded by TAILOR, AGAUR SGR and Spanish grant PID2019-108141GB-I00.

References

- Andre, D. and Russell, S. J. Programmable Reinforcement Learning Agents. In *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS) Conference*, pp. 1019–1025, 2000.
- Andreas, J., Klein, D., and Levine, S. Modular Multitask Reinforcement Learning with Policy Sketches. In *Proceedings of the International Conference on Machine Learning (ICML)*, pp. 166–175, 2017.
- Angluin, D. Inductive Inference of Formal Languages from Positive Data. *Inf. Control.*, 45(2):117–135, 1980.
- Araki, B., Li, X., Vodrahalli, K., DeCastro, J. A., Fry, M. J., and Rus, D. The Logical Options Framework. In *Proceedings of the International Conference on Machine Learning (ICML)*, pp. 307–317, 2021.
- Ardon, L., Furelos-Blanco, D., and Russo, A. Learning Reward Machines in Cooperative Multi-Agent Tasks. In *Proceedings of the Neuro-Symbolic AI for Agent and Multi-Agent Systems (NeSyMAS) Workshop at the International Conference on Autonomous Agents and Multiagent Systems (AAMAS)*, 2023.
- Bacon, P., Harb, J., and Precup, D. The Option-Critic Architecture. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, pp. 1726–1734, 2017.
- Barto, A. G. and Mahadevan, S. Recent Advances in Hierarchical Reinforcement Learning. *Discrete Event Dynamic Systems*, 13(4):341–379, 2003.
- Bellemare, M. G., Srinivasan, S., Ostrovski, G., Schaul, T., Saxton, D., and Munos, R. Unifying Count-Based Exploration and Intrinsic Motivation. In *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS) Conference*, pp. 1471–1479, 2016.
- Bengio, Y., Louradour, J., Collobert, R., and Weston, J. Curriculum Learning. In *Proceedings of the International Conference on Machine Learning (ICML)*, pp. 41–48, 2009.
- Bourel, H., Jonsson, A., Maillard, O.-A., and Sadegh Talebi, M. Exploration in Reward Machines with Low Regret. In *Proceedings of the International Conference on Artificial Intelligence and Statistics (AISTATS)*, pp. 4114–4146, 2023.
- Camacho, A., Toro Icarte, R., Klassen, T. Q., Valenzano, R. A., and McIlraith, S. A. LTL and Beyond: Formal Languages for Reward Function Specification in Reinforcement Learning. In *Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI)*, pp. 6065–6073, 2019.
- Camacho, A., Varley, J., Zeng, A., Jain, D., Iscen, A., and Kalashnikov, D. Reward Machines for Vision-Based Robotic Manipulation. In *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, pp. 14284–14290, 2021.
- Chevalier-Boisvert, M., Willems, L., and Pal, S. Minimalistic Gridworld Environment for OpenAI Gym. <https://github.com/maximecb/gym-minigrid>, 2018.
- Chevalier-Boisvert, M., Bahdanau, D., Lahlou, S., Willems, L., Saharia, C., Nguyen, T. H., and Bengio, Y. BabyAI: A Platform to Study the Sample Efficiency of Grounded Language Learning. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2019.
- Christoffersen, P. J. K., Li, A. C., Toro Icarte, R., and McIlraith, S. A. Learning Symbolic Representations for Reinforcement Learning of Non-Markovian Behavior. In *Proceedings of the Knowledge Representation and Reasoning Meets Machine Learning (KR2ML) Workshop at the Advances in Neural Information Processing Systems (NeurIPS) Conference*, 2020.
- Corazza, J., Gavran, I., and Neider, D. Reinforcement Learning with Stochastic Reward Machines. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, pp. 6429–6436, 2022.
- Dabney, W., Ostrovski, G., and Barreto, A. Temporally-Extended ϵ -Greedy Exploration. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021.
- Dann, M., Yao, Y., Alechina, N., Logan, B., and Thangarajah, J. Multi-Agent Intention Progression with Reward Machines. In *Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI)*, pp. 215–222, 2022.
- De Giacomo, G., Favorito, M., Iocchi, L., Patrizi, F., and Ronca, A. Temporal Logic Monitoring Rewards via Transducers. In *Proceedings of the International Conference on Principles of Knowledge Representation and Reasoning (KR)*, pp. 860–870, 2020.
- Dietterich, T. G. Hierarchical Reinforcement Learning with the MAXQ Value Function Decomposition. *J. Artif. Intell. Res.*, 13:227–303, 2000.
- Dietterich, T. G., Domingos, P. M., Getoor, L., Muggleton, S., and Tadepalli, P. Structured machine learning: the next ten years. *Mach. Learn.*, 73(1):3–23, 2008.
- Dohmen, T., Topper, N., Atia, G. K., Beckus, A., Trivedi, A., and Velasquez, A. Inferring Probabilistic Reward

- Machines from Non-Markovian Reward Signals for Reinforcement Learning. In *Proceedings of the International Conference on Automated Planning and Scheduling (ICAPS)*, pp. 574–582, 2022.
- Eiter, T. and Gottlob, G. On the Computational Cost of Disjunctive Logic Programming: Propositional Case. *Ann. Math. Artif. Intell.*, 15(3-4):289–323, 1995.
- Furelos-Blanco, D., Law, M., Jonsson, A., Broda, K., and Russo, A. Induction and Exploitation of Subgoal Automata for Reinforcement Learning. *J. Artif. Intell. Res.*, 70:1031–1116, 2021.
- Gaon, M. and Brafman, R. I. Reinforcement Learning with Non-Markovian Rewards. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, pp. 3980–3987, 2020.
- Gelfond, M. and Kahl, Y. *Knowledge Representation, Reasoning, and the Design of Intelligent Agents: The Answer-Set Programming Approach*. Cambridge University Press, 2014.
- Gelfond, M. and Lifschitz, V. The Stable Model Semantics for Logic Programming. In *Proceedings of the International Conference and Symposium on Logic Programming (ICLP/SLP)*, pp. 1070–1080, 1988.
- Hasanbeig, M., Jeppu, N. Y., Abate, A., Melham, T., and Kroening, D. DeepSynth: Automata Synthesis for Automatic Task Segmentation in Deep Reinforcement Learning. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, pp. 7647–7656, 2021.
- Hinton, G., Srivastava, N., and Swersky, K. Neural Networks for Machine Learning - Lecture 6e - RMSprop: Divide the Gradient by a Running Average of its Recent Magnitude. https://www.cs.toronto.edu/~tijmen/csc321/slides/lecture_slides_lec6.pdf, 2012.
- Igl, M., Ciosek, K., Li, Y., Tschitschek, S., Zhang, C., Devlin, S., and Hofmann, K. Generalization in Reinforcement Learning with Selective Noise Injection and Information Bottleneck. In *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS) Conference*, pp. 13956–13968, 2019.
- Illanes, L., Yan, X., Toro Icarte, R., and McIlraith, S. A. Symbolic Plans as High-Level Instructions for Reinforcement Learning. In *Proceedings of the International Conference on Automated Planning and Scheduling (ICAPS)*, pp. 540–550, 2020.
- Jiang, M., Grefenstette, E., and Rocktäschel, T. Prioritized Level Replay. In *Proceedings of the International Conference on Machine Learning (ICML)*, pp. 4940–4950, 2021.
- Jinnai, Y., Park, J. W., Abel, D., and Konidaris, G. D. Discovering Options for Exploration by Minimizing Cover Time. In *Proceedings of the International Conference on Machine Learning (ICML)*, pp. 3130–3139, 2019.
- Jothimurugan, K., Alur, R., and Bastani, O. A Composable Specification Language for Reinforcement Learning Tasks. In *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS) Conference*, pp. 13021–13030, 2019.
- Karpathy, A. REINFORCEjs: WaterWorld demo. <http://cs.stanford.edu/people/karpathy/reinforcejs/waterworld.html>, 2015.
- Kulkarni, T. D., Narasimhan, K., Saeedi, A., and Tenenbaum, J. Hierarchical Deep Reinforcement Learning: Integrating Temporal Abstraction and Intrinsic Motivation. In *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS) Conference*, pp. 3675–3683, 2016.
- Law, M. *Inductive Learning of Answer Set Programs*. PhD thesis, Imperial College London, UK, 2018.
- Law, M., Russo, A., and Broda, K. The ILASP System for Learning Answer Set Programs, 2015. URL <https://www.ilasp.com>.
- Law, M., Russo, A., and Broda, K. Iterative Learning of Answer Set Programs from Context Dependent Examples. *Theory Pract. Log. Program.*, 16(5-6):834–848, 2016.
- Law, M., Russo, A., and Broda, K. The Meta-program Injection Feature in ILASP. Technical report, Imperial College London, June 2018. URL <https://www.doc.ic.ac.uk/~ml1909/ILASP/inject.pdf>.
- León, B. G., Shanahan, M., and Belardinelli, F. Systematic Generalisation through Task Temporal Logic and Deep Reinforcement Learning. *arXiv preprint, arXiv:2006.08767*, 2020.
- Leonetti, M., Iocchi, L., and Patrizi, F. Automatic Generation and Learning of Finite-State Controllers. In *Proceedings of the International Conference on Artificial Intelligence: Methodology, Systems, Applications (AIMSA)*, pp. 135–144, 2012.
- Levy, A., Konidaris, G. D., Jr., R. P., and Saenko, K. Learning Multi-Level Hierarchies with Hindsight. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2019.

- Li, A. C., Chen, Z., Vaezipoor, P., Klassen, T. Q., Toro Icarte, R., and McIlraith, S. A. Noisy Symbolic Abstractions for Deep RL: A case study with Reward Machines. In *Proceedings of the Deep Reinforcement Learning Workshop at the Advances in Neural Information Processing Systems (NeurIPS) Conference*, 2022.
- Machado, M. C., Bellemare, M. G., and Bowling, M. H. A Laplacian Framework for Option Discovery in Reinforcement Learning. In *Proceedings of the International Conference on Machine Learning (ICML)*, pp. 2295–2304, 2017.
- Matiisen, T., Oliver, A., Cohen, T., and Schulman, J. Teacher-Student Curriculum Learning. *IEEE Trans. Neural Networks Learn. Syst.*, 31(9):3732–3740, 2020.
- McGovern, A. and Barto, A. G. Automatic Discovery of Subgoals in Reinforcement Learning using Diverse Density. In *Proceedings of the International Conference on Machine Learning (ICML)*, pp. 361–368, 2001.
- Minton, S. Quantitative Results Concerning the Utility of Explanation-Based Learning. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, pp. 564–569, 1988.
- Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M. A., Fidjeland, A., Ostrovski, G., Petersen, S., Beattie, C., Sadik, A., Antonoglou, I., King, H., Kumaran, D., Wierstra, D., Legg, S., and Hassabis, D. Human-level control through deep reinforcement learning. *Nature*, 518(7540): 529–533, 2015.
- Nangue Tasse, G., James, S. D., and Rosman, B. A Boolean Task Algebra for Reinforcement Learning. In *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS) Conference*, pp. 9497–9507, 2020.
- Neary, C., Xu, Z., Wu, B., and Topcu, U. Reward Machines for Cooperative Multi-Agent Reinforcement Learning. In *Proceedings of the International Conference on Autonomous Agents and Multiagent Systems (AAMAS)*, pp. 934–942, 2021.
- Parr, R. and Russell, S. J. Reinforcement Learning with Hierarchies of Machines. In *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS) Conference*, pp. 1043–1049, 1997.
- Pierrot, T., Ligner, G., Reed, S. E., Sigaud, O., Perrin, N., Laterre, A., Kas, D., Beguir, K., and de Freitas, N. Learning Compositional Neural Programs with Recursive Tree Search and Planning. In *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS) Conference*, pp. 14646–14656, 2019.
- Sidor, S. Reinforcement Learning with Natural Language Signals. Master’s thesis, Massachusetts Institute of Technology, 2016.
- Sipser, M. *Introduction to the Theory of Computation*. PWS Publishing Company, 1997.
- Stolle, M. and Precup, D. Learning Options in Reinforcement Learning. In *Proceedings of the International Symposium on Abstraction, Reformulation and Approximation (SARA)*, pp. 212–223, 2002.
- Sutton, R. S. and Barto, A. G. *Reinforcement Learning: An Introduction*. MIT Press, 2018.
- Sutton, R. S., Precup, D., and Singh, S. P. Intra-Option Learning about Temporally Abstract Actions. In *Proceedings of the International Conference on Machine Learning (ICML)*, pp. 556–564, 1998.
- Sutton, R. S., Precup, D., and Singh, S. P. Between MDPs and Semi-MDPs: A Framework for Temporal Abstraction in Reinforcement Learning. *Artif. Intell.*, 112(1-2):181–211, 1999.
- Toro Icarte, R., Klassen, T. Q., Valenzano, R. A., and McIlraith, S. A. Using Reward Machines for High-Level Task Specification and Decomposition in Reinforcement Learning. In *Proceedings of the International Conference on Machine Learning (ICML)*, pp. 2112–2121, 2018.
- Toro Icarte, R., Waldie, E., Klassen, T. Q., Valenzano, R. A., Castro, M. P., and McIlraith, S. A. Learning Reward Machines for Partially Observable Reinforcement Learning. In *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS) Conference*, pp. 15497–15508, 2019.
- Toro Icarte, R., Klassen, T. Q., Valenzano, R., and McIlraith, S. A. Reward Machines: Exploiting Reward Function Structure in Reinforcement Learning. *J. Artif. Intell. Res.*, 73:173–208, 2022.
- van Hasselt, H., Guez, A., and Silver, D. Deep Reinforcement Learning with Double Q-Learning. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, pp. 2094–2100, 2016.
- Verginis, C. K., Köprülü, C., Chinchali, S., and Topcu, U. Joint Learning of Reward Machines and Policies in Environments with Partially Known Semantics. *arXiv preprint, arXiv:2204.11833*, 2022.
- Wang, G., Trimbach, C., Lee, J. K., Ho, M. K., and Littman, M. L. Teaching a Robot Tasks of Arbitrary Complexity via Human Feedback. In *Proceedings of the ACM/IEEE International Conference on Human-Robot Interaction (HRI)*, pp. 649–657, 2020.

- Xu, Z., Gavran, I., Ahmad, Y., Majumdar, R., Neider, D., Topcu, U., and Wu, B. Joint Inference of Reward Machines and Policies for Reinforcement Learning. *arXiv preprint*, arXiv:1909.05912, 2019.
- Xu, Z., Gavran, I., Ahmad, Y., Majumdar, R., Neider, D., Topcu, U., and Wu, B. Joint Inference of Reward Machines and Policies for Reinforcement Learning. In *Proceedings of the International Conference on Automated Planning and Scheduling (ICAPS)*, pp. 590–598, 2020.
- Xu, Z., Wu, B., Ojha, A., Neider, D., and Topcu, U. Active Finite Reward Automaton Inference and Reinforcement Learning Using Queries and Counterexamples. In *Proceedings of the International Cross-Domain Conference for Machine Learning and Knowledge Extraction (CD-MAKE)*, pp. 115–135, 2021.
- Zhou, W. and Li, W. A Hierarchical Bayesian Approach to Inverse Reinforcement Learning with Symbolic Reward Machines. In *Proceedings of the International Conference on Machine Learning (ICML)*, pp. 27159–27178, 2022.

A. Formalism Details

In this appendix, we extend Example 4 by showing all the intermediate steps (Appendix A.1), and provide the proofs for Theorems 3.6 and 3.7 (Appendix A.2).

A.1. Hierarchy Traversal Example

The HRM in Figure 1c accepts trace $\lambda = \langle \{\ddagger\}, \{\heartsuit\}, \{\}, \{\clubsuit\}, \{\spadesuit\}, \{\clubsuit\}, \{\heartsuit\} \rangle$, whose traversal is $H(\lambda) = \langle v_0, v_1, v_2, v_3, v_4, v_5, v_6 \rangle$, where:

$$\begin{aligned}
 v_0 &= \langle M_0, u_0^0, \top, [] \rangle, \\
 v_1 &= \delta_H(v_0, \{\ddagger\}) \\
 &= \delta_H(\langle M_0, u_0^0, \top, [] \rangle, \{\ddagger\}) \\
 &= \delta_H(\langle M_1, u_1^0, \neg\clubsuit, [\langle u_0^0, u_0^1, M_0, M_1, \neg\clubsuit, \top \rangle] \rangle, \{\ddagger\}) \\
 &= \delta_H(\langle M_\top, u_\top^0, \neg\clubsuit \wedge \ddagger, [\langle u_0^0, u_0^1, M_0, M_1, \neg\clubsuit, \top \rangle, \langle u_1^0, u_1^1, M_1, M_\top, \ddagger, \neg\clubsuit \rangle] \rangle, \{\ddagger\}) \\
 &= \delta_H(\langle M_1, u_1^1, \top, [\langle u_0^0, u_0^1, M_0, M_1, \neg\clubsuit, \top \rangle] \rangle, \perp) \\
 &= \langle M_1, u_1^1, \top, [\langle u_0^0, u_0^1, M_0, M_1, \neg\clubsuit, \top \rangle] \rangle, \\
 v_2 &= \delta_H(v_1, \{\heartsuit\}) \\
 &= \delta_H(\langle M_1, u_1^1, \top, [\langle u_0^0, u_0^1, M_0, M_1, \neg\clubsuit, \top \rangle] \rangle, \{\heartsuit\}) \\
 &= \delta_H(\langle M_\top, u_\top^0, \heartsuit, [\langle u_0^0, u_0^1, M_0, M_1, \neg\clubsuit, \top \rangle, \langle u_1^1, u_1^A, M_1, M_\top, \heartsuit, \top \rangle] \rangle, \{\heartsuit\}) \\
 &= \delta_H(\langle M_1, u_1^A, \top, [\langle u_0^0, u_0^1, M_0, M_1, \neg\clubsuit, \top \rangle] \rangle, \perp) \\
 &= \delta_H(\langle M_0, u_0^1, \top, [] \rangle, \perp) \\
 &= \langle M_0, u_0^1, \top, [] \rangle, \\
 v_3 &= \delta_H(v_2, \{\}) \\
 &= \delta_H(\langle M_0, u_0^1, \top, [] \rangle, \{\}) \\
 &= \langle M_0, u_0^1, \top, [] \rangle, \\
 v_4 &= \delta_H(v_3, \{\clubsuit\}) \\
 &= \delta_H(\langle M_0, u_0^1, \top, [] \rangle, \{\clubsuit\}) \\
 &= \delta_H(\langle M_2, u_2^0, \top, [\langle u_0^1, u_0^3, M_0, M_2, \top, \top \rangle] \rangle, \{\clubsuit\}) \\
 &= \delta_H(\langle M_\top, u_\top^0, \clubsuit, [\langle u_0^1, u_0^3, M_0, M_2, \top, \top \rangle, \langle u_2^0, u_2^1, M_2, M_\top, \clubsuit, \top \rangle] \rangle, \{\clubsuit\}) \\
 &= \delta_H(\langle M_2, u_2^1, \top, [\langle u_0^1, u_0^3, M_0, M_2, \top, \top \rangle] \rangle, \perp) \\
 &= \langle M_2, u_2^1, \top, [\langle u_0^1, u_0^3, M_0, M_2, \top, \top \rangle] \rangle, \\
 v_5 &= \delta_H(v_4, \{\heartsuit\}) \\
 &= \delta_H(\langle M_2, u_2^1, \top, [\langle u_0^1, u_0^3, M_0, M_2, \top, \top \rangle] \rangle, \{\heartsuit\}) \\
 &= \delta_H(\langle M_\top, u_\top^0, \heartsuit, [\langle u_0^1, u_0^3, M_0, M_2, \top, \top \rangle, \langle u_2^1, u_2^A, M_2, M_\top, \heartsuit, \top \rangle] \rangle, \{\heartsuit\}) \\
 &= \delta_H(\langle M_2, u_2^A, \top, [\langle u_0^1, u_0^3, M_0, M_2, \top, \top \rangle] \rangle, \perp) \\
 &= \delta_H(\langle M_0, u_0^3, \top, [] \rangle, \perp) \\
 &= \langle M_0, u_0^3, \top, [] \rangle, \\
 v_6 &= \delta_H(v_5, \{\clubsuit\}) \\
 &= \delta_H(\langle M_0, u_0^3, \top, [] \rangle, \{\clubsuit\}) \\
 &= \delta_H(\langle M_\top, u_\top^0, \clubsuit, [\langle u_0^3, u_0^A, M_0, M_\top, \clubsuit, \top \rangle] \rangle, \{\clubsuit\}) \\
 &= \delta_H(\langle M_0, u_0^A, \top, [] \rangle, \perp) \\
 &= \langle M_0, u_0^A, \top, [] \rangle.
 \end{aligned}$$

A.2. Equivalence to Flat Hierarchies of Reward Machines

In this section, we prove the theorems introduced in Section 3 regarding the equivalence of an arbitrary HRM to a flat HRM.

A.2.1. PROOF OF THEOREM 3.6

We formally show that any HRM can be transformed into an equivalent one consisting of a single non-leaf RM. The latter HRM type is called *flat* since there is a single hierarchy level.

Definition A.1. Given an HRM $H = \langle \mathcal{M}, M_r, \mathcal{P} \rangle$, a constituent RM $M_i \in \mathcal{M}$ is *flat* if its height h_i is 1.

Definition A.2. An HRM $H = \langle \mathcal{M}, M_r, \mathcal{P} \rangle$ is *flat* if the root RM M_r is flat.

We now define what it means for two HRMs to be equivalent. This definition is based on that used in automaton theory (Sipser, 1997).

Definition A.3. Given a set of propositions \mathcal{P} and a labeling function l , two HRMs $H = \langle \mathcal{M}, M_r, \mathcal{P} \rangle$ and $H' = \langle \mathcal{M}', M'_r, \mathcal{P} \rangle$ are *equivalent* if for any label trace λ one of the following conditions holds: (i) both HRMs accept λ , (ii) both HRMs reject λ , or (iii) neither of the HRMs accepts or rejects λ .

We now have all the required definitions to prove Theorem 3.6, which is restated below.

Theorem 3.6. *Given an HRM H , there exists an equivalent flat HRM \bar{H} .*

To prove the theorem, we introduce an algorithm for flattening any HRM. Without loss of generality, we work on the case of an HRM with two hierarchy levels; that is, an HRM consisting of a root RM that calls flat RMs. Note that an HRM with an arbitrary number of levels can be flattened by considering the RMs in two levels at a time. We start flattening RMs in the second level (i.e., with height 2), which use RMs in the first level (by definition, these are already flat), and once the second level RMs are flat, we repeat the process with the levels above until the root is reached. This process is applicable since, by assumption, the hierarchies do not have cyclic dependencies (including recursion). For simplicity, we use the MDP reward assumption made in Section 2, i.e. the reward transition function of any RM M_i is $r_i(u, u') = \mathbb{1}[u \notin \mathcal{U}_i^A \wedge u' \in \mathcal{U}_i^A]$ like in Section 4. However, the proof below could be adapted to arbitrary definitions of $r_i(u, u')$.

Preliminary Transformation Algorithm. Before proving Theorem 3.6, we introduce an intermediate step that transforms a flat HRM into an equivalent one that takes contexts with which it may be called into account. Remember that a call to an RM is associated with a context. In the case of two-level HRMs such as the ones we are considering in this flattening process, the context and the exit condition from the called flat RM must be satisfied. Crucially, the context must only be satisfied at the time of the call; that is, it only lasts for a single transition. Therefore, if we revisit the initial state of the called RM by taking an edge to it, the context should not be checked anymore.

To make the need for this transformation clearer, we use the HRM illustrated in Figure 5a. The flattening algorithm described later embeds the called RM into the calling one; crucially, the context of the call is taken into account by putting it in conjunction with the outgoing edges from the initial state of the called RM.³ Figure 5b is a flat HRM obtained using the flattening algorithm; however, it does not behave like the HRM in Figure 5a. Following the definition of the hierarchical transition function δ_H , the context of a call only lasts for a single transition in the called RM in Figure 5a (i.e., $a \wedge \neg c$ is only checked when M_1 is started), but the context is kept permanently in Figure 5b, which is problematic if we go back to the initial state at some point. We later come back to this example after presenting the transformation algorithm.

To deal with the situation above, we need to transform an RM to ensure that contexts are only checked once from the initial state. We describe this transformation as follows. Given a flat HRM $H = \langle \mathcal{M}, M_r, \mathcal{P} \rangle$ with root $M_r = \langle \mathcal{U}_r, \mathcal{P}, \varphi_r, r_r, u_r^0, \mathcal{U}_r^A, \mathcal{U}_r^R \rangle$, we construct a new HRM $H' = \langle \mathcal{M}', M'_r, \mathcal{P} \rangle$ with root $M'_r = \langle \mathcal{U}'_r, \mathcal{P}, \varphi'_r, r'_r, u_r^0, \mathcal{U}_r^A, \mathcal{U}_r^R \rangle$ such that:

- $\mathcal{U}'_r = \mathcal{U}_r \cup \{\hat{u}_r^0\}$, where \hat{u}_r^0 plays the role of the initial state after the first transition is taken.
- The state transition function φ'_r is built by copying φ_r and applying the following changes:
 1. Remove the edges to the actual initial state from any state $v \in \mathcal{U}'_r$: $\varphi'_r(v, u_r^0, M_\top) = \perp$. Note that since the RM is flat, the only callable RM is the leaf M_\top .

³We refer the reader to the ‘Flattening Algorithm’ description introduced later for specific details.

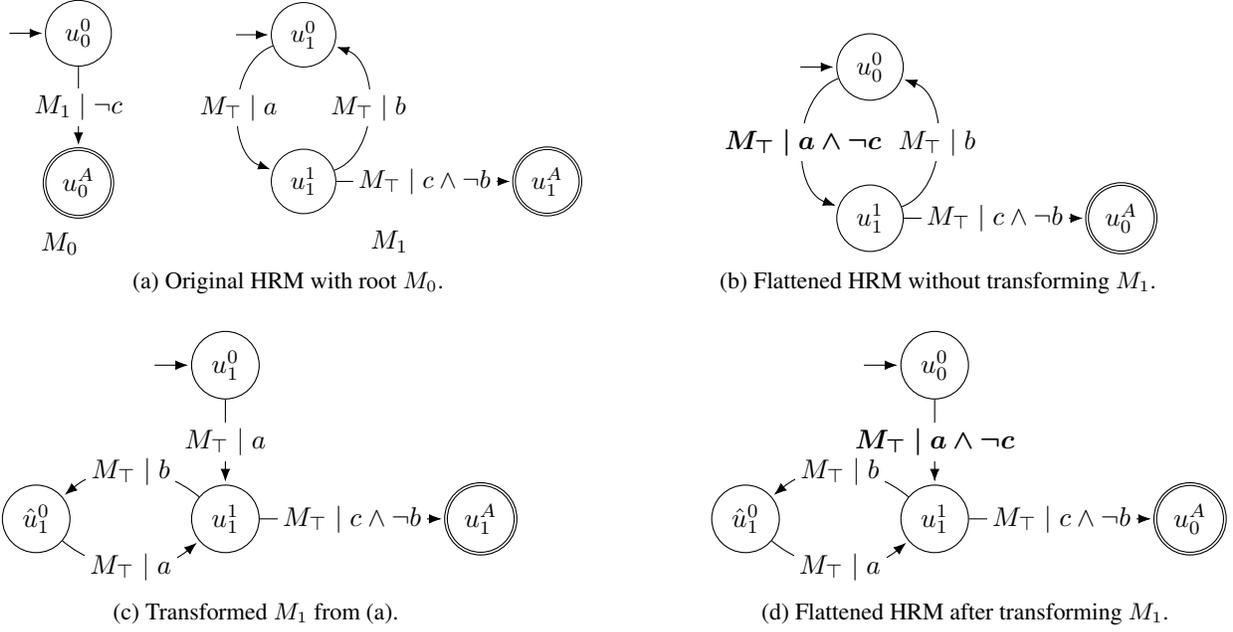


Figure 5: Example to justify the need for the preliminary transformation algorithm.

2. Add edges to the dummy initial state \hat{u}_r^0 from all states $v \in \mathcal{U}_r'$ that had an edge to the actual initial state: $\varphi_r'(v, \hat{u}_r^0, M_\top) = \varphi_r(v, u_r^0, M_\top)$.
3. Add edges from the dummy initial state \hat{u}_r^0 to all those states $v \in \mathcal{U}_r'$ that the actual initial state u_r^0 points to: $\varphi_r'(\hat{u}_r^0, v, M_\top) = \varphi_r(u_r^0, v, M_\top)$.

- The reward transition function $r_r'(u, u') = \mathbb{1}[u \notin \mathcal{U}_r^A \wedge u' \in \mathcal{U}_r^A]$ is defined as stated at the beginning of the section.

The HRM H' is such that $\mathcal{M}' = \{M_r', M_\top\}$. Note that this transformation is only required in HRMs where the RMs have initial states with incoming edges.

We now prove that this transformation is correct; that is, the HRMs are equivalent. There are two cases depending on whether the initial state has incoming edges or not. First, if the initial state u_r^0 does not have incoming edges, step 1 does not remove any edges going to u_r^0 , and step 2 does not add any edges going to \hat{u}_r^0 , making it unreachable. Even though edges from \hat{u}_r^0 to other states may be added, it is irrelevant since it is unreachable. Therefore, we can safely say that in this case, the transformed HRM is equivalent to the original one. Second, if the initial state has incoming edges, we prove equivalence by examining the traversals $H(\lambda)$ and $H'(\lambda)$ for the original HRM $H = \langle \mathcal{M}, M_r, \mathcal{P} \rangle$ and the transformed one $H' = \langle \mathcal{M}', M_r', \mathcal{P} \rangle$ given a generic label trace $\lambda = \langle \mathcal{L}_0, \dots, \mathcal{L}_n \rangle$. By construction, both $H(\lambda)$ and $H'(\lambda)$ will be identical until reaching a state w with an outgoing transition to u_r^0 in the case of H and the dummy initial state \hat{u}_r^0 in the case of H' . More specifically, upon reaching w and satisfying an outgoing formula to the aforementioned states, the traversals are:

$$\begin{aligned} H(\lambda) &= \langle \langle M_r, u_r^0, \top, [] \rangle, \dots, \langle M_r, w, \top, [] \rangle \rangle, \\ H'(\lambda) &= \langle \langle M_r', u_r^0, \top, [] \rangle, \dots, \langle M_r', w, \top, [] \rangle \rangle. \end{aligned}$$

By construction, state w is in both HRMs, and both of the aforementioned transitions from this state are associated with the same formula, i.e. $\varphi_r(w, u_r^0, M_\top) = \varphi_r'(w, \hat{u}_r^0, M_\top)$. Therefore, if one of them is satisfied, the other will be too, and the traversals will become:

$$\begin{aligned} H(\lambda) &= \langle \langle M_r, u_r^0, \top, [] \rangle, \dots, \langle M_r, w, \top, [] \rangle, \langle M_r, u_r^0, \top, [] \rangle \rangle, \\ H'(\lambda) &= \langle \langle M_r', u_r^0, \top, [] \rangle, \dots, \langle M_r', w, \top, [] \rangle, \langle M_r', \hat{u}_r^0, \top, [] \rangle \rangle. \end{aligned}$$

We stay in u_r^0 and \hat{u}_r^0 until a transition to a state w' is satisfied. By construction, w' is in both HRMs and the same formula

is satisfied, i.e., $\varphi_r(u_r^0, w', M_\top) = \varphi'_r(\hat{u}^0, w', M_\top)$. The hierarchy traversals then become:

$$\begin{aligned} H(\lambda) &= \langle \langle M_r, u_r^0, \top, [] \rangle, \dots, \langle M_r, w, \top, [] \rangle, \langle M_r, u_r^0, \top, [] \rangle, \dots, \langle M_r, u_r^0, \top, [] \rangle, \langle M_r, w', \top, [] \rangle \rangle, \\ H'(\lambda) &= \langle \langle M'_r, u_r^0, \top, [] \rangle, \dots, \langle M'_r, w, \top, [] \rangle, \langle M'_r, \hat{u}_r^0, \top, [] \rangle, \dots, \langle M'_r, \hat{u}_r^0, \top, [] \rangle, \langle M'_r, w', \top, [] \rangle \rangle. \end{aligned}$$

From here both traversals will be the same until transitions to u_r^0 and \hat{u}_r^0 are respectively satisfied again (if any) in H and H' . Clearly, the only change in $H(\lambda)$ with respect to $H'(\lambda)$ (except for the different roots) is that the hierarchy states of the form $\langle M'_r, \hat{u}_r^0, \top, [] \rangle$ in the latter appear as $\langle M_r, u_r^0, \top, [] \rangle$ in the former. We now check if the equivalence conditions from Definition A.3 hold:

- If $H(\lambda)$ ends with state u_r^0 , $H'(\lambda)$ ends with state \hat{u}_r^0 following the reasoning above. By construction, neither of these states is accepting or rejecting; therefore, neither of these HRMs accepts or rejects λ .
- If $H(\lambda)$ ends with state w , $H'(\lambda)$ will also end with this state following the reasoning above. Therefore, if w is an accepting state, both HRMs accept λ ; if w is a rejecting state, both HRMs reject λ ; and if w is not an accepting or rejecting state, neither of the HRMs accepts or rejects λ .

Since all equivalence conditions are satisfied for any trace λ , H and H' are equivalent.

Figure 5c exemplifies the output of the transformation algorithm given M_1 in Figure 5a as input, whereas Figure 5d is the output of the flattening algorithm discussed next, which properly handles the context unlike the HRM in Figure 5b.

Flattening Algorithm. We describe the algorithm for flattening an HRM. As previously stated, we assume without loss of generality that the HRM to be flattened consists of two hierarchy levels (i.e., the root calls flat RMs). We also assume that the flat RMs have the form produced by the previously presented transformation algorithm.

Given an HRM $H = \langle \mathcal{M}, M_r, \mathcal{P} \rangle$ with root $M_r = \langle \mathcal{U}_r, \mathcal{P}, \varphi_r, r_r, u_r^0, \mathcal{U}_r^A, \mathcal{U}_r^R \rangle$, we build a flat RM $\bar{M}_r = \langle \bar{\mathcal{U}}_r, \mathcal{P}, \bar{\varphi}_r, \bar{r}_r, \bar{u}_r^0, \bar{\mathcal{U}}_r^A, \bar{\mathcal{U}}_r^R \rangle$ using the following steps:

1. Copy the sets of states and initial state from M_r (i.e., $\bar{\mathcal{U}}_r = \mathcal{U}_r$, $\bar{u}_r^0 = u_r^0$, $\bar{\mathcal{U}}_r^A = \mathcal{U}_r^A$, $\bar{\mathcal{U}}_r^R = \mathcal{U}_r^R$).
2. Loop through the non-false entries of the transition function φ_r and decide what to copy. That is, for each triplet (u, u', M_j) where $u, u' \in \mathcal{U}_r$ and $M_j \in \mathcal{M}$ such that $\varphi_r(u, u', M_j) \neq \perp$:
 - (a) If $M_j = M_\top$ (i.e., the called RM is the leaf), we copy the transition: $\bar{\varphi}_r(u, u', M_\top) = \varphi_r(u, u', M_\top)$.
 - (b) If $M_j \neq M_\top$, we embed the transition function of $M_j = \langle \mathcal{U}_j, \mathcal{P}, \varphi_j, r_j, u_j^0, \mathcal{U}_j^A, \mathcal{U}_j^R \rangle$ into \bar{M}_r . Remember that M_j is flat. To do so, we run the following steps:
 - i. Update the set of states by adding all non-initial and non-accepting states from M_j . Similarly, the set of rejecting states is also updated by adding all rejecting states of the called RM. The initial and accepting states from M_j are unimportant: their roles are played by u and u' respectively. In contrast, the rejecting states are important since, by assumption, they are global. Note that the added states v are renamed to $v_{u,u',j}$ in order to take into account the edge being embedded: if the same state v was reused for another edge, then we would not be able to distinguish them.

$$\begin{aligned} \bar{\mathcal{U}}_r &= \bar{\mathcal{U}}_r \cup \{v_{u,u',j} \mid v \in (\mathcal{U}_j \setminus (\{u_j^0\} \cup \mathcal{U}_j^A))\}, \\ \bar{\mathcal{U}}_r^R &= \bar{\mathcal{U}}_r^R \cup \{v_{u,u',j} \mid v \in \mathcal{U}_j^R\}. \end{aligned}$$

- ii. Embed the transition function φ_j of M_j into $\bar{\varphi}_r$. Since M_j is flat, we can make copies of the transitions straightaway: the only important thing is to check whether these transitions involve initial or accepting states which, as stated before, are going to be replaced by u and u' accordingly. Given a triplet (v, w, M_\top) such that $v, w \in \mathcal{U}_j$ and for which $\varphi_j(v, w, M_\top) = \phi$ and $\phi \neq \perp$ we update $\bar{\varphi}_r$ as follows:⁴
 - A. If $v = u_j^0$ and $w \notin \mathcal{U}_j^A$, then $\bar{\varphi}_r(u, w_{u,u',j}, M_\top) = \text{DNF}(\phi \wedge \varphi_r(u, u', M_j))$. The initial state of M_j has been substituted by u , we use the clone of w associated with the call $(w_{u,u',j})$, and append the context of the call to M_j to the formula ϕ .

⁴We do not cover the case where v is an accepting state since, by assumption, there are no outgoing transitions from it. In the case of rejecting states, we keep all of them as explained in the previous case and, therefore, there are no substitutions to be made. We also do not cover the case where $w = u_j^0$ since the input flat machines never have edges to their initial states, but to the dummy initial state.

- B. If $v = u_j^0$ and $w \in \mathcal{U}_j^A$, then $\bar{\varphi}_r(u, u', M_\top) = \text{DNF}(\phi \wedge \varphi_r(u, u', M_j))$. Like the previous case but performing two substitutions: u replaces v and u' replaces w . The context is appended since it is a transition from the initial state of M_j .
- C. If $v \neq u_j^0$ and $w \in \mathcal{U}_j^A$, then $\bar{\varphi}_r(v_{u,u',j}, u', M_\top) = \phi$. We substitute the accepting state w by u' , and use the clone of v associated with the call $(v_{u,u',j})$. This time the call's context is not added since v is not the initial state of M_j .
- D. If none of the previous cases holds, there are no substitutions to be made nor contexts to be taken into account. Hence, $\bar{\varphi}_r(v_{u,u',j}, w_{u,u',j}, M_\top) = \phi$. We just use the clones of v and w corresponding to the call $(v_{u,u',j}$ and $w_{u,u',j})$.
3. We apply the transformation algorithm we described before, and form a new flat HRM $\bar{H} = \langle \{\bar{M}_r, M_\top\}, \bar{M}_r, \mathcal{P} \rangle$ with the flattened (and transformed) root \bar{M}_r .

The reward transition function $r'_r(u, u') = \mathbb{1}[u \notin \bar{\mathcal{U}}_r^A \wedge u' \in \bar{\mathcal{U}}_r^A]$ is defined as stated at the beginning of the section. Note that u might not necessarily be a state of the non-flat root, but derived from an RM with lower height.

We now have everything to prove the previous theorem. Without loss of generality and for simplicity, we assume that the transformation algorithm has not been applied over the flattened root (we have already shown that the transformation produces an equivalent flat machine).

Theorem 3.6. *Given an HRM H , there exists an equivalent flat HRM \bar{H} .*

Proof. Let us assume that an HRM $\bar{H} = \langle \bar{\mathcal{M}}, \bar{M}_r, \mathcal{P} \rangle$, where $\bar{M}_r = \langle \bar{\mathcal{U}}_r, \mathcal{P}, \bar{\varphi}_r, \bar{r}_r, \bar{u}_r^0, \bar{\mathcal{U}}_r^A, \bar{\mathcal{U}}_r^R \rangle$, is a flat HRM that results from applying the flattening algorithm on an HRM $H = \langle \mathcal{M}, M_r, \mathcal{P} \rangle$, where $M_r = \langle \mathcal{U}_r, \mathcal{P}, \varphi_r, r_r, u_r^0, \mathcal{U}_r^A, \mathcal{U}_r^R \rangle$. For these HRMs to be equivalent, any label trace $\lambda = \langle \mathcal{L}_0, \dots, \mathcal{L}_n \rangle$ must satisfy one of the conditions in Definition A.3. To prove the equivalence, we examine the hierarchy traversals $H(\lambda)$ and $\bar{H}(\lambda)$ given a generic label trace λ .

Let $u \in \mathcal{U}_r$ be a state in the root M_r of H and let $\varphi_r(u, u', M_\top)$ be a satisfied transition from that state. By construction, u is also in the root \bar{M}_r of the flat hierarchy \bar{H} , and \bar{M}_r has an identical transition $\bar{\varphi}_r(u, u', M_\top)$, which must also be satisfied. If the hierarchy states are $\langle M_r, u, \top, [] \rangle$ and $\langle \bar{M}_r, u, \top, [] \rangle$ for H and \bar{H} respectively, then the next hierarchy states upon application of δ_H will be $\langle M_r, u', \top, [] \rangle$ and $\langle \bar{M}_r, u', \top, [] \rangle$. Therefore, both HRMs behave equivalently when calls to the leaf RM are made.

We now examine what occurs when a non-leaf RM is called in H . Let $\varphi_r(u, u', M_j)$ be a satisfied transition in M_r , and let $\varphi_j(u_j^0, w, M_\top)$ be a satisfied transition from M_j 's initial state. By construction, \bar{M}_r contains a transition whose associated formula is the conjunction of the previous two, i.e. $\varphi_r(u, u', M_j) \wedge \varphi_j(u_j^0, w, M_\top)$. Now, the hierarchy traversals will be different depending on w :

- If $w \notin \mathcal{U}_j^A$ (i.e., w is not an accepting state of M_j), by construction, \bar{M}_r contains the transition $\bar{\varphi}_r(u, w_{u,u',j}, M_\top) = \varphi_r(u, u', M_j) \wedge \varphi_j(u_j^0, w, M_\top)$. If the current hierarchy states are (the equivalent) $\langle M_r, u, \top, [] \rangle$ and $\langle \bar{M}_r, u, \top, [] \rangle$ for H and \bar{H} , then the next hierarchy states upon application of δ_H are $\langle M_j, w, \top, [(u, u', M_r, M_j, \varphi_r(u, u', M_j), \top)] \rangle$ and $\langle \bar{M}_r, w_{u,u',j}, \top, [] \rangle$. These hierarchy states are equivalent since $w_{u,u',j}$ is a clone of w that saves all the call information (i.e., a call to machine M_j for transitioning from u to u').
- If $w \in \mathcal{U}_j^A$ (i.e., w is an accepting state of M_j), by construction, \bar{M}_r contains the transition $\bar{\varphi}_r(u, u', M_\top) = \varphi_r(u, u', M_j) \wedge \varphi_j(u_j^0, w, M_\top)$. If the current hierarchy states are (the equivalent) $\langle M_r, u, \top, [] \rangle$ and $\langle \bar{M}_r, u, \top, [] \rangle$ for H and \bar{H} , then the next hierarchy states upon application of δ_H are $\langle M_r, u', \top, [] \rangle$ and $\langle \bar{M}_r, u', \top, [] \rangle$. These hierarchy states are clearly equivalent since the machine states are exactly the same.

We now check the case in which we are inside a called RM. Let $\varphi_r(u, u', M_j)$ be the transition that caused H to start running M_j , and let $\varphi_j(v, w, M_\top)$ be a satisfied transition within M_j such that $v \neq u_j^0$. By construction, \bar{M}_r contains a transition associated with the same formula $\varphi_j(v, w, M_\top)$. The hierarchy traversals vary depending on w :

- If $w \notin \mathcal{U}_j^A$ (i.e., w is not an accepting state of M_j), by construction, \bar{M}_r contains the transition $\bar{\varphi}_r(v_{u,u',j}, w_{u,u',j}, M_\top) = \varphi_j(v, w, M_\top)$. For the transition to be taken in H , the hierarchy state must be $\langle M_j, v, \top, [(u, u', M_r, M_j, \varphi_r(u, u', M_j), \top)] \rangle$, whereas in \bar{H} it will be $\langle \bar{M}_r, v_{u,u',j}, \top, [] \rangle$. These hierarchy states

are clearly equivalent: $v_{u,u',j}$ is a clone of v that saves all information related to the call being made (the called machine, and the starting and resulting states in the transition). Upon application of δ_H , the hierarchy states will remain equivalent: $\langle M_j, w, \top, [\langle u, u', M_r, M_j, \varphi_r(u, u', M_j), \top \rangle] \rangle$ and $\langle \bar{M}_r, w_{u,u',j}, \top, [] \rangle$ (again $w_{u,u',j}$ saves all the call information, just like the stack).

- If $w \in \mathcal{U}_j^A$ (i.e., w is an accepting state of M_j), by construction, \bar{M}_r contains the transition $\bar{\varphi}_r(v_{u,u',j}, u', M_\top) = \varphi_j(v, w, M_\top)$. This case corresponds to that where control is returned to the calling RM. Like in the previous case, for the transition to be taken in H , the hierarchy state must be $\langle M_j, v, \top, [\langle u, u', M_r, M_j, \varphi_r(u, u', M_j), \top \rangle] \rangle$, whereas in \bar{H} it will be $\langle \bar{M}_r, v_{u,u',j}, \top, [] \rangle$. The resulting hierarchy states then become $\langle M_r, u', \top, [] \rangle$ and $\langle \bar{M}_r, u', \top, [] \rangle$ respectively, which are clearly equivalent (the state is exactly the same and both come from equivalent hierarchy states).

We have shown both HRMs have equivalent traversals for any given trace, implying that both will accept, reject, or not accept nor reject a trace. Therefore, the HRMs are equivalent. \square

Figure 6a shows the result of applying the flattening algorithm on the BOOK HRM shown in Figure 1c. Note that the resulting HRM can be compressed: there are two states having an edge with the same label to a specific state. Indeed, the presented algorithm might not produce the smallest possible flat equivalent. Figure 6b shows the resulting compressed HRM, which is like Figure 1b but naming the states following the algorithm for clarity. Estimating how much a flat HRM (or any HRM) can be compressed and designing an algorithm to perform such compression are left as future work.

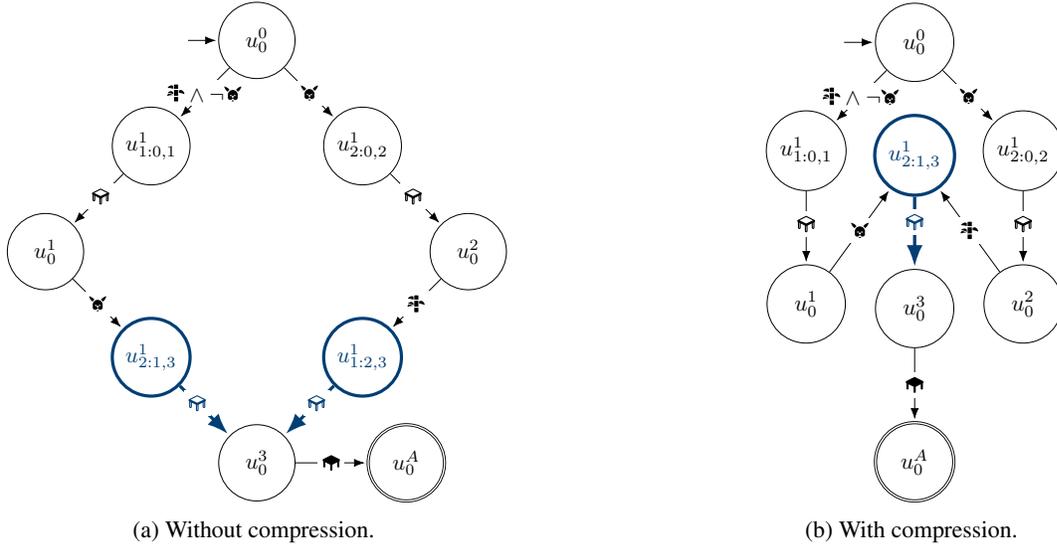


Figure 6: Results of flattening the HRM in Figure 1c. The notation $u_{j:x,y}^i$ denotes the i -th state of RM j in the call between states x and y in the parent RM. Note that x and y appear only if that state comes from a called RM. The blue states and edges in (a) can be compressed as shown in (b).

A.2.2. PROOF OF THEOREM 3.7

We prove Theorem 3.7 by first characterizing an HRM H using a set of abstract parameters. Then, we describe how the number of states and edges in an HRM and its corresponding flat equivalent are computed, and use these quantities to give an example for which the theorem holds. The parameters are the following:

- The height of the root h_r .
- The number of RMs with height i , $N^{(i)}$.
- The number of states in RMs with height i , $U^{(i)}$.
- The number of edges from each state in RMs with height i , $E^{(i)}$.

We assume that (i) RMs with height i only call RMs with height $i - 1$; (ii) all RMs have a single accepting state and no rejecting states; (iii) all RMs except for the root are called; and (iv) the HRM is well-formed (i.e., it behaves deterministically and there are no cyclic dependencies). Note that $N^{(h_r)} = 1$ since there is a single root. Assumption (i) can be made since for the root to have height h_r we need it to call at least one RM with height $h_r - 1$. Considering that all called RMs have the same height simplifies the analysis since we can characterize the RMs at each height independently. Assumption (ii) is safe to make since a single accepting state is enough, and helps simplify the counting since only some RMs could have rejecting states. Assumption (iii) ensures that the flat HRM will comprise all RMs in the original HRM. This is also a fair assumption: if a given RM is not called by any RM in the hierarchy, we could remove it beforehand.

The number of states $|H|$ in the HRM H is obtained by summing the number of states of each RM:

$$|H| = \sum_{i=1}^{h_r} N^{(i)} U^{(i)}.$$

The number of states $|\bar{H}|$ in the flat HRM \bar{H} is given by the number of states in the flattened root RM

$$|\bar{H}| = \bar{U}^{(h_r)},$$

where $\bar{U}^{(i)}$ is the number of states in the flattened representation of an RM with height i , which is recursively defined as:

$$\bar{U}^{(i)} = \begin{cases} U^{(i)} & \text{if } i = 1, \\ U^{(i)} + (\bar{U}^{(i-1)} - 2) (U^{(i)} - 1) E^{(i)} & \text{if } i > 1. \end{cases}$$

That is, the number of states in a flattened RM with height i has all states that the non-flat HRM had. In addition, for each of the $U^{(i)} - 1$ non-accepting states in the non-flat RM, there are $E^{(i)}$ edges, each of which calls an RM with height $i - 1$ whose number of states is $\bar{U}^{(i-1)}$. These edges are replaced by the called RM except for the initial and accepting states, whose role is now played by the states involved in the substituted edge (hence the -2). This construction process corresponds to the one used to prove Theorem 3.6.

The total of number of edges in an HRM is given by:

$$\sum_{i=1}^{h_r} N^{(i)} (U^{(i)} - 1) E^{(i)},$$

where $(U^{(i)} - 1) E^{(i)}$ is the total number of edges in an RM with height i (the -1 is because the accepting state is discarded), so $N^{(i)} (U^{(i)} - 1) E^{(i)}$ determines how many edges there are across RMs with height i .

The total number of edges in the flat HRM is given by the total number of edges in the flattened root RM, $\bar{E}^{(h_r)}$, where $\bar{E}^{(i)}$ is the total number of edges in the flattened representation of an RM with height i , which is recursively defined as follows:

$$\bar{E}^{(i)} = \begin{cases} (U^{(i)} - 1) E^{(i)} & \text{if } i = 1, \\ (U^{(i)} - 1) E^{(i)} \bar{E}^{(i-1)} & \text{if } i > 1. \end{cases}$$

That is, each of the $(U^{(i)} - 1) E^{(i)}$ edges in an RM with height i is replaced by $\bar{E}^{(i-1)}$ edges given by an RM with height $i - 1$ (if any).

The key intuition is that an HRM with root height $h_r > 1$ is beneficial representation-wise if the number of calls across RMs with height i is higher than the number of RMs with height $i - 1$; in other words, RMs with lower heights are being reused. Numerically, the total number of edges/calls in an RM with height i is $(U^{(i)} - 1) E^{(i)}$ and, therefore, the total number of calls across RMs with height i is $(U^{(i)} - 1) E^{(i)} N^{(i)}$. If this quantity is higher than $N^{(i-1)}$, then RMs with lower heights are reused, and therefore having RMs with different heights is beneficial.

Theorem 3.7. *Let $H = \langle \mathcal{M}, M_r, \mathcal{P} \rangle$ be an HRM and h_r be the height of its root M_r . The number of states and edges in an equivalent flat HRM \bar{H} can be exponential in h_r .*

Proof. By example. Let $H = \langle \mathcal{M}, M_r, \mathcal{P} \rangle$ be an HRM whose root M_r has height h_r and is parameterized by $N^{(i)} = 1$, $U^{(i)} = 3$, $E^{(i)} = 1$ for $i = 1, \dots, h_r$. Figure 7 shows an instance of this hierarchy. Let us write the number of states in the

flat RMs of each level:

$$\begin{aligned}
 \bar{U}^{(1)} &= U^{(1)} = 3, \\
 \bar{U}^{(2)} &= U^{(2)} + (\bar{U}^{(1)} - 2)(U^{(2)} - 1)E^{(2)} = 3 + (3 - 2)(3 - 1)1 = 5, \\
 \bar{U}^{(3)} &= U^{(3)} + (\bar{U}^{(2)} - 2)(U^{(3)} - 1)E^{(3)} = 3 + (5 - 2)(3 - 1)1 = 9, \\
 &\vdots \\
 \bar{U}^{(i)} &= 2\bar{U}^{(i-1)} - 1 = 2^i + 1.
 \end{aligned}$$

Hence, the number of states in the flat HRM is $|\bar{H}| = \bar{U}^{(h_r)} = 2^{h_r} + 1$, showing that the number of states in the flat HRM grows exponentially with the height of the root. In contrast, the number of states in the HRM grows linearly with the height of the root, $|H| = \sum_{i=1}^{h_r} N^{(i)}U^{(i)} = \sum_{i=1}^{h_r} 1 \cdot 3 = 3h_r$.

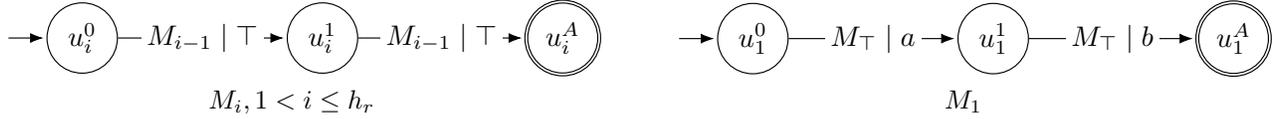


Figure 7: Example of an HRM whose root has height h_r used in the proof of Theorem 3.7.

In the case of the total number of edges, we again write some iterations to derive a general expression:

$$\begin{aligned}
 \bar{E}^{(1)} &= (U^{(1)} - 1)E^{(1)} = (3 - 1)1 = 2, \\
 \bar{E}^{(2)} &= (U^{(2)} - 1)E^{(2)}\bar{E}^{(1)} = (3 - 1) \cdot 1 \cdot 2 = 4, \\
 \bar{E}^{(3)} &= (U^{(3)} - 1)E^{(3)}\bar{E}^{(2)} = (3 - 1) \cdot 1 \cdot 4 = 8, \\
 &\vdots \\
 \bar{E}^{(i)} &= 2\bar{E}^{(i-1)} = 2^i.
 \end{aligned}$$

Therefore, the total number of edges in the flat HRM is $\bar{E}^{(h_r)} = 2^{h_r}$. In contrast, the total number of edges in the HRM grows linearly: $\sum_{i=1}^{h_r} N^{(i)}(U^{(i)} - 1)E^{(i)} = \sum_{i=1}^{h_r} 1(3 - 1)1 = 2h_r$.

Finally, we emphasize that the resulting flat HRM cannot be compressed, unlike the HRM in Figure 6: each state has at most one incoming edge, so there are not multiple paths that can be merged. We have thus shown that there are HRMs whose equivalent flat HRM has a number of states and edges that grows exponentially with the height of the root. \square

Using the aforementioned intuition, we observe that the hierarchical structure is actually expected to be useful: the number of calls across RMs with height i is $(U^{(i)} - 1)E^{(i)} = (3 - 1)1 = 2$, which is greater than the number of RMs with height $i - 1$ (only 1).

There are cases where having a multi-level hierarchy (i.e., with $h_r > 1$) is not beneficial. For instance, given an HRM whose root has height h_r and parameterized by $N^{(i)} = 1$, $U^{(i)} = 2$ and $E^{(i)} = 1$, the number of states in the equivalent flat HRM is constant (2), whereas in the HRM itself it grows linearly with h_r . The same occurs with the number of edges. By checking the previously introduced intuition, we observe that $(U^{(i)} - 1)E^{(i)}N^{(i)} = (2 - 1) \cdot 1 \cdot 1 = 1 \not\geq N^{(i-1)} = 1$, which verifies that having non-reused RMs with multiple heights is not useful.

B. Policy Learning Implementation Details

In this appendix, we describe some implementation details that were omitted in Section 4 for simplicity. First, we start by describing some methods used in policy learning. Second, we explain the option selection algorithm step-by-step and provide examples to ease its understanding.

B.1. Policies

Deep Q-networks (DQNs; Mnih et al., 2015). We use Double DQNs (van Hasselt et al., 2016) for both formula and call options. The DQNs associated with formula options simply take an MDP state and output a Q-value for each action. In contrast, the DQNs associated with call options also take an RM state and a context, which are encoded as follows:

- The RM state is encoded using a one-hot vector. The size of the vector is given by the number of states in the RM.
- The context, which is either \top or a DNF formula with a single disjunct/conjunction, is encoded using a vector whose size is the number of propositions $|\mathcal{P}|$. Each vector position corresponds to a proposition $p \in \mathcal{P}$ whose value depends on how p appears in the context: (i) +1 if p appears positively, (ii) -1 if p appears negatively, or (iii) 0 if p does not appear. Note that if the context is \top , the vector solely consists of zeros.

These DQNs output a value for each possible call in the RM; however, some of these values must be masked if the corresponding calls are not available from the RM state-context used as input. For instance, the DQN for M_0 in Figure 1c outputs a value for $\langle M_1, \neg \text{A} \rangle$, $\langle M_2, \top \rangle$, $\langle M_1, \top \rangle$, and $\langle M_\top, \text{A} \rangle$. If the RM state was u_0^0 and the context was \top , only the values for the first two calls are relevant. Just like unavailable calls, we also mask unsatisfiable calls (i.e., calls whose context cannot be satisfied in conjunction with the accumulated context used as input).

To speed up learning, a subset of the Q-functions associated with formula options is updated after each step. Updating all the Q-functions after each step is costly and we observed that similar performance could be obtained with this strategy. To determine the subset, we keep an update counter c_ϕ for each Q-function q_ϕ , and a global counter c (i.e., the total number of times Q-functions have been updated). The probability of updating q_ϕ is:

$$p_\phi = \frac{s_\phi}{\sum_{\phi'} s_{\phi'}}, \text{ where } s_\phi = c - c_\phi - 1.$$

A subset of Q-functions is chosen using this probability distribution without replacement.

Exploration. During training, the formula and call option policies are ϵ -greedy. In the case of formula options, akin to Q-functions, each option $\omega_{i,u,\Phi}^{j,\phi}$ performs exploration with an exploration factor $\epsilon_{\phi \wedge \Phi}$, which linearly decreases with the number of steps performed using the policy induced by $q_{\phi \wedge \Phi}$. Likewise, Kulkarni et al. (2016) keep an exploration factor for each subgoal, but vary it depending on the option’s success rather than on the number of performed steps. In the case of call options, each RM state-context pair is associated with its own exploration factor, which linearly decreases as options started from that pair terminate.

The Formula Tree. As explained in Section 4, each formula option’s policy is induced by a Q-function associated with a formula. In domains where certain proposition sets cannot occur, it is unnecessary to consider formulas that cover some of these sets. For instance, in a domain where two propositions a and b cannot be simultaneously observed (i.e., it is impossible to observe $\{a, b\}$), formulas such as $a \wedge \neg b$ or $b \wedge \neg a$ could instead be represented by the more abstract formulas a or b ; therefore, $a \wedge \neg b$ and a could be both associated with a Q-function q_a , whereas $b \wedge \neg a$ and b could be both associated with a Q-function q_b . By reducing the number of Q-functions, learning naturally becomes more efficient.

We represent relationships between formulas using a *formula tree* which, as the name suggests, arranges a set of formulas in a tree structure. Formally, given a set of propositions \mathcal{P} , a formula tree is a tuple $\langle \mathcal{F}, F_r, \mathbb{L} \rangle$, where \mathcal{F} is a set of nodes, each associated with a formula; $F_r \in \mathcal{F}$ is the root of the tree and it is associated with the formula \top ; and $\mathbb{L} \subseteq (2^{\mathcal{P}})^*$ is a set of labels. All the nodes in the tree except for the root are associated with conjunctions. Let $\nu(X) \subseteq 2^{\mathcal{P}}$ denote the set of literals of a formula X , e.g. if $X = a \wedge \neg b$, then $\nu(X) = \{a, \neg b\}$. A formula X *subsumes* a formula Y if (1) $X = \top$, or (2.i) $\nu(X) \subseteq \nu(Y)$ and (2.ii) for all labels $\mathcal{L} \in \mathbb{L}$, either $\mathcal{L} \models X$ and $\mathcal{L} \models Y$, or $\mathcal{L} \not\models X$ and $\mathcal{L} \not\models Y$. Case (2) indicates that Y is a special case of X (it adds literals but it is satisfied by exactly the same labels). The tree is organized such that the formula at a given node subsumes all its descendants. The set of Q-functions is determined by the children of the root.

During the agent-environment interaction, the formula tree is updated if (i) a new formula appears in the learned HRMs, or (ii) a new label is observed. Algorithm 1 contains the pseudo-code for updating the tree in these two cases. When a new formula is added (line 1), we create a node for the formula (line 2) and add it to the tree. The insertion place is determined by exploring the tree top-down from the root F_r (lines 3–19). First, we check whether a child of the current node subsumes the new node (line 7). If such a node exists, then we go down this path (lines 8–9); otherwise, the new node is going to be a child of the current node (lines 16–17). In the latter case, in addition, all those children nodes of the current node that

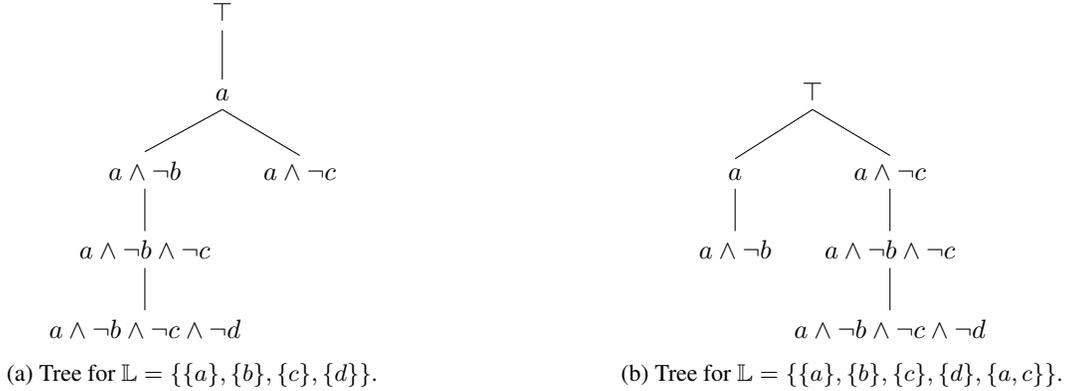


Figure 8: Examples of formula trees for different sets of literals. Note that the node $a \wedge \neg b \wedge \neg c$ in (a) could also be a child of $a \wedge \neg c$ (the parent depends on the insertion order).

are subsumed by the new node need to become children of the new node (lines 11–15). The other core case in which the tree may need an update occurs when a new label is observed (lines 20–25) since we need to make sure that parenting relationships comply with the set of labels \mathbb{L} . First, we find nodes inconsistent with the new label: a parenting relationship is broken (line 39) when the formula of the parent non-root node is satisfied by the label but the formula of the child node is not (or vice versa). Once the inconsistent nodes are found, we remove their current parenting relationship (lines 45–46) and reinsert them in the tree (line 47). Figure 8 shows two simple examples of formula trees, where the Q-functions are q_a in (a), and q_a and $q_{a \wedge \neg c}$ in (b).

B.2. Option Selection Algorithm

Algorithm 2 shows how options are selected, updated, and interrupted during an episode. Lines 1–3 correspond to the algorithm’s initialization. The initial state is that of the environment, while the initial hierarchy state is formed by the root RM M_r , its initial state u_r^0 , an empty context (i.e., $\Phi = \top$), and an empty call stack. The option stack Ω_H contains the options we are currently running, where options at the front are the shallowest ones (e.g., the first option in the list is taken in the root RM). The steps taken during an episode are shown in lines 4–14, which are grouped as follows:

1. The agent fills the option stack Ω_H by selecting options in the HRM from the current hierarchy state until a formula option is chosen (lines 15–25). The context is propagated and augmented through the HRM (i.e., the context of the calls is conjuncted with the propagating context and converted into DNF form). Note that the context is initially \top (true), and not that of the hierarchy state. It is possible that no new options are selected if the formula option chosen in a previous step has not terminated yet.
2. The agent chooses an action according to the last option in the option stack (line 6), which will always be a formula option whose policy maps states into actions. The action is applied, and the agent observes the next state and label (line 7). The next hierarchy state is obtained by applying the hierarchical transition function δ_H using the observed label (line 8). The Q-functions associated with formula options’ policies are updated after this step (line 9).
3. The option stack Ω_H is updated by removing those options that have terminated (lines 10, 26–45). The terminated options are saved in a different list Ω_β to update the Q-functions of the RMs where they were initiated later on (line 11). The termination of the options is performed as described in Section 4. All options terminate if a terminal state is reached (lines 27–28). Otherwise, we check options in Ω_H from deeper to shallower levels. The first checked option is always a formula option, which terminates if the hierarchy state has changed (line 40). In contrast, a call option terminates if it does not appear in the stack (lines 33, 46–51).⁵ When an option is found to terminate, it is added to Ω_β and removed from Ω_H (lines 35–36, 41–42). If a non-terminating option is found (lines 37, 43), we stop checking for termination (no higher level options can have terminated in this case).

⁵We denote by $\phi_1 \subseteq \phi_2$, where $\phi_1, \phi_2 \in \text{DNF}_{\mathcal{P}}$, the fact that all the disjuncts of ϕ_1 appear in ϕ_2 . This containment relationship also holds if both formulas are \top . For instance, $(a \wedge \neg c) \subseteq (a \wedge \neg c) \vee d$.

Algorithm 1 Formula tree operations

Input: a formula tree $\langle \mathcal{F}, F_r, \mathbb{L} \rangle$, where \mathcal{F} is a set of nodes, $F_r \in \mathcal{F}$ is the root node (associated with the formula \top), and \mathbb{L} is a set of labels.

```

1: function ADDFORMULA( $f$ )
2:   ADDNODE(CREATENODE( $f$ ))
3: function ADDNODE(new_node)
4:   current_node  $\leftarrow F_r$ 
5:   added_node  $\leftarrow \perp$ 
6:   while added_node =  $\perp$  do
7:     child_node  $\leftarrow$  FINDSUBSUMINGCHILD(current_node, new_node)
8:     if child_node  $\neq$  nil then {Keep exploring down this path}
9:       current_node  $\leftarrow$  child_node
10:    else {Insert the node}
11:      subsumed_children  $\leftarrow$  GETSUBSUMEDCHILDREN(current_node, new_node)
12:      new_node.children  $\leftarrow$  new_node.children  $\cup$  subsumed_children
13:      for child  $\in$  subsumed_children do
14:        current_node.children  $\leftarrow$  current_node.children  $\setminus$  {child}
15:        child.parent  $\leftarrow$  new_node
16:        current_node.children  $\leftarrow$  current_node.children  $\cup$  {new_node}
17:        new_node.parent  $\leftarrow$  current_node
18:      added_node  $\leftarrow \top$ 
19:    $\mathcal{F} \leftarrow \mathcal{F} \cup \{\text{new\_node}\}$ 
20: function ONLABEL( $\mathcal{L}$ )
21:    $\mathbb{L} \leftarrow \mathbb{L} \cup \{\mathcal{L}\}$ 
22:   inconsistent_nodes  $\leftarrow \{\}$ 
23:   for child  $\in F_r$ .children do
24:     FINDINCONSISTENTNODES(child,  $\mathcal{L}$ , inconsistent_nodes)
25:   REINSERTINCONSISTENTNODES(inconsistent_nodes)
26: function FINDSUBSUMINGCHILD(current_node, new_node)
27:   for child  $\in$  current_node.children do
28:     if child.formula subsumes new_node.formula then
29:       return child
30:   return nil
31: function GETSUBSUMEDCHILDREN(current_node, new_node)
32:   subsumed_children  $\leftarrow \{\}$ 
33:   for child  $\in$  current_node.children do
34:     if new_node.formula subsumes child.formula then
35:       subsumed_children  $\leftarrow$  subsumed_children  $\cup$  {new_node}
36:   return subsumed_children
37: function FINDINCONSISTENTNODES(current_node,  $\mathcal{L}$ , inconsistent_nodes)
38:   for child  $\in$  current_node.children do
39:     if  $\mathcal{L} \models$  current_node.formula  $\oplus$   $\mathcal{L} \models$  child.formula then
40:       inconsistent_nodes  $\leftarrow$  inconsistent_nodes  $\cup$  {child}
41:     else
42:       FINDINCONSISTENTNODES(child,  $\mathcal{L}$ , inconsistent_nodes)
43: function REINSERTINCONSISTENTNODES(inconsistent_nodes)
44:   for node  $\in$  inconsistent_nodes do
45:     node.parent.children  $\leftarrow$  node.parent.children  $\setminus$  {node}
46:     node.parent  $\leftarrow$  nil
47:   ADDNODE(node)

```

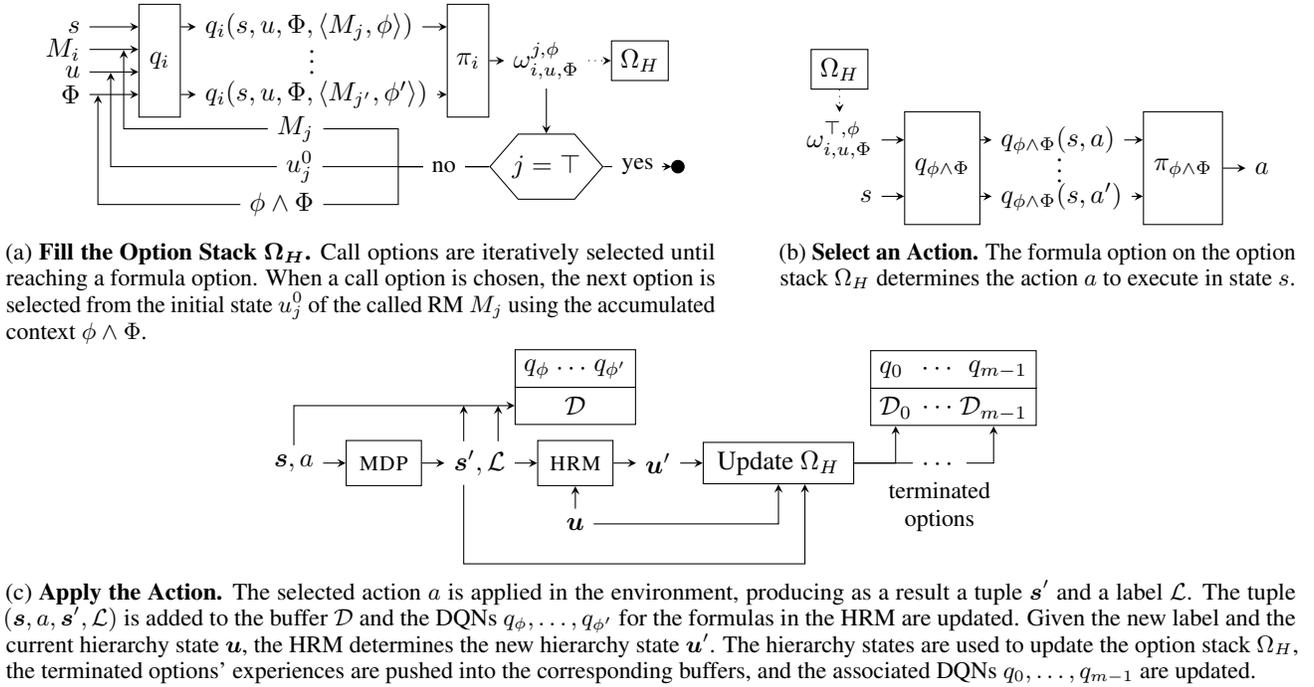


Figure 9: The core procedures involved in the policy learning algorithm that exploits HRMs.

- If at least one option has terminated (line 12), the option stack is updated such that it contains all options appearing in the call stack (lines 13, 52–70). Options are derived for the full stack if Ω_H is empty (lines 53–54), or for the part of the stack not appearing in Ω_H (lines 56–59). The new derived options (lines 61–70) from the call stack are assumed to start in the same state as the last terminated option (i.e., the shallowest terminated option, line 63) and to have been run for the same number of steps too. Crucially, the contexts should be propagated accordingly, starting from the context of the last terminated option (line 69).

As a result of the definition of the hierarchical transition function δ_H , the contexts in the stack may be DNF formulas with more than one disjunct. In contrast, the contexts associated with options are either \top or DNFs with a single disjunct (remember that an option is formed for each disjunct). For instance, this occurs if the context is $a \vee b$ and $\{a, b\}$ is observed: since both disjuncts are satisfied, the context shown in the call stack will be the full disjunction $a \vee b$. In the simplest case, the derived option (which as said before is associated with a DNF with a single disjunct or \top) can include one of these disjuncts chosen uniformly at random (line 67). Alternatively, we could memorize all the derived options and perform identical updates for both later on once terminated.

Figure 9 illustrates the core procedures that constitute the option selection algorithm: (i) filling the option stack, (ii) selecting an action using the formula option in the option stack, and (iii) applying the action and updating the Q-functions and the option stack accordingly.

Examples. We briefly describe some examples of how policy learning is performed in the HRM of Figure 1c. We first enumerate the options in the hierarchy. The formula options are $\omega_{1,0,\neg\mathcal{A}}^{\top,\mathcal{A}}$, $\omega_{2,0,\top}^{\top,\mathcal{A}}$, $\omega_{1,0,\top}^{\top,\mathcal{B}}$, $\omega_{1,1,\top}^{\top,\mathcal{B}}$, $\omega_{2,1,\top}^{\top,\mathcal{B}}$, and $\omega_{0,3,\top}^{\top,\mathcal{B}}$. The first option should lead the agent to observe the label $\{\mathcal{A}\}$ to satisfy $\mathcal{A} \wedge \neg\mathcal{A}$. The Q-functions associated with this set of options are $q_{\mathcal{A} \wedge \neg\mathcal{A}}$, $q_{\mathcal{A}}$, $q_{\mathcal{B}}$, $q_{\mathcal{B}}$ and $q_{\mathcal{B}}$. Note that $\omega_{1,1,\top}^{\top,\mathcal{B}}$ and $\omega_{2,1,\top}^{\top,\mathcal{B}}$ are both associated with $q_{\mathcal{B}}$. Conversely, the call options are $\omega_{0,0,\top}^{1,\neg\mathcal{A}}$, $\omega_{0,0,\top}^{2,\top}$, $\omega_{0,1,\top}^{2,\top}$, and $\omega_{0,2,\top}^{1,\top}$, where the first one achieves its local goal if formula options $\omega_{1,0,\neg\mathcal{A}}^{\top,\mathcal{A}}$ and $\omega_{1,1,\top}^{\top,\mathcal{B}}$ sequentially achieve theirs. The associated Q-functions are q_0 , q_1 and q_2 . Note that $\omega_{0,0,\top}^{2,\top}$ and $\omega_{0,1,\top}^{2,\top}$ are both associated with q_2 .

We now describe a few steps of the aforementioned option selection algorithm in two scenarios. First, we consider the scenario where all chosen options are run to completion (i.e., until their local goals are achieved):

Algorithm 2 Episode execution using an HRM (continues on the next page)

Input: an HRM $H = \langle \mathcal{M}, M_r, \mathcal{P} \rangle$ and an environment $\text{ENV} = \langle \mathcal{S}, \mathcal{A}, p, r, \gamma, \mathcal{P}, l, \tau \rangle$.

- 1: $s_0 \leftarrow \text{ENV.INIT}() \{ \text{Initial MDP tuple} \}$
- 2: $\langle M_i, u, \Phi, \Gamma \rangle \leftarrow \langle M_r, u_r^0, \top, \emptyset \rangle \{ \text{Initial hierarchy state} \}$
- 3: $\Omega_H \leftarrow \emptyset \{ \text{Initial option stack} \}$
- 4: **for** each step $t = 0, \dots$, **do**
- 5: $\Omega_H \leftarrow \text{FILLOPTIONSTACK}(s_t, \langle M_i, u, \Phi, \Gamma \rangle, \Omega_H) \{ \text{Expand the option stack} \}$
- 6: $a \leftarrow \text{SELECTACTION}(s_t, \Omega_H) \{ \text{Choose } a \text{ according to the last option in } \Omega_H \}$
- 7: $s_{t+1}, \mathcal{L}_{t+1} \leftarrow \text{ENV.APPLYACTION}(a)$
- 8: $\langle M_j, u', \Phi', \Gamma' \rangle \leftarrow \delta_H(\langle M_i, u, \Phi, \Gamma \rangle, \mathcal{L}_{t+1}) \{ \text{Apply transition function} \}$
- 9: $\text{UPDATEFORMULAQFUNCTIONS}(s_t, a, s_{t+1}, \mathcal{L}_{t+1})$
- 10: $\Omega_\beta, \Omega_H \leftarrow \text{TERMINATEOPTIONS}(\Omega_H, s, \langle M_i, u, \Phi, \Gamma \rangle, \langle M_j, u', \Phi', \Gamma' \rangle)$
- 11: $\text{UPDATECALLQFUNCTIONS}(\Omega_\beta, s_{t+1}, \mathcal{L}_{t+1})$
- 12: **if** $|\Omega_\beta| > 0$ **then**
- 13: $\Omega_H \leftarrow \text{ALIGNOPTIONSTACK}(\Omega_H, \Gamma', \Omega_\beta)$
- 14: $\langle M_i, u, \Phi, \Gamma \rangle \leftarrow \langle M_j, u', \Phi', \Gamma' \rangle$
- 15: **function** $\text{FILLOPTIONSTACK}(s, \langle M_i, u, \cdot, \Gamma \rangle, \Omega_H)$
- 16: $\Omega'_H \leftarrow \Omega_H$
- 17: $\Phi \leftarrow \top \{ \text{The context is initially true} \}$
- 18: $M_j \leftarrow M_i; v \leftarrow u \{ \text{The state-automaton pair in which an option is selected} \}$
- 19: **while** the last option in Ω'_H is not a formula option **do**
- 20: $\omega_{j,v,\Phi}^{x,\phi} \leftarrow \text{SELECTOPTION}(s, M_j, v, \Phi) \{ \text{Select an option (e.g., with } \epsilon\text{-greedy)} \}$
- 21: **if** $x \neq \top$ **then** $\{ \text{If the option is a call option} \}$
- 22: $M_j \leftarrow M_x; v \leftarrow u_x^0 \{ \text{Next option is chosen on the called RM's initial state} \}$
- 23: $\Phi \leftarrow \text{DNF}(\Phi \wedge \phi) \{ \text{Update the context} \}$
- 24: $\Omega'_H \leftarrow \Omega'_H \oplus \omega_{j,v,\Phi}^{x,\phi} \{ \text{Update the option stack (concatenate new option)} \}$
- 25: **return** Ω'_H
- 26: **function** $\text{TERMINATEOPTIONS}(\Omega_H, s, \langle M_i, u, \Phi, \Gamma \rangle, \langle M_j, u', \Phi', \Gamma' \rangle)$
- 27: **if** $s^T = \top$ **then**
- 28: **return** $\Omega_H, \emptyset \{ \text{All options terminate} \}$
- 29: $\Omega_\beta \leftarrow \emptyset; \Omega'_H \leftarrow \Omega_H \{ \text{Initialize structures} \}$
- 30: **while** $|\Omega'_H| > 0$ **do** $\{ \text{While the option stack is not empty} \}$
- 31: $\omega_{k,v,\Psi}^{x,\phi} \leftarrow \text{last option in } \Omega'_H$
- 32: **if** $x \neq \top$ **then** $\{ \text{If the option is a call option} \}$
- 33: $\text{in_stack}, _ \leftarrow \text{OPTIONINSTACK}(\omega_{k,v,\Psi}^{x,\phi}, \Gamma')$
- 34: **if** $\neg \text{in_stack}$ **then**
- 35: $\Omega_\beta \leftarrow \Omega_\beta \oplus \omega_{k,v,\Psi}^{x,\phi} \{ \text{Update the list of terminated options} \}$
- 36: $\Omega'_H \leftarrow \Omega'_H \ominus \omega_{k,v,\Psi}^{x,\phi} \{ \text{Remove the last option from the option stack} \}$
- 37: **else**
- 38: **break** $\{ \text{Stop terminating} \}$
- 39: **else**
- 40: **if** $\langle M_i, u, \Phi, \Gamma \rangle \neq \langle M_j, u', \Phi', \Gamma' \rangle$ **then** $\{ \text{If the hierarchy state has changed...} \}$
- 41: $\Omega_\beta \leftarrow \Omega_\beta \oplus \omega_{k,v,\Psi}^{x,\phi} \{ \text{Update the list of terminated options} \}$
- 42: $\Omega'_H \leftarrow \Omega'_H \ominus \omega_{k,v,\Psi}^{x,\phi} \{ \text{Remove the last option from the option stack} \}$
- 43: **else**
- 44: **break** $\{ \text{Stop terminating} \}$
- 45: **return** Ω_β, Ω'_H
- 46: **function** $\text{OPTIONINSTACK}(\omega_{k,v,\Phi}^{x,\phi}, \Gamma)$
- 47: **for** $l = 0 \dots |\Gamma| - 1$ **do**
- 48: $\langle u_f, \cdot, M_i, M_j, \phi', \Phi' \rangle \leftarrow \Gamma_l$
- 49: **if** $u_f = v \wedge i = k \wedge j = x \wedge \phi \subseteq \phi' \wedge \Phi \subseteq \Phi'$ **then** $\{ \text{The call option is in the call stack} \}$
- 50: **return** $\top, l \{ \text{Return whether it appears in the stack and the index} \}$
- 51: **return** $\perp, -1$

```

52: function ALIGNOPTIONSTACK( $\Omega_H, \Gamma, \Omega_\beta$ )
53:   if  $|\Omega_H| = 0$  then
54:     return ALIGNOPTIONSTACKHELPER( $\Omega_H, \Gamma, \Omega_\beta, 0$ )
55:   else
56:      $\omega_{k,v,\Phi}^{x,\phi} \leftarrow$  last option in  $\Omega_H$ 
57:     in_stack, stack_index  $\leftarrow$  OPTIONINSTACK( $\omega_{k,v,\Phi}^{x,\phi}, \Gamma$ )
58:     if in_stack then
59:       return ALIGNOPTIONSTACKHELPER( $\Omega_H, \Gamma, \Omega_\beta, \text{stack\_index}$ )
60:   return  $\Omega_H$ 
61: function ALIGNOPTIONSTACKHELPER( $\Omega_H, \Gamma, \Omega_\beta, \text{stack\_index}$ )
62:    $\Omega'_H \leftarrow \Omega_H$ 
63:    $\omega_{\cdot,\cdot,\Phi}^{j,\cdot} \leftarrow$  last option in  $\Omega_\beta$  {Shallowest terminated option}
64:    $\Phi' \leftarrow \Phi$  {Context initialized from last option}
65:   for  $l = \text{stack\_index} \dots |\Gamma| - 1$  do
66:      $\langle u_f, \cdot, M_i, M_j, \phi, \cdot \rangle \leftarrow \Gamma_l$ 
67:      $\phi_{sel} \leftarrow$  Select disjunct from  $\phi$  (e.g., randomly)
68:      $\Omega'_H \leftarrow \Omega'_H \oplus \omega_{i,u_f,\Phi'}^{j,\phi_{sel}}$  {Append new option to the option stack}
69:      $\Phi' \leftarrow \text{DNF}(\Phi' \wedge \phi_{sel})$ 
70:   return  $\Omega'_H$ 

```

1. The initial hierarchy state is $\langle M_0, u_0^0, \top, [] \rangle$ and the option stack Ω_H is empty. We select options to fill Ω_H . The first option is chosen from u_0^0 in M_0 using a policy induced by q_0 . At this state, the available options are $\omega_{0,0,\top}^{1,-\forall}$ and $\omega_{0,0,\top}^{2,\top}$. Let us assume that the former is chosen. Then an option from the initial state of M_1 under context $\neg\forall$ is chosen, which can only be $\omega_{1,0,-\forall}^{\top,\ddagger}$. Since this option is a formula option (the call is made to M_\top), we do not select any more options and the option stack is $\Omega_H = \langle \omega_{0,0,\top}^{1,-\forall}, \omega_{1,0,-\forall}^{\top,\ddagger} \rangle$.
2. The agent selects options according to the formula option in Ω_H , $\omega_{1,0,-\forall}^{\top,\ddagger}$, whose policy is induced by $q_{\ddagger \wedge \neg\forall}$. Let us assume that the policy tells the agent to turn right. Since the label at this location is empty, the hierarchy state remains the same; therefore, no options terminate, and the option stack does not change.
3. Let us assume that the agent moves forward twice, thus observing $\{\ddagger\}$. The hierarchy state then becomes $\langle M_1, u_1^1, \top, [\langle u_0^0, u_0^1, M_0, M_1, \neg\forall, \top \rangle] \rangle$ (see Appendix A.1 for a step-by-step application of the hierarchical transition function). We check which options in Ω_H have terminated starting from the last chosen one. The formula option $\omega_{1,0,-\forall}^{\top,\ddagger}$ terminates because the hierarchy state has changed. In contrast, the call option $\omega_{0,0,\top}^{1,-\forall}$ does not terminate since there is an item in the call stack, $\langle u_0^0, u_0^1, M_0, M_1, \neg\forall, \top \rangle$ that can be mapped into it (meaning that the option is running).
4. An experience $(s, \omega_{1,0,-\forall}^{\top,\ddagger}, s')$ is formed for the terminated option, where s and s' are the observed tuples on initiation and termination respectively. This tuple is added to the replay buffer associated with the RM where the option appears, \mathcal{D}_1 , since it achieved its goal (i.e., a label that satisfied $\ddagger \wedge \neg\forall$ was observed).
5. We align Ω_H with the new stack. In this case, Ω_H remains unchanged since its only option can be mapped into an item of the new stack.
6. We start a new step. Since the option stack does not contain a formula option, we select new options from the current hierarchy state according to a policy induced by q_1 . In this case, there is a single eligible option: $\omega_{1,1,\top}^{\top,\forall}$.

In the second scenario, we observe what occurs when the HRM traversal differs from the options chosen by the agent:

1. The initial step is like the one in the previous scenario, but we assume $\omega_{0,0,\top}^{2,\top}$ is selected instead. Then, since this is a call option, an option from the initial state of M_2 under context \top is chosen, which can only be $\omega_{2,0,\top}^{\top,\forall}$. The option stack thus becomes $\Omega_H = \langle \omega_{0,0,\top}^{2,\top}, \omega_{2,0,\top}^{\top,\forall} \rangle$.

2. Let us assume that by taking actions according to $\omega_{2,0,\top}^{\top,\heartsuit}$ we end up observing $\{\clubsuit\}$. Like in the previous scenario, the hierarchy state becomes $\langle M_1, u_1^1, \top, [\langle u_0^0, u_0^1, M_0, M_1, \neg\heartsuit, \top \rangle] \rangle$. We check which options in Ω_H have terminated. The formula option $\omega_{2,0,\top}^{\top,\heartsuit}$ terminates since the hierarchy state has changed, and the call option $\omega_{0,0,\top}^2$ also terminates since it cannot be mapped into an item of the call stack. Note that these options should intuitively finish since the HRM is being traversed through a path different from that chosen by the agent.
3. The replay buffers are not updated for these options since they have not achieved their local goals.
4. We align Ω_H with the new stack. The only item of the stack $\langle u_0^0, u_0^1, M_0, M_1, \neg\heartsuit, \top \rangle$ can be mapped into option $\omega_{0,0,\top}^{1,\heartsuit}$. We assume that this option starts on the same tuple s and that it has run for the same number of steps as the last terminated option $\omega_{0,0,\top}^2$.

C. HRM Learning Implementation Details

In this appendix, we present some implementation details omitted in Section 5. First, we explain the specifics of our curriculum learning mechanism (Appendix C.1). Second, we describe how an HRM is learned from traces using ILASP (Appendix C.2). Finally, we describe additional details of the algorithm that interleaves RL and HRM learning (Appendix C.3).

C.1. Curriculum Learning

We here describe the details of the curriculum learning method described in Section 5. When an episode is completed for \mathbb{M}_{ij} , R_{ij} is updated using the episode’s undiscounted return r as $R_{ij} \leftarrow \beta R_{ij} + (1 - \beta)r$, where $\beta \in [0, 1]$ is a hyperparameter. A score $c_{ij} = 1 - R_{ij}$ is computed from the return and used to determine the probability of selecting tasks and instances. Note that this scoring function, also used in the curriculum method by Andreas et al. (2017), assumes that the undiscounted return ranges between 0 and 1 (see Section 2). The probability of choosing task i is $\max_j c_{ij} / \sum_k \max_l c_{kl}$; that is, the task for which an instance is performing very poorly has a higher probability. Having selected task i , the probability of choosing instance j is $c_{ij} / \sum_k c_{ik}$, i.e. instances where performance is worse have a higher probability of being chosen. The average undiscounted returns R_{ij} for each task-instance pair are periodically updated using the undiscounted return obtained by the greedy policies in a single evaluation episode.

C.2. Learning an HRM from Traces with ILASP

We formalize the task of learning an HRM using ILASP (Law et al., 2015), an inductive logic programming system that learns answer set programs (ASP) from examples. We refer the reader to Gelfond & Kahl (2014) for an introduction to ASP, and to Law (2018) for ILASP. Our formalization is close to that by Furelos-Blanco et al. (2021) for flat finite-state machines. Without loss of generality, as stated in Section 5, we assume that each RM has exactly one accepting and one rejecting state.

We first describe how HRMs are represented in ASP (Appendix C.2.1), and then explain the encoding of the HRM learning task in ILASP (Appendix C.2.2). Finally, we detail the version of ILASP and the flags we use to run it (Appendix C.2.3).

C.2.1. REPRESENTATION OF AN HRM IN ANSWER SET PROGRAMMING

In this section, we explain how HRMs are represented using Answer Set Programming (ASP). First, we describe how traces are represented. Then, we present how HRMs themselves are represented and also introduce the general rules that describe the behavior of these hierarchies. Finally, we prove the correctness of the representation. We use $\mathbb{A}(X)$ to denote the ASP representation of X (e.g., a trace).

Definition C.1 (ASP representation of a label trace). Given a label trace $\lambda = \langle \mathcal{L}_0, \dots, \mathcal{L}_n \rangle$, $M(\lambda)$ denotes the set of ASP facts that describe it:

$$\mathbb{A}(\lambda) = \{\text{label}(p, t). \mid 0 \leq t \leq n, p \in \mathcal{L}_t\} \cup \{\text{step}(t). \mid 0 \leq t \leq n\} \cup \{\text{last}(n).\}.$$

The $\text{label}(p, t)$ fact indicates that proposition $p \in \mathcal{P}$ is observed in step t , $\text{step}(t)$ states that t is a step of the trace, and $\text{last}(n)$ indicates that the trace ends in step n .

Example 6. The set of ASP facts for the label trace $\lambda = \langle \{\heartsuit\}, \{\}, \{\clubsuit\} \rangle$ is $\mathbb{A}(\lambda) = \{\text{label}(\heartsuit, 0)., \text{label}(\clubsuit, 2)., \text{step}(0)., \text{step}(1)., \text{step}(2)., \text{last}(2).\}$.

Definition C.2 (ASP representation of an HRM). Given an HRM $H = \langle \mathcal{M}, M_r, \mathcal{P} \rangle$, $\mathbb{A}(H) = \bigcup_{M_i \in \mathcal{M} \setminus \{M_\top\}} \mathbb{A}(M_i)$, where:

$$\begin{aligned} \mathbb{A}(M_i) &= \mathbb{A}_{\mathcal{U}}(M_i) \cup \mathbb{A}_{\varphi}(M_i), \\ \mathbb{A}_{\mathcal{U}}(M_i) &= \{\text{state}(u, M_i) \mid u \in \mathcal{U}_i\}, \\ \mathbb{A}_{\varphi}(M_i) &= \left\{ \begin{array}{l} \text{call}(u, u', x + e, M_i, M_j). \\ \bar{\varphi}(u, u', x + e, M_i, \mathbf{T}) :- \text{not label}(p_1, \mathbf{T}), \text{step}(\mathbf{T}). \\ \vdots \\ \bar{\varphi}(u, u', x + e, M_i, \mathbf{T}) :- \text{not label}(p_n, \mathbf{T}), \text{step}(\mathbf{T}). \\ \bar{\varphi}(u, u', x + e, M_i, \mathbf{T}) :- \text{label}(p_{n+1}, \mathbf{T}), \text{step}(\mathbf{T}). \\ \vdots \\ \bar{\varphi}(u, u', x + e, M_i, \mathbf{T}) :- \text{label}(p_m, \mathbf{T}), \text{step}(\mathbf{T}). \end{array} \right\} \end{aligned}$$

$$\left. \begin{array}{l} M_j \in \mathcal{M}, u, u' \in \mathcal{U}_i, \\ \varphi_i(u, u', M_j) \neq \perp, \\ x = \sum_{k=0}^{j-1} |\varphi_i(u, u', M_k)|, \\ e \in [1, |\varphi_i(u, u', M_j)|], \\ \phi_e \in \varphi_i(u, u', M_j), \\ \phi_e = p_1 \wedge \dots \wedge p_n \\ \wedge \neg p_{n+1} \wedge \dots \wedge \neg p_m \end{array} \right\}.$$

Note that each non-leaf RM M_i in the hierarchy is associated with its own set of rules $\mathbb{A}(M_i)$, which are described as follows:

- Facts $\text{state}(u, M_i)$ indicate that u is a state of RM M_i .
- Facts $\text{call}(u, u', e, M_i, M_j)$ indicate that edge e between states u and u' in RM M_i is labeled with a call to RM M_j .
- Normal rules whose *head* is of the form $\bar{\varphi}(u, u', e, M_i, \mathbf{T})$ indicate that the transition from state u to u' with edge e in RM M_i does not hold at step \mathbf{T} . The *body* of these rules consists of a single $\text{label}(p, \mathbf{T})$ literal and a $\text{step}(\mathbf{T})$ atom indicating that \mathbf{T} is a step. Commonly, variables are represented using upper case letters in ASP, which is the case of steps \mathbf{T} here.

There are some important things to take into account regarding the encoding:

- There is no leaf RM M_\top . We later introduce the ASP rules to emulate it.
- The edge identifiers e between a given pair of states (u, u') range from 1 to the total number of conjunctions/disjuncts between them. Note that in \mathbb{A}_{φ} we assume that the leaf RM has an index, just like the other RMs in the HRM. The index could be n since the rest are numbered from 0 to $n - 1$.

Example 7. The following rules represent the HRM in Figure 1c:

$$\left\{ \begin{array}{l} \text{state}(u_0^0, M_0). \text{state}(u_0^1, M_0). \text{state}(u_0^2, M_0). \\ \text{call}(u_0^0, u_0^1, 1, M_0, M_1). \text{call}(u_0^0, u_0^2, 1, M_0, M_2). \\ \text{call}(u_0^3, u_0^A, 1, M_0, M_\top). \\ \bar{\varphi}(u_0^3, u_0^A, 1, M_0, \mathbf{T}) :- \text{not label}(\heartsuit, \mathbf{T}), \text{step}(\mathbf{T}). \end{array} \right\} \cup \left\{ \begin{array}{l} \text{state}(u_0^3, M_0). \text{state}(u_0^A, M_0). \\ \text{call}(u_0^3, u_0^3, 1, M_0, M_2). \text{call}(u_0^3, u_0^3, 1, M_0, M_1). \\ \bar{\varphi}(u_0^0, u_0^1, 1, M_0, \mathbf{T}) :- \text{label}(\heartsuit, \mathbf{T}), \text{step}(\mathbf{T}). \end{array} \right\}$$

$$\left\{ \begin{array}{l} \text{state}(u_1^0, M_1). \text{state}(u_1^1, M_1). \text{state}(u_1^A, M_1). \\ \bar{\varphi}(u_1^0, u_1^1, 1, M_1, \mathbf{T}) :- \text{not label}(\heartsuit, \mathbf{T}), \text{step}(\mathbf{T}). \\ \text{state}(u_2^0, M_2). \text{state}(u_2^1, M_2). \text{state}(u_2^A, M_2). \\ \bar{\varphi}(u_2^0, u_2^1, 1, M_2, \mathbf{T}) :- \text{not label}(\heartsuit, \mathbf{T}), \text{step}(\mathbf{T}). \end{array} \right\} \cup \left\{ \begin{array}{l} \text{call}(u_1^0, u_1^1, 1, M_1, M_\top). \text{call}(u_1^1, u_1^A, 1, M_1, M_\top). \\ \bar{\varphi}(u_1^1, u_1^A, 1, M_1, \mathbf{T}) :- \text{not label}(\heartsuit, \mathbf{T}), \text{step}(\mathbf{T}). \\ \text{call}(u_2^0, u_2^1, 1, M_2, M_\top). \text{call}(u_2^1, u_2^A, 1, M_2, M_\top). \\ \bar{\varphi}(u_2^1, u_2^A, 1, M_2, \mathbf{T}) :- \text{not label}(\heartsuit, \mathbf{T}), \text{step}(\mathbf{T}). \end{array} \right\}.$$

General Rules. The following sets of rules, whose union is denoted by $\mathcal{R} = \bigcup_{i=0}^5 \mathcal{R}_i$, represent how an HRM functions (e.g., how transitions are taken or the acceptance/rejection criteria). For simplicity, all initial, accepting and rejecting states are denoted by u^0 , u^A and u^R respectively.

The rule below is the inversion of the negation of the state transition function $\bar{\varphi}$. Note that the predicate for φ includes the called RM M_2 as an argument.

$$\mathcal{R}_0 = \{ \varphi(X, Y, E, M, M_2, \mathbf{T}) :- \text{not } \bar{\varphi}(X, Y, E, M, \mathbf{T}), \text{call}(X, Y, E, M, M_2), \text{step}(\mathbf{T}). \}.$$

The rule set \mathcal{R}_1 introduces the $\text{pre_sat}(X, M, \mathbf{T})$ predicate, which encodes the exit condition presented in Section 3 and indicates whether a call from state X of RM M can be started at time \mathbf{T} . The first rule corresponds to the base case and

indicates that if the leaf M_\top is called then the condition is satisfied if the associated formula is satisfied. The second rule applies to calls to non-leaf RMs, where we need to satisfy the context of the call (like in the base case), and also check whether a call from the initial state of the potentially called RM can be started.

$$\mathcal{R}_1 = \left\{ \begin{array}{l} \text{pre_sat}(X, M, T) : - \varphi(X, \rightarrow, M, M_\top, T). \\ \text{pre_sat}(X, M, T) : - \varphi(X, \rightarrow, M, M_2, T), \text{pre_sat}(u^0, M_2, T), M_2 \neq M_\top. \end{array} \right\}.$$

The rule set \mathcal{R}_2 introduces the $\text{reachable}(X, M, T_0, T_2)$ predicate, which indicates that state X of RM M is reached between steps T_0 and T_2 . The latter step can also be seen as the step we are currently at. The first fact indicates that the initial state of the root RM is reached from step 0 to step 0. The second rule indicates that the initial state of a non-root RM is reached from step T to step T (i.e., it is reached anytime). The third rule represents the loop transition in the initial state of the root M_r : we stay there if no call can be started at T (i.e., we are not moving in the HRM). The fourth rule is analogous to the third but for the accepting state of the root instead of the initial state. Remember this is the only accepting state in the HRM that does not return control to the calling RM. The fifth rule is also similar to the previous ones: it applies to states reached after T_0 that are non-accepting, which excludes looping in initial states of non-root RMs at the time of starting them (i.e., loops are permitted in the initial state of a non-root RM if we can reach it afterwards by going back to it). The last rule indicates that Y is reached at step T_2 in RM M started at T_0 if there is an outgoing transition from the current state X to Y at time T that holds between T and T_2 , and state X has been reached between T_0 and T . We will later see how δ is defined.

$$\mathcal{R}_2 = \left\{ \begin{array}{l} \text{reachable}(u^0, M_r, 0, 0). \\ \text{reachable}(u^0, M, T, T) : - \text{state}(u^0, M), M \neq M_r, \text{step}(T). \\ \text{reachable}(X, M, T_0, T+1) : - \text{reachable}(X, M, T_0, T), \text{not pre_sat}(X, M, T), \\ \quad \text{step}(T), X = u^0, M = M_r. \\ \text{reachable}(X, M, T_0, T+1) : - \text{reachable}(X, M, T_0, T), \text{not pre_sat}(X, M, T), \\ \quad \text{step}(T), X = u^A, M = M_r. \\ \text{reachable}(X, M, T_0, T+1) : - \text{reachable}(X, M, T_0, T), \text{not pre_sat}(X, M, T), \\ \quad \text{step}(T), T_0 < T, X \neq u^A. \\ \text{reachable}(Y, M, T_0, T_2) : - \text{reachable}(X, M, T_0, T), \delta(X, Y, M, T, T_2). \end{array} \right\}.$$

The rule set \mathcal{R}_3 introduces two predicates. The predicate $\text{satisfied}(M, T_0, T_E)$ indicates that RM M is satisfied if its accepting state u^A is reached between steps T_0 and T_E . Likewise, the predicate $\text{failed}(M, T_0, T_E)$ indicates that RM M fails if its rejecting state u^R is reached between steps T_0 and T_E . These two descriptions correspond to the first and third rules. The second rule applies to the leaf RM M_\top , which always returns control immediately; thus, it is always satisfied between any two consecutive steps.

$$\mathcal{R}_3 = \left\{ \begin{array}{l} \text{satisfied}(M, T_0, T_E) : - \text{reachable}(u^A, M, T_0, T_E). \\ \text{satisfied}(M_\top, T, T+1) : - \text{step}(T). \\ \text{failed}(M, T_0, T_E) : - \text{reachable}(u^R, M, T_0, T_E). \end{array} \right\}$$

The following set, \mathcal{R}_4 , encodes multi-step transitions within an RM. The predicate $\delta(X, Y, M, T, T_2)$ expresses that the transition from state X to state Y in RM M is satisfied between steps T and T_2 . The first rule indicates that this occurs if the context labeling a call to an RM M_2 is satisfied and that RM is also satisfied (i.e., its accepting state is reached) between these two steps. In contrast, the second rule is used for the case in which the rejecting state of the called RM is reached between those steps. In the latter case, we transition to the local rejecting state u^R of M (i.e., the state we would have transitioned to does not matter). This follows from the assumption that rejecting states are global rejectors (see Section 3). The idea of this rule is that rejection is propagated bottom-up in the HRM.

$$\mathcal{R}_4 = \left\{ \begin{array}{l} \delta(X, Y, M, T, T_2) : - \varphi(X, Y, \rightarrow, M, M_2, T), \text{satisfied}(M_2, T, T_2). \\ \delta(X, u^R, M, T, T_2) : - \varphi(X, \rightarrow, M, M_2, T), \text{failed}(M_2, T, T_2). \end{array} \right\}.$$

The last set, \mathcal{R}_5 , encodes the accepting/rejecting criteria. Remember that the $\text{last}(T)$ predicate indicates that T is the last step of a trace. Therefore, the trace is accepted if the root RM is satisfied from the initial step 0 to step $T + 1$ (the step after the last step of the trace, once the final label has been processed). In contrast, the trace is rejected if a rejecting state in the hierarchy is reached between these two same steps.

$$\mathcal{R}_5 = \left\{ \begin{array}{l} \text{accept} : - \text{last}(T), \text{satisfied}(M_r, 0, T+1). \\ \text{reject} : - \text{last}(T), \text{failed}(M_r, 0, T+1). \end{array} \right\}$$

Unlike the formalism introduced in Section 3, this encoding does not use stacks, which would be costly to do. Here we know the trace to be processed and, therefore, the RMs can be evaluated bottom-up; that is, we start evaluating the lowest level RMs first on different substraces, and the result of this evaluation is used in higher level RMs.

We now prove the correctness of the ASP encoding. To do so, we first introduce what means for an HRM to be valid with respect to a trace, as well as a definition and a theorem due to Gelfond & Lifschitz (1988) that will help us derive the proof.

Definition C.3. Given a label trace λ^* , where $*$ \in $\{G, D, I\}$, an HRM H is valid with respect to λ^* if H accepts λ^* and $*$ = G (i.e., λ^* is a goal trace), or H rejects λ^* and $*$ = D (i.e., λ^* is a dead-end trace), or H does not accept nor reject λ^* and $*$ = I (i.e., λ^* is an incomplete trace).

Definition C.4. An ASP program P is stratified when there is a partition

$$P = P_0 \cup P_1 \cup \dots \cup P_n \quad (P_i \text{ and } P_j \text{ disjoint for all } i \neq j)$$

such that, (1) for every predicate p , the definition of p (all clauses with p in the head) is contained in one of the partitions P_i and, (2) for each $1 \leq i \leq n$, if a predicate occurs positively in a clause of P_i then its definition is contained within $\bigcup_{j < i} P_j$, and if a predicate occurs negatively in a clause of P_i then its definition is contained within $\bigcup_{j < i} P_j$.

Theorem C.5. If an ASP program P is stratified, then it has a unique answer set.

Proposition C.6 (Correctness of the ASP encoding). Given a finite label trace λ^* , where $*$ \in $\{G, D, I\}$, and an HRM $H = \langle \mathcal{M}, M_r, \mathcal{P} \rangle$ that is valid with respect to λ^* , the program $P = \mathbb{A}(H) \cup \mathcal{R} \cup \mathbb{A}(\lambda^*)$ has a unique answer set AS and (1) `accept` $\in AS$ if and only if $*$ = G , and (2) `reject` $\in AS$ if and only if $*$ = D .

Proof. First, we prove that the program $P = \mathbb{A}(H) \cup \mathcal{R} \cup \mathbb{A}(\lambda^*)$, where $\mathcal{R} = \bigcup_{i=0}^5 \mathcal{R}_i$, has a unique answer set. By Theorem C.5, if P is stratified then it has a unique answer set. We show there is a way of partitioning P following the constraints in Definition C.4. A possible partition is $P = P_0 \cup P_1 \cup P_2 \cup P_3$, where $P_0 = \mathbb{A}(\lambda^*)$, $P_1 = \mathbb{A}(H)$, $P_2 = \mathcal{R}_0 \cup \mathcal{R}_1$, $P_3 = \mathcal{R}_2 \cup \mathcal{R}_3 \cup \mathcal{R}_4 \cup \mathcal{R}_5$. The unique answer set $AS = AS_0 \cup AS_1 \cup AS_2 \cup AS_3$, where AS_i corresponds to partition P_i , is shown in Figure 10. For simplicity, $\lambda^*[t]$ denotes the t -th label in trace λ^* , $\lambda^*[t:]$ denotes the substrace starting from the t -th label onwards, and $M_i(\lambda^*)$ denotes the hierarchy traversal using RM M_i as the root.

We now prove that `accept` $\in AS$ if and only if $*$ = G (i.e., the trace achieves the goal). If $*$ = G then, since the hierarchy is valid with respect to λ^* (see Definition C.3), the hierarchy traversal $H(\lambda^*)$ finishes in the accepting state u^A of the root; that is, $H(\lambda^*)[n+1] = \langle M_r, u_r^A, \cdot, \cdot \rangle$. This holds if and only if `accept` $\in AS$.

The proof showing that `reject` $\in AS$ if and only if $*$ = D (i.e., the trace reaches a dead-end) is similar to the previous one. If $*$ = D then, since the hierarchy is valid with respect to λ^* , the hierarchy traversal $H(\lambda^*)$ finishes in a rejecting state u^R ; that is, $H(\lambda^*)[n+1] = \langle M_k, u^R, \cdot, \cdot \rangle$, where $M_k \in \mathcal{M}$. This holds if and only if `reject` $\in AS$. \square

C.2.2. REPRESENTATION OF THE HRM LEARNING TASK IN ILASP

We here formalize the learning of an HRM and its mapping to a general ILASP learning task. We start by defining the HRM learning task introduced in Section 5.

Definition C.7. An HRM learning task is a tuple $T_H = \langle r, \mathcal{U}, \mathcal{P}, \mathcal{M}, \mathcal{M}_C, u^0, u^A, u^R, \Lambda, \kappa \rangle$, where r is the index of the root RM in the HRM; $\mathcal{U} \supseteq \{u^0, u^A, u^R\}$ is a set of states of the root RM always containing an initial state u^0 , an accepting state u^A , and a rejecting state u^R ; \mathcal{P} is a set of propositions; $\mathcal{M} \supseteq \{M_\top\}$ is a set of RMs; $\mathcal{M}_C \subseteq \mathcal{M}$ is a set of callable RMs; $\Lambda = \Lambda^G \cup \Lambda^D \cup \Lambda^I$ is a set of label traces; and κ is the maximum number of conjunctions/disjuncts in each formula. An HRM $H = \langle \mathcal{M} \cup \{M_r\}, M_r, \mathcal{P} \rangle$ is a solution of T_H if and only if it is valid with respect to all the traces in Λ .

We make some assumptions about the sets of RMs \mathcal{M} : (i) all RMs reachable from RMs in \mathcal{M}_C must be in \mathcal{M} , (ii) all RMs in \mathcal{M} are deterministic, and (iii) all RMs in \mathcal{M} are defined over the same set of propositions \mathcal{P} (or a subset of it).

For completeness, we provide the definition of an ILASP task introduced by Law et al. (2016). The first definition corresponds to the form of the examples taken by ILASP, while the second corresponds to the ILASP tasks themselves.

Definition C.8. A context-dependent partial interpretation (CDPI) is a pair $\langle \langle e^{inc}, e^{exc} \rangle, e^{ctx} \rangle$, where $\langle e^{inc}, e^{exc} \rangle$ is a pair of sets of atoms, called a partial interpretation, and e^{ctx} is an ASP program called a context. A program P accepts a CDPI $\langle \langle e^{inc}, e^{exc} \rangle, e^{ctx} \rangle$ if and only if there is an answer set AS of $P \cup e^{ctx}$ such that $e^{inc} \subseteq AS$ and $e^{exc} \cap AS = \emptyset$.

$$\begin{aligned}
 AS_0 &= \{\text{label}(p, t). \mid 0 \leq t \leq n, p \in \mathcal{L}_t\} \cup \{\text{step}(t). \mid 0 \leq t \leq n\} \cup \{\text{last}(n).\}, \\
 AS_1 &= \left\{ \begin{array}{l} \{\text{state}(u, M_i). \mid M_i \in \mathcal{M} \setminus \{M_\top\}, u \in \mathcal{U}_i\} \cup \\ \{\text{call}(u, u', x + e, M_i, M_j). \mid M_i \in \mathcal{M} \setminus \{M_\top\}, M_j \in \mathcal{M}, u, u' \in \mathcal{U}_i, \varphi_i(u, u', M_j) \neq \perp, \\ x = \sum_{k=0}^{j-1} |\varphi_i(u, u', M_k)|, e \in [1, |\varphi_i(u, u', M_j)|]\} \cup \\ \{\bar{\varphi}(u, u', x + e, M_i, t). \mid 0 \leq t \leq n, M_i \in \mathcal{M} \setminus \{M_\top\}, M_j \in \mathcal{M}, u, u' \in \mathcal{U}_i, \\ \varphi_i(u, u', M_j) \neq \perp, x = \sum_{k=0}^{j-1} |\varphi_i(u, u', M_k)|, \\ e \in [1, |\varphi_i(u, u', M_j)|], \lambda^*[t] \neq \varphi_i(u, u', M_j)[e]\} \end{array} \right\} \cup, \\
 AS_2 &= \left\{ \begin{array}{l} \{\varphi(u, u', x + e, M_i, t). \mid 0 \leq t \leq n, M_i \in \mathcal{M} \setminus \{M_\top\}, M_j \in \mathcal{M}, u, u' \in \mathcal{U}_i, \\ \varphi_i(u, u', M_j) \neq \perp, x = \sum_{k=0}^{j-1} |\varphi_i(u, u', M_k)|, \\ e \in [1, |\varphi_i(u, u', M_j)|], \lambda^*[t] \models \varphi_i(u, u', M_j)[e]\} \cup, \\ \{\text{pre_sat}(u, M_i, t). \mid 0 \leq t \leq n, M_i \in \mathcal{M} \setminus \{M_\top\}, u \in \mathcal{U}_i, \lambda^*[t] \models \xi_{i, u, \top}\} \end{array} \right\}, \\
 AS_3 &= \left\{ \begin{array}{l} \{\text{reachable}(u^0, M_r, 0, 0).\} \cup \\ \{\text{reachable}(u^0, M_i, t, t). \mid 0 \leq t \leq n, M_i \in \mathcal{M} \setminus \{M_\top, M_r\}, u^0 \in \mathcal{U}_i\} \cup \\ \{\text{reachable}(u, M_r, t_1, t_2). \mid 0 \leq t_1 < t_2 \leq n + 1, u \in \mathcal{U}_r, \\ H(\lambda^*[t_1 :])[t_2 - t_1] = \langle M_r, u, \cdot, \cdot \rangle\} \cup \\ \{\text{reachable}(u, M_i, t_1, t_2). \mid 0 \leq t_1 < t_2 \leq n + 1, M_i \in \mathcal{M} \setminus \{M_r, M_\top\}, u \in \mathcal{U}_i, \\ \lambda^*[t_1] \models \xi_{i, u^0, \top}, \\ M_i(\lambda^*[t_1 :])[t_2 - t_1] = \langle M_i, u, \cdot, \cdot \rangle, \\ M_i(\lambda^*[t_1 :])[t_2 - t_1 - 1] \neq \langle M_i, u^A, \cdot, \cdot \rangle\} \cup \\ \{\text{satisfied}(M_r, t_1, t_2) \mid 0 \leq t_1 < t_2 \leq n + 1, H(\lambda^*[t_1 :])[t_2 - t_1] = \langle M_r, u^A, \cdot, \cdot \rangle\} \\ \{\text{satisfied}(M_i, t_1, t_2). \mid 0 \leq t_1 < t_2 \leq n + 1, M_i \in \mathcal{M} \setminus \{M_r, M_\top\}, \\ \lambda^*[t_1] \models \xi_{i, u^0, \top}, \\ M_i(\lambda^*[t_1 :])[t_2 - t_1] = \langle M_i, u^A, \cdot, \cdot \rangle, \\ M_i(\lambda^*[t_1 :])[t_2 - t_1 - 1] \neq \langle M_i, u^A, \cdot, \cdot \rangle\} \cup \\ \{\text{satisfied}(M_\top, t, t + 1) \mid 0 \leq t \leq n\} \cup \\ \{\text{failed}(M_r, t_1, t_2) \mid 0 \leq t_1 < t_2 \leq n + 1, H(\lambda^*[t_1 :])[t_2 - t_1] = \langle \cdot, u^R, \cdot, \cdot \rangle\} \cup \\ \{\text{failed}(M_i, t_1, t_2). \mid 0 \leq t_1 < t_2 \leq n + 1, M_i \in \mathcal{M} \setminus \{M_r, M_\top\}, \\ \lambda^*[t_1] \models \xi_{i, u^0, \top}, \\ M_i(\lambda^*[t_1 :])[t_2 - t_1] = \langle \cdot, u^R, \cdot, \cdot \rangle\} \cup \\ \{\delta(u, u', M_i, t, t + 1). \mid 0 \leq t \leq n, M_i \in \mathcal{M} \setminus \{M_\top\}, u, u' \in \mathcal{U}_i, \\ \lambda^*[t_1] \models \varphi_i(u, u', M_\top)\} \cup \\ \{\delta(u, u', M_i, t_1, t_2). \mid 0 \leq t_1 < t_2 \leq n + 1, M_i \in \mathcal{M} \setminus \{M_\top\}, u, u' \in \mathcal{U}_i, \\ \exists M_j \in \mathcal{M} \setminus \{M_\top\} \text{ s.t. } \phi = \varphi_i(u, u', M_j), \lambda^*[t_1] \models \xi_{j, u^0, \phi}, \\ M_j(\lambda^*[t_1 :])[t_2 - t_1] = \langle M_j, u^A, \cdot, \cdot \rangle, \\ M_j(\lambda^*[t_1 :])[t_2 - t_1 - 1] \neq \langle M_j, u^A, \cdot, \cdot \rangle\} \cup \\ \{\delta(u, u^R, M_i, t_1, t_2). \mid M_i \in \mathcal{M} \setminus \{M_\top\}, u \in \mathcal{U}_i, 0 \leq t_1 < t_2 \leq n + 1, \\ \exists M_j \in \mathcal{M} \setminus \{M_\top\} \text{ s.t. } \phi = \varphi_i(u, u', M_j), \lambda^*[t_1] \models \xi_{j, u^0, \phi}, \\ M_j(\lambda^*[t_1 :])[t_2 - t_1] = \langle M_k, u^R, \cdot, \cdot \rangle, M_k \in \mathcal{M}\} \cup \\ \{\text{accept} \mid H(\lambda^*)[n + 1] = \langle M_r, u^A, \cdot, \cdot \rangle\} \cup \\ \{\text{reject} \mid H(\lambda^*)[n + 1] = \langle M_k, u^R, \cdot, \cdot \rangle, M_k \in \mathcal{M} \setminus \{M_\top\}\} \end{array} \right\}.
 \end{aligned}$$

Figure 10: Answer sets for each of the partitions in the program $P = \mathbb{A}(H) \cup \mathcal{R} \cup \mathbb{A}(\lambda^*)$, where H is an HRM, \mathcal{R} is the set of general rules and λ^* is a label trace.

Definition C.9. An *ILASP task* is a tuple $T = \langle \mathcal{B}, \mathcal{S}_M, \langle \mathcal{E}^+, \mathcal{E}^- \rangle \rangle$ where \mathcal{B} is the ASP background knowledge, which describes a set of known concepts before learning; \mathcal{S}_M is the set of ASP rules allowed in the hypotheses; and \mathcal{E}^+ and \mathcal{E}^- are sets of CDPIs called, respectively, the positive and negative examples. A hypothesis $\mathcal{H} \subseteq \mathcal{S}_M$ is an *inductive solution* of T if and only if (i) $\forall e \in \mathcal{E}^+, \mathcal{B} \cup \mathcal{H}$ accepts e , and (ii) $\forall e \in \mathcal{E}^-, \mathcal{B} \cup \mathcal{H}$ does not accept e .

Given an HRM learning task T_H , we map it into an ILASP learning task $\mathbb{A}(T_H) = \langle \mathcal{B}, \mathcal{S}_M, \langle \mathcal{E}^+, \emptyset \rangle \rangle$ and use the ILASP system (Law et al., 2015) to find an inductive solution $\mathbb{A}_\varphi(H) \subseteq \mathcal{S}_M$ that covers the examples. Note that we do not use *negative examples* ($\mathcal{E}^- = \emptyset$). We define the components of $\mathbb{A}(T_H)$ below.

Background Knowledge. The background knowledge $\mathcal{B} = \mathcal{B}_U \cup \mathcal{B}_M \cup \mathcal{R}$ is a set of rules that describe the behavior of the HRM. The set \mathcal{B}_U consists of $\text{state}(u, M_r)$ facts for each state $u \in \mathcal{U}$ of the root RM with index r we aim to induce, whereas $\mathcal{B}_M = \bigcup_{M_i \in \mathcal{M} \setminus \{M_r\}} \mathbb{A}(M_i)$ contains the ASP representations of all RMs. Finally, \mathcal{R} is the set of general rules introduced in Appendix C.2.1 that defines how HRMs process label traces. Importantly, the index of the root r in these rules must correspond to the one used in T_H .

Hypothesis Space. The hypothesis space \mathcal{S}_M contains all ed and φ rules that characterize a transition from a non-terminal state $u \in \mathcal{U} \setminus \{u^A, u^R\}$ to a different state $u' \in \mathcal{U} \setminus \{u\}$ using edge $i \in [1, \kappa]$. Formally, it is defined as

$$\mathcal{S}_M = \left\{ \begin{array}{l} \text{call}(u, u', i, M). \\ \varphi(u, u', i, M, T) :- \text{label}(p, T), \text{step}(T). \\ \varphi(u, u', i, M, T) :- \text{not label}(p, T), \text{step}(T). \end{array} \mid \begin{array}{l} u \in \mathcal{U} \setminus \{u^A, u^R\}, \\ u' \in \mathcal{U} \setminus \{u\}, i \in [1, \kappa], \\ M \in \mathcal{M}_C, p \in \mathcal{P} \end{array} \right\}.$$

Example Sets. Given a set of traces $\Lambda = \Lambda^G \cup \Lambda^D \cup \Lambda^I$, the set of *positive examples* is defined as

$$\mathcal{E}^+ = \{ \langle e^*, \mathbb{A}(\lambda) \rangle \mid * \in \{G, D, I\}, \lambda \in \Lambda^* \},$$

where $e^G = \langle \{\text{accept}\}, \{\text{reject}\} \rangle$, $e^D = \langle \{\text{reject}\}, \{\text{accept}\} \rangle$, and $e^I = \langle \{\}, \{\text{accept}, \text{reject}\} \rangle$ are the partial interpretations for goal, dead-end and incomplete traces. The `accept` and `reject` atoms express whether a trace is accepted or rejected by the HRM; hence, goal traces must only be accepted, dead-end traces must only be rejected, and incomplete traces cannot be accepted or rejected. Note that the context of each example is the set of ASP facts $\mathbb{A}(\lambda)$ that represents the corresponding trace (see Definition C.1).

Correctness of the Learning Task. The following theorem captures the correctness of the HRM learning task.

Theorem C.10. *Given an HRM learning task $T_H = \langle r, \mathcal{U}, \mathcal{P}, \mathcal{M}, \mathcal{M}_C, u^0, u^A, u^R, \Lambda, \kappa \rangle$, an HRM $H = \langle \mathcal{M} \cup \{M_r\}, M_r, \mathcal{P} \rangle$ is a solution of T_H if and only if $\mathbb{A}_\varphi(M_r)$ is an inductive solution of $\mathbb{A}(T_H) = \langle \mathcal{B}, \mathcal{S}_M, \langle \mathcal{E}^+, \emptyset \rangle \rangle$.*

Proof. Assume H is a solution of T_H .

$\iff H$ is valid with respect to all traces in Λ (i.e., H accepts all traces in Λ^G , rejects all traces in Λ^D and does not accept nor reject any trace in Λ^I).

\iff By Proposition C.6, for each trace $\lambda^* \in \Lambda^*$ where $* \in \{G, D, I\}$, $\mathbb{A}(H) \cup \mathcal{R} \cup \mathbb{A}(\lambda^*)$ has a unique answer set AS and (1) `accept` $\in AS$ if and only if $* = G$, and (2) `reject` $\in AS$ if and only if $* = D$.

\iff For each example $e \in \mathcal{E}^+$, $\mathcal{R} \cup \mathbb{A}(H)$ accepts e .

\iff For each example $e \in \mathcal{E}^+$, $\mathcal{B} \cup \mathbb{A}_\varphi(M_r)$ accepts e (the two programs are identical).

$\iff \mathbb{A}_\varphi(M_r)$ is an inductive solution of $\mathbb{A}(T_H)$. □

Constraints. We introduce several constraints encoding structural properties of the HRMs we want to learn. Some of these constraints are expressed in terms of facts $\text{pos}(u, u', e, m, p)$ and $\text{neg}(u, u', e, m, p)$, which indicate that proposition $p \in \mathcal{P}$ appears positively (resp. negatively) in edge e from state u to state u' in RM M_m . These facts are derived from φ rules in $\mathbb{A}(H)$ and injected in the ILASP tasks using meta-program injection (Law et al., 2018).

The following set of constraints ensures that the learned root RM is *deterministic* using the *saturation* technique (Eiter & Gottlob, 1995). The idea is to check determinism top-down by selecting two edges from a given state in the root, each

- Rule out inductive solutions containing states different from the accepting and rejecting states without outgoing edges. In general, these states are not interesting.

$$\left\{ \begin{array}{l} \text{has_outgoing_edges}(X, M) :- \text{ed}(X, _, _, M). \\ :- \text{state}(X, M), \text{not has_outgoing_edges}(X, M), X \neq u^A, X \neq u^R. \end{array} \right\}$$

- Rule out inductive solutions containing cycles; that is, solutions where two states can be reached from each other. The $\text{path}(X, Y, M)$ predicate indicates there is a directed path (i.e., a sequence of directed edges) from X to Y in RM M . The first rule states that there is a path from X to Y if there is an edge from X to Y . The second rule indicates that there is a path from X to Y if there is an edge from X to an intermediate state Z from which there is a path to Y . Finally, the third rule discards the solutions where X and Y can be reached from each other through directed edges.

$$\left\{ \begin{array}{l} \text{path}(X, Y, M) :- \text{ed}(X, Y, _, M). \\ \text{path}(X, Y, M) :- \text{ed}(X, Z, _, M), \text{path}(Z, Y, M). \\ :- \text{path}(X, Y, M), \text{path}(Y, X, M). \end{array} \right\}$$

C.2.3. ILASP FLAGS

We use ILASP2 (Law et al., 2015) to learn the HRMs. For efficiency, the default calls to the underlying ASP solver are modified to be made with the flag `---opt-mode=ignore`, meaning that non-minimal solutions might be obtained (i.e., solutions involving more rules than needed), so the learned root might contain some unnecessary edges. In practice, the solutions produced by ILASP rarely contain such edges and, if they do, these edges eventually disappear by observing an appropriate counterexample. We hypothesize that using this flag helps since no optimization is made every time ILASP is called. We highlight that this notion of minimality is not related to that of a minimal RM (i.e., an RM with the fewest number of states) described in Section 5.

C.3. Interleaving Algorithm

Akin to some methods for learning RMs (Toro Icarte et al., 2019; Furelos-Blanco et al., 2021), we compress label traces by merging consecutive equal labels into a single one; e.g., $\langle \{\}, \{\blacksquare\}, \{\blacksquare\}, \{\}, \{\}, \{\spadesuit\}, \{\spadesuit\} \rangle$ becomes $\langle \{\}, \{\blacksquare\}, \{\}, \{\spadesuit\}, \{\spadesuit\} \rangle$. Our method does not require traces to be compressed, but performance is enhanced since traces usually become shorter.

D. Experimental Details

In this section, we describe the details of the experiments introduced in Section 6. First, we discuss the implementation of the domains, the network architecture we used for each of them, and provide HRM examples for the different tasks (Appendix D.1). Second, we provide the list of hyperparameters used in the methods evaluated in this paper (Appendix D.2). Finally, we provide all specific results summarized in Section 6 (Appendix D.3). All timed experiments ran on 3.40GHz Intel® Core™ i7-6700 processors, while non-timed experiments have also run on 2.90GHz Intel® Core™ i7-10700, 4.20GHz Intel® Core™ i7-7700K, and 3.20GHz Intel® Core™ i7-8700 processors.

D.1. Domains

We here describe how the domains we use in our experiments are implemented, the architecture of the DQNs, and provide some example HRMs for the tasks we have considered.

D.1.1. CRAFTWORLD

Implementation. This domain is based on MiniGrid (Chevalier-Boisvert et al., 2018), thus inheriting many of its features. At each step, the agent observes a $W \times H \times 3$ tensor, where W and H are the width and height of the grid. The three channels contain the object IDs, the color IDs, and object state IDs (including the orientation of the agent) respectively. Each of the objects we define (except for the lava \spadesuit , which already existed in MiniGrid) has its own object and color IDs. Before providing the agent with the state, the content of all matrices is scaled between -1 and 1. Note that even though the agent gets a full view of the grid, it is still unaware of the completion degree of a task. Other works have previously used the full view of the grid (Igl et al., 2019; Jiang et al., 2021).

The grids are randomly generated. In all settings (OP, OPL, FR, FRL), the agent and the objects are randomly assigned an unoccupied position. In the case of FR and FRL, no object occupies a position between rooms or its adjoining positions.

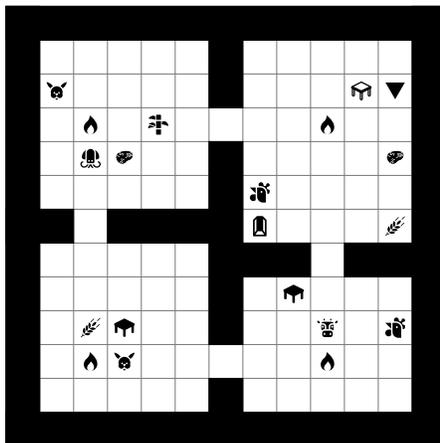


Figure 11: An instance of the CRAFTWORLD grid in the FRL setting.

There is a single object per object type (i.e., proposition) in OP and OPL, whereas there can be one or two per type in FR and FRL. Finally, there is a single lava location in OPL, which is randomly assigned (like the rest of the propositions), whereas in FRL there are four fixed lava locations placed in the intersections between doors as shown in Figure 11.

Network Architecture. The DQNs for CRAFTWORLD consist of a 3-layer convolutional neural network (CNN) with 16, 32, and 32 filters respectively. All kernels are 2×2 and use a stride of 1. In the FR and FRL settings, there is a max pooling layer with kernel size 2×2 after the first convolutional layer. This part of the architecture is based on that by Igl et al. (2019) and Jiang et al. (2021), who also work on MiniGrid using the full view of the grid. In DQNs associated with formulas, the CNN’s output is fed to a 3-layer multilayer perceptron (MLP) where the hidden layer has 256 rectifier units and the output layer has a single output for each action. In the case of DQNs for RMs, the output of the CNN is extended with the encoding of the RM state and the context (as discussed in Appendix B) before being fed to a 3-layer MLP where the hidden layer has 256 rectifier units and the output layer has a single output for each call in the RM.

Examples of HRMs. Figure 12 shows minimal root RMs for the CRAFTWORLD tasks listed in Table 1. Note that (i) since two or more propositions can never occur simultaneously, the mutual exclusivity between formulas could be enforced differently, and (ii) these RMs correspond to the settings without dead-ends (thus, they do not include rejecting states).

D.1.2. WATERWORLD

Implementation. This domain (cf. Figure 13) has a continuous state space. The states are vectors containing the absolute position and velocity of the agent, and the relative positions and velocities of the other balls. The agent does not know the color of each ball. In all settings (WOD and WD), a WATERWORLD instance is created by assigning a random position and direction to each ball. Like in CRAFTWORLD, the agent does not know the degree of completion of a task.

Network Architecture. The architecture for WATERWORLD is a simple modification of the one introduced by Toro Icarte et al. (2018). The formula DQNs consist of a 5-layer MLP, where each of the 3 hidden layers has 512 rectifier units. The DQN for the RMs share the same architecture and, like in CRAFTWORLD, the state from the environment is extended with the state and context encodings.

Examples of HRMs. Figure 14 shows minimal root RMs for the tasks listed in Figure 3 (right). Note that these RMs correspond to the settings without dead-ends; thus, they do not include rejecting states.

D.2. Hyperparameters

Policy and HRM Learning. Table 2 lists the hyperparameters used in the experiments with our approach. We also provide the hyperparameters used for CRM (Toro Icarte et al., 2022) to learn policies in flat HRMs. The DQNs for CRM are like those associated with formulas in our approach.

Flat HRM Learning Baselines. We briefly describe the methods used to learn flat HRMs in Section 6.2. Each run consists of 150,000 episodes, and the set of instances is exactly the same across methods. The core difference with respect to learning

non-flat HRMs is that there is a single task for which the HRM is learned. Our method, LHRM, is therefore not able to reuse previously learned HRMs for other tasks; however, it still uses the same hyperparameters (see Table 2). In the case of DeepSynth (Hasanbeig et al., 2021), JIRP (Xu et al., 2020) and LRM (Toro Icarte et al., 2019), we exclusively evaluate their RM learning components using traces collected through random walks.⁶ For a fair comparison against LHRM (both in the non-flat and flat learning cases), we (i) compress the traces using the methodology described in Appendix C.3, and (ii) use the OP and WOD settings of CRAFTWORLD and WATERWORLD respectively, where observing goal traces by randomly exploring the environment is relatively easy (especially for simple tasks such as MILKBUCKET). In these approaches, a different instance is selected at each episode following a cyclic order (i.e., $1, 2, \dots, I - 1, I, 1, 2, \dots$). The proposition set in these approaches includes a proposition covering the case where none of the original propositions are observed (if needed). In the case of LRM, one of the parameters is the maximum number of RM states, which we set to that of the minimal RM. Akin to other approaches, we modify DeepSynth to call the learner only when a counterexample trace is observed instead of calling it periodically, which repeatedly produced the same RM and resulted in avoidable timeouts.

D.3. Extended Results

We here present the tables and figures on which the discussion in Section 6 is based.

D.3.1. LEARNING OF NON-FLAT HRMS

We present tables containing the results for the HRM learning component of LHRM. The content of the columns is the following left-to-right: (1) task name; (2) number of runs in which at least one goal trace was observed; (3) number of runs in which at least one HRM was learned; (4) time spent to learn the HRMs; (5) number of calls made to ILASP to learn the HRMs; (6) number of states of the final HRM; (7) number of edges of the final HRM; (8) number of episodes between the learning of the first HRM and the activation of the task’s level; (9) number of example traces of a given type ($G = \text{goal}$, $D = \text{dead-end}$, $I = \text{incomplete}$); and (10) length of the example traces of a given type. In addition, the bottom of the tables contains the number of completed runs (i.e., the number of runs that have not timed out), the total time spent on learning the HRMs, and the total number of calls made to ILASP. In the case of CRAFTWORLD, Table 3 shows the results for the default case (all lower level RMs are callable and options are used for exploration), Table 4 shows the results when the set of callable RMs contains only those actually needed, and Table 5 shows the results using primitive actions for exploration instead of options. Analogous results are shown for WATERWORLD in Tables 6, 7 and 8.

The performance decay for RGB&CMY observed in Figure 3 is due to a new RM for RG&BC being learned, which is indicated by a vertical line for the latter occurring exactly at the time of the decay. Following our curriculum method (see Section 5 and Appendix C.1), the average return for RG&BC is reset to 0 and the current level is set to 2; hence, the agent stops performing RGB&CMY (level 3), which causes the performance decay (the reward is 0 while a task is not active). When the average return for RG&BC is again above the threshold, the agent continues learning RGB&CMY.

D.3.2. LEARNING OF FLAT HRMS

Table 9 shows the results of learning a non-flat HRM using LHRM, and the results of learning a flat HRM using several approaches (LHRM, DeepSynth, JIRP and LRM). An extended discussion of these results can be found in Section 6.2.

D.3.3. POLICY LEARNING IN HANDCRAFTED HRMS

Figure 15 shows the plots for the settings omitted in the main paper (the remaining CRAFTWORLD setting, FRL, is shown in Figure 4). As mentioned in Section 6.3 and shown in the figure, the efficacy of non-flat HRMs is less evident in two scenarios. First, when the task’s goal is reachable regardless of the chosen options (e.g., if there are no rejecting states, like in OP and FR) because the policies over options become irrelevant. Second, when the reward is not sparse, like in OPL (the grid is small) or WATERWORLD (the balls easily get near the agent, so even a poor agent can achieve the goal if the number of steps per episode is large). In addition, Figure 15b shows a case where the convergence in the non-flat case is delayed with respect to the flat one. Remember that in DQNs (Mnih et al., 2015), learning does not start until the buffers contain a certain number of experiences. In our approach, as described in Section 4, there is a DQN and a replay buffer for each RM; thus, in the flat case, there is a single DQN and buffer, while in the non-flat case there are several. Filling the buffers in

⁶The codebases for DeepSynth (<https://github.com/grockious/deepsynth>) and LRM (<https://bitbucket.org/RToroIcarte/lrm>) are linked in the papers, whereas the one for JIRP (<https://github.com/corazza/stochastic-reward-machines>) was referred to us by one of the authors through personal communication.

the non-flat case is slower since there are higher-level options (i.e., call options) that do not occur as often as others (i.e., formula options), hence explaining the convergence delay. Nevertheless, we emphasize that in more complex scenarios this does not occur, as shown in Figure 4.

We observe CRM performs closely to the HRL methods when the reward is not very sparse; for instance, when the CRAFTWORLD grid is small (OP, OPL). However, as the reward becomes sparser, CRM struggles more to converge since it does not decompose tasks into independently solvable subtasks and instead relies on a single non-zero reward signal (i.e., the one coming from the transitions to the accepting state).

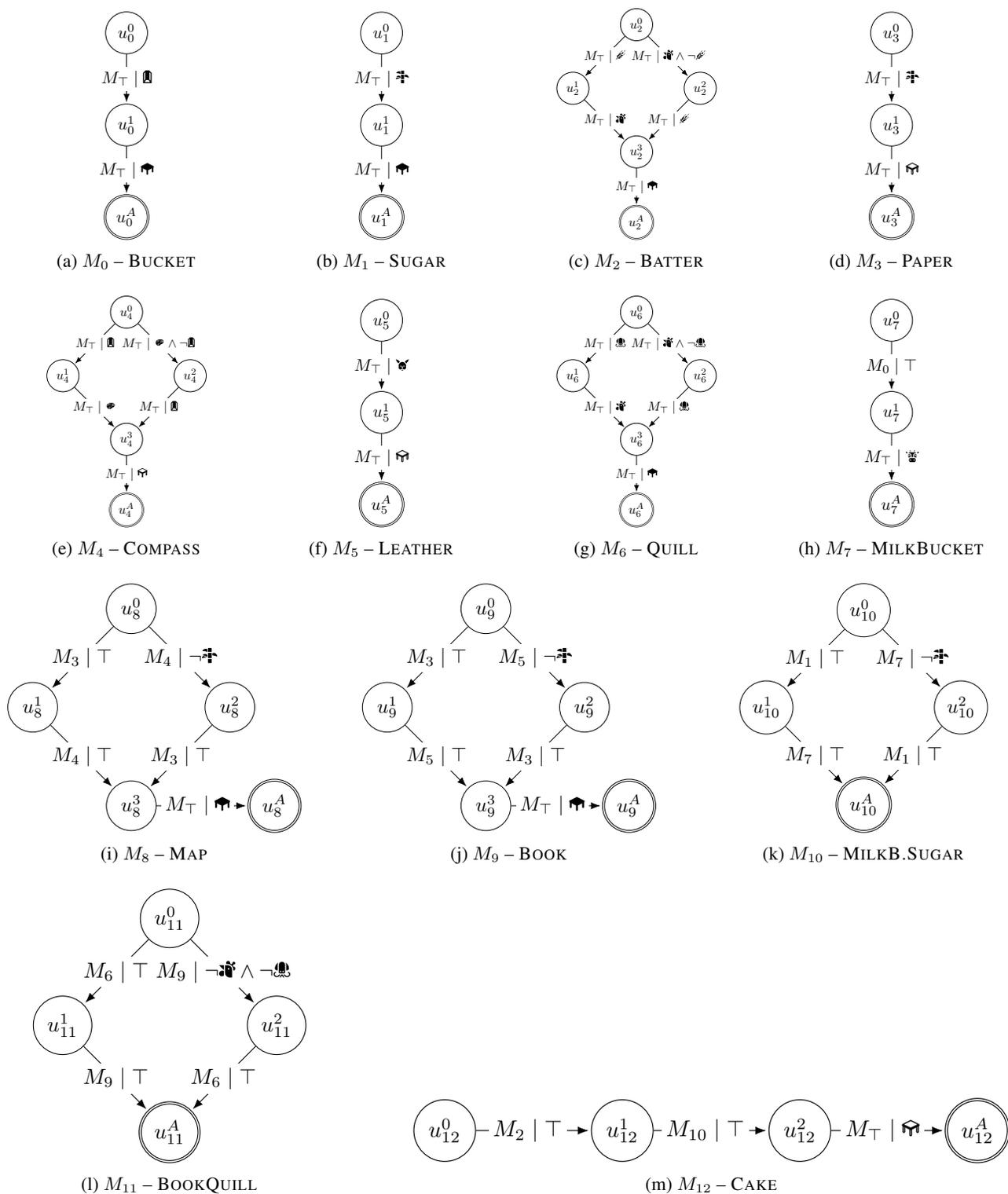


Figure 12: Root reward machines for each of the CRAFTWORLD tasks.

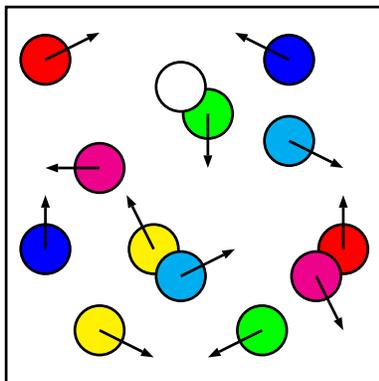


Figure 13: An instance of the WATERWORLD (Toro Icarte et al., 2018) in the WOD setting. Image taken from (Furelos-Blanco et al., 2021).

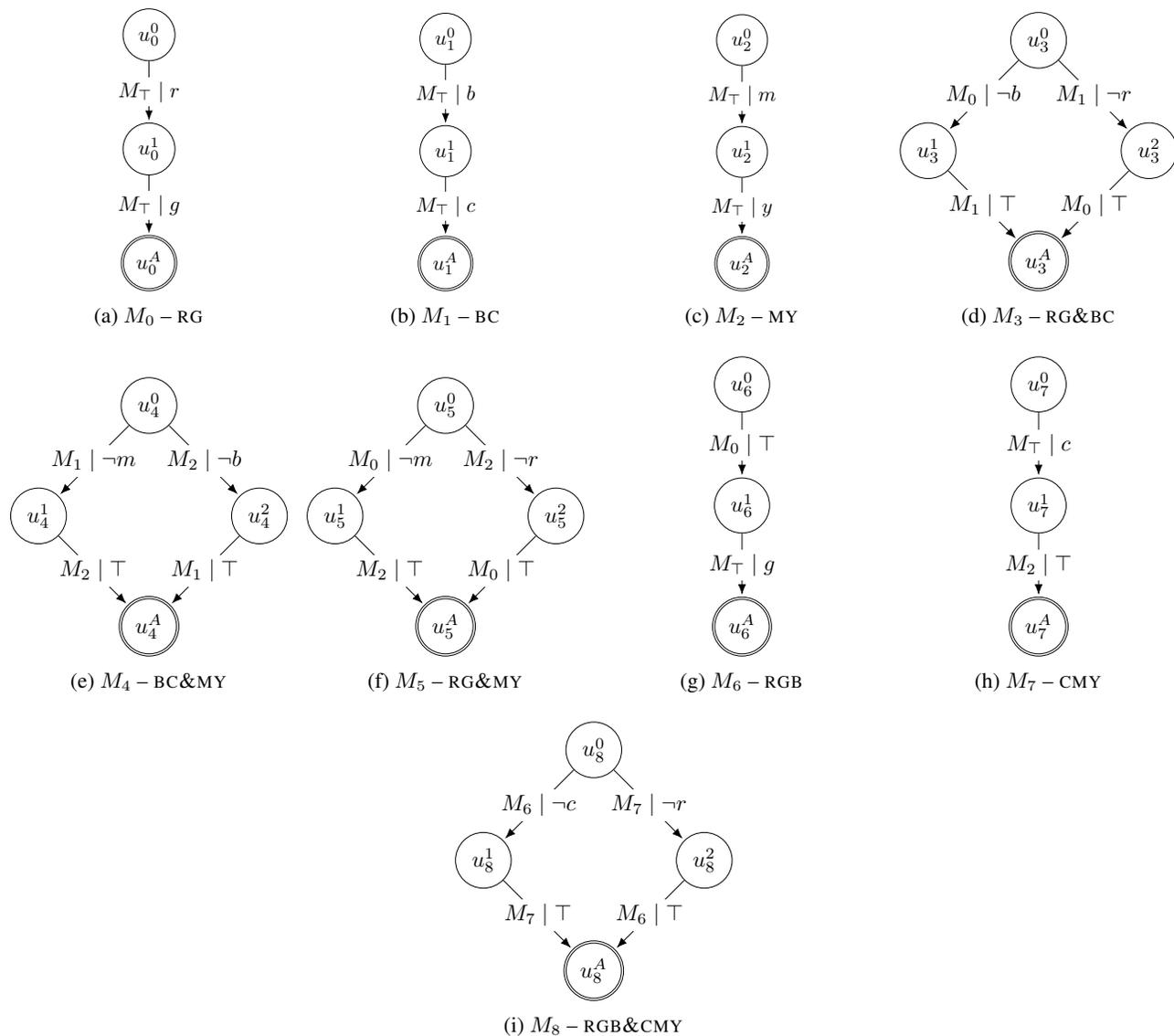


Figure 14: Root reward machines for each of the WATERWORLD tasks.

Table 2: List of hyperparameters and their values.

Parameter	CRAFTWORLD	WATERWORLD
<i>General</i>		
Episodes		
Without HRM learning	100,000 (OP, OPL); 200,000 (FR, FRL)	100,000 (WOD); 200,000 (WD)
With HRM learning	150,000 (OP, OPL); 300,000 (FR, FRL)	150,000 (WOD); 300,000 (WD)
Maximum episode length	1,000	1,000
Num. of instances I	10	10
<i>HRM policy learning (Section 4, Appendix B)</i>		
Learning rate α	5×10^{-4}	1×10^{-5}
Learning rate (SMDP) α	5×10^{-4}	1×10^{-3}
Optimizer	RMSprop (Hinton et al., 2012)	RMSprop (Hinton et al., 2012)
Discount γ	0.9	0.9
Discount (SMDP) γ	0.99	0.99
Updated formula Q-functions per step	4	4
Replay memory size	500,000	500,000
Replay start size	100,000	100,000
Target network update frequency	1,500	1,500
Replay memory size (SMDP)	10,000	10,000
Replay start size (SMDP)	1,000	1,000
Target network update frequency (SMDP)	500	500
Minibatch size	32	32
Initial exploration	1.0	1.0
Final exploration	0.1	0.1
Annealing steps	2,000,000	5,000,000
Annealing steps (SMDP)	10,000	10,000
<i>HRM learning (Section 5, Appendix C)</i>		
Curriculum weight β	0.99	0.99
Curriculum threshold	0.85	0.75
Curriculum update frequency (# episodes)	100	100
ILASP time budget	2 hours	2 hours
Num. collected goal traces ρ (height 1)	25	25
Num. collected goal traces ρ (height ≥ 2)	150	150
Num. goal traces ρ_s to learn first HRM	10	10
<i>CRM (Toro Icarte et al., 2022)</i>		
Learning rate α	5×10^{-4}	1×10^{-5}
Optimizer	RMSprop (Hinton et al., 2012)	RMSprop (Hinton et al., 2012)
Discount	0.99	0.99
Replay memory size	1,000,000	1,000,000
Replay start size	100,000	100,000
Target network update frequency	1,500	1,500
Minibatch size	32	32
Initial exploration	1.0	1.0
Final exploration	0.1	0.1
Annealing steps	100,000,000	2,000,000

Hierarchies of Reward Machines

Table 3: Results of LHRM in CRAFTWORLD for the default case.

(a) OP

Task	# G	# L	Time (s.)	Calls	States	Edges	Ep. First HRM ($\times 10^2$)	# Examples		Example Length		
								G	I	G	I	
BATTER	5	5	11.1 (1.7)	17.8 (1.9)	5.0 (0.0)	5.2 (0.2)	1.8 (0.1)	12.2 (0.7)	11.6 (1.4)	26.5 (2.1)	24.2 (3.2)	
BUCKET	5	5	0.9 (0.0)	3.6 (0.2)	3.0 (0.0)	2.0 (0.0)	1.7 (0.1)	10.0 (0.0)	1.6 (0.2)	19.4 (1.1)	19.3 (5.7)	
COMPASS	5	5	135.4 (73.3)	18.6 (1.6)	5.0 (0.0)	5.2 (0.2)	1.8 (0.2)	11.8 (0.6)	12.8 (1.4)	28.7 (1.9)	20.3 (2.8)	
LEATHER	5	5	0.9 (0.0)	3.8 (0.2)	3.0 (0.0)	2.0 (0.0)	1.8 (0.1)	10.0 (0.0)	1.8 (0.2)	16.7 (1.7)	17.9 (4.4)	
PAPER	5	5	0.8 (0.1)	3.4 (0.2)	3.0 (0.0)	2.0 (0.0)	1.6 (0.1)	10.0 (0.0)	1.4 (0.2)	19.8 (2.0)	40.6 (27.0)	
QUILL	5	5	18.0 (3.5)	19.8 (1.2)	5.0 (0.0)	5.2 (0.2)	2.1 (0.1)	13.2 (0.4)	12.6 (1.1)	29.6 (2.5)	24.4 (3.2)	
SUGAR	5	5	0.8 (0.1)	3.2 (0.2)	3.0 (0.0)	2.0 (0.0)	1.7 (0.2)	10.0 (0.0)	1.2 (0.2)	17.7 (1.6)	17.5 (3.2)	
BOOK	5	5	191.2 (36.4)	22.8 (2.6)	5.0 (0.0)	5.8 (0.2)	6.0 (0.2)	11.4 (0.7)	17.4 (2.2)	20.5 (1.8)	24.8 (1.5)	
MAP	5	5	549.4 (149.5)	33.4 (3.2)	5.0 (0.0)	5.6 (0.2)	6.0 (0.2)	12.2 (0.6)	27.2 (2.9)	29.5 (3.2)	28.7 (1.7)	
MILKBUCKET	5	5	1.5 (0.2)	4.6 (0.4)	3.0 (0.0)	2.0 (0.0)	6.8 (0.5)	10.0 (0.0)	2.6 (0.4)	11.6 (0.7)	15.3 (4.3)	
BOOKQUILL	5	5	17.9 (1.4)	19.6 (1.1)	4.0 (0.0)	4.0 (0.0)	3.8 (0.1)	10.0 (0.0)	16.6 (1.1)	27.2 (1.3)	20.8 (1.4)	
MILKB.SUGAR	5	5	7.3 (1.2)	12.4 (1.2)	4.0 (0.0)	4.0 (0.0)	3.8 (0.1)	10.2 (0.2)	9.2 (1.2)	16.9 (0.8)	14.3 (1.7)	
CAKE	5	5	74.5 (25.7)	26.4 (3.7)	4.0 (0.0)	3.2 (0.2)	2.1 (0.1)	10.2 (0.2)	23.2 (3.6)	38.4 (0.9)	22.7 (1.6)	
Completed Runs = 5			Total Time (s.) = 1009.8 (122.3)					Total Calls = 189.4 (4.1)				

(b) OPL

Task	# G	# L	Time (s.)	Calls	States	Edges	Ep. First HRM ($\times 10^2$)	# Examples			Example Length		
								G	D	I	G	D	I
BATTER	5	5	13.7 (2.9)	23.0 (3.0)	6.0 (0.0)	9.2 (0.2)	12.0 (1.0)	11.4 (0.4)	7.0 (1.2)	10.6 (1.6)	20.4 (1.1)	18.7 (1.6)	12.1 (1.7)
BUCKET	5	5	1.8 (0.2)	7.2 (0.6)	4.0 (0.0)	4.0 (0.0)	8.0 (0.5)	10.2 (0.2)	2.2 (0.2)	2.8 (0.4)	10.2 (0.5)	13.4 (1.9)	6.8 (1.7)
COMPASS	5	5	13.1 (1.7)	22.0 (1.7)	6.0 (0.0)	9.2 (0.2)	10.4 (1.4)	11.0 (0.6)	6.8 (1.0)	10.2 (1.0)	17.2 (1.6)	20.9 (1.9)	14.3 (0.8)
LEATHER	5	5	1.9 (0.2)	7.0 (0.5)	4.0 (0.0)	4.0 (0.0)	6.9 (0.5)	10.0 (0.0)	2.4 (0.2)	2.6 (0.4)	11.1 (0.9)	16.9 (5.6)	8.9 (3.3)
PAPER	5	5	2.0 (0.2)	7.6 (0.6)	4.0 (0.0)	4.0 (0.0)	7.7 (1.1)	10.0 (0.0)	3.0 (0.3)	2.6 (0.4)	10.1 (0.9)	18.9 (3.3)	5.6 (0.8)
QUILL	5	5	11.3 (1.2)	22.0 (1.2)	6.0 (0.0)	9.2 (0.2)	12.8 (1.5)	10.6 (0.2)	6.4 (0.7)	11.0 (0.9)	15.3 (1.3)	13.5 (1.0)	12.1 (1.4)
SUGAR	5	5	1.7 (0.1)	6.4 (0.4)	4.0 (0.0)	4.0 (0.0)	6.5 (0.7)	10.0 (0.0)	2.4 (0.2)	2.0 (0.3)	9.6 (0.6)	15.3 (3.6)	16.6 (9.2)
BOOK	5	5	427.8 (201.6)	32.6 (4.2)	6.0 (0.0)	6.6 (0.2)	5.6 (0.2)	12.0 (0.3)	3.6 (0.7)	23.0 (3.4)	21.6 (1.5)	25.9 (3.4)	23.7 (1.3)
MAP	5	5	647.9 (110.7)	38.6 (3.6)	6.0 (0.0)	6.4 (0.2)	5.6 (0.2)	11.2 (0.4)	3.8 (0.9)	29.6 (3.5)	23.1 (1.0)	27.8 (4.6)	26.1 (0.4)
MILKBUCKET	5	5	2.1 (0.2)	5.4 (0.4)	4.0 (0.0)	3.0 (0.0)	7.6 (0.5)	10.0 (0.0)	1.4 (0.4)	2.0 (0.0)	11.1 (0.5)	26.3 (6.5)	15.2 (5.8)
BOOKQUILL	5	5	18.7 (2.3)	16.6 (1.3)	4.0 (0.0)	4.0 (0.0)	3.7 (0.2)	10.0 (0.0)	0.4 (0.2)	13.2 (1.4)	29.0 (1.1)	6.2 (5.5)	27.8 (1.4)
MILKB.SUGAR	5	5	7.7 (0.7)	12.2 (0.9)	4.0 (0.0)	4.0 (0.0)	3.8 (0.2)	10.0 (0.0)	0.2 (0.2)	9.0 (0.9)	16.0 (0.9)	1.6 (1.6)	16.3 (1.3)
CAKE	5	5	472.9 (216.6)	36.0 (6.0)	5.0 (0.0)	4.6 (0.2)	2.1 (0.0)	10.0 (0.0)	1.6 (0.4)	31.4 (5.7)	39.5 (1.2)	41.5 (8.6)	26.9 (0.8)
Completed Runs = 5			Total Time (s.) = 1622.6 (328.7)					Total Calls = 236.6 (9.3)					

(c) FR

Task	# G	# L	Time (s.)	Calls	States	Edges	Ep. First HRM ($\times 10^2$)	# Examples		Example Length		
								G	I	G	I	
BATTER	5	5	12.3 (1.7)	17.6 (1.3)	5.0 (0.0)	5.4 (0.2)	9.2 (1.2)	11.6 (0.4)	12.0 (1.2)	30.3 (2.3)	27.8 (2.0)	
BUCKET	5	5	1.2 (0.1)	3.8 (0.2)	3.0 (0.0)	2.0 (0.0)	6.7 (0.9)	10.0 (0.0)	1.8 (0.2)	16.6 (2.5)	28.5 (4.1)	
COMPASS	5	5	14.1 (1.6)	20.2 (1.7)	5.0 (0.0)	5.2 (0.2)	9.8 (0.7)	11.6 (0.6)	14.6 (1.2)	26.5 (0.8)	26.5 (2.1)	
LEATHER	5	5	1.1 (0.1)	3.6 (0.2)	3.0 (0.0)	2.0 (0.0)	4.5 (0.7)	10.0 (0.0)	1.6 (0.2)	13.4 (1.3)	16.7 (3.6)	
PAPER	5	5	1.2 (0.0)	4.0 (0.0)	3.0 (0.0)	2.0 (0.0)	4.9 (0.9)	10.0 (0.0)	2.0 (0.0)	12.4 (1.1)	10.9 (2.5)	
QUILL	5	5	8.9 (0.9)	16.0 (0.8)	5.0 (0.0)	5.2 (0.2)	9.4 (1.7)	10.6 (0.2)	11.4 (0.6)	25.4 (0.3)	25.5 (2.7)	
SUGAR	5	5	1.1 (0.1)	3.8 (0.2)	3.0 (0.0)	2.0 (0.0)	5.2 (0.3)	10.0 (0.0)	1.8 (0.2)	15.3 (1.7)	21.0 (10.1)	
BOOK	5	5	220.2 (83.3)	25.2 (3.4)	5.0 (0.0)	5.6 (0.2)	6.1 (0.2)	10.2 (0.2)	21.0 (3.4)	21.9 (1.0)	18.4 (0.7)	
MAP	5	5	628.3 (85.4)	37.8 (3.7)	5.0 (0.0)	5.6 (0.2)	5.8 (0.1)	10.0 (0.0)	33.8 (3.7)	26.4 (1.0)	21.4 (0.7)	
MILKBUCKET	5	5	1.9 (0.2)	5.0 (0.3)	3.0 (0.0)	2.0 (0.0)	9.8 (0.7)	10.0 (0.0)	3.0 (0.3)	13.2 (0.7)	12.8 (3.2)	
BOOKQUILL	5	5	12.9 (2.2)	15.6 (1.7)	4.0 (0.0)	4.0 (0.0)	3.9 (0.1)	10.0 (0.0)	12.6 (1.7)	29.0 (1.5)	13.3 (0.8)	
MILKB.SUGAR	5	5	7.2 (0.6)	12.0 (0.7)	4.0 (0.0)	4.0 (0.0)	3.9 (0.2)	10.0 (0.0)	9.0 (0.7)	18.9 (0.9)	10.1 (1.0)	
CAKE	5	5	121.1 (41.1)	34.0 (4.8)	4.0 (0.0)	3.0 (0.0)	2.2 (0.0)	10.0 (0.0)	31.0 (4.8)	42.2 (1.7)	16.2 (1.1)	
Completed Runs = 5			Total Time (s.) = 1031.6 (150.3)					Total Calls = 198.6 (11.3)				

(d) FRL

Task	# G	# L	Time (s.)	Calls	States	Edges	Ep. First HRM ($\times 10^2$)	# Examples			Example Length		
								G	D	I	G	D	I
BATTER	5	5	11.3 (1.4)	23.4 (2.5)	6.0 (0.0)	9.2 (0.2)	468.4 (121.9)	10.4 (0.2)	7.6 (0.9)	11.4 (1.9)	11.9 (0.6)	10.1 (1.3)	9.9 (0.4)
BUCKET	5	5	2.3 (0.2)	7.0 (0.3)	4.0 (0.0)	4.0 (0.0)	129.5 (69.4)	10.2 (0.2)	2.8 (0.2)	2.0 (0.3)	7.8 (0.5)	9.9 (1.7)	6.4 (2.1)
COMPASS	5	5	13.0 (1.9)	24.6 (2.2)	6.0 (0.0)	9.4 (0.2)	550.8 (156.4)	10.4 (0.2)	7.8 (1.0)	12.4 (1.2)	12.5 (1.6)	9.4 (1.0)	8.4 (0.5)
LEATHER	5	5	2.5 (0.3)	7.8 (0.7)	4.0 (0.0)	4.0 (0.0)	89.0 (18.0)	10.0 (0.0)	3.2 (0.4)	2.6 (0.4)	7.3 (0.4)	9.3 (1.7)	3.7 (0.4)
PAPER	5	5	2.2 (0.1)	7.0 (0.3)	4.0 (0.0)	4.0 (0.0)	82.7 (18.8)	10.0 (0.0)	3.0 (0.0)	2.0 (0.3)	6.9 (0.7)	10.2 (1.8)	4.7 (2.7)
QUILL	5	5	11.6 (1.1)	23.8 (1.5)	6.0 (0.0)	9.6 (0.2)	458.9 (61.0)	10.6 (0.2)	8.0 (0.9)	11.2 (1.2)	11.9 (0.6)	13.1 (2.7)	9.2 (0.8)
SUGAR	5	5	2.7 (0.2)	8.4 (0.7)	4.0 (0.0)	4.0 (0.0)	103.5 (39.5)	10.0 (0.0)	3.6 (0.4)	2.8 (0.5)	8.2 (0.7)	10.1 (1.9)	5.0 (1.1)
BOOK	5	5	301.7 (98.1)	36.4 (1.9)	6.0 (0.0)	6.8 (0.2)	5.3 (0.1)	10.2 (0.2)	5.0 (0.7)	27.2 (1.9)	21.7 (1.1)	18.8 (2.2)	16.1 (0.6)
MAP	5	5	754.1 (158.2)	44.6 (2.6)	6.0 (0.0)	7.0 (0.0)	5.5 (0.2)	10.2 (0.2)	5.2 (0.4)	35.2 (2.3)	25.6 (0.5)	20.4 (2.9)	18.7 (0.6)
MILKBUCKET	5	5	2.8 (0.1)	6.6 (0.2)	4.0 (0.0)	3.0 (0.0)	6.9 (0.4)	10.0 (0.0)	2.0 (0.0)	2.6 (0.2)	12.5 (0.8)	13.1 (3.7)	7.4 (2.2)
BOOKQUILL	5	5	19.8 (2.9)	19.6 (1.6)	4.0 (0.0)	4.0 (0.0)	4.3 (0.1)	10.0 (0.0)	0.8 (0.4)	15.8 (1.2)	28.4 (1.1)	2.7 (1.3)	13.5 (0.9)
MILKB.SUGAR	5	5	8.8 (0.9)	12.6 (1.0)	4.0 (0.0)	4.0 (0.0)	4.0 (0.1)	10.0 (0.0)	1.2 (0.5)	8.4 (0.7)	19.3 (1.3)	3.7 (2.0)	10.7 (2.0)
CAKE	5	5	344.0 (87.7)	46.2 (4.9)	5.0 (0.0)	4.8 (0.2)	2.8 (0.1)	10.0 (0.0)	2.8 (0.7)	40.4 (4.5)	44.5 (2.3)	21.8 (2.2)	17.3 (1.0)
Completed Runs = 5			Total Time (s.) = 1476.8 (175.3)					Total Calls = 268.0 (6.5)					

Hierarchies of Reward Machines

Table 4: Results of LHRM in CRAFTWORLD with a restricted set of callable RMs.

(a) OP

Task	# G	# L	Time (s.)	Calls	States	Edges	Ep. First HRM ($\times 10^2$)	# Examples		Example Length		
								G	I	G	I	
BATTER	5	5	11.2 (1.6)	17.8 (1.9)	5.0 (0.0)	5.2 (0.2)	1.8 (0.1)	12.2 (0.7)	11.6 (1.4)	26.5 (2.1)	24.2 (3.2)	
BUCKET	5	5	0.9 (0.0)	3.6 (0.2)	3.0 (0.0)	2.0 (0.0)	1.7 (0.1)	10.0 (0.0)	1.6 (0.2)	19.4 (1.1)	19.3 (5.7)	
COMPASS	5	5	15.5 (4.2)	18.6 (1.6)	5.0 (0.0)	5.2 (0.2)	1.8 (0.2)	11.8 (0.6)	12.8 (1.4)	28.7 (1.9)	20.3 (2.8)	
LEATHER	5	5	0.9 (0.0)	3.8 (0.2)	3.0 (0.0)	2.0 (0.0)	1.8 (0.1)	10.0 (0.0)	1.8 (0.2)	16.7 (1.7)	17.9 (4.4)	
PAPER	5	5	0.9 (0.0)	3.4 (0.2)	3.0 (0.0)	2.0 (0.0)	1.6 (0.1)	10.0 (0.0)	1.4 (0.2)	19.8 (2.0)	40.6 (27.0)	
QUILL	5	5	18.2 (3.5)	19.8 (1.2)	5.0 (0.0)	5.2 (0.2)	2.1 (0.1)	13.2 (0.4)	12.6 (1.1)	29.6 (2.5)	24.4 (3.2)	
SUGAR	5	5	0.8 (0.0)	3.2 (0.2)	3.0 (0.0)	2.0 (0.0)	1.7 (0.2)	10.0 (0.0)	1.2 (0.2)	17.7 (1.6)	17.5 (3.2)	
BOOK	5	5	45.8 (4.5)	19.6 (0.9)	5.0 (0.0)	5.6 (0.2)	6.0 (0.2)	11.2 (1.0)	14.4 (0.9)	21.6 (1.8)	21.0 (1.7)	
MAP	5	5	64.1 (10.6)	22.0 (2.6)	5.0 (0.0)	5.2 (0.2)	6.1 (0.2)	10.8 (0.4)	17.2 (2.7)	22.5 (1.6)	23.0 (1.2)	
MILKBUCKET	5	5	1.2 (0.1)	4.4 (0.4)	3.0 (0.0)	2.0 (0.0)	6.8 (0.3)	10.0 (0.0)	2.4 (0.4)	12.1 (0.7)	15.3 (1.6)	
BOOKQUILL	5	5	4.5 (0.8)	10.2 (1.4)	4.0 (0.0)	4.0 (0.0)	3.9 (0.1)	10.0 (0.0)	7.2 (1.4)	26.1 (0.8)	22.4 (0.9)	
MILKB.SUGAR	5	5	3.5 (0.5)	9.6 (1.3)	4.0 (0.0)	4.0 (0.0)	3.9 (0.1)	10.2 (0.2)	6.4 (1.2)	17.4 (0.5)	12.5 (0.8)	
CAKE	5	5	9.1 (0.9)	17.0 (0.9)	4.0 (0.0)	3.2 (0.2)	2.1 (0.1)	10.0 (0.0)	14.0 (0.9)	37.5 (1.9)	18.0 (1.9)	
Completed Runs = 5			Total Time (s.) = 176.6 (13.1)					Total Calls = 153.0 (3.6)				

(b) OPL

Task	# G	# L	Time (s.)	Calls	States	Edges	Ep. First HRM ($\times 10^2$)	# Examples			Example Length		
								G	D	I	G	D	I
BATTER	5	5	13.9 (3.0)	23.0 (3.0)	6.0 (0.0)	9.2 (0.2)	12.0 (1.0)	11.4 (0.4)	7.0 (1.2)	10.6 (1.6)	20.4 (1.1)	18.7 (1.6)	12.1 (1.7)
BUCKET	5	5	1.8 (0.1)	7.2 (0.6)	4.0 (0.0)	4.0 (0.0)	8.0 (0.5)	10.2 (0.2)	2.2 (0.2)	2.8 (0.4)	10.2 (0.5)	13.4 (1.9)	6.8 (1.7)
COMPASS	5	5	13.2 (1.7)	22.0 (1.7)	6.0 (0.0)	9.2 (0.2)	10.4 (1.4)	11.0 (0.6)	6.8 (1.0)	10.2 (1.0)	17.2 (1.6)	20.9 (1.9)	14.3 (0.8)
LEATHER	5	5	1.9 (0.1)	7.0 (0.5)	4.0 (0.0)	4.0 (0.0)	6.9 (0.5)	10.0 (0.0)	2.4 (0.2)	2.6 (0.4)	11.1 (0.9)	16.9 (5.6)	8.9 (3.3)
PAPER	5	5	2.0 (0.2)	7.6 (0.6)	4.0 (0.0)	4.0 (0.0)	7.7 (1.1)	10.0 (0.0)	3.0 (0.3)	2.6 (0.4)	10.1 (0.9)	18.9 (3.3)	5.6 (0.8)
QUILL	5	5	11.5 (1.3)	22.0 (1.2)	6.0 (0.0)	9.2 (0.2)	12.8 (1.5)	10.6 (0.2)	6.4 (0.7)	11.0 (0.9)	15.3 (1.3)	13.5 (1.0)	12.1 (1.4)
SUGAR	5	5	1.6 (0.1)	6.4 (0.4)	4.0 (0.0)	4.0 (0.0)	6.5 (0.7)	10.0 (0.0)	2.4 (0.2)	2.0 (0.3)	9.6 (0.6)	15.3 (3.6)	16.6 (9.2)
BOOK	5	5	69.0 (20.5)	21.8 (2.2)	6.0 (0.0)	6.2 (0.2)	5.5 (0.1)	10.4 (0.2)	5.2 (0.9)	12.2 (1.9)	20.4 (1.3)	21.2 (2.0)	20.8 (1.7)
MAP	5	5	76.5 (6.0)	24.2 (1.3)	6.0 (0.0)	6.4 (0.2)	5.7 (0.3)	11.6 (0.8)	4.0 (0.3)	14.6 (1.0)	24.8 (3.0)	21.4 (2.1)	25.7 (0.8)
MILKBUCKET	5	5	1.7 (0.2)	6.0 (0.6)	4.0 (0.0)	3.0 (0.0)	7.5 (0.7)	10.2 (0.2)	1.4 (0.2)	2.4 (0.2)	11.7 (0.7)	25.4 (6.4)	14.2 (3.4)
BOOKQUILL	5	5	5.3 (0.9)	10.8 (1.4)	4.0 (0.0)	4.0 (0.0)	3.7 (0.1)	10.0 (0.0)	1.0 (0.5)	6.8 (0.9)	27.7 (1.0)	11.2 (5.4)	21.1 (1.7)
MILKB.SUGAR	5	5	4.0 (0.9)	9.8 (1.7)	4.0 (0.0)	4.0 (0.0)	3.8 (0.1)	10.0 (0.0)	1.6 (0.7)	5.2 (1.2)	18.4 (0.7)	8.3 (2.9)	15.6 (1.7)
CAKE	5	5	16.2 (0.4)	20.8 (0.2)	5.0 (0.0)	4.0 (0.0)	2.1 (0.1)	10.0 (0.0)	3.2 (0.2)	14.6 (0.2)	38.1 (0.9)	22.5 (3.3)	25.8 (1.7)
Completed Runs = 5			Total Time (s.) = 218.6 (21.1)					Total Calls = 188.6 (5.4)					

(c) FR

Task	# G	# L	Time (s.)	Calls	States	Edges	Ep. First HRM ($\times 10^2$)	# Examples		Example Length		
								G	I	G	I	
BATTER	5	5	12.6 (1.8)	17.6 (1.3)	5.0 (0.0)	5.4 (0.2)	9.2 (1.2)	11.6 (0.4)	12.0 (1.2)	30.3 (2.3)	27.8 (2.0)	
BUCKET	5	5	1.2 (0.1)	3.8 (0.2)	3.0 (0.0)	2.0 (0.0)	6.7 (0.9)	10.0 (0.0)	1.8 (0.2)	16.6 (2.5)	28.5 (4.1)	
COMPASS	5	5	14.1 (1.5)	20.2 (1.7)	5.0 (0.0)	5.2 (0.2)	9.8 (0.7)	11.6 (0.6)	14.6 (1.2)	26.5 (0.8)	26.5 (2.1)	
LEATHER	5	5	1.1 (0.1)	3.6 (0.2)	3.0 (0.0)	2.0 (0.0)	4.5 (0.7)	10.0 (0.0)	1.6 (0.2)	13.4 (1.3)	16.7 (3.6)	
PAPER	5	5	1.2 (0.1)	4.0 (0.0)	3.0 (0.0)	2.0 (0.0)	4.9 (0.9)	10.0 (0.0)	2.0 (0.0)	12.4 (1.1)	10.9 (2.5)	
QUILL	5	5	9.3 (0.8)	16.0 (0.8)	5.0 (0.0)	5.2 (0.2)	9.4 (1.7)	10.6 (0.2)	11.4 (0.6)	25.4 (0.3)	25.5 (2.7)	
SUGAR	5	5	1.4 (0.2)	3.8 (0.2)	3.0 (0.0)	2.0 (0.0)	5.2 (0.3)	10.0 (0.0)	1.8 (0.2)	15.3 (1.7)	21.0 (1.0)	
BOOK	5	5	43.8 (13.0)	20.0 (1.9)	5.0 (0.0)	5.4 (0.2)	6.0 (0.1)	10.0 (0.0)	16.0 (1.9)	21.9 (1.0)	14.7 (1.4)	
MAP	5	5	85.2 (13.4)	22.2 (2.5)	5.0 (0.0)	5.2 (0.2)	5.9 (0.1)	10.2 (0.2)	18.0 (2.6)	26.5 (0.9)	18.2 (1.2)	
MILKBUCKET	5	5	1.4 (0.1)	4.4 (0.2)	3.0 (0.0)	2.0 (0.0)	10.2 (0.9)	10.0 (0.0)	2.4 (0.2)	13.0 (0.8)	12.2 (2.8)	
BOOKQUILL	5	5	6.3 (0.9)	13.2 (1.7)	4.0 (0.0)	4.0 (0.0)	3.8 (0.1)	10.0 (0.0)	10.2 (1.7)	30.6 (2.0)	11.9 (1.2)	
MILKB.SUGAR	5	5	4.8 (0.6)	11.8 (1.3)	4.0 (0.0)	4.0 (0.0)	3.8 (0.1)	10.0 (0.0)	8.8 (1.3)	19.8 (0.7)	8.6 (1.0)	
CAKE	5	5	12.5 (1.8)	20.8 (2.5)	4.0 (0.0)	3.0 (0.0)	2.3 (0.1)	10.0 (0.0)	17.8 (2.5)	44.2 (2.6)	13.1 (1.2)	
Completed Runs = 5			Total Time (s.) = 194.9 (17.6)					Total Calls = 161.4 (7.0)				

(d) FRL

Task	# G	# L	Time (s.)	Calls	States	Edges	Ep. First HRM ($\times 10^2$)	# Examples			Example Length		
								G	D	I	G	D	I
BATTER	5	5	11.2 (1.4)	23.4 (2.5)	6.0 (0.0)	9.2 (0.2)	468.4 (121.9)	10.4 (0.2)	7.6 (0.9)	11.4 (1.9)	11.9 (0.6)	10.1 (1.3)	9.9 (0.4)
BUCKET	5	5	2.4 (0.1)	7.0 (0.3)	4.0 (0.0)	4.0 (0.0)	129.5 (69.4)	10.2 (0.2)	2.8 (0.2)	2.0 (0.3)	7.8 (0.5)	9.9 (1.7)	6.4 (2.1)
COMPASS	5	5	13.1 (1.9)	24.6 (2.2)	6.0 (0.0)	9.4 (0.2)	550.8 (156.4)	10.4 (0.2)	7.8 (1.0)	12.4 (1.2)	12.5 (1.6)	9.4 (1.0)	8.4 (0.5)
LEATHER	5	5	2.5 (0.4)	7.8 (0.7)	4.0 (0.0)	4.0 (0.0)	89.0 (18.0)	10.0 (0.0)	3.2 (0.4)	2.6 (0.4)	7.3 (0.4)	9.3 (1.7)	3.7 (0.4)
PAPER	5	5	2.1 (0.1)	7.0 (0.3)	4.0 (0.0)	4.0 (0.0)	82.7 (18.8)	10.0 (0.0)	3.0 (0.0)	2.0 (0.3)	6.9 (0.7)	10.2 (1.8)	4.7 (2.7)
QUILL	5	5	11.6 (1.2)	23.8 (1.5)	6.0 (0.0)	9.6 (0.2)	458.9 (61.0)	10.6 (0.2)	8.0 (0.9)	11.2 (1.2)	11.9 (0.6)	13.1 (2.7)	9.2 (0.8)
SUGAR	5	5	2.6 (0.2)	8.4 (0.7)	4.0 (0.0)	4.0 (0.0)	103.5 (39.5)	10.0 (0.0)	3.6 (0.4)	2.8 (0.5)	8.2 (0.7)	10.1 (1.9)	5.0 (1.1)
BOOK	5	5	62.2 (13.2)	27.4 (2.2)	6.0 (0.0)	6.6 (0.2)	5.3 (0.1)	10.2 (0.2)	5.6 (0.6)	17.6 (1.7)	23.0 (1.0)	16.5 (2.0)	13.4 (1.0)
MAP	5	5	131.3 (28.0)	34.0 (3.0)	6.0 (0.0)	6.6 (0.2)	5.5 (0.2)	10.2 (0.2)	6.8 (0.7)	23.0 (2.4)	26.2 (0.7)	16.9 (1.6)	14.5 (0.5)
MILKBUCKET	5	5	2.7 (0.7)	6.6 (0.6)	4.0 (0.0)	3.0 (0.0)	6.8 (0.3)	10.0 (0.0)	2.2 (0.2)	2.4 (0.4)	12.0 (0.8)	9.4 (1.0)	9.9 (2.3)
BOOKQUILL	5	5	6.8 (0.6)	12.6 (0.7)	4.0 (0.0)	4.0 (0.0)	4.4 (0.2)	10.0 (0.0)	1.6 (0.5)	8.0 (0.3)	32.3 (2.3)	4.9 (1.5)	11.6 (1.1)
MILKB.SUGAR	5	5	5.4 (0.6)	12.2 (1.0)	4.0 (0.0)	4.0 (0.0)	4.0 (0.1)	10.2 (0.2)	1.0 (0.4)	8.0 (0.4)	20.4 (1.7)	2.9 (1.4)	10.3 (1.2)
CAKE	5	5	16.3 (1.2)	21.2 (1.0)	5.0 (0.0)	4.0 (0.0)	2.8 (0.0)	10.0 (0.0)	2.6 (0.2)	15.6 (1.2)	47.7 (1.9)	15.0 (0.7)	16.0 (0.9)
Completed Runs = 5			Total Time (s.) = 270.1 (34.6)					Total Calls = 216.0 (5.1)					

Hierarchies of Reward Machines

Table 5: Results of LHRM in CRAFTWORLD without exploration using options.

(a) OP

Task	# G	# L	Time (s.)	Calls	States	Edges	Ep. First HRM ($\times 10^2$)	# Examples		Example Length		
								G	I	G	I	
BATTER	5	5	11.2 (1.7)	17.8 (1.9)	5.0 (0.0)	5.2 (0.2)	1.8 (0.1)	12.2 (0.7)	11.6 (1.4)	26.5 (2.1)	24.2 (3.2)	
BUCKET	5	5	0.9 (0.0)	3.6 (0.2)	3.0 (0.0)	2.0 (0.0)	1.7 (0.1)	10.0 (0.0)	1.6 (0.2)	19.4 (1.1)	19.3 (5.7)	
COMPASS	5	5	15.6 (4.1)	18.6 (1.6)	5.0 (0.0)	5.2 (0.2)	1.8 (0.2)	11.8 (0.6)	12.8 (1.4)	28.7 (1.9)	20.3 (2.8)	
LEATHER	5	5	0.9 (0.1)	3.8 (0.2)	3.0 (0.0)	2.0 (0.0)	1.8 (0.1)	10.0 (0.0)	1.8 (0.2)	16.7 (1.7)	17.9 (4.4)	
PAPER	5	5	0.9 (0.1)	3.4 (0.2)	3.0 (0.0)	2.0 (0.0)	1.6 (0.1)	10.0 (0.0)	1.4 (0.2)	19.8 (2.0)	40.6 (27.0)	
QUILL	5	5	18.3 (3.6)	19.8 (1.2)	5.0 (0.0)	5.2 (0.2)	2.1 (0.1)	13.2 (0.4)	12.6 (1.1)	29.6 (2.5)	24.4 (3.2)	
SUGAR	5	5	0.9 (0.0)	3.2 (0.2)	3.0 (0.0)	2.0 (0.0)	1.7 (0.2)	10.0 (0.0)	1.2 (0.2)	17.7 (1.6)	17.5 (3.2)	
BOOK	5	5	529.0 (164.2)	21.2 (1.4)	5.0 (0.0)	5.8 (0.2)	6.8 (0.2)	10.2 (0.2)	17.0 (1.5)	33.0 (2.6)	23.7 (1.3)	
MAP	5	5	1924.2 (443.5)	28.0 (3.8)	5.0 (0.0)	5.4 (0.2)	7.8 (0.4)	10.4 (0.2)	23.6 (3.7)	40.1 (1.0)	29.4 (1.3)	
MILKBUCKET	5	5	1.6 (0.2)	4.4 (0.4)	3.0 (0.0)	2.0 (0.0)	6.1 (0.3)	10.0 (0.0)	2.4 (0.4)	16.0 (1.0)	14.2 (1.3)	
BOOKQUILL	5	5	42.7 (10.1)	24.6 (3.9)	4.0 (0.0)	4.0 (0.0)	6.8 (0.2)	10.0 (0.0)	21.6 (3.9)	55.8 (2.7)	21.2 (1.1)	
MILKB.SUGAR	5	5	8.1 (0.8)	11.8 (1.0)	4.0 (0.0)	4.0 (0.0)	4.9 (0.1)	10.2 (0.2)	8.6 (1.2)	31.1 (0.7)	13.1 (0.8)	
CAKE	5	5	198.3 (47.5)	43.0 (5.3)	4.0 (0.0)	3.8 (0.2)	5.5 (0.2)	10.0 (0.0)	40.0 (5.3)	65.0 (0.9)	22.0 (0.9)	
Completed Runs = 5			Total Time (s.) = 2752.8 (503.2)					Total Calls = 203.2 (11.8)				

(b) OPL

Task	# G	# L	Time (s.)	Calls	States	Edges	Ep. First HRM ($\times 10^2$)	# Examples			Example Length		
								G	D	I	G	D	I
BATTER	5	5	14.1 (3.2)	23.0 (3.0)	6.0 (0.0)	9.2 (0.2)	12.0 (1.0)	11.4 (0.4)	7.0 (1.2)	10.6 (1.6)	20.4 (1.1)	18.7 (1.6)	12.1 (1.7)
BUCKET	5	5	1.8 (0.1)	7.2 (0.6)	4.0 (0.0)	4.0 (0.0)	8.0 (0.5)	10.2 (0.2)	2.2 (0.2)	2.8 (0.4)	10.2 (0.5)	13.4 (1.9)	6.8 (1.7)
COMPASS	5	5	13.5 (1.8)	22.0 (1.7)	6.0 (0.0)	9.2 (0.2)	10.4 (1.4)	11.0 (0.6)	6.8 (1.0)	10.2 (1.0)	17.2 (1.6)	20.9 (1.9)	14.3 (0.8)
LEATHER	5	5	1.8 (0.1)	7.0 (0.5)	4.0 (0.0)	4.0 (0.0)	6.9 (0.5)	10.0 (0.0)	2.4 (0.2)	2.6 (0.4)	11.1 (0.9)	16.9 (5.6)	8.9 (3.3)
PAPER	5	5	2.0 (0.2)	7.6 (0.6)	4.0 (0.0)	4.0 (0.0)	7.7 (1.1)	10.0 (0.0)	3.0 (0.3)	2.6 (0.4)	10.1 (0.9)	18.9 (3.3)	5.6 (0.8)
QUILL	5	5	11.8 (1.3)	22.0 (1.2)	6.0 (0.0)	9.2 (0.2)	12.8 (1.5)	10.6 (0.2)	6.4 (0.7)	11.0 (0.9)	15.3 (1.3)	13.5 (1.0)	12.1 (1.4)
SUGAR	5	5	1.6 (0.1)	6.4 (0.4)	4.0 (0.0)	4.0 (0.0)	6.5 (0.7)	10.0 (0.0)	2.4 (0.2)	2.0 (0.3)	9.6 (0.6)	15.3 (3.6)	16.6 (9.2)
BOOK	5	5	224.8 (71.6)	27.0 (1.9)	6.0 (0.0)	6.4 (0.2)	139.7 (21.8)	11.6 (0.4)	3.2 (0.4)	18.2 (1.4)	22.0 (1.6)	24.7 (6.5)	23.5 (1.2)
MAP	5	5	339.9 (33.6)	33.0 (2.8)	6.0 (0.0)	6.4 (0.2)	204.8 (27.1)	10.6 (0.2)	2.8 (0.5)	25.6 (2.5)	25.4 (0.8)	21.8 (3.1)	25.2 (1.1)
MILKBUCKET	5	5	3.5 (0.3)	8.2 (0.6)	4.0 (0.0)	3.0 (0.0)	47.6 (3.7)	10.2 (0.2)	2.6 (0.4)	3.4 (0.4)	10.3 (0.7)	16.2 (1.7)	14.2 (1.8)
BOOKQUILL	5	5	19.0 (2.2)	15.4 (1.5)	4.0 (0.0)	4.0 (0.0)	383.4 (83.7)	10.0 (0.0)	1.0 (0.3)	11.4 (1.3)	38.2 (1.6)	14.1 (4.9)	23.8 (1.0)
MILKB.SUGAR	5	5	11.4 (2.1)	14.4 (1.7)	4.0 (0.0)	4.0 (0.0)	87.4 (8.9)	10.4 (0.2)	1.0 (0.4)	10.0 (1.3)	19.7 (1.2)	8.7 (4.9)	17.6 (1.2)
CAKE	4	1	277.4 (0.0)	33.0 (0.0)	5.0 (0.0)	4.0 (0.0)	264.1 (0.0)	10.0 (0.0)	2.0 (0.0)	28.0 (0.0)	46.7 (0.0)	36.0 (0.0)	22.9 (0.0)
Completed Runs = 5			Total Time (s.) = 701.0 (111.2)					Total Calls = 199.8 (6.9)					

(c) FRL

Task	# G	# L	Time (s.)	Calls	States	Edges	Ep. First HRM ($\times 10^2$)	# Examples			Example Length		
								G	D	I	G	D	I
BATTER	5	5	11.1 (1.4)	23.4 (2.5)	6.0 (0.0)	9.2 (0.2)	468.4 (121.9)	10.4 (0.2)	7.6 (0.9)	11.4 (1.9)	11.9 (0.6)	10.1 (1.3)	9.9 (0.4)
BUCKET	5	5	2.2 (0.1)	7.0 (0.3)	4.0 (0.0)	4.0 (0.0)	129.5 (69.4)	10.2 (0.2)	2.8 (0.2)	2.0 (0.3)	7.8 (0.5)	9.9 (1.7)	6.4 (2.1)
COMPASS	5	5	12.9 (1.9)	24.6 (2.2)	6.0 (0.0)	9.4 (0.2)	550.8 (156.4)	10.4 (0.2)	7.8 (1.0)	12.4 (1.2)	12.5 (1.6)	9.4 (1.0)	8.4 (0.5)
LEATHER	5	5	2.8 (0.4)	7.8 (0.7)	4.0 (0.0)	4.0 (0.0)	89.0 (18.0)	10.0 (0.0)	3.2 (0.4)	2.6 (0.4)	7.3 (0.4)	9.3 (1.7)	3.7 (0.4)
PAPER	5	5	2.1 (0.1)	7.0 (0.3)	4.0 (0.0)	4.0 (0.0)	82.7 (18.8)	10.0 (0.0)	3.0 (0.0)	2.0 (0.3)	6.9 (0.7)	10.2 (1.8)	4.7 (2.7)
QUILL	5	5	11.6 (1.1)	23.8 (1.5)	6.0 (0.0)	9.6 (0.2)	458.9 (61.0)	10.6 (0.2)	8.0 (0.9)	11.2 (1.2)	11.9 (0.6)	13.1 (2.7)	9.2 (0.8)
SUGAR	5	5	2.6 (0.3)	8.4 (0.7)	4.0 (0.0)	4.0 (0.0)	103.5 (39.5)	10.0 (0.0)	3.6 (0.4)	2.8 (0.5)	8.2 (0.7)	10.1 (1.9)	5.0 (1.1)
BOOK	5	0	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)
MAP	3	0	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)
MILKBUCKET	5	2	4.7 (0.5)	11.0 (1.0)	4.0 (0.0)	3.0 (0.0)	885.6 (142.3)	10.0 (0.0)	2.0 (0.0)	7.0 (1.0)	8.2 (0.4)	10.2 (1.7)	8.8 (0.2)
BOOKQUILL	0	0	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)
MILKB.SUGAR	0	0	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)
CAKE	0	0	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)
Completed Runs = 5			Total Time (s.) = 47.1 (0.9)					Total Calls = 106.4 (2.9)					

Hierarchies of Reward Machines

Table 6: Results of LHRM in WATERWORLD for the default case.

(a) WOD

Task	# G	# L	Time (s.)	Calls	States	Edges	Ep. First HRM	# Examples		Example Length		
								$(\times 10^2)$		G	I	G
RG	5	5	0.9 (0.0)	4.0 (0.0)	3.0 (0.0)	2.0 (0.0)	0.9 (0.1)	10.0 (0.0)	2.0 (0.0)	11.2 (1.0)	5.8 (1.1)	
BC	5	5	0.9 (0.1)	3.8 (0.2)	3.0 (0.0)	2.0 (0.0)	0.8 (0.1)	10.0 (0.0)	1.8 (0.2)	10.8 (0.8)	11.9 (3.4)	
MY	5	5	0.9 (0.0)	3.6 (0.2)	3.0 (0.0)	2.0 (0.0)	0.7 (0.0)	10.0 (0.0)	1.6 (0.2)	8.7 (0.8)	6.6 (1.9)	
RG&BC	5	5	4.5 (0.3)	13.4 (0.4)	4.0 (0.0)	4.0 (0.0)	8.8 (0.3)	11.8 (0.6)	8.6 (0.7)	12.2 (0.9)	14.8 (1.2)	
BC&MY	5	5	5.8 (1.0)	15.6 (2.1)	4.0 (0.0)	4.0 (0.0)	8.1 (0.2)	12.8 (1.3)	9.8 (1.5)	13.2 (1.7)	17.1 (1.6)	
RG&MY	5	5	4.7 (0.5)	13.2 (1.0)	4.0 (0.0)	4.0 (0.0)	8.5 (0.2)	10.8 (0.2)	9.4 (0.9)	12.2 (0.7)	18.6 (1.2)	
RGB	5	5	1.2 (0.1)	4.8 (0.5)	3.0 (0.0)	2.0 (0.0)	8.6 (0.2)	10.0 (0.0)	2.8 (0.5)	7.8 (0.2)	7.0 (1.4)	
CMY	5	5	1.4 (0.2)	5.4 (0.7)	3.0 (0.0)	2.0 (0.0)	8.8 (0.5)	10.0 (0.0)	3.4 (0.7)	8.0 (0.3)	10.2 (1.3)	
RGB&CMY	5	5	15.1 (1.7)	21.6 (1.7)	4.0 (0.0)	4.0 (0.0)	2.3 (0.0)	11.0 (0.4)	17.6 (1.7)	17.3 (0.4)	22.6 (1.6)	
Completed Runs = 5			Total Time (s.) = 35.4 (2.0)					Total Calls = 85.4 (3.1)				

(b) WD

Task	# G	# L	Time (s.)	Calls	States	Edges	Ep. First HRM	# Examples			Example Length		
								$(\times 10^2)$			G	D	I
RG	5	5	1.9 (0.2)	7.8 (1.1)	4.0 (0.0)	4.0 (0.0)	3.0 (0.3)	10.0 (0.0)	3.0 (0.3)	2.8 (0.9)	7.0 (0.7)	10.3 (1.6)	4.2 (0.8)
BC	5	5	1.8 (0.3)	7.2 (1.0)	4.0 (0.0)	4.0 (0.0)	2.7 (0.3)	10.2 (0.2)	2.4 (0.2)	2.6 (0.7)	8.4 (0.5)	6.9 (1.5)	6.3 (1.7)
MY	5	5	1.4 (0.1)	5.8 (0.4)	4.0 (0.0)	4.0 (0.0)	2.9 (0.3)	10.0 (0.0)	2.4 (0.2)	1.4 (0.2)	6.9 (0.4)	5.8 (1.5)	4.6 (1.6)
RG&BC	5	5	11.7 (2.5)	24.0 (3.8)	4.8 (0.2)	4.8 (0.2)	12.0 (0.4)	13.0 (0.7)	6.2 (1.6)	11.8 (1.9)	11.7 (0.8)	6.9 (1.5)	11.8 (0.8)
BC&MY	5	5	9.5 (1.5)	20.8 (2.4)	4.8 (0.2)	4.8 (0.2)	11.5 (0.4)	11.2 (0.6)	5.2 (0.7)	11.4 (1.6)	10.6 (0.9)	8.8 (1.4)	13.2 (0.9)
RG&MY	5	5	5.4 (0.5)	14.0 (1.3)	4.2 (0.2)	4.2 (0.2)	11.8 (0.6)	10.6 (0.2)	3.2 (0.7)	7.2 (1.1)	9.8 (0.3)	6.2 (1.9)	11.6 (0.7)
RGB	5	5	2.5 (0.3)	8.2 (0.7)	4.0 (0.0)	3.0 (0.0)	11.9 (0.6)	10.2 (0.2)	3.0 (0.3)	3.0 (0.5)	7.9 (0.4)	14.0 (1.8)	10.6 (2.0)
CMY	5	5	3.6 (0.4)	11.4 (1.2)	4.0 (0.0)	3.0 (0.0)	10.8 (0.3)	10.0 (0.0)	4.2 (0.7)	5.2 (0.6)	7.9 (0.2)	8.8 (1.6)	10.5 (1.1)
RGB&CMY	5	5	29.0 (4.1)	31.4 (2.8)	4.4 (0.2)	4.4 (0.2)	4.9 (0.3)	11.2 (0.4)	5.4 (1.6)	21.8 (1.3)	16.6 (0.8)	7.4 (0.9)	17.2 (0.6)
Completed Runs = 5			Total Time (s.) = 67.0 (6.2)					Total Calls = 130.6 (6.0)					

Table 7: Results of LHRM in WATERWORLD with a restricted set of callable RMs.

(a) WOD

Task	# G	# L	Time (s.)	Calls	States	Edges	Ep. First HRM	# Examples		Example Length		
								$(\times 10^2)$		G	I	G
RG	5	5	0.9 (0.0)	4.0 (0.0)	3.0 (0.0)	2.0 (0.0)	0.9 (0.1)	10.0 (0.0)	2.0 (0.0)	11.2 (1.0)	5.8 (1.1)	
BC	5	5	0.9 (0.1)	3.8 (0.2)	3.0 (0.0)	2.0 (0.0)	0.8 (0.1)	10.0 (0.0)	1.8 (0.2)	10.8 (0.8)	11.9 (3.4)	
MY	5	5	0.9 (0.0)	3.6 (0.2)	3.0 (0.0)	2.0 (0.0)	0.7 (0.0)	10.0 (0.0)	1.6 (0.2)	8.7 (0.8)	6.6 (1.9)	
RG&BC	5	5	5.3 (0.4)	15.2 (0.9)	4.0 (0.0)	4.0 (0.0)	8.6 (0.3)	12.4 (0.2)	9.8 (0.8)	14.7 (1.3)	16.0 (0.8)	
BC&MY	5	5	3.9 (0.1)	12.4 (0.2)	4.0 (0.0)	4.0 (0.0)	8.3 (0.4)	11.8 (0.7)	7.6 (0.7)	11.2 (0.8)	13.2 (1.0)	
RG&MY	5	5	4.6 (0.3)	13.8 (0.9)	4.0 (0.0)	4.0 (0.0)	8.5 (0.2)	10.2 (0.2)	10.6 (0.9)	10.7 (0.5)	15.8 (1.6)	
RGB	5	5	1.2 (0.1)	4.8 (0.7)	3.0 (0.0)	2.0 (0.0)	8.7 (0.2)	10.2 (0.2)	2.6 (0.7)	8.3 (0.5)	16.2 (3.8)	
CMY	5	5	1.6 (0.2)	6.2 (0.7)	3.0 (0.0)	2.0 (0.0)	8.6 (0.5)	10.0 (0.0)	4.2 (0.7)	8.0 (0.3)	10.8 (1.2)	
RGB&CMY	5	5	5.7 (0.8)	15.0 (1.6)	4.0 (0.0)	4.0 (0.0)	2.6 (0.1)	10.4 (0.4)	11.6 (1.6)	17.0 (1.1)	15.9 (1.3)	
Completed Runs = 5			Total Time (s.) = 24.9 (0.9)					Total Calls = 78.8 (2.7)				

(b) WD

Task	# G	# L	Time (s.)	Calls	States	Edges	Ep. First HRM	# Examples			Example Length		
								$(\times 10^2)$			G	D	I
RG	5	5	1.9 (0.3)	7.8 (1.1)	4.0 (0.0)	4.0 (0.0)	3.0 (0.3)	10.0 (0.0)	3.0 (0.3)	2.8 (0.9)	7.0 (0.7)	10.3 (1.6)	4.2 (0.8)
BC	5	5	1.8 (0.3)	7.2 (1.0)	4.0 (0.0)	4.0 (0.0)	2.7 (0.3)	10.2 (0.2)	2.4 (0.2)	2.6 (0.7)	8.4 (0.5)	6.9 (1.5)	6.3 (1.7)
MY	5	5	1.4 (0.1)	5.8 (0.4)	4.0 (0.0)	4.0 (0.0)	2.9 (0.3)	10.0 (0.0)	2.4 (0.2)	1.4 (0.2)	6.9 (0.4)	5.8 (1.5)	4.6 (1.6)
RG&BC	5	5	6.9 (0.7)	17.6 (1.5)	4.6 (0.2)	4.6 (0.2)	12.0 (0.4)	10.8 (0.2)	4.6 (0.5)	9.2 (1.0)	10.4 (0.5)	9.4 (1.8)	12.6 (0.7)
BC&MY	5	5	9.3 (1.8)	21.4 (2.9)	4.8 (0.2)	4.8 (0.2)	11.7 (0.5)	12.2 (1.0)	5.8 (1.1)	10.4 (1.8)	11.4 (0.7)	6.9 (1.3)	12.1 (0.7)
RG&MY	5	5	7.8 (1.1)	18.8 (1.9)	4.8 (0.2)	4.8 (0.2)	11.8 (0.5)	11.0 (0.3)	4.8 (0.4)	10.0 (2.0)	9.8 (0.2)	8.6 (0.8)	13.0 (0.8)
RGB	5	5	2.1 (0.1)	7.6 (0.2)	4.0 (0.0)	3.0 (0.0)	11.9 (0.5)	10.0 (0.0)	2.6 (0.2)	3.0 (0.0)	7.6 (0.5)	11.8 (1.7)	10.7 (1.7)
CMY	5	5	2.3 (0.2)	8.2 (0.8)	4.0 (0.0)	3.0 (0.0)	10.7 (0.2)	10.0 (0.0)	2.2 (0.4)	4.0 (0.5)	7.8 (0.2)	9.9 (1.6)	8.9 (0.5)
RGB&CMY	5	5	9.6 (1.5)	20.6 (2.6)	5.0 (0.0)	5.0 (0.0)	5.0 (0.6)	10.2 (0.2)	6.4 (0.9)	11.0 (1.8)	15.0 (0.5)	12.3 (0.8)	14.2 (1.2)
Completed Runs = 5			Total Time (s.) = 42.9 (3.7)					Total Calls = 115.0 (7.5)					

Table 8: Results of LHRM in WATERWORLD without exploration using options.

(a) WOD

Task	# G	# L	Time (s.)	Calls	States	Edges	Ep. First HRM ($\times 10^2$)	# Examples		Example Length		
								G	I	G	I	
RG	5	5	0.9 (0.0)	4.0 (0.0)	3.0 (0.0)	2.0 (0.0)	0.9 (0.1)	10.0 (0.0)	2.0 (0.0)	11.2 (1.0)	5.8 (1.1)	
BC	5	5	0.9 (0.1)	3.8 (0.2)	3.0 (0.0)	2.0 (0.0)	0.8 (0.1)	10.0 (0.0)	1.8 (0.2)	10.8 (0.8)	11.9 (3.4)	
MY	5	5	0.9 (0.0)	3.6 (0.2)	3.0 (0.0)	2.0 (0.0)	0.7 (0.0)	10.0 (0.0)	1.6 (0.2)	8.7 (0.8)	6.6 (1.9)	
RG&BC	5	5	4.2 (0.4)	12.2 (0.9)	4.0 (0.0)	4.0 (0.0)	9.5 (0.3)	10.6 (0.2)	8.6 (1.1)	13.8 (0.2)	15.3 (1.6)	
BC&MY	5	5	4.3 (0.3)	11.8 (0.7)	4.0 (0.0)	4.0 (0.0)	9.8 (0.1)	11.6 (0.2)	7.2 (0.6)	15.3 (0.9)	16.7 (1.7)	
RG&MY	5	5	4.6 (0.3)	12.6 (0.7)	4.0 (0.0)	4.0 (0.0)	9.5 (0.1)	11.2 (0.4)	8.4 (0.7)	14.2 (0.9)	14.8 (0.8)	
RGB	5	5	1.2 (0.2)	4.6 (0.7)	3.0 (0.0)	2.0 (0.0)	9.0 (0.1)	10.0 (0.0)	2.6 (0.7)	9.4 (0.4)	8.7 (1.7)	
CMY	5	5	1.4 (0.1)	5.0 (0.5)	3.0 (0.0)	2.0 (0.0)	8.8 (0.2)	10.0 (0.0)	3.0 (0.5)	8.8 (0.2)	10.6 (1.5)	
RGB&CMY	5	5	16.1 (1.1)	19.8 (1.1)	4.0 (0.0)	4.0 (0.0)	4.1 (0.1)	11.2 (0.6)	15.6 (1.3)	26.0 (1.2)	21.6 (0.9)	
Completed Runs = 5			Total Time (s.) = 34.4 (1.4)					Total Calls = 77.4 (2.0)				

(b) WD

Task	# G	# L	Time (s.)	Calls	States	Edges	Ep. First HRM ($\times 10^2$)	# Examples			Example Length		
								G	D	I	G	D	I
RG	5	5	1.9 (0.3)	7.8 (1.1)	4.0 (0.0)	4.0 (0.0)	3.0 (0.3)	10.0 (0.0)	3.0 (0.3)	2.8 (0.9)	7.0 (0.7)	10.3 (1.6)	4.2 (0.8)
BC	5	5	1.8 (0.2)	7.2 (1.0)	4.0 (0.0)	4.0 (0.0)	2.7 (0.3)	10.2 (0.2)	2.4 (0.2)	2.6 (0.7)	8.4 (0.5)	6.9 (1.5)	6.3 (1.7)
MY	5	5	1.4 (0.1)	5.8 (0.4)	4.0 (0.0)	4.0 (0.0)	2.9 (0.3)	10.0 (0.0)	2.4 (0.2)	1.4 (0.2)	6.9 (0.4)	5.8 (1.5)	4.6 (1.6)
RG&BC	5	5	8.1 (1.4)	18.2 (2.2)	4.6 (0.2)	4.6 (0.2)	97.4 (4.2)	10.8 (0.4)	5.2 (0.6)	9.2 (1.4)	10.5 (0.5)	10.3 (1.7)	13.7 (1.5)
BC&MY	5	5	6.2 (0.5)	15.6 (0.7)	4.6 (0.2)	4.6 (0.2)	91.5 (5.8)	10.6 (0.2)	4.6 (0.6)	7.4 (0.7)	9.7 (0.2)	7.0 (1.2)	11.3 (1.0)
RG&MY	5	5	8.6 (1.8)	19.2 (2.7)	4.4 (0.2)	4.4 (0.2)	90.3 (5.3)	11.2 (0.8)	5.6 (0.9)	9.4 (1.3)	10.4 (0.7)	7.7 (0.7)	13.3 (0.9)
RGB	5	5	2.3 (0.1)	7.6 (0.2)	4.0 (0.0)	3.0 (0.0)	65.3 (1.6)	10.2 (0.2)	2.6 (0.4)	2.8 (0.4)	7.6 (0.3)	11.1 (3.0)	8.8 (1.7)
CMY	5	5	4.4 (0.6)	13.2 (1.5)	3.8 (0.2)	3.0 (0.0)	59.2 (2.9)	10.2 (0.2)	3.8 (0.6)	7.2 (1.0)	6.9 (0.5)	8.2 (1.0)	7.6 (0.8)
RGB&CMY	5	5	32.1 (5.0)	31.4 (3.3)	4.4 (0.2)	4.4 (0.2)	125.7 (9.9)	11.4 (0.5)	5.8 (1.2)	21.2 (2.2)	17.1 (0.6)	8.4 (1.6)	17.8 (1.0)
Completed Runs = 5			Total Time (s.) = 66.7 (6.6)					Total Calls = 126.0 (6.3)					

Table 9: Results of learning non-flat and flat HRMs using different methods. The columns are the following: the number of completed runs without timing out, the amount of time needed to learn the HRMs or RMs, and the number of states and edges of the RM.

Task	LHRM (Non-Flat)			LHRM (Flat)			DeepSynth			JRP			LRM				
	C	Time (s.)	States	C	Time (s.)	States	C	Time (s.)	States	C	Time (s.)	States	C	Time (s.)	States	Edges	
MILKBUCKET	5	1.5 (0.2)	3.0 (0.0)	5	3.2 (0.6)	4.0 (0.0)	5	325.6 (29.7)	13.4 (0.4)	93.2 (1.7)	5	17.1 (5.5)	4.0 (0.0)	5	347.5 (64.5)	4.0 (0.0)	14.0 (1.0)
BOOK	5	191.2 (36.4)	5.0 (0.0)	0	-	-	5	288.9 (31.7)	16.6 (3.1)	119.0 (19.4)	0	-	-	5	2261.0 (552.2)	8.0 (0.0)	31.2 (2.0)
BOOKQUILL	5	17.9 (1.4)	4.0 (0.0)	0	-	-	5	308.6 (52.6)	12.8 (0.5)	92.8 (2.3)	0	-	-	0	-	-	-
CAKE	5	74.5 (25.7)	4.0 (0.0)	0	-	-	4	290.6 (36.4)	17.2 (2.5)	110.2 (11.6)	0	-	-	0	-	-	-
RG	5	0.9 (0.0)	3.0 (0.0)	5	0.9 (0.0)	3.0 (0.0)	0	-	-	-	5	32.3 (7.9)	3.8 (0.2)	0	-	-	-
RG&BC	5	4.5 (0.3)	4.0 (0.0)	0	-	-	0	-	-	-	0	-	-	0	-	-	-
RGB&CMY	5	15.1 (1.7)	4.0 (0.0)	0	-	-	0	-	-	-	0	-	-	0	-	-	-

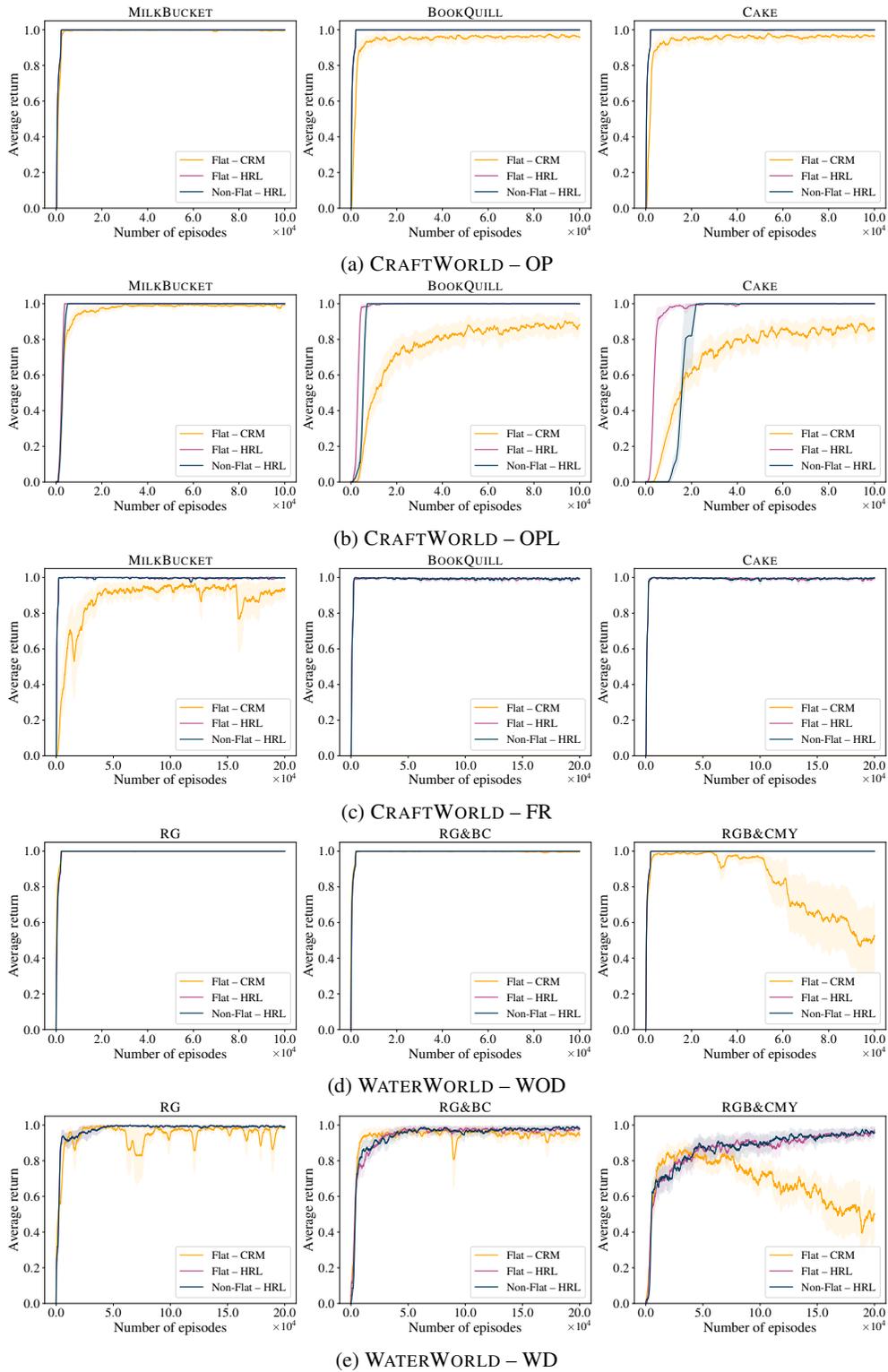


Figure 15: Learning curves comparing the performance of policy learning algorithms exploiting handcrafted non-flat and flat HRMs.