

Machine Learning

Transferable Dynamics Models for Efficient Object-Oriented Reinforcement Learning --Manuscript Draft--

Manuscript Number:	MACH-D-20-00735
Full Title:	Transferable Dynamics Models for Efficient Object-Oriented Reinforcement Learning
Article Type:	Manuscript
Keywords:	reinforcement learning; transfer learning; object-oriented representations; efficient learning
Abstract:	<p>The Reinforcement Learning (RL) framework offers a general paradigm for constructing autonomous agents that can make effective decisions when solving tasks. While in principle RL can be utilised to solve a multitude of important tasks in the field of Artificial Intelligence (AI), in practice the framework struggles to scale to large real-world applications. One of the major reasons for this is that standard RL representations are not suitable for transfer, and so each task typically needs to be solved independently. Previous research has shown that object-oriented representations are, in many cases, more suitable for the purposes of transfer. Such representations leverage the notion of object classes to learn lifted rules that apply to grounded object instantiations. In this paper, we extend previous research on object-oriented representations and introduce two formalisms: the first is based on deictic predicates, and is used to learn a transferable transition dynamics model; the second is based on propositions, and is used to learn a transferable reward dynamics model. In addition, we extend previously introduced efficient learning algorithms for object-oriented representations to our proposed formalisms. Our frameworks can be combined into a single efficient algorithm that learns transferable transition and reward dynamics models across a domain of related tasks. We illustrate our proposed algorithm empirically on an extended version of the Taxi domain as well as the difficult Sokoban domain showing the benefits of our approach with regards to efficient learning and transfer.</p>

Noname manuscript No.
(will be inserted by the editor)

Transferable Dynamics Models for Efficient Object-Oriented Reinforcement Learning

Ofir Marom · Benjamin Rosman

Received: date / Accepted: date

Abstract The Reinforcement Learning (RL) framework offers a general paradigm for constructing autonomous agents that can make effective decisions when solving tasks. While in principle RL can be utilised to solve a multitude of important tasks in the field of Artificial Intelligence (AI), in practice the framework struggles to scale to large real-world applications. One of the major reasons for this is that standard RL representations are not suitable for transfer, and so each task typically needs to be solved independently. Previous research has shown that object-oriented representations are, in many cases, more suitable for the purposes of transfer. Such representations leverage the notion of object classes to learn **lifted rules** that apply to grounded object instantiations. In this paper, we extend previous research on object-oriented representations and introduce two formalisms: the first is based on deictic predicates, and is used to learn a transferable transition dynamics model; the second is based on propositions, and is used to learn a transferable reward dynamics model. In addition, we extend previously introduced efficient learning algorithms for object-oriented representations to our proposed formalisms. Our frameworks can be combined into a single efficient algorithm that learns transferable transition and reward dynamics models across a domain of related tasks. We illustrate our proposed algorithm empirically on an extended version of the Taxi domain as well as the difficult Sokoban domain showing the benefits of our approach with regards to efficient learning and transfer.

Keywords reinforcement learning · transfer learning · object-oriented representations · efficient learning

1 Introduction

A longstanding goal in the field of Artificial Intelligence (AI) is the construction of autonomous agents that make effective decisions in the real-world to complete tasks. The Reinforcement Learning (RL) framework [26] offers a general paradigm for constructing

Ofir Marom
University of the Witwatersrand, Johannesburg, South Africa, E-mail: ofiremarom@gmail.com

Benjamin Rosman
University of the Witwatersrand, Johannesburg, South Africa, E-mail: benjamin.Rosman1@wits.ac.za

autonomous agents with these capabilities. In the RL setting, an agent interacts with its environment with the aim of learning a policy that instructs the agent on which action to take in any given state to achieve an optimal payoff. For reasons of computational tractability, an RL task is typically modelled as a Markov Decision Process (MDP). Given an MDP description of an RL task, a set of fairly simple algorithms exist with convergence guarantees to optimal policies. **RL has had many successes, achieving master-level performance on a variety of complex tasks such as Backgammon, a suite of Atari 2600 games and chess** [28, 16, 21, 22].

While such progress is encouraging, RL has struggled to scale to real-world tasks. A major factor of this limitation is that standard RL algorithms are highly inefficient because these algorithms start with no prior knowledge and so have to learn everything from scratch when solving a new task. Such a blank-slate approach is at odds with human cognition that transfers vast amounts of prior knowledge gained from previous experiences when solving new tasks [18]. Therefore, an intuitive approach to scaling up autonomous agents under the RL framework is to transfer prior knowledge between tasks in the same way that humans do.

While such an approach has intuitive appeal, transfer learning in RL is a challenging research area with many open questions [27]. In particular, it is not always clear what, when or how to transfer in order to achieve improved learning performance. Furthermore, the research area of transfer learning has tended to rely on empirical methods that perform well in limited settings and that lack theoretical foundations.

One concept that has shown much promise for effective transfer is based on relational representations in RL [6, 17]. These representations rely on expressing the components of an MDP in terms of logical statements over variables. Transfer is then achieved by expressing lifted rules about groups of variables, and then grounding these lifted rules for specific instantiations of variable values.

Object-oriented representations are a form of relational representation that draws inspiration from the structure of the physical world [8, 7, 11, 5, 4, 19]. Object-oriented representations define a set of lifted object classes from which grounded objects are then instantiated for a particular task. Such representations aim to learn a set of lifted rules that apply to the object classes. These lifted rules then generalise to grounded instantiations of objects, and are therefore transferable between related tasks of a domain.

In particular, previous research has introduced the Propositional Object-Oriented MDP (Propositional OO-MDP) framework [5, 4]. Propositional OO-MDPs operate by formulating the transition dynamics of an MDP in terms of propositional preconditions over lifted object classes that map to effects over object class attributes. Since the preconditions and effects that represent the transition dynamics are fully lifted, they transfer across all tasks of a domain. Propositional OO-MDPs have limited expressive power because they require that all preconditions be propositional. However, they have the advantage that probably efficient learning algorithms exist under the formalism. In particular, the research that introduces Propositional OO-MDPs also introduces a provably efficient learning algorithm called *DOORMAX* that operates under the Propositional OO-MDP formalism in domains with deterministic transition dynamics.

Unfortunately, the core restriction of Propositional OO-MDPs that the transition dynamics be described only in terms of propositional preconditions is a strong one, and precludes efficient learning for a large class of domains. Such domains include those where it is required to distinguish between different objects of the same object class for tasks of the domain. As a specific example of this, and further elaborated in subsection 3.2, consider the Sokoban domain where a person attempts to push a box but cannot do so if that box is

adjacent to a wall. In this case, there is no way of tying the box that is adjacent to the person with the box that is adjacent to the wall with propositions.

To accommodate this, our previous research introduced the Deictic Object-Oriented MDP (Deictic OO-MDP) framework [15]. Deictic OO-MDPs use deictic predicates for their preconditions. A deictic predicate is a predicate that is grounded only with respect to a central deictic object, therefore that **object may relate itself to non-grounded object classes**, but not to other grounded objects. Returning to the Sokoban domain, a deictic predicate over boxes allows a specific box to ascertain whether *any* wall is adjacent to it, but not whether a *specific* wall is adjacent to it. Deictic OO-MDPs are more general than Propositional OO-MDPs, and so are able to represent a larger class of domains. Furthermore, the research that introduces Deictic OO-MDPs extends *DOORMAX* with an algorithm called *DOORMAX_D* so that efficient learning is possible under the Deictic OO-MDP framework as well [15].

Both Propositional OO-MDPs and Deictic OO-MDPs assume known reward dynamics and focus on learning transition dynamics. In the most general model-based RL setting, the agent must learn both of these components. In this paper we combine Deictic OO-MDPs and Propositional OO-MDPs by leveraging the former to efficiently learn transferable transition dynamics and the latter to efficiently learn transferable reward dynamics. We then extend *DOORMAX_D* to accommodate learning of both of these components.

We illustrate our framework empirically on an extended version of the Taxi domain, called the All-Passenger Any-Destination Taxi domain, as well as the Sokoban domain. In both these domains we show that it is possible to learn transferable models of both the transition and reward dynamics from some simple source tasks of the domain and then transfer them to more complex target tasks of the domain. In fact, we show that zero-shot transfer is possible under our proposed frameworks thus greatly improving sample efficiency.

Throughout this paper we will use the Taxi and Sokoban domains to illustrate key points related to our frameworks. This paper is organised as follows: in section 2 we discuss background information for efficient learning and object-oriented representations in RL; in section 3 we introduce the Deictic OO-MDP formalism and associated efficient learning algorithms; in section 4 we extend previous work on Propositional OO-MDPs to accommodate a formalism and associated efficient learning algorithms for reward dynamics; in section 5 we introduce the complete *DOORMAX_D* algorithm for learning transferable models of the transition and reward dynamics under the formalisms of sections 3 and 4 respectively, while also providing formal proofs of efficiency; in section 7 we run experiments on the All-Passenger Any-Destination Taxi and Sokoban domains illustrating the benefits of our proposed framework; and in section 8 we conclude with final remarks.

2 Background

2.1 Markov Decision Process

For reasons of computational tractability, an RL task is typically formulated as a Markov Decision Process (MDP) [26]. An MDP is described by a tuple $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma, \rho)$, where:

- \mathcal{S} is a finite set of states called the state-space;
- \mathcal{A} is a finite set of all actions called the action-space;
- \mathcal{P} is the transition dynamics for which $\mathcal{P}(s'|s, a)$ returns the transition probability for transitioning to state $s' \in \mathcal{S}$ conditional on the agent being in state s and taking action $a \in \mathcal{A}$;

- \mathcal{R} is the reward dynamics for which $\mathcal{R}(s, a, s')$ returns the (stationary, finite variance) distribution that governs the reward signal that the agent receives when in state s , takes action a and transitions to state s' ;
- $\gamma \in [0, 1)$ is a discount rate that controls the trade-off between the importance of immediate and future rewards; and
- ρ is a distribution over start states for which $\rho(s_0)$ returns the probability of a task starting in state $s_0 \in \mathcal{S}$.

In this paper, we further restrict attention to episodic tasks, so that the MDP has a set of terminal states $\mathcal{S}_{terminal}$ and a task terminates when one of these states is reached.

In RL the agent has no control over the environment dynamics, \mathcal{P} and \mathcal{R} . What the agent does control is the policy it follows in the environment. Formally, a policy is defined by $\pi(a|s)$ and represents the probability of taking action a conditional on the agent being in state s . The goal of an RL agent is to find an optimal policy π_* that maximises the expected future discounted rewards given any state $s \in \mathcal{S}$. Broadly speaking, algorithms that learn an optimal policy for an MDP can be categorised into model-free and model-based methods. Model-free methods go straight from experience to a policy, and so do not learn models of \mathcal{P} and \mathcal{R} ; model-based methods learn models of \mathcal{P} and \mathcal{R} in conjunction with learning an optimal policy. Given an MDP formulation of an RL task, a host of model-free and model-based algorithms exist with convergence guarantees to an optimal policy.

2.2 Efficient Learning

Unfortunately, algorithms that only have convergence guarantees to an optimal policy are not sufficient. In practice, such algorithms may run for an infeasible amount of time before convergence. Therefore, in addition to convergence guarantees, we also care about the performance of RL algorithms. In particular, we want algorithms to be *sample efficient*. Sample efficient algorithms learn an optimal policy in as few interactions as possible.

One well-known sample efficient RL algorithm is R_{max} [2]. The R_{max} algorithm is a model-based algorithm that operates under the *optimism in the face of uncertainty principle*. Under this principle, the agent initially assumes optimistic models of the dynamics, such that taking any action in any state leads to a terminal state with optimal payoff. The agent then applies an exact planning algorithm to obtain an optimal policy under these incorrect models. This optimal policy drives exploration to the parts of the environment where the agent expects to receive these optimal payoffs. From there, the agent learns the underlying truth and updates its models with this knowledge. By applying this methodology, the agent is able to construct iteratively more accurate models of the dynamics every time it interacts with its environment. As a result of this exploration strategy, it can be shown that R_{max} has polynomial sample complexity with respect to $|\mathcal{S}|$ and $|\mathcal{A}|$ under the PAC (probably approximately correct) terms [10] – i.e. the policy produced by R_{max} is near-optimal with high probability.

The R_{max} algorithm was designed for the case of finite MDPs. Unfortunately, even with the polynomial sample complexity guarantees of R_{max} , learning a near-optimal policy may be slow. This is because R_{max} has to run an exact planning algorithm at every iteration and, when the state-space or action-space is large, such planning algorithms tend to have high computational complexity.

An alternative way to improve the sample efficiency of RL algorithms is with compact representations. Such representations do not learn the dynamics models in a tabular way, but

rather use a representation that compresses the size of the state-space or action-space so that the models can be learned more efficiently.

The KWIK (knows what it knows) framework extends R_{max} with the KWIK- R_{max} algorithm [14, 13]. The KWIK- R_{max} algorithm generalises R_{max} so that it works with a broader range of representations on the condition that they respect the KWIK protocol. Under the KWIK protocol, when the agent is in some state s and takes some action a , then rather than making a next-state prediction s' the agent may instead return \perp , indicating that they are not yet able to make an accurate prediction. The KWIK protocol restricts the number of times that the agent may return \perp , called the KWIK-bound, to be a polynomial bound. Under the restriction that the KWIK-bound is a polynomial bound, KWIK- R_{max} can be shown to have polynomial sample complexity under the PAC terms as well. Examples of compact representations that have been used in conjunction with KWIK- R_{max} include dynamic Bayesian networks [23, 24], linear models [25] and relations [5, 4].

2.3 Object-Oriented MDPs

Object-oriented MDPs (OO-MDPs) form part of the broader field of relational RL, and were first introduced under the Relational MDP framework for planning domains [8, 7]. OO-MDPs view the state-space of an MDP in terms of objects, where each object belongs to some object class (or simply class) that has a set of attributes. The transition and reward dynamics are then expressed through the relational structure between object class attributes. The main advantage of object-oriented representations is that they allow for a compact description of an MDP that is fully lifted. Therefore, such an MDP can be used to describe an entire domain, rather than a specific task. Such a lifted MDP is called a *schema*, from which grounded MDPs can then be instantiated.

As a result, OO-MDPs are favourable for the purposes of transfer learning. In the transfer learning research problem, the agent aims to learn some knowledge from a set of source tasks that can then be used to accelerate solving a related target task [27]. Since with OO-MDPs knowledge is represented at the class-level, it is possible to learn some knowledge on source tasks of the schema and then transfer that knowledge to target tasks of the schema.

Formally, the state-space for such a schema consists of a set of object classes $\mathfrak{C} = \{C_i\}_{i=1}^{N_C}$. Each object class $C \in \mathfrak{C}$ has a set of attributes $Att(C) = \{C.\alpha_i\}_{i=1}^{N_C}$ and each attribute $C.\alpha \in Att(C)$ of an object class has a domain of possible values $Dom(C.\alpha)$. Given a schema, a grounded state-space is instantiated by first selecting a grounded object set which consists of n objects $O = \{o_i\}_{i=1}^n$ where each $o \in O$ is an instance of some object class C . The value of attribute $C.\alpha$ for object o is denoted by $o.\alpha$. Then the grounded state-space, denoted \mathcal{S}_O , is an assignment of each $o.\alpha$ for all objects in O . The schema state-space, denoted \mathcal{S} , is the set of all states for all possible object sets O .

To make the notion of a schema state-space concrete, consider the classical Taxi domain [3] as show in figure 1. The original Taxi task takes place in a 5×5 gridworld where a taxi has to pick up a passenger that is at one of four pickup locations and drop them off at one of four destination locations. The possible pickup and destination locations as shown in figure 1 and are fixed upfront for the task.

The set of actions available to the agent are *North*, *East*, *South* and *West*, that control the taxi's navigation in the gridworld, as well as *Pickup* and *Dropoff*. The action *Pickup* picks up the passenger if the taxi is at the pickup location and the passenger is not already in the taxi, while the action *Dropoff* drops off the passenger if the taxi has already picked up the passenger and the taxi is on the destination location. Walls in the gridworld limit the taxi's

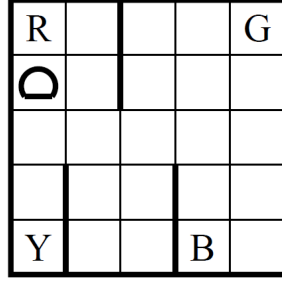


Fig. 1: The original Taxi task. Letters mark possible pickup and destination locations; \odot marks the taxi; thicker lines mark walls.

movements. The reward dynamics for the Taxi task is -1 for each navigation step, -10 for applying the *Pickup* or *Dropoff* actions incorrectly and 0 for reaching a terminal state that drops off the passenger at the correct destination location.

Under an object-oriented representation, we would define the schema state-space with four object classes: *Taxi*, *Wall*, *Passenger* and *Destination*. Each object class has attributes x and y for their location in the grid. Object class *Wall* has an additional attribute *pos* to mark one of four positions in a square, while *Passenger* has an additional Boolean attribute *in-taxi* to indicate if the passenger is in the taxi. Given the schema state-space for this domain, we can instantiate a set of grounded objects from the object classes and a resulting grounded state.

2.4 Propositional OO-MDPs

OO-MDPs were first introduced under the Relational MDP framework to learn transferable class-level value functions for planning domains [8, 7]. Thereafter, they were used conjunction with KWIK- R_{max} for model-based RL under the Propositional OO-MDP framework [5, 4]. The key insight that underlies Propositional OO-MDPs is that, for certain domains, it is possible to compactly represent the transition dynamics with conjunctions of propositional statements over object classes that map to effects over object class attributes as illustrated in table 1 for the Taxi domain.

The statement

$$Touch_E(Taxi, Wall) = 0 \implies Taxi.x \leftarrow Taxi.x + 1$$

for action *East* can be read as: if an object of class *Taxi* has an object of class *Wall* one square east of it is false, then all objects of class *Taxi* have their x attributes increased by 1. The same logic applies to the $Touch_N$, $Touch_S$ and $Touch_W$ expressions for the *North*, *South* and *West* actions respectively. The statement

$$On(Taxi, Passenger) = 1 \implies Passenger.in-taxi \leftarrow 1$$

for can be read as as: if an object of class *Taxi* is on the same square as an object of class *Passenger* is true, then every object of class *Passenger* has its *in-taxi* attribute set to 1.

Since the preconditions in table 1 only refer to object class attributes, the transition dynamics described in this way can be transferred to different Taxi tasks. For example, the

Action	Precondition	Effect
<i>North</i>	$Touch_N(Taxi, Wall) = 0$	$Taxi.y \leftarrow Taxi.y + 1$
<i>East</i>	$Touch_E(Taxi, Wall) = 0$	$Taxi.x \leftarrow Taxi.x + 1$
<i>South</i>	$Touch_S(Taxi, Wall) = 0$	$Taxi.y \leftarrow Taxi.y - 1$
<i>West</i>	$Touch_W(Taxi, Wall) = 0$	$Taxi.x \leftarrow Taxi.x - 1$
<i>Pickup</i>	$On(Taxi, Passenger) = 1$	$Passenger.in-taxi \leftarrow 1$
<i>Dropoff</i>	$Passenger.in-taxi = 1 \wedge On(Taxi, Destination) = 1$	$Passenger.in-taxi \leftarrow 0$

Table 1: Transition dynamics of the Taxi domain under the Propositional OO-MDP formalism. Under the formalism, any other assignment of truth values to the propositions leads to a global *failure condition* that leaves all object class attributes unchanged.

dynamics described in his way transfer to a Taxi task with different gridworld dimensions, pickup, dropoff and wall locations.

An alternative view of transition dynamics under the Propositional OO-MDP formalism is through binary trees. For each action a and attribute C_α , we can represent the transition dynamics as a full binary tree with propositions at the non-leaf nodes and effects at the leaf nodes. This is shown in figure 2a for the action *East* and attribute *Taxi.x* and in figure 2b for the action *Dropoff* and attribute *Passenger.in-taxi*.

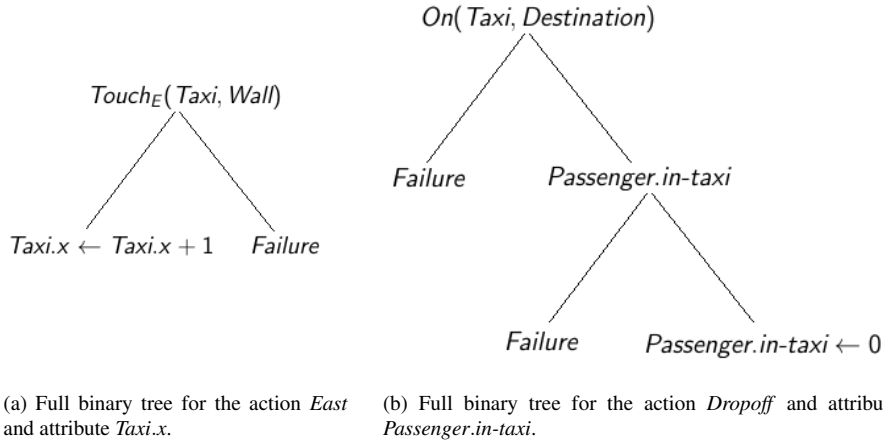


Fig. 2: Examples of full binary tree structures under the Propositional OO-MDP formalism. The leaf node ‘Failure’ refers to a failure condition. Right branches represent a truth value of 1.

The Propositional OO-MDP framework also introduces a KWIK- R_{max} based algorithm called *DOORMAX* (Deterministic Object-Oriented R_{max}) [5, 4]. The *DOORMAX* algorithm operations in domains with deterministic transition dynamics that are representable under the formalism. The **algorithm learns the binary trees** that represent the transition dynamics as shown in figure 2. Under *DOORMAX*, for each action a and attribute $C.\alpha$, a set of propo-

sitions and an effect type ¹ are hypothesised. The algorithm then learns which conjunctions over the hypothesised propositions map to which effect of the effect type. The main assumption required for *DOORMAX* to be correct is that each effect can occur at most at one non-leaf node in the tree. All leaf nodes that do not produce a unique effect map to a global *failure condition* that leaves the state unchanged.

The Propositional OO-MDP framework extends *DOORMAX* to allow for the learner to hypothesise $N > 0$ effect types under the assumption that exactly one effect type is true and, furthermore, that the maximum number of effects per effect type is bounded by some $K > 0$. Under these additional assumptions one can prove that *DOORMAX* has a KIWK-bound $K(D + 1) + 1$ to learn a single effect type for an action a and attribute C_α , while learning given N effect types has a KWIK-bound of $N(K(D + 1) + 1) + (N - 1)$ [5]. These bounds exclude the learning of failure conditions, which are learned through memorisation [14, 13]. **The *DOORMAX* algorithm assumes known reward dynamics and focuses on learning transition dynamics.**

3 Transferable Transition Dynamics with Deictic OO-MDPs

The Propositional OO-MDP formalism has the benefits of efficient learning, however, it lacks expressive power. In particular, Propositional OO-MDPs cannot compactly represent the transition dynamics of domains where it is required to distinguish between different objects of the same object class. To illustrate this, consider a simple extension to the Taxi domain which we call the All-Passenger Any-Destination Taxi domain. This domain is similar to the original Taxi domain, except that the taxi is tasked to pick up multiple passengers and drop each of them off at one of any destination locations. The taxi can only pick up one passenger at a time, so if a passenger is already in the taxi and the *Pickup* action is taken while the taxi is at the pickup location of another passenger, the state does not change. The state-space of the schema for this domain under an object-oriented representation is identical to that of the original Taxi domain described in section 2.3 except that we add an additional Boolean attribute to the *Passenger* object class *at-destination*, to indicate if a passenger has already been dropped off at a destination. Figure 3 shows sample states of the schema.

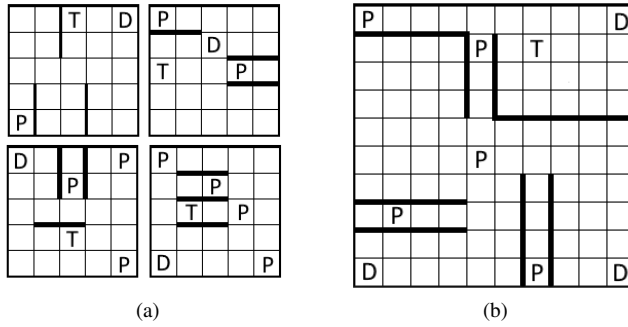


Fig. 3: ‘P’ marks a passenger; ‘D’ marks a destination; ‘T’ marks a taxi; thicker lines mark walls. Five possible states from the All-Passenger Any-Destination Taxi domain schema.

¹ An example of an effect type is the relative effect type: $Rel_i(C.\alpha) = C.\alpha + i$.

Propositions over object classes are insufficient to compactly represent the transition dynamics for this version of the Taxi domain. To see why, suppose we have a task with two passenger objects and the proposition $On(Taxi, Passenger)$ with a truth value 1. This can be read as: an object of class *Taxi* is on the same square as an object of class *Passenger* is true. Clearly, this information is insufficient to determine which passenger object's *in-taxi* attribute should change given the *Pickup* action. To overcome this ambiguity under a propositional approach, we must resort to propositions over the grounded passenger objects, $passenger_1$ and $passenger_2$, of the form $On(Taxi, passenger_1)$ and $On(Taxi, passenger_2)$. Note that while this resolves the ambiguity, the number of propositions needed for the precondition now changes as we change the number of *passenger* objects in the task. This complicates both learning and transfer procedures.

As a further example, consider the Sokoban domain where a warehouse keeper is tasked to push boxes to storage locations in a warehouse, but cannot do so if a box is against a wall or in front of another box.² Suppose we have the propositions $Touch_W(Box, Person)$ and $Touch_E(Box, Wall)$ both with truth value 1. The first proposition can be read as: an object of class *Box* has an object of class *Person* one square west of it is true, while the second can be read as: an object of class *Box* has an object of class *Wall* one square east of it is true. Then the conjunction $Touch_W(Box, Person) = 1 \wedge Touch_E(Box, Wall) = 1$ is insufficient to determine the transition dynamics of any *box* object's *x* attribute when taking action *East* since there is no way to know if the statements are referring to the same box. See figure 4 for an illustration. A similar example of ambiguity that is not resolvable under a propositional approach is also described by the original authors of the Propositional OO-MDP framework [4].

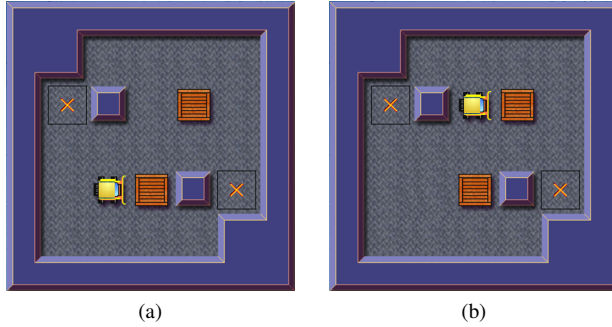


Fig. 4: For action *East* and the box adjacent to the warehouse keeper, in figure (4a) the effect is $box.x \leftarrow box.x + 0$; in figure (4b) the effect is $box.x \leftarrow box.x + 1$ while the conjunction $Touch_W(Box, Person) = 1 \wedge Touch_E(Box, Wall) = 1$ is true in both cases.

An example of a more expressive formalism that can be used to resolve this type of ambiguity is the First-Order MDP (FO MDP) [1] that can use first order predicates over grounded objects. Unfortunately, such an expressive formalism complicates efficient learning and transfer procedures.

In this section we introduce the **Deictic Object-Oriented MDP** (Deictic OO-MDP) framework. Deictic OO-MDPs are more expressive than Propositional OO-MDPs, thus allow-

² All Sokoban images in this paper are taken from JSoko: <https://www.sokoban-online.de/jsoko/credits>.

ing for a broader range of domains to be represented under this formalism. While being less expressive than FO MDPs, Deictic OO-MDPs still have the probably efficient learning guarantees that underlies Propositional OO-MDPs and which are not held by FO MDPs.

The core idea behind the Deictic OO-MDP formalism is the notion of a deictic predicate. A deictic predicate is grounded only with respect to a single grounded reference object that must relate itself to non-grounded object classes. Therefore, deictic predicates are more expressive than propositions that may only depend on object classes, while being less expressive than first-order predicates that may depend on an arbitrary number of grounded objects.

3.1 Formalism

The Deictic OO-MDP framework uses the same schema state-space \mathcal{S} as described in section 2.3 while deictic predicate preconditions are used to define the schema transition dynamics as described below. Let \mathcal{A} be a set of actions. Then for each action $a \in \mathcal{A}$ and attribute $C.\alpha$ of a class C define a set of effects of size $K_{a,C,\alpha}$:

$$\mathcal{E}_{a,C,\alpha} = \{e_i : \text{Dom}(C.\alpha) \rightarrow \text{Dom}(C.\alpha)\}_{i=1}^{K_{a,C,\alpha}},$$

Define a set of deictic predicates of size $D_{a,C,\alpha}$:

$$\mathcal{F}_{a,C,\alpha} = \{f_i : O[C] \times \mathcal{S} \rightarrow \mathfrak{B}\}_{i=1}^{D_{a,C,\alpha}},$$

where $O[C]$ is a set that contains objects with all possible attribute value assignments that are instances of C , and $\mathfrak{B} = \{0, 1\}$.

Then the probabilistic transition dynamics for a and $C.\alpha$ are defined by

$$\mathcal{P}_{a,C,\alpha} : \mathfrak{B}^{D_{a,C,\alpha}} \times \mathcal{E}_{a,C,\alpha} \rightarrow [0, 1].$$

The schema transition dynamics \mathcal{P} is the set of transition dynamics for all actions and attributes,

$$\mathcal{P} = \{\mathcal{P}_{a,C,\alpha} | a \in \mathcal{A}, C \in \mathfrak{C}, C.\alpha \in \text{Att}(C)\}.$$

The schema reward dynamics are defined by

$$\mathcal{R} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}.$$

Given an object set O we can instantiate a grounded MDP $\mathcal{M}_{O,p} = (\mathcal{S}_O, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma, \rho)$. Then if the agent is currently in state $s \in \mathcal{S}_O$ and takes action a , the transition dynamics for $\mathcal{M}_{O,p}$ operate as follows: for each object o in s that is an instance of C and each attribute $C.\alpha$ we compute the Boolean truth values $\mathcal{B} = \{f_i(o, s)\}_{i=1}^{D_{a,C,\alpha}}$ for the deictic predicates in $\mathcal{F}_{a,C,\alpha}$. Then for an effect $e \in \mathcal{E}_{a,C,\alpha}$ we compute $\mathcal{P}_{a,C,\alpha}(\mathcal{B}, e)$ which returns the probability of e occurring given \mathcal{B} . This implies a distribution over effects which in turn implies a distribution over the attribute values of o by applying the effect to $o.\alpha$ in s and obtaining $e(o.\alpha) = o.\alpha'$ in s' .

As an example, consider attribute *Taxi.x* and action *East* for the All-Passenger Any-Destination Taxi domain. We can define a set of relative effects $\text{Rel}_i(x) = x + i$ that produce a shift of i squares from the current location x , as well as a deictic predicate $\text{Touch}_E(\text{taxi}, s)$

that returns 1 if the *taxi* object has an object of class *Wall* one square east of it in *s*, otherwise 0. Then the transition dynamics can be described for any *taxi* object as

$$Touch_E(taxi, s) = 1 \implies taxi.x \leftarrow Rel_0(taxi.x)$$

with probability one and

$$Touch_E(taxi, s) = 0 \implies taxi.x \leftarrow Rel_1(taxi.x)$$

with probability one.

The key insight with Deictic OO-MDPs is that the parameters we pass to each deictic predicate in $\mathcal{P}_{a,C,\alpha}$ are a grounded deictic object *o* that must be an instance of *C* and *s* which is a state of the schema, not a grounded state. As a result, these deictic predicates may not refer to specific objects in *s*; however, they may relate *o* to object classes of the schema. For example, with $Touch_E(taxi, s)$ as defined above only *taxi* is grounded while we never refer to a grounded object in *s*.

Similarly to the Propositional OO-MDP formalism, a requirement for the definition of $\mathcal{P}_{a,C,\alpha}$ to be correct is that the set $\mathcal{E}_{a,C,\alpha}$ must be *invertible* so that if the current value assignment for some attribute *C.α* of a grounded object *o* in state *s* is *o.α* and we take action *a* to subsequently observe *o.α'* in *s'* then there must be a unique $e \in \mathcal{E}_{a,C,\alpha}$ such that $e(o.\alpha) = o.\alpha'$. Clearly, if this requirement is not met then either the effects are not able to correctly capture the true transition dynamics or there are duplicate effects that lead to the same *o.α'* which creates ambiguity over which effect occurred.

Note that the Deictic OO-MDP formalism is more general than the Propositional OO-MDP formalism since a deictic predicate that does not make reference to its deictic object is simply a proposition. Consequently, any domain that is representable under the Propositional OO-MDP formalism is also representable under the Deictic OO-MDP formalism.

3.2 Resolving the Limitations of Propositional OO-MDPs

Deictic OO-MDPs are more expressive than Propositional OO-MDPs allowing us to compactly represent the transition dynamics for the All-Passenger Any-Destination Taxi and Sokoban domains. In particular, define the following effect sets:

- $Rel_i(x) = x + i$ for $i \in \{-1, 0, 1\}$ where *x* is an integer.
- $SetBool_i(x) = x\mathbb{1}(i = 0) + (1 - x)\mathbb{1}(i = 1)$ for $i \in \{0, 1\}$ where $x \in \{0, 1\}$, where $\mathbb{1}$ is the indicator function.

For the All-Passenger Any-Destination Taxi domain define the following deictic predicates³:

- $AnyWallNorthOfTaxi(taxi, s)$: returns 1 if there is any object of class *Wall* one square north of *taxi* in *s*, and 0 otherwise. (f_1)
- $AnyWallEastOfTaxi(taxi, s)$: returns 1 if there is any object of class *Wall* one square east of *taxi* in *s*, and 0 otherwise. (f_2)

³ The statements $AnyPassengerInAnyTaxi(passenger, s)$ and $AnyTaxiOnAnyDestination(destination, s)$ are actually a proposition as they do not refer to the grounded *passenger* object. We could simplify notation by writing, for example, $AnyPassengerInAnyTaxi(s)$. However, we include the *passenger* object to remain consistent with the notation described in section 3.1.

- *AnyWallSouthOfTaxi*(*taxi*, *s*): returns 1 if there is any object of class *Wall* one square south of *taxi* in *s*, and 0 otherwise. (f_3)
- *AnyWallWestOfTaxi*(*taxi*, *s*): returns 1 if there is object of class *Wall* one square west of *taxi* in *s*, and 0 otherwise. (f_4)
- *AnyTaxiOnPassenger*(*passenger*, *s*): returns 1 if there is any object of class *Taxi* on the same square as *passenger* in *s*, and 0 otherwise. (f_5)
- *PassengerAtAnyDestination*(*passenger*, *s*): returns 1 if *passenger.at-destination* is set to 1 in *s*, and 0 otherwise. (f_6)
- *AnyPassengerInAnyTaxi*(*passenger*, *s*): returns 1 if any object of class *Passenger* has its *in-taxi* attribute is set to 1 in *s*, and 0 otherwise. (f_7)
- *PassengerInAnyTaxi*(*passenger*, *s*): returns 1 if *passenger.in-taxi* is set to 1 in *s*, and 0 otherwise. (f_8)
- *AnyTaxiOnAnyDestination*(*destination*, *s*): returns 1 if any object of class *Taxi* is on the same square as any object of class *Destination* in *s*, and 0 otherwise. (f_9)

Then table 2 shows for each attribute and action the relevant preconditions and effects that describe the transition dynamics of the domain under the Deictic OO-MDP formalism.

Action	Attribute	Precondition	Effect
<i>North</i>	<i>Taxi.y</i>	$f_1 = 0$	Rel_1
<i>North</i>	<i>Taxi.y</i>	$f_1 = 1$	Rel_0
<i>East</i>	<i>Taxi.x</i>	$f_2 = 0$	Rel_1
<i>East</i>	<i>Taxi.x</i>	$f_2 = 1$	Rel_0
<i>South</i>	<i>Taxi.y</i>	$f_3 = 0$	Rel_{-1}
<i>South</i>	<i>Taxi.y</i>	$f_3 = 1$	Rel_0
<i>West</i>	<i>Taxi.x</i>	$f_4 = 0$	Rel_{-1}
<i>West</i>	<i>Taxi.x</i>	$f_4 = 1$	Rel_0
<i>Pickup</i>	<i>Passenger.in-taxi</i>	$f_5 = 1 \wedge f_6 = 0 \wedge f_7 = 0$	$SetBool_1$
<i>Pickup</i>	<i>Passenger.in-taxi</i>	$f_5 = 0 \vee f_6 = 1 \vee f_7 = 1$	$SetBool_0$
<i>Dropoff</i>	<i>Passenger.in-taxi</i>	$f_8 = 1 \wedge f_9 = 1$	$SetBool_1$
<i>Dropoff</i>	<i>Passenger.in-taxi</i>	$f_8 = 0 \vee f_9 = 0$	$SetBool_0$
<i>Dropoff</i>	<i>Passenger.at-destination</i>	$f_8 = 1 \wedge f_9 = 1$	$SetBool_1$
<i>Dropoff</i>	<i>Passenger.at-destination</i>	$f_8 = 0 \vee f_9 = 0$	$SetBool_0$

Table 2: Preconditions and effects for each action and attribute in the All-Passenger Any-Destination Taxi domain. We exclude *Wall* attributes because they never change.

Deictic OO-MDPs are also able to compactly represent the transition dynamics of the Sokoban domain without ambiguity. For a deictic *person* object there are four deictic predicates required when taking the action *East* that resolve the following questions: is there any object of class *Box* one square east of *person*? Is there any object of class *Box* two squares east of *person*? Is there any object of class *Wall* one square east of *person*? Is there any object of class *Wall* two squares east of *person*? Meanwhile for a deictic *box* object there are three deictic predicates required to resolve the questions: is there any object of class *Person* one square west of *box*? Is there any object of class *Wall* one square east of *box*? Is there any object of class *Box* one square east of *box*? More specifically, define the following deictic predicates:

- *AnyWallEastOfPerson1*(*person*, *s*): returns 1 if there is any object of class *Wall* one square east of *person* in state *s*, and 0 otherwise. (f_1)

- $\text{AnyWallEastOfPerson2}(\text{person}, s)$: returns 1 if there is any object of class *Wall* two squares east of *person* in state *s*, and 0 otherwise. (\hat{f}_2)
- $\text{AnyBoxEastOfPerson1}(\text{person}, s)$: returns 1 if there is any object of class *Box* one square east of *person* in state *s*, and 0 otherwise. (\hat{f}_3)
- $\text{AnyBoxEastOfPerson2}(\text{person}, s)$: returns 1 if there is any object of class *Box* two squares east of *person* in state *s*, and 0 otherwise. (\hat{f}_4)
- $\text{AnyPersonWestOfBox}(\text{box}, s)$: returns 1 if there is any object of class *Person* one square west of *box* in *s*, and otherwise 0. (\hat{f}_5)
- $\text{AnyBoxEastOfBox}(\text{box}, s)$: returns 1 if there is any object of class *Box* one square east of *box* in *s*, and otherwise 0. (\hat{f}_6)
- $\text{AnyWallEastOfBox}(\text{box}, s)$: returns 1 if there is any object of class *Wall* one square east of *box* in *s*, and otherwise 0. (\hat{f}_7)

Then table 3 shows for each attribute and action the relevant preconditions and effects that describe the transition dynamics of the domain under the Deictic OO-MDP formalism for action *East*.

Action	Attribute	Precondition	Effect
<i>East</i>	<i>Person.x</i>	$\hat{f}_1 = 0 \wedge \hat{f}_3 = 0 \vee \hat{f}_3 = 1$ $\wedge \hat{f}_2 = 0 \wedge \hat{f}_4 = 0$	Rel_1
<i>East</i>	<i>Person.x</i>	$\hat{f}_1 = 1 \vee \hat{f}_3 = 1 \wedge \hat{f}_4 = 1$ $\vee \hat{f}_3 = 1 \wedge \hat{f}_2 = 1$	Rel_0
<i>East</i>	<i>Box.x</i>	$\hat{f}_5 = 1 \wedge \hat{f}_6 = 0 \wedge \hat{f}_7 = 0$	Rel_1
<i>East</i>	<i>Box.x</i>	$\hat{f}_5 = 0 \vee \hat{f}_6 = 1 \vee \hat{f}_7 = 1$	Rel_0
<i>Reset</i>	<i>Person.reset</i>	$\hat{f}_8 = 1$	$SetBool_1$

Table 3: Precondition and effects for action *East* and *Reset* in the Sokoban domain. We exclude *Wall* attributes since they cannot change.

The actions *North*, *South* and *West* are analogous. In addition, we include a *reset* attribute for the *Person* class that activates when a *Reset* action is taken, along with the following deictic predicate.

- $\text{ResetActivated}(\text{person}, s)$: returns 1 if *person.reset* = 1 in *s*, otherwise 0. (\hat{f}_8).

The *Reset* action immediately terminates the task, and its inclusion is necessary because in Sokoban it is possible to reach a deadlock state from which the task is no longer solvable.

3.3 Learning

Given a set of *D* deictic predicates we want to learn the transition dynamics for each action *a* and attribute *C.α*. If the transition dynamics are deterministic this can be done using memoisation with 2^D unique observations. However, this is prohibitive if *D* is large. As discussed in section 2.4 Propositional OO-MDPs introduce a learning algorithm called *DOORMAX* for deterministic transition dynamics that, under certain assumptions, has a KWIK-bound that is linear in *D*.

For *DOORMAX* to be correct, the transition dynamics for each action and attribute must be representable as a full binary tree with propositions at the non-leaf nodes and effects at

the leaf nodes. Furthermore, each possible effect of an effect type can occur at most at one leaf node of the tree, except for a special effect called a *failure condition* that may occur at multiple leaf nodes. A failure condition implies that globally no attribute changes when an action was taken i.e. $s = s'$ when a is taken. See figure 5a for how this is represented for the $Taxi.x$ attribute with action *East*.

The intuition behind *DOORMAX* is that, in many cases, **the number of propositions that an effect depends on is much smaller than D** . Furthermore, since an effect can occur at most once in the tree, we can invalidate multiple propositions with a single observation. To illustrate the core learning mechanism of *DOORMAX*, consider the original Taxi domain with action *East*, attribute $Taxi.x$ and effect type Rel_i . We can propose four propositional statements $Touch_N(Taxi, Wall)$, $Touch_E(Taxi, Wall)$, $Touch_S(Taxi, Wall)$ and $Touch_W(Taxi, Wall)$. The aim of the KWIK-learner then is to learn that

$$Touch_E(Taxi, Wall) = 0 \implies Taxi.x \leftarrow Taxi.x + 1$$

when action *East* is taken. Suppose that agent interacts with its environment and takes action *East* to observe the effect: $Taxi.x + 1$. Prior to taking the action *East*, the agent's state determined the following assignments to the propositions:

$$Touch_N(Taxi, Wall) = 0 \wedge Touch_E(Taxi, Wall) = 0 \wedge Touch_S(Taxi, Wall) = 0 \wedge Touch_W(Taxi, Wall) = 0.$$

Suppose the next time the agent takes the action *East* and observes the effect $Taxi.x + 1$ and the agent's previous state determined the following assignments to the propositions:

$$Touch_N(Taxi, Wall) = 1 \wedge Touch_E(Taxi, Wall) = 0 \wedge Touch_S(Taxi, Wall) = 1 \wedge Touch_W(Taxi, Wall) = 1.$$

Since the *DOORMAX* algorithm assumes that each effect can occur at most once in the tree, we can infer that the propositions $Touch_N(Taxi, Wall)$, $Touch_S(Taxi, Wall)$ and $Touch_W(Taxi, Wall)$ are irrelevant to the precondition of the effect. As a result the agent has learned

$$Touch_E(Taxi, Wall) = 0 \implies Taxi \leftarrow Taxi.x + 1$$

with only two unique observations. Of course, this is a best-case scenario. A worst-case scenario requires 4 unique observations. In general, given $D \geq 0$ propositions, learning requires at most $D + 1$ unique observations.⁴

We adapt the *DOORMAX* algorithm to Deictic OO-MDPs, which we call *DOORMAX_D* (algorithm 7). The *DOORMAX_D* algorithm requires two sub-algorithms for the transition dynamics: one to learn and one to make predictions. These sub-algorithms are presented in this section while the full *DOORMAX_D* algorithm that calls these sub-algorithms is then presented in section 5. The main difference between *DOORMAX_D* and *DOORMAX* is that we remove the notion of a global failure condition. Instead we require that all effects apply to a single attribute. See figure 5b for how this is represented for the $Taxi.x$ attribute with action *East*. To achieve this, we extend the notion of an effect type to include a partition function over effects that groups them into those that can occur at most at one leaf node and those that can occur at multiple leaf nodes. For example, as shown in figure 5b, the effects $taxi.x \leftarrow taxi.x + 1 (Rel_1)$ and $taxi.x \leftarrow taxi.x + 0 (Rel_0)$ are both unique in the tree. Therefore, by using an appropriate partition function, we can efficiently learn both these branches. Meanwhile, the design choice used by *DOORMAX* of a global failure conditions implies that $Taxi.x \leftarrow Taxi.x + 1 (Rel_1)$ is efficiently learned, while $Taxi.x \leftarrow Taxi.x + 0 (Rel_0)$ is learned by memorisation. A trade-off exists between these design choices. The *DOORMAX_D*

⁴ It is $D + 1$ and not D because an effect can depend on no propositions.

algorithm requires a partition function for each action, attribute and effect set that improves learning efficiency. However, constructing such a partition function requires additional prior knowledge about a domain that is not required by *DOORMAX*.

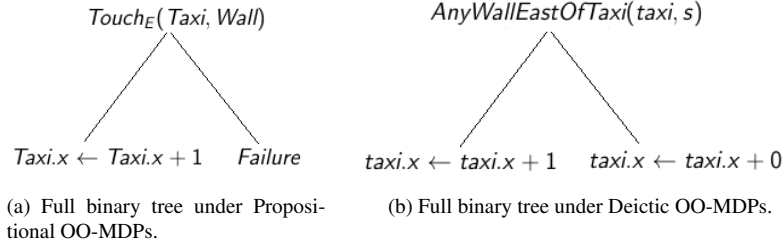


Fig. 5: Full binary tree structure for the transition dynamics of $Taxi.x$ attribute and action *East*. Right branches represent a truth value of 1.

In this section we introduce two algorithms that are called by *DOORMAX_D* to KWIK-learn the transition dynamics of an MDP under the Deictic OO-MDP formalism:

- *UpdateTree1* (algorithm 1) updates the binary tree for action a and attribute $C.\alpha$.
- *Predict1* (algorithm 2) makes a prediction for action a and attribute $C.\alpha$.

Before introducing the algorithms, we require some definitions. Let \mathcal{F} be a set of deictic predicates and \mathcal{E} be a set of effects.

Definition 1 A term is a tuple (f, b) where $f \in \mathcal{F}$ and $b \in \mathcal{B}$. A set of terms is denoted by \mathcal{T} . A set that contains sets of terms is denoted by \mathcal{T} .

Definition 2 Term (f_1, b_1) mismatches term (f_2, b_2) if $f_1 = f_2$ and $b_1 \neq b_2$.

Definition 3 $\Pi : \mathcal{E} \rightarrow \mathcal{B}$ is a binary partition function over \mathcal{E} and assigns each effect in \mathcal{E} to one of two partitions, 0 or 1. We call the tuple $g = (\mathcal{E}, \Pi)$ an effect type. Denote by K_0^g and K_1^g the number of effects in partition 0 and 1 respectively. We use \cdot notation to refer to an element in a tuple g , so for example $g.\mathcal{E}$ refers to \mathcal{E} in g . We denote sets of effect types with \mathcal{G} .

Definition 4 Let g be an effect type. Let $M > 1$ be a constant. Then $Tree(g, \mathcal{F}, M)$ is the set of all full binary trees such that non-leaf nodes are elements of \mathcal{F} and leaf nodes are elements of $g.\mathcal{E}$. Furthermore, if $g.\Pi$ assigns an effect in $g.\mathcal{E}$ to partition 1 then that effect can occur at most at one leaf node and we call that effect conjunctive, otherwise that effect may occur at most at M leaf nodes and we call that effect disjunctive.

We now introduce the *UpdateTree1* (algorithm 1) and *Predict1* (algorithm 2) algorithms. Note that $\hat{\mathcal{F}}(a, C.\alpha)$, $\hat{\mathcal{G}}(a, C.\alpha)$, $\mathfrak{T}_e^g(a, C.\alpha)$ and M are initialised globally in *DOORMAX_D* as discussed in section 5. The set $\hat{\mathcal{F}}(a, C.\alpha)$ contains hypothesised deictic predicates for action a and attribute $C.\alpha$; the set $\hat{\mathcal{G}}(a, C.\alpha)$ contains hypothesised effect types for action a and attribute $C.\alpha$; while the set $\mathfrak{T}_e^g(a, C.\alpha)$ contains the sets of terms that map to effect e of effect type g for action a and attribute $C.\alpha$. The $\hat{\cdot}$ notation used for $\hat{\mathcal{F}}(a, C.\alpha)$ and

Algorithm 1: *UpdateTreeI*: update binary tree for action a and attribute $C.\alpha$.

Input: $C.\alpha \in \text{Att}(C)$, $s \in \mathcal{S}$, $o \in O[C]$, $a \in \mathcal{A}$, $o.\alpha' \in \text{Dom}(C.\alpha)$

```

1  pass  $o$  and  $s$  to the deictic predicates in  $\hat{\mathcal{F}}(a, C.\alpha)$  to retrieve a set of terms  $\mathcal{T}$ 
2  foreach  $g \in \mathcal{G}(a, C.\alpha)$  do
3      foreach  $e \in g.\mathcal{E}$  do
4          if  $e(o.\alpha) = o.\alpha'$  then
5              if  $\mathfrak{T}_e^g(a, C.\alpha) = \emptyset$  then
6                  if  $\exists e' \in g.\mathcal{E}$  with  $\mathcal{T}_{e'} \in \mathfrak{T}_{e'}^g(a, C.\alpha)$  such that  $\mathcal{T}_{e'} \subseteq \mathcal{T}$  then
7                      remove  $g$  from  $\mathcal{G}(a, C.\alpha)$ 
8                  else
9                      add  $\mathcal{T}$  to  $\mathfrak{T}_e^g(a, C.\alpha)$ 
10                 end
11             else
12                 if  $g.\Pi(e) = 0$  then
13                     add  $\mathcal{T}$  to  $\mathfrak{T}_e^g(a, C.\alpha)$ 
14                 else
15                      $\mathcal{T}_{\text{temp}} \leftarrow \mathcal{T}$ 
16                      $\mathcal{T} \leftarrow$  the only element in  $\mathfrak{T}_e^g$ 
17                     remove from  $\mathcal{T}$  any terms that mismatch terms in  $\mathcal{T}_{\text{temp}}$ 
18                      $\mathfrak{T}_e^g(a, C.\alpha) \leftarrow \emptyset$ 
19                     add  $\mathcal{T}$  to  $\mathfrak{T}_e^g(a, C.\alpha)$ 
20                 end
21                 if  $(\exists e' \in (g.\mathcal{E} - \{e\})$  with  $\mathcal{T}_{e'} \in \mathfrak{T}_{e'}^g(a, C.\alpha)$  such that  $(\mathcal{T} \subseteq \mathcal{T}_{e'}$  or  $\mathcal{T}_{e'} \subseteq \mathcal{T})$ 
22                     or  $|\mathfrak{T}_e^g(a, C.\alpha)| > M$  then
23                     remove  $g$  from  $\mathcal{G}(a, C.\alpha)$ 
24                 end
25             end
26         end
27     end
28 end

```

$\mathcal{G}(a, C.\alpha)$ indicates that these sets are sets of hypotheses, of which the correct hypotheses must be determined by the algorithm.

At a high-level, the algorithm *UpdateTreeI* ensures that for all $g \in \mathcal{G}(a, C.\alpha)$ the constraint $\text{Tree}(g, \mathcal{F}(a, C.\alpha), M)$ is maintained. That is, the set $\{\mathfrak{T}_e^g(a, C.\alpha)\}_{e \in g.\mathcal{E}}$ induces a binary tree subject to the constraints of the partition function $g.\Pi$. Each time we observe a set of terms \mathcal{T} and an associated effect e we update $\mathfrak{T}_e^g(a, C.\alpha)$, and in doing so may discover that the resulting set $\{\mathfrak{T}_e^g(a, C.\alpha)\}_{e \in g.\mathcal{E}}$ can no longer induce an appropriate binary tree at which point we remove g from $\mathcal{G}(a, C.\alpha)$.

At the lower-level, the algorithm *UpdateTreeI* updates the binary tree for action a and attribute $C.\alpha$ given an object o from state s and the resulting attribute value $o.\alpha'$ in s' . The algorithm starts at line 1 where it computes \mathcal{T} , the terms for $\hat{\mathcal{F}}(a, C.\alpha)$. In lines 2-4 the algorithm iterates over each effect type $g \in \mathcal{G}(a, C.\alpha)$ and each effect $e \in g.\mathcal{E}$ where it checks if e applied to $o.\alpha$ is equal to $o.\alpha'$. If it is, then it proceeds to update $\mathfrak{T}_e^g(a, C.\alpha)$ in lines 5-25.

If $\mathfrak{T}_e^g(a, C.\alpha)$ is empty, then as per lines 5-11, the algorithm checks if \mathcal{T} already maps to different effect in g . If it does, then the effect type g is invalidated because adding \mathcal{T} to $\mathfrak{T}_e^g(a, C.\alpha)$ would result in an invalid tree; otherwise, \mathcal{T} is added to $\mathfrak{T}_e^g(a, C.\alpha)$.

If $\mathfrak{T}_e^g(a, C.\alpha)$ is not empty then, as per lines 12-20, it first checks if e is disjunctive or conjunctive based on the partition function $g.\Pi$. If e is disjunctive then it adds \mathcal{T} to $\mathfrak{T}_e^g(a, C.\alpha)$. If e is conjunctive then it updates the existing element in $\mathfrak{T}_e^g(a, C.\alpha)$ by remov-

ing from it any terms that mismatch terms in \mathcal{T} - this refines the set of terms that e depends on in the tree, and is the core mechanism for efficient learning in $DOORMAX_D$.

Finally, after the update to $\mathfrak{T}_e^g(a, C, \alpha)$, in lines 21-23 the algorithm checks if \mathcal{T} maps to a different existing effect in g or if the number of effects in $\mathfrak{T}_e^g(a, C, \alpha)$ exceeds the limit M . In either of these cases, g is invalidated.

To provide a better intuition on the invalidation logic of the algorithm, consider a case where $\mathcal{F} = \{f_1, f_2, f_3\}$ and there are two effects $\mathcal{E} = \{e_1, e_2\}$ where e_1 is conjunctive and e_2 is disjunctive. Consider the following examples:

- Figure 6a: we currently have $\mathfrak{T}_{e_1} = \{\mathcal{T}_{e_1} = \{(f_1, 1), (f_2, 1), (f_3, 1)\}\}$ and $\mathfrak{T}_{e_2} = \emptyset$. Suppose we then observe e_2 with $\mathcal{T} = \{(f_1, 1), (f_2, 1), (f_3, 1)\}$. Then \mathfrak{T}_{e_2} is empty and $\mathcal{T}_{e_1} \subseteq \mathcal{T}$.
- Figure 6b: we currently have $\mathfrak{T}_{e_1} = \{\mathcal{T}_{e_1} = \{(f_1, 1), (f_2, 1), (f_3, 1)\}\}$ and $\mathfrak{T}_{e_2} = \{\mathcal{T}_{e_2} = \{(f_1, 1), (f_2, 1), (f_3, 0)\}\}$. Suppose we then observe e_1 with $\mathcal{T} = \{(f_1, 1), (f_2, 0), (f_3, 0)\}$. As e_1 is conjunctive and \mathfrak{T}_{e_1} is not empty we first remove mismatching terms. Then $\mathfrak{T}_{e_1} = \{\mathcal{T} = \{(f_1, 1)\}\}$ and now $\mathcal{T} \subseteq \mathfrak{T}_{e_2}$.
- Figure 6c: we currently have $\mathfrak{T}_{e_1} = \{\mathcal{T}_{e_1} = \{(f_1, 1)\}\}$ and $\mathfrak{T}_{e_2} = \{\mathcal{T}_{e_2} = \{(f_1, 0), (f_2, 0), (f_3, 0)\}\}$. Suppose we then observe e_2 with $\mathcal{T} = \{(f_1, 1), (f_2, 0), (f_3, 1)\}$. We add \mathcal{T} to \mathfrak{T}_{e_2} and now $\mathcal{T}_{e_1} \subseteq \mathcal{T}$.

In all the above cases, we conclude that the effect type is invalid since the observed data can no longer induce a binary tree subject to the specified constraints. Note that we do not place any restrictions on the order in which the preconditions may appear in the tree, but there is no reordering that can recover an appropriate binary tree given the data.

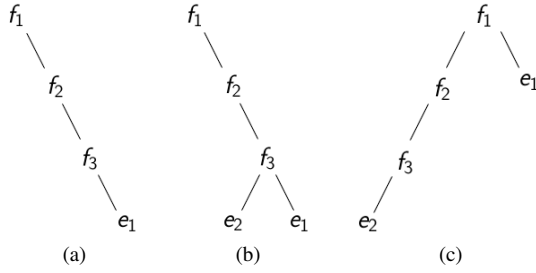


Fig. 6: Binary trees induced by \mathfrak{T}_{e_1} and \mathfrak{T}_{e_2} . Right branches represent a truth value of 1.

The *Predict1* procedure is analogous to that used in *DOOMAX*. At a high-level the algorithm makes a prediction for attribute $C.\alpha$ of an object o of class C given a schema state s and action a . The algorithm requires that for each effect type $g \in \mathcal{G}(a, C, \alpha)$ a prediction is made that is not \perp , and furthermore that all effect types make the same prediction; otherwise the algorithm returns \perp . This ensures that the algorithm only makes a correct prediction if it does not return \perp . Once the transition dynamics are learned, then given a state s and action a the resulting state s' can be determined by calling *Predict1* for each object attribute $o.\alpha$ in s to obtain $o.\alpha'$ in s' .

In line 2 of the algorithm it computes \mathcal{T} , the terms for $\hat{\mathcal{F}}(a, C, \alpha)$. In line 3 the algorithm iterates over the effect types $g \in \mathcal{G}(a, C, \alpha)$. In lines 4-10 the algorithm checks whether any effect $e \in g.\mathcal{E}$ is mapped to by \mathcal{T} in $\mathfrak{T}_e^g(a, C, \alpha)$. If it is, then the prediction $e(o.\alpha)$ is recorded

Algorithm 2: *Predict1*: prediction procedure for action a and attribute $C.\alpha$.

Input: $C.\alpha \in \text{Att}(C)$, $s \in S$, $o \in O[C]$, $a \in \mathcal{A}$

- 1 initialise an empty set $\mathcal{V} \langle \text{Dom}(C.\alpha) \rangle$
- 2 pass o and s to the deictic predicates in $\hat{\mathcal{F}}(a, C.\alpha)$ to retrieve a set of terms \mathcal{T}
- 3 **foreach** $g \in \mathcal{G}(a, C.\alpha)$ **do**
- 4 $\text{pred} \leftarrow \perp$
- 5 **foreach** $e \in g.\mathcal{E}$ **do**
- 6 **if** $\exists \mathcal{T}_e \in \mathcal{T}_e^g(a, C.\alpha)$ such that $\mathcal{T}_e \subseteq \mathcal{T}$ **then**
- 7 $\text{pred} \leftarrow e(o.\alpha)$
- 8 **exit loop**
- 9 **end**
- 10 **end**
- 11 **if** $\text{pred} = \perp$ **then**
- 12 **return** \perp
- 13 **else**
- 14 add pred to \mathcal{V}
- 15 **if** $|\mathcal{V}| > 1$ **then**
- 16 **return** \perp
- 17 **end**
- 18 **end**
- 19 **end**
- 20 **return** only element in \mathcal{V}

and added to a set \mathcal{V} in line 14; otherwise the algorithm returns \perp in line 12. In lines 15-17 the algorithm checks if $|\mathcal{V}| > 1$. If it is, then there are two effect types that make different predictions for $o.\alpha$ and the algorithm returns \perp ; otherwise, all effect types record the same prediction for $o.\alpha$ and the algorithm returns this prediction in line 20.

4 Transferable Reward Dynamics with Propositional OO-MDPs

Propositional OO-MDPs introduce a KWIK-learning algorithm to learn the transition dynamics of a domain, while assuming known reward dynamics. In this section we introduce an algorithm to KWIK-learn a family of reward dynamics under a propositional object-oriented approach.

This family consists of reward dynamics whereby for most transitions the agent receives a default reward signal, while a small number of transition groups lead to different reward signals. An example of a domain that exhibits such reward dynamics is the Taxi domain. In the Taxi domain, the agent receives a default reward of -1 for all transitions except for the groups of transitions that: apply an illegal pickup operation, which produce a reward signal of -10 ; apply an illegal dropoff operation, which produce a reward signal of -10 ; or lead to terminal state, which produces a reward signal of 0 . Furthermore, it is required that these groups can be described under an object-oriented approach through propositional statements over object class attributes.

We note that while it is possible to extend the propositional formalism described in this section to a deictic formalism, the structure of reward dynamics in most MDPs makes a deictic approach unnecessary. This is because most RL tasks have global goals such as ‘get all boxes to some storage location’ in the Sokoban domain or ‘get all passengers to some destination location’ in the All-Passenger Any-Destination Taxi domain.

4.1 Refactoring Reward Dynamics

We note that for any finite MDP, the reward dynamics can be refactored as a sum of $L + 1$ terms:

$$\mathcal{R}(s, a, s') = \sum_{i=1}^L z_i(s, a, s') U_i + (1 - \sum_{i=1}^L z_i(s, a, s')) U_{L+1}, \quad (1)$$

where the $z_i(s, a, s') \in \mathfrak{B}$ are indicator variables such that $\sum_{i=1}^L z_i(s, a, s') \in \mathfrak{B}$ - that is, for any (s, a, s') at most one z_i has value one and all others have value zero - and U_i is a *reward token* that maps to a stationary reward distribution (or a scalar in the case of deterministic reward dynamics). We call U_{L+1} the *default reward token*.

Rewriting the reward dynamics in this way does not sacrifice generality since any arbitrary reward dynamics can be mapped to equation 1 as follows:

- Set $L = |\mathcal{S}|^2 |\mathcal{A}| - 1$.
- Pick one (s, a, s') arbitrarily and map $U_{L+1} = \mathcal{R}(s, a, s')$.
- For all other L transitions map some previously unmapped $i \in [1 : L]$ to $z_j(s, a, s') = 1$ if $i = j$, otherwise $z_j(s, a, s') = 0$; and map $U_i = \mathcal{R}(s, a, s')$.

However, for many tasks we can group transitions together and therefore define the reward dynamics with $L \ll |\mathcal{S}|^2 |\mathcal{A}| - 1$. This is particularly evident in tasks with sparse rewards where the agent gets a constant default reward for almost all transitions in an MDP. For example, consider the original Taxi task. In this task the agent gets a default reward of -1 for all steps except for an illegal *Pickup* action, an illegal *Dropoff* action or for reaching a terminal state. Therefore, we can represent the reward dynamics for the entire Taxi domain with $L = 3$ given the following propositions as per table 4:

- P_a : (s, a, s') is a transition where the *Pickup* action is applied illegally.
- P_b : (s, a, s') is a transition where the *Dropoff* action is applied illegally.
- P_c : (s, a, s') is a transition which reaches a terminal state.

Proposition	z_1	z_2	z_3	Reward Token
$P_a = 1$	1	0	0	$U_1 = -10$
$P_b = 1$	0	1	0	$U_2 = -10$
$P_c = 1$	0	0	1	$U_3 = 0$
Default	0	0	0	$U_4 = -1$

Table 4: Expressing the reward dynamics for the Taxi domain as per equation 1 with $L = 3$.

Furthermore, under an object-oriented approach, each of the propositions P_a , P_b and P_c can be constructed from propositions over object class attributes as shown below:⁵

- P_a : *Pickup* action with $On(Taxi, Passenger) = 0 \vee On(Taxi, Passenger) = 1 \wedge Passenger.in-taxi = 1$.
- P_b : *Dropoff* action with $Passenger.in-taxi = 0 \vee Passenger.in-taxi = 1 \wedge On(Taxi, Destination) = 0$.
- P_c : *Dropoff* action with $On(Taxi, Destination) = 1 \wedge Passenger.in-taxi = 1$.

⁵ These propositions are extracted from the state s of the transition tuple (s, a, s') . In general, the propositional statements for the reward dynamics can be extracted from s' as well. However, for this domain only (s, a) is needed to capture the reward dynamics.

4.2 Formalism

We enhance the Propositional OO-MDP formalism to include the representation of reward dynamics under equation 1. Given $L \in \mathbb{N}$, $a \in \mathcal{A}$ and $i \in [1 : L]$, define a set of propositions of size $D_{a,i}$:

$$\mathcal{F}_{a,i} = \{f_{a,i} : S \times S \rightarrow \mathcal{B}\}_{i=1}^{D_{a,i}},$$

and a binary mapping function that is used determine whether reward token mapped to index i triggers given the truth values of the propositions in $\mathcal{F}_{a,i}$:

$$\mathcal{Z}_{a,i} : \mathcal{B}^{D_{a,i}} \rightarrow \mathcal{B}.$$

Define a set of $L + 1$ reward tokens:

$$\mathcal{U} = \{U_j\}_{j=1}^{L+1}.$$

Then the schema reward dynamics are defined by the set:

$$\mathcal{R} = \{\mathcal{F}_{a,i}, \mathcal{Z}_{a,i}, U_j | j \in [1 : L + 1], i \in [1 : L], a \in \mathcal{A}\}.$$

Then given a grounded MDP $\mathcal{M}_{O,\rho} = (\mathcal{S}_O, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma, \rho)$ the reward dynamics operate as follows: the agent is currently in state $s \in \mathcal{S}_O$, takes action a and transitions to state $s' \in \mathcal{S}_O$. For each $i \in [1 : L]$ compute the set of truth values for the propositional statements in $\mathcal{F}_{a,i}$ to get a set of binary values $\mathcal{B}_i = \{f_k(s, s')\}_{k=1}^{D_{a,i}}$. Compute $\mathcal{Z}_{a,i}(\mathcal{B}_i)$ to obtain some indicator variable z_i . The set of indicator variables is then $\mathcal{Z} = \{z_i\}_{i=1}^L$. The elements in the sets \mathcal{Z} and \mathcal{U} are then passed to equation 1 to compute a reward. Recall that equation 1 is restricted so that $\sum_{i=1}^L z_i \in \mathcal{B}$. Therefore, if $\sum_{i=1}^L z_i = 1$ then the agent receives a reward generated from whichever reward token U_j is activated by $z_j = 1$, and if $\sum_{i=1}^L z_i = 0$ then the agent receives a reward generated from the default reward token U_{L+1} .

As a practical example of the formalism, consider the action *Dropoff* with $\mathcal{U} = \{U_1, U_2, U_3, U_4\}$ as described in table 4. Suppose we have $\mathcal{F}_{Dropoff,2} = \mathcal{F}_{Dropoff,3} = \{AnyTaxiOnAnyDestination(s, s'), AnyPassengerInAnyTaxi(s, s')\}$ that are defined as:⁶

- *AnyTaxiOnAnyDestination(s, s')*: returns 1 if any objects of class *Taxi* is on the same square as any object of class *Destination* in s , and otherwise 0.
- *AnyPassengerInAnyTaxi(s, s')*: returns 1 if any objects of class *Passenger* has their *in-taxi* attribute set to 1 in s , and otherwise 0.

Given assignments $\{b_1, b_2\}$ to these propositional statements the mapping functions would be defined as:

- $\mathcal{Z}_{Dropoff,2}$: return 1 if $(b_2 = 0)$ or $(b_2 = 1 \text{ and } b_1 = 0)$, else return 0.
- $\mathcal{Z}_{Dropoff,3}$: return 1 if $(b_1 = 1 \text{ and } b_2 = 1)$, else return 0.

Note that we would simply define $\mathcal{F}_{Dropoff,1} = \emptyset$ and $\mathcal{Z}_{Dropoff,1}$ would always return a value of 0. This is because an illegal pickup is not possible when the action *Dropoff* is taken.

Suppose that the agent is currently in state s , takes action *Dropoff* and the resulting state is s' . The states s and s' are passed to *AnyTaxiOnAnyDestination(s, s')* and *AnyPassengerInAnyTaxi(s, s')* to obtain their respective truth values. Suppose that the resulting truth values are:

⁶ We can actually simplify *AnyTaxiOnAnyDestination(s, s')* and *AnyPassengerInAnyTaxi(s, s')* to *AnyTaxiOnAnyDestination(s)* and *AnyPassengerInAnyTaxi(s)* respectively because these propositions do not depend on s' . However, we include s' to remain consistent with the notation of the formalism.

- $\{1, 0\}$: then $z_2 = 1$ and $z_3 = 0$, so the agent receives a reward from the reward token $U_2 = -10$.
- $\{1, 1\}$: then $z_2 = 0$ and $z_3 = 1$, so the agent receives a reward from the reward token $U_3 = 0$.

4.3 Examples

The reward dynamics for the All-Passenger Any-Destination Taxi domain require $L = 3$ reward tokens in addition to the default reward token. **Define the following propositions:**

- $AnyTaxiOnAnyPassenger(s, s')$: returns 1 if any objects of class *Taxi* is on the same square as any object of class *Passenger* in s , and otherwise 0. (p_1)
- $AnyPassengerInAnyTaxi(s, s')$: returns 1 if any objects of class *Passenger* has their *in-taxi* attribute set to 1 in s , and otherwise 0. (p_2)
- $AnyTaxiOnAnyDestination(s, s')$: returns 1 if any objects of class *Taxi* is on the same square as any object of class *Destination* in s , and otherwise 0. (p_3)
- $AllPassengersAtAnyDestination(s, s')$: returns 1 if all objects of class *Passenger* have their *at-destination* attributes set to 1 in s' , and otherwise 0. (p_4)

Define the following reward tokens:

- $\bar{U}_1 = -10$
- $\bar{U}_2 = -10$
- $\bar{U}_3 = 0$
- $\bar{U}_4 = -1$

Then table 5 shows the reward token that activates for each action and precondition.

Action	Precondition	Reward Token
<i>Pickup</i>	$p_1 = 0 \vee p_1 = 1 \wedge p_2 = 1$	\bar{U}_1
<i>Dropoff</i>	$p_2 = 0 \vee p_2 = 1 \wedge p_3 = 0$	\bar{U}_2
<i>Dropoff</i>	$p_4 = 1$	\bar{U}_3

Table 5: Reward tokens that activate for each action and precondition in the All-Passenger Any-Destination Taxi domain. Any other combination leads to the default reward token \bar{U}_4 .

For the Sokoban domain we set $L = 2$. We define the following propositions for the reward dynamics:

- $AllBoxesAtAnyStorage(s, s')$: returns 1 if all objects of class *Box* are on the same square as any object of class *Storage* in s' , and otherwise 0. (\acute{p}_1)
- $ResetActivated(s, s')$: returns 1 if any object of class *Person* has their *reset* attribute set to 1 in s' , otherwise 0. (\acute{p}_2)

Define the following reward tokens ⁷:

⁷ These reward dynamics encourage the agent to get all boxes to a storage location in as few steps as possible. The high reward of 300 for achieving this is necessary so that the agent does not learn an optimal policy that applies the *Reset* action prematurely.

- $\dot{U}_1 = 300$
- $\dot{U}_2 = -1$
- $\dot{U}_3 = -1$

Then table 6 shows the reward token that activates for each action and precondition.

Action	Precondition	Reward Token
<i>North</i>	$\dot{p}_1 = 1$	\dot{U}_1
<i>East</i>	$\dot{p}_1 = 1$	\dot{U}_1
<i>South</i>	$\dot{p}_1 = 1$	\dot{U}_1
<i>West</i>	$\dot{p}_1 = 1$	\dot{U}_1
<i>Reset</i>	$\dot{p}_2 = 1$	\dot{U}_2

Table 6: Reward tokens that activate for each action and precondition in the Sokoban domain. Any other combination leads to the default reward token \dot{U}_3 .

4.4 Learning

In this section we propose an algorithm to KWIK-learn the mapping functions $\mathcal{L}_{a,i}$ described in section 4.2 given a set of propositions $\mathcal{F}_{a,i}$ for $i \in [1 : L]$. The procedure for learning these mapping functions is analogous to learning transition dynamics as described in section 3.3.

When learning transition dynamics, **we use a full binary tree structure for each action and attribute with logical statements at the leaf nodes and effects at the non-leaf nodes.** For learning the mapping functions, we will also use a full binary tree structure with logical statements, in the form of propositions, at the leaf nodes; but now the non-leaf nodes are binary indicators. Each such tree is used to indicate if a particular reward token is activated for a given transition. Therefore, we learn a full binary tree for each action $a \in \mathcal{A}$ and $i \in [1 : L]$.

Learning the mapping functions no longer requires the notion of an effect type, since the trees have binary leaf nodes. Therefore, in this setting, the partition function is $\Pi : \mathcal{B} \rightarrow \mathcal{B}$ that maps which of the binary indicator values is conjunctive and disjunctive in the tree. In a similar way to learning transition dynamics, the key to achieving efficiency under this framework is through the notion of conjunctive and disjunctive indicator values. That is, if an indicator value for action a and index i is known to occur at most at one non-leaf node in a tree then that branch can be KWIK-learned with at most $D_{a,i} + 1$ unique observations.

Returning to the example of the Taxi task with reward dynamics as described in table 4, $z_3 = 1$ can be efficiently learned because the precondition that maps to the activation of U_3 is the conjunction: $On(Taxi, Destination) = 1 \wedge Passenger.in-taxi = 1$. Meanwhile, $z_1 = 1$ and $z_2 = 1$ must be learned through memorisation because they occur at two leaf nodes in the tree. However, $z_1 = 0$ and $z_2 = 0$ can be efficiently learned since they occur when

$$On(Taxi, Passenger) = 1 \wedge Passenger.in-taxi = 0$$

and

$$On(Taxi, Destination) = 1 \wedge Passenger.in-taxi = 1$$

respectively.

We further assume that each U_j is KWIK-learnable with KWIK-bound B_j for $j \in [1 : L + 1]$. Under this assumption we can then KWIK-learn the complete reward dynamics with the algorithms described in this section. We introduce two algorithms that are called by the $DOORMAX_D$ algorithm presented in section 5 to KWIK-learn the reward dynamics for an MDP under the Propositional OO-MDP formalism.

- *UpdateTree2* (algorithm 3) updates the binary tree for action a and reward token mapped to index i .
- *Predict2* (algorithm 4) makes a binary prediction for action a and reward token mapped to index i .

Note that $\hat{\mathcal{F}}(a, i)$, $\mathcal{T}_b(a, i)$ and $\Pi(a, i)$ are initialised globally in $DOORMAX_D$. The set $\hat{\mathcal{F}}(a, i)$ contains hypothesised propositions for action a and reward token mapped to index i ; the set $\mathcal{T}_b(a, i)$ contains sets of terms that map to the binary indicator value b for action a and reward token mapped to index i ; while $\Pi(a, i)$ is a partition function for action a and reward token mapped to index i that determines if an indicator value is conjunctive or disjunctive in the tree induced by $\{\mathcal{T}_b(a, i)\}_{b \in \mathcal{B}}$. The $\hat{\cdot}$ notation used for $\hat{\mathcal{F}}(a, i)$ indicates that these are sets of hypotheses of which the correct hypotheses must be determined by algorithm.

The algorithm *UpdateTree2* operates in an analogous manner to the *UpdateTree1* algorithm presented in section 3.3. At a high-level, the algorithm ensures that the set $\{\mathcal{T}_b^g(a, i)\}_{b \in \mathcal{B}}$ at all times induces a binary tree subject to the constraints of the partition function $\Pi(a, i)$. Each time it observes a set of terms \mathcal{T} and an associated indicator value z it updates $\mathcal{T}_z^g(a, i)$ to refine tree.

At the lower-level the algorithm *UpdateTree2* updates the binary tree for action a and reward token mapped on index i given a schema state s , a resulting schema state s' and a binary indicator value z . The algorithm starts at line 1 where it computes \mathcal{T} , the terms for $\hat{\mathcal{F}}(a, i)$. In lines 2-3 the algorithm checks if $\mathcal{T}_z(a, i)$ is empty, and if it is it adds \mathcal{T} to $\mathcal{T}_z(a, i)$. If \mathcal{T}_z is not empty then, as per lines 5-13, it checks if z is disjunctive or conjunctive based on the partition function $\Pi(a, i)$. If z is disjunctive then it adds \mathcal{T} to $\mathcal{T}_z(a, i)$; otherwise, if z is conjunctive then it updates the only element in $\mathcal{T}_z(a, i)$ by removing from it any terms that mismatch terms in \mathcal{T} .

Algorithm 3: *UpdateTree2*: update binary tree for action a and reward token mapped to index i .

Input: $i \in [1 : L]$, $s \in \mathcal{S}$, $a \in \mathcal{A}$, $s' \in \mathcal{S}$, $z \in \mathcal{B}$

```

1 pass  $s$  and  $s'$  to the propositions in  $\hat{\mathcal{F}}(a, i)$  to retrieve a set of terms  $\mathcal{T}$ 
2 if  $\mathcal{T}_z(a, i) = \emptyset$  then
3   | add  $\mathcal{T}$  to  $\mathcal{T}_z(a, i)$ 
4 else
5   | if  $\Pi(a, i)(z) = 0$  then
6     | add  $\mathcal{T}$  to  $\mathcal{T}_z(a, i)$ 
7   | else
8     |  $\mathcal{T}_{temp} \leftarrow \mathcal{T}$ 
9     |  $\mathcal{T} \leftarrow$  the only element in  $\mathcal{T}_z(a, i)$ 
10    | remove from  $\mathcal{T}$  any terms that mismatch terms in  $\mathcal{T}_{temp}$ 
11    |  $\mathcal{T}_z(a, i) \leftarrow \emptyset$ 
12    | add  $\mathcal{T}$  to  $\mathcal{T}_z(a, i)$ 
13   | end
14 end
```

The algorithm *Predict2* takes as input a reward token index i , schema state s , action a and resulting schema state s' . The algorithm then returns a binary value indicating if the reward token mapped to index i activates, or \perp if the algorithm cannot yet make an accurate prediction. Once the reward dynamics are learned, then given a state s , action a and resulting state s' , the reward token that activates can be determined by calling *Predict2* for each $i \in [1 : L]$ and observing which of these returns a value of 1, or if they are all 0 then activating the default reward token. In line 1 the algorithm computes \mathcal{T} , the terms for $\hat{\mathcal{F}}(a, i)$. In lines 2-6 it checks whether any binary indicator value b is mapped to by \mathcal{T} in $\mathcal{T}_b(a, i)$. If it is, then b is returned; otherwise, the algorithm returns \perp in line 7.

Algorithm 4: *Predict2*: prediction procedure for action a and reward token mapped to index i .

Input: $i \in [1 : L], s \in S, a \in \mathcal{A}, s' \in S$
1 pass s and s' to the propositions in $\hat{\mathcal{F}}(a, i)$ to retrieve a set of terms \mathcal{T}
2 **foreach** $b \in \mathcal{B}$ **do**
3 **if** $\exists \mathcal{T}_b \in \mathcal{T}_b(a, i)$ such that $\mathcal{T}_b \subseteq \mathcal{T}$ **then**
4 return b
5 **end**
6 **end**
7 return \perp

An important point to emphasise with regards to the *UpdateTree2* and *Predict2* algorithms is that they take the index of reward token i as one of their inputs. This information is assumed to come from the environment that the agent interacts with as per line 9 of the *DOORMAX_D* algorithm in section 5. This additional knowledge is uncharacteristic of standard RL, where the agent only receives a scalar reward at each timestep. This assumption is necessary because the algorithm needs to know which of the trees to update with an indicator value of 1 and 0 when learning reward dynamics.

This additional knowledge inherently assumes that the agent has a built-in notion of ‘reward categorisation’. The idea of expanding the environment reward in RL to more than just a scalar has previously been introduced and shown to have benefits both in terms of learning performance and interpretability [12, 9, 20]. For example, by assuming that the agent receives a reward vector from its environment, where each element in the vector represents a reward type that encourages a specific behaviour (such as speed or safety), it is possible for an RL agent to learn a multi-dimensional policy that can execute any of the desired behaviours [12]. Such an approach assumes that the agent can categorise rewards, as it is able to perceive a reward vector that comprises of different reward types from its environment.

5 The *DOORMAX_D* Algorithm

In this section we combine the algorithms in sections 3 and 4 to present the full *DOORMAX_D* algorithm, while also providing formal guarantees of *DOORMAX_D* efficiency for the algorithm.

5.1 Algorithms

This section introduces the following algorithms / procedures:

- *DefineGlobal* (algorithm 5) defines the global data structures needed by the procedures called by *DOORMAX_D*.
- *BuildFullModels* (algorithm 6) builds the models \mathcal{P}_{model} and \mathcal{R}_{model} .
- *DOORMAX_D* (algorithm 7) is the full *DOORMAX_D* algorithm.

Note that the inputs to *DOORMAX_D* are assumed to be globally accessible to all the procedures called by the main algorithm.

The procedure *DefineGlobal* initialises the data structures and variables required by *DOORMAX_D*. In lines 1-8 the procedure initialises the $\mathcal{F}(a, C, \alpha)$, $\mathcal{G}(a, C, \alpha)$ and $\mathcal{T}_e^g(a, C, \alpha)$ data structures as well as the variable M required for learning transition dynamics. In lines 9-14 it initialises the $\mathcal{F}(a, i)$ and $\mathcal{T}_b(a, i)$ data structures as well as the partition functions $\Pi(a, i)$ and reward tokens \hat{U}_j for learning reward dynamics.

Algorithm 5: *DefineGlobal*: define global variables and data structures required for *DOORMAX_D*.

```

1  foreach  $C \in \mathcal{C}$  do
2    foreach  $(a, C, \alpha) \in \mathcal{A} \times \text{Att}(C)$  do
3      define  $\mathcal{F}(a, C, \alpha)$  globally as the hypothesised deictic predicates for each action  $a$  and
      attribute  $C, \alpha$ 
4      define  $\mathcal{G}(a, C, \alpha)$  globally as the hypothesised effect types for each action  $a$  and attribute
       $C, \alpha$ 
5      foreach  $g \in \mathcal{G}(a, C, \alpha)$  do
6        foreach  $e \in g, \mathcal{E}$  do
7          define  $\mathcal{T}_e^g(a, C, \alpha)$  globally as the set of sets for each action  $a$ , attribute  $C, \alpha$  and
          effect  $e$ 
8  define  $M > 1$  such that  $M - 1$  is the maximum number of elements allowed in any set  $\mathcal{T}_e^g(a, C, \alpha)$ 
9  foreach  $(a, i) \in \mathcal{A} \times [1 : L]$  do
10   define  $\mathcal{F}(a, i)$  globally as the hypothesised propositions for each action  $a$  and index  $i$  mapped
      to  $z_i$ 
11   define  $\mathcal{T}_b(a, i)$  globally as a set of sets for each action  $a$ , index  $i$  mapped to  $z_i$  and Boolean  $b$ 
12   define  $\Pi(a, i)$  globally as a partition function for each action  $a$  and index  $i$  mapped to  $z_i$ 
13 foreach  $j \in [1 : L + 1]$  do
14   define  $\hat{U}_j$  globally as the estimated reward token mapped to index  $j$ 

```

The algorithm *BuildFullModels* builds the full models \mathcal{P}_{model} and \mathcal{R}_{model} . The algorithm starts at lines 1-2 where it initialises the models to the most optimistic case where every transition returns r_{max} . In line 3 the algorithm iterates over all tuples $(s, a) \in \mathcal{S}_O \times \mathcal{A}$ of the grounded MDP. In lines 4-15 the algorithm calls *Predict1* for every object attribute value o, α in s and checks whether it can accurately predict the object attribute value o, α' when taking action a in order to construct the next state s' . If it can, then the algorithm proceeds to lines 16-25 where it calls *Predict2* on the tuple (s, a, s') for every reward token index $i \in [1 : L]$. The algorithm checks if it can make an accurate prediction for every $i \in [1 : L]$. If it can, the algorithm proceeds to lines 26-31 where it selects the activated reward token U_j . If U_j has been KWIK-learned, then the algorithm sets $\mathcal{P}_{model}(s, a)$ to s' and $\mathcal{R}_{model}(s, a, s')$ to \hat{U}_j .

The *DOORMAX_D* algorithm is the root algorithm of this paper. In line 1 it calls *DefineGlobal* to initialise the required data structures and variables. A start state s of the MDP is initialised in line 3 and the agent-environment interaction starts at line 4 with a while loop

Algorithm 6: *BuildFullModels*: build \mathcal{P}_{model} and \mathcal{R}_{model} for $DOORMAX_D$.

```

1  initialise  $\mathcal{P}_{model}(s, a) \leftarrow \perp$  for all  $(s, a) \in \mathcal{S}_O \times \mathcal{A}$ 
2  initialise  $\mathcal{R}_{model} \leftarrow r_{max}$  for all  $(s, a, s') \in \mathcal{S}_O \times \mathcal{A} \times \mathcal{S}_O$ 
3  foreach  $(s, a) \in \mathcal{S}_O \times \mathcal{A}$  do
4       $s' \leftarrow s$ 
5       $allPredKnown1 \leftarrow 1$ 
6      foreach object  $o$  in  $s'$  do
7          foreach attribute  $C.\alpha \in Att(C)$  from the class  $C$  of object  $o$  do
8               $pred \leftarrow Predict1(C.\alpha, s, o, a)$ 
9              if  $pred = \perp$  then
10                   $allPredKnown1 \leftarrow 0$ 
11              else
12                  replace attribute  $C.\alpha$  of object  $o$  in  $s'$  with  $pred$ 
13              end
14          end
15      end
16      if  $allPredKnown1 = 1$  then
17           $allPredKnown2 \leftarrow 1$ 
18           $j \leftarrow L + 1$ 
19          foreach  $i \in [1 : L]$  do
20               $pred \leftarrow Predict2(i, s, a, s')$ 
21              if  $pred = \perp$  then
22                   $allPredKnown2 \leftarrow 0$ 
23              else if  $pred = 1$  then
24                   $j \leftarrow i$ 
25          end
26          if  $allPredKnown2 = 1$  then
27              if  $\hat{U}_j$  with KWIK-bound  $B_j$  has been KWIK-learned then
28                   $\mathcal{P}_{model}(s, a) \leftarrow s'$ 
29                   $\mathcal{R}_{model}(s, a, s') \leftarrow \hat{U}_j$ 
30              end
31          end
32      end
33  end
34  return  $(\mathcal{P}_{model}, \mathcal{R}_{model})$ 

```

that ends when a terminal state is reached. In line 5 the models \mathcal{P}_{model} and \mathcal{R}_{model} are built by calling *BuildFullModels*. In line 6 the algorithm constructs an MDP \hat{M} with transition dynamics \mathcal{P}_{model} and reward dynamics \mathcal{R}_{model} . In line 7 the algorithm computes an optimal policy $\hat{\pi}_*$ for \hat{M} by running an exact planning algorithm. In lines 8-9 the algorithm selects an action from $\hat{\pi}_*$ to observe a next state s' , reward r and reward token index j from the environment. In lines 10-15 the algorithm calls *UpdateTree1* to update the appropriate binary trees for every $o.\alpha'$ in s' . In lines 16-22 the algorithm calls *UpdateTree2* to update the appropriate binary trees for every reward token index $i \in [1 : L]$. In lines 23-25 the algorithm updates \hat{U}_j with r if \hat{U}_j is not yet KWIK-learned. The loop ends after line 26 by setting the current state s to the next state s' for the start of the next iteration.

6 Theory

In this section we present a proof that the learning algorithm presented in section 3.3 KIWK-learns the transition dynamics for action a and attribute $C.\alpha$, as well as derive the corresponding KWIK-bound. The proof is analogous to that of Propositional OO-MDPs [5, 4].

Algorithm 7: $DOORMAX_D$: the $DOORMAX_D$ algorithm for episodic tasks.

Input: $\mathcal{C}, L, \mathcal{S}_O, \mathcal{A}, \gamma, r_{max}$

```

1 DefineGlobal()
2 repeat
3   start episode at initial state  $s$ 
4   while  $s \notin \mathcal{S}_{terminal}$  do
5      $(\mathcal{P}_{model}, \mathcal{R}_{model}) \leftarrow BuildFullModels()$ 
6     define an MDP  $\hat{M} = (\mathcal{S}_O, \mathcal{A}, \mathcal{P}_{model}, \mathcal{R}_{model}, \gamma)$ 
7     compute an optimal policy  $\hat{\pi}_*$  for  $\hat{M}$  using exact planning algorithm
8     choose action  $a \in \mathcal{A}(s)$  from  $\hat{\pi}_*$ 
9     observe some next state  $s'$  and reward  $r$  generated from reward token  $j$ 
10    foreach object  $o$  in  $s$  do
11      foreach attribute  $C.\alpha \in Att(C)$  from the class  $C$  of object  $o$  do
12         $o.\alpha' \leftarrow$  value of attribute  $C.\alpha$  for object  $o$  in  $s'$ 
13         $UpdateTree1(C.\alpha, s, o, a, o.\alpha')$ 
14      end
15    end
16    foreach  $i \in [1 : L]$  do
17       $z \leftarrow 0$ 
18      if  $i = j$  then
19         $z \leftarrow 1$ 
20      end
21       $UpdateTree2(i, s, a, s', z)$ 
22    end
23    if  $\hat{U}_j$  with KWIK-bound  $B_j$  has not been KWIK-learned then
24      update  $\hat{U}_j$  with  $r$ 
25    end
26     $s \leftarrow s'$ 
27  end
28 until forever

```

Theorem 1 (KWIK-bound under $DOORMAX_D$) For action a and attribute $C.\alpha$ let \hat{F} be a set of size D that contains hypothesised deictic predicate preconditions on which the transition dynamics of $C.\alpha$ when taking action a may depend. Let $\mathcal{G} = \{g_i\}_{i=1}^N$ be a set of size N that contains hypothesised effect types where each $e \in g_i.E$ has domain $Dom(C.\alpha)$. Let $K_0 = \max_{g \in \mathcal{G}} K_0^g$ and $K_1 = \max_{g \in \mathcal{G}} K_1^g$. Let $\mathcal{H} = \{Tree(g, \hat{F}, M) | g \in \mathcal{G}\}$ for some constant $M > 1$. Then if exactly one $h^* \in \mathcal{H}$ is true, the transition dynamics for a and $C.\alpha$ can be KWIK-learned under the $DOORMAX_D$ algorithm with KWIK-bound $N(K_0M + K_1(D + 1) + 1) + N - 1$.

Proof Consider an effect type $g \in \mathcal{G}$. If this is the correct effect type then it can be learned with KWIK bounds $K_0^gM + K_1^g(D + 1)$. This is because the K_1^g conjunctive effects require at most $D + 1$ observations each to learn the terms they depend on while the terms for the K_0^g disjunctive effects can be memorised M times each. If we then consider all N effect types in \mathcal{G} , an upper bound on the number of observations required so that all effect types are either removed or return some prediction is $N(K_0M + K_1(D + 1) + 1)$.

Now when we call *Predict1* if two effect types provide a prediction that is not the same we remove one of them on the subsequent run of the *UpdateTree1* procedure and this can occur at most $N - 1$ times. This gives a total KWIK-bound of $N(K_0M + K_1(D + 1) + 1) + N - 1$.

An analogous proof can be constructed for the learning algorithm of reward dynamics presented in section 4.4 showing that if a binary value occurs at most once in the tree for

action a and index i mapped to z_i , then that branch can be learned with KWIK-bound $D_{a,i} + 1$.

7 Experiments

In this section we conduct experiments to illustrate the **benefits of the object-oriented frameworks presented in this paper**. In section 7.1 we conduct experiments on the All-Passenger Any-Destination Taxi domain to efficiently learn transition dynamics, under the assumption of known reward dynamics. In section 7.2 we conduct experiments on the All-Passenger Any-Destination Taxi domain to efficiently learn reward dynamics, under the assumption of known transition dynamics. In section 7.3 we demonstrate that the reward dynamics and transition dynamics can be efficiently learned together for the Sokoban domain.

7.1 All-Passenger Any-Destination Taxi Domain: Learning Transition Dynamics

We conduct two sets of experiments on the All-Passenger Any-Destination Taxi domain under the assumption that reward dynamics are known while transition dynamics need to be learned. In the first set of experiments we have one destination and we fix the number of passengers, n . We generate a grounded MDP with an initial state by randomly sampling n passenger locations and one destination location from one of six pre-specified locations and we also sample a random taxi start location together with one of four wall configurations as shown in figure 3a.

We apply 20 independent runs of the following procedure: we sample 10 test MDPs with random initial states. We then randomly sample a training MDP and run $DOORMAX_D$ on it for one episode until we reach a terminal state. Upon termination, we test performance by running $DOORMAX_D$ for one episode on each of the 10 test MDPs, stopping an episode early if we exceed 500 steps. We repeat this for 100 training MDPs. Since all the MDPs come from the same schema we can share transition dynamics between our MDPs - but we only update the transition dynamics on training MDPs.

In our experiments we start with $n = 1$ passenger and incrementally increase to $n = 4$ passengers. We run our experiments for Propositional OO-MDPs and two versions of Deictic OO-MDPs. In the first, without transfer, we relearn the transition dynamics for each n while for the second, with transfer, we transfer the previously learned transition dynamics each time we increase n . We report results in figure 7 that averages over the 20 independent runs the average number of steps for the 10 test MDPs with error bars included.

We use the hypothesis space as per table 7 in our experiments, where P_1, P_2, P_3, P_4 and P_5 are propositions as defined as below. This hypothesis space is chosen so as to mimic the experimental setup of Propositional OO-MDPs on the original Taxi domain as closely as possible [5].

- $Touch_N(Taxi, Wall)$: returns 1 if any object of class *Taxi* has an object of class *Wall* one square north of it, and 0 otherwise. (P_1)
- $Touch_E(Taxi, Wall)$: returns 1 if any object of class *Taxi* has an object of class *Wall* one square east of it, and 0 otherwise. (P_2)
- $Touch_S(Taxi, Wall)$: returns 1 if any object of class *Taxi* has an object of class *Wall* one square south of it, and 0 otherwise. (P_3)
- $Touch_W(Taxi, Wall)$: returns 1 if any object of class *Taxi* has an object of class *Wall* one square west of it, and 0 otherwise. (P_4)

- $On(Taxi, Passenger)$: returns 1 if any object of class *Taxi* is on the same square as any object of class *Passenger*, and 0 otherwise. (P_5)
- $On(Taxi, Destination)$: returns 1 if any object of class *Taxi* is on the same square as any object of class *Destination*, and 0 otherwise. ($P_6 = f_9$)
- $InTaxi(Passenger)$: returns 1 if any object of class *Passenger* has its *in-taxi* attribute set to 1, and 0 otherwise. ($P_7 = f_7$)

Action	Attribute	Hypothesis	Conjunctive
Any	<i>Taxi.x</i>	$\{f_1, f_2, f_3, f_4, P_5, P_6, P_7\}$	$\{-1, 0, 1\}$
Any	<i>Taxi.y</i>	$\{f_1, f_2, f_3, f_4, P_5, P_6, P_7\}$	$\{-1, 0, 1\}$
Pickup	<i>Passenger.in-taxi</i>	$\{f_5, f_6, P_1, P_2, P_3, P_4, P_7\}$	$\{1\}$
Dropoff	<i>Passenger.in-taxi</i>	$\{f_8, P_1, P_2, P_3, P_4, P_5, P_6\}$	$\{1\}$
Dropoff	<i>Passenger.at-destination</i>	$\{f_8, P_1, P_2, P_3, P_4, P_5, P_6\}$	$\{1\}$

Table 7: Hypothesised deictic predicates / propositions for each action and attribute as well as conjunctive effects. Any other action and attribute uses the hypothesis $\{P_1, P_2, P_3, P_4, P_5, P_6, P_7\}$ with all effects being conjunctive. Integer attributes use *Rel* effects while Boolean attributes use *SetBool* effects.

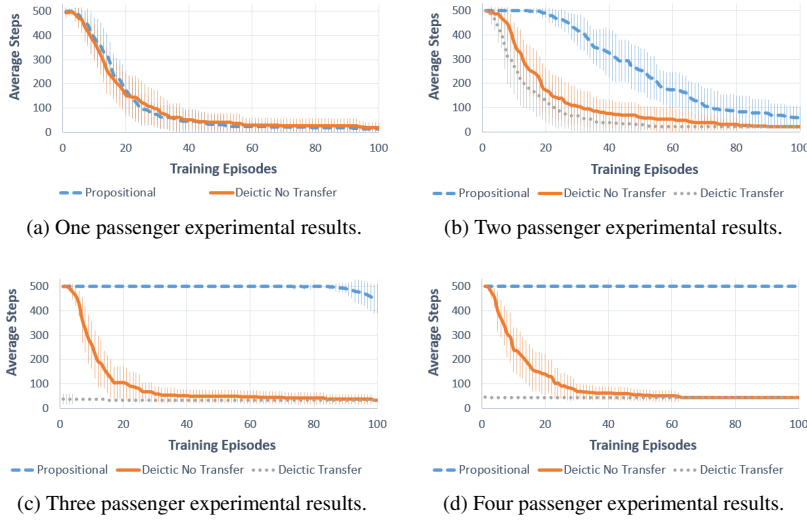


Fig. 7: Experimental results for learning transitions dynamics in the All-Passenger Any-Destination Taxi domain with different number of passengers.

We see from the results that for this domain, Deictic OO-MDPs outperform Propositional OO-MDPs as we increase the number of passengers. This is because with Propositional OO-MDPs we need to add more propositions to the preconditions as we increase the

number of passengers - in fact the propositional representation is unable to learn the task when $n = 4$ even after 100 training episodes.

In particular, the classical Taxi domain with a single passenger and destination requires 7 propositions $\{P_1, P_2, P_3, P_4, P_5, P_6, P_7\}$ [5]. Under a propositional approach, the All-Passenger Any-Destination Taxi Domain requires an additional proposition: *PassengerAtDestination* (P_8) that returns 1 if any object of class *Passenger* is on the same square as any object of class *Destination*. Furthermore, P_5 , P_7 and P_8 need to be grounded for each object of class *Passenger* in order to avoid ambiguity, as further discussed in section 3.2. Therefore, in totality the All-Passenger Any-Destination Taxi Domain requires $5 + 3 \times n$ propositions for a task of n passengers.

Meanwhile, the Deictic OO-MPD framework requires 8 deictic predicates for any task of the domain. Furthermore, as the MDPs belong to the same schema it is beneficial to transfer the previous transition dynamics under the deictic representation. Specifically, once we get to $n = 4$ passengers the transition dynamics we transfer from the $n = 3$ experiment have learned the schema transition dynamics completely and we have zero-shot transfer of the transition dynamics.

We observe from figure 7 that using the deictic representation without transfer is actually learning slightly faster as we add more passengers. This is somewhat misleading. What is actually happening is that as the tasks become more complex, the agent is able to learn more about the transition dynamics over a single training episode, but that episode will require many more steps to complete. To illustrate this, and also highlight the robustness of the deictic representation with transfer methodology, we conduct a second set of experiments. These experiments are similar to those conducted for the first set but we now use a larger 10×10 gridworld with five passengers and three destinations as in figure 3b. In these experiments we stop after 100 steps for each episode of the training MDPs. Furthermore, for the deictic representation with transfer we simply transfer the learned transition dynamics of $n = 4$ passengers and do no additional learning on the new larger gridworld.

In figure 8 we plot for all the experiments run the average number of steps relative to optimal number of steps. We see that the deictic representation with transfer is able to solve the larger gridworld optimally with no additional learning of the transition dynamics. Meanwhile the deictic representation with no transfer, which was decreasing up to $n = 4$, now has a jump between at $n = 5$ because the agent does not have the benefit of learning for a full training episode. We cap the graph’s y axis at 10 to make it more readable, but remark that Propositional OO-MDPs exhibits exponentially worse performance relative to optimal as the tasks become more complex.

To provide more intuition on how efficient learning is achieved with conjunctive effects, consider tables 8 and 9. In these tables we show results from a simulation of running *DOORMAX_D* for action *East*, attribute *Taxi.x* and effect *Rel₁* as well as action *Pickup*, attribute *Passenger.in-taxi* and effect *SetBool₁* respectively.

For example, in table 8 the first time we observe the effect *Rel₁* for action *East* and attribute *Taxi.x* the following assignment to the truth values are observed: $\{f_1 = 0, f_2 = 0, f_3 = 0, f_4 = 1, P_5 = 0, P_6 = 0, P_7 = 0\}$ (row 1). The next time, the following assignment to the truth values are observed: $\{f_1 = 1, f_2 = 0, f_3 = 0, f_4 = 1, P_5 = 0, P_6 = 1, P_7 = 0\}$ (row 2). Since the effect is conjunctive, f_1 and P_6 can be invalidated. As a result, the remaining possible deictic predicates that make up the precondition is $\{f_2, f_3, f_4, P_5, P_7\}$. This process continues until the true precondition is learned in line 7 i.e. $f_2 = 0$. Note that in table 8 we are able to learn the correct precondition with only 7 observations, whereas learning by memorisation would require 2^6 unique observations.

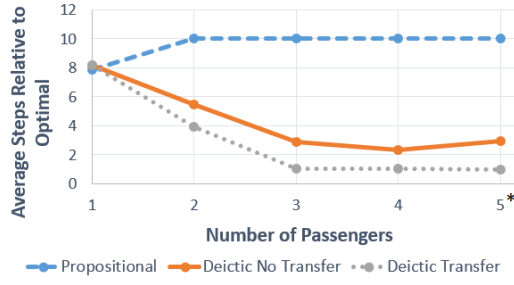


Fig. 8: Average number of steps relative to optimal number of steps as we add more passengers - for $n = 5$ we also increase the gridworld size and add more destinations hence it is marked with a *. The y axis is capped at 12 to make the graph more readable.

observation	f_1	f_2	f_3	f_4	P_5	P_6	P_7
1	0	0	0	1	0	0	0
2	-	0	0	1	0	-	0
3	-	0	0	-	0	-	0
4	-	0	0	-	0	-	-
5	-	0	0	-	0	-	-
6	-	0	0	-	0	-	-
7	-	0	-	-	-	-	-

Table 8: Learning the precondition for action *East*, attribute *Taxi.x* and effect Rel_1 in the All-Passenger Any-Destination Taxi domain.

observation	P_1	P_2	P_3	P_4	f_5	f_6	P_7
1	0	1	1	0	1	0	0
2	-	1	-	0	1	0	0
3	-	-	-	-	1	0	0

Table 9: Learning the precondition for action *Pickup*, attribute *Passenger.in-taxi* and effect $SetBool_1$ in the All-Passenger Any-Destination Taxi domain.

7.2 All-Passenger Any-Destination Taxi Domain: Learning Reward Dynamics

In this section, we run a similar experiment to that of section 7.1, but we now assume known transition dynamics and focus on learning and transfer of reward dynamics. In addition to p_1 , p_2 , p_3 and p_4 defined in section 4.3 we define:

- *AnyWallNorthOfAnyTaxi*(s, s'): returns 1 if any objects of class *Taxi* has an object of class *Wall* one square north of it in s' , and otherwise 0. (p_5)
- *AnyWallEastOfAnyTaxi*(s, s'): returns 1 if any objects of class *Taxi* has an object of class *Wall* one square east of it in s' , and otherwise 0. (p_6)
- *AnyWallSouthOfAnyTaxi*(s, s'): returns 1 if any objects of class *Taxi* has an object of class *Wall* one square south of it in s' , and otherwise 0. (p_7)
- *AnyWallWestOfAnyTaxi*(s, s'): returns 1 if any objects of class *Taxi* has an object of class *Wall* one square west of it in s' , and otherwise 0. (p_8)

Then in our experiments we use the we use the hypothesis space as per table 10.

Action	Reward Token	Hypothesis	Conjunctive
<i>Pickup</i>	\bar{U}_1	$\{p_1, p_2, p_3, p_4, p_5, p_6, p_7, p_8\}$	$\{0\}$
<i>Dropoff</i>	\bar{U}_2	$\{p_1, p_2, p_3, p_4, p_5, p_6, p_7, p_8\}$	$\{0\}$
<i>Dropoff</i>	\bar{U}_3	$\{p_1, p_2, p_3, p_4, p_5, p_6, p_7, p_8\}$	$\{0, 1\}$

Table 10: Hypothesised propositions for each action and reward token as well as their conjunctive indicator values. All other actions and reward token pairs use the hypothesis $\{p_1, p_2, p_3, p_4, p_5, p_6, p_7, p_8\}$ with conjunctive indicators $\{0, 1\}$. The default reward token is \bar{U}_4 .

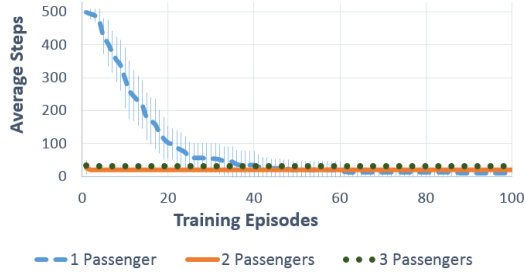


Fig. 9: Experimental results for learning reward dynamics in the All-Passenger Any-Destination Taxi domain with different number of passengers.

In figure 9 we see that, as in the case of learning transition dynamics, we can transfer the model of the reward dynamics between tasks of the domain as more passengers are added. The reward dynamics are completely learned after completing the $n = 2$ passenger MDPs.⁸ From there, the reward dynamics are zero-shot transferred to the $n = 3$ passenger MDPs. In fact, the model of the reward dynamics can be zero-shot transferred to any MDP of the All-Passenger Any-Destination Taxi domain schema after learning on the $n = 2$ passenger MDPs.

7.3 Sokoban Domain

To demonstrate the benefits of Deictic OO-MDPs, we conduct an experiment on a more challenging Sokoban domain. In this experiment we first learn both the transition dynamics and reward dynamics by randomly sampling MDPs with start states as shown in figure 10a and continue learning on these until we have no \perp predictions in these MDPs. These

⁸ While it appears from the figure that the reward dynamics are completely learned after $n = 1$, there is a small amount of learning required for $n = 2$ as the agent needs to learn the dynamics when picking up a passenger while another passenger is already in the taxi. It is not clearly visible on the figure because this learning happens in the first few training episodes and does not have a material impact on the optimal number of steps to solve the test tasks.

simple MDPs, each with approximately 8000 states, are enough to completely learn the schema transition and reward dynamics of this domain under the object-oriented frameworks presented in this paper.

Once learned we zero-shot transfer the transition dynamics to a more complex Sokoban task as shown in 10b. This task comes from the ‘Micro-Cosmos’ level pack and has approximately 10^6 states while the optimal number of steps to solve this task 209. With no additional learning we run value iteration and solve for an optimal policy. Note that the ability to transfer here is critical. The larger MDP has approximately 125 times more states than the toy MDPs. Running R_{max} based algorithms directly on the larger MDP is very slow because at each step it is required to compute a policy with an exact planning algorithm and this has high computational complexity. By transferring the transition dynamics learned in the toy MDPs we can solve the larger MDP with only a *single run* of value iteration.

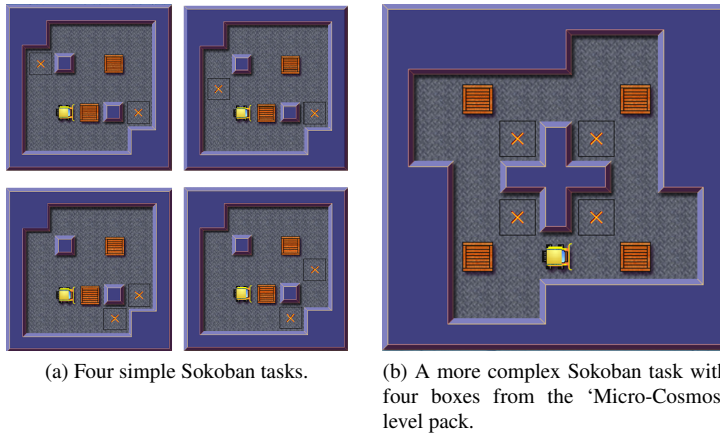


Fig. 10: Training tasks and test task for Sokoban.

8 Conclusion

This paper has introduced and integrated two object-oriented frameworks for efficient model-based RL: Deictic OO-MDPs for transition dynamics and Propositional OO-MDPs for reward dynamics. These frameworks apply to a domain as described by a schema, and so generalise across all instantiated tasks from the schema.

The Deictic OO-MDP framework is based on deictic predicates. A deictic predicate is a predicate that is grounded with respect to a single reference object that relates itself to lifted object classes, and can be used to compactly represent the transition dynamics of domains not possible under a propositional approach. **The Propositional OO-MDP framework extend previous work that focuses only on learning transition dynamics under a propositional approach to the case of learning reward dynamics.**

Both of these frameworks are then combined into a KWIK- R_{max} based algorithm called $DOORMAX_D$ that efficiently learns the full dynamics models. Since the models are schema-based, they transfer across all tasks of the domains. To illustrate the benefits of our proposed

frameworks we run experiments on a modified version of the Taxi domain, the All-Passenger Any-Destination Taxi domain, as well as the Sokoban domain. In both these domains we illustrate that the dynamics models can be learned on a simple set of source task from the domain and then zero-shot transferred to a more complex target task of the domain for efficient RL.

9 Declarations

9.1 Funding

Not applicable.

9.2 Conflicts of interest/Competing interests

Not applicable.

9.3 Availability of data and material

Not applicable.

9.4 Code availability

An implementation of the Deictic OO-MDP framework can be found here.

References

1. Boutilier, C., Reiter, R., Price, B.: Symbolic dynamic programming for first-order mdps. Proceedings of the 17th International Joint Conference on Artificial Intelligence pp. 690–697 (2001)
2. Brafman, R.I., Tennenholtz, M.: R-max - a general polynomial time algorithm for near-optimal reinforcement learning. Journal of Machine Learning Research **3**, 213–231 (2002)
3. Dietterich, T.: The maxq method for hierarchical reinforcement learning. Proceedings of the 15th International Conference on Machine Learning (1998)
4. Diuk, C.: An object-oriented representation for efficient reinforcement learning. Ph.D. thesis, Rutgers The State University of New Jersey-New Brunswick (2010)
5. Diuk, C., Cohen, A., Littman, M.L.: An object-oriented representation for efficient reinforcement learning. Proceedings of the 25th International Conference on Machine learning pp. 240–247 (2008)
6. Džeroski, S., DeRaedt, L., Driessens, K.: Relational reinforcement learning. Machine Learning Journal **43**, 7–52 (2001)
7. Guestrin, C.: Planning under uncertainty in complex structured environments. Ph.D. thesis, Stanford University (2003)
8. Guestrin, C., Koller, D., Gearhart, C., Kanodia, N.: Generalizing plans to new environments in relational mdps. Proceedings of the 18th Joint Conference on Artificial Intelligence pp. 1003–1010 (2003)
9. Juozapaitis, Z., Koul, A., Ferm, A., Erwig, M., Doshi-Velez, F.: Explainable reinforcement learning via reward decomposition. Workshop on Explainable Artificial Intelligence at International Joint Conference on Artificial Intelligence pp. 47–53 (2019)
10. Kakade, M.S.: On the sample complexity of reinforcement learning. Ph.D. thesis, Gatsby Computational Neuroscience Unit, University College London (2003)
11. Kanksy, K., Silver, T., Mèly, D.A., Eldawy, M., Lázaro-Gredilla, M., Lou, X., Dorfman, N., Sidor, S., Phoenix, S., George, D.: Schema networks: Zero-shot transfer with a generative causal model of intuitive physics. Proceedings of the 34th International Conference on Machine Learning pp. 1809–1818 (2017)

12. K.Kiguchi, T.Nanayakkara, K.Watanabe, T.Fukuda: Multi-dimensional reinforcement learning using a vector q-net - application to mobile robots. *Internal Journal of Control, Automation and Systems* **1**, 142–148 (2003)
13. Li, L.: A unifying framework for computational reinforcement learning theory. Ph.D. thesis, Rutgers The State University of New Jersey-New Brunswick (2009)
14. Li, L., Littman, M.L., Walsh, T.J.: Knows what it knows: A framework for self-aware learning. *Proceedings of the 25th International Conference on Machine Learning* pp. 568–575 (2008)
15. Marom, O., Rosman, B.S.: Zero-shot transfer with deictic object-oriented representation in reinforcement learning. *Proceedings of the 32nd Conference on Neural Information Processing Systems* (2018)
16. Mnih, V., Kavukcuoglu, K., Silver, D., Graves, A., Antonoglou, I., Wierstra, D., Riedmiller, M.: Playing atari with deep reinforcement learning. *CoRR abs/1312.5602* (2013)
17. van Otterlo, M.: A survey of reinforcement learning in relational domains. *CTIT Technical Report Series* pp. 1381–3625 (2005)
18. Perkins, D.N., Salomon, G.: Transfer of learning. *International Encyclopedia of Education* (1992)
19. Scholz, J., Levihn, M., Isbell, C.L., Wingate, D.: A physics-based model prior for object-oriented mdps. *Proceedings of the 31st International Conference on Machine Learning* pp. 1089–1097 (2014)
20. van Seijen, H., Fatemi, M., Romoff, J., Larocche, R., Barnes, T., Tsang, J.: Hybrid reward architecture for reinforcement learning. *CoRR abs/1706.04208* (2017)
21. Silver, D., Huang, A., Maddison, C.J., Guez, A., Sifre, L., van den Driessche, G., Schrittwieser, J., Antonoglou, I., Panneershelvam, V., Lanctot, M., Dieleman, S., Grewe, D., Nham, J., Kalchbrenner, N., Sutskever, I., Lillicrap, T., Leach, M., Kavukcuoglu, K., Graepel, T., Hassabis, D.: Mastering the game of go with deep neural networks and tree search. *Nature* **529**, 484–489 (2016)
22. Silver, D., Hubert, T., Schrittwieser, J., Antonoglou, I., Lai, M., Guez, A., Lanctot, M., Sifre, L., Kumar, D., Graepel, T., Lillicrap, T.P., Simonyan, K., Hassabis, D.: Mastering chess and shogi by self-play with a general reinforcement learning algorithm. *CoRR/abs/1712.01815* (2017)
23. Strehl, A.L.: Model-based reinforcement learning in factored-state mdps. *Proceedings of the 2007 IEEE International Symposium on Approximate Dynamic Programming and Reinforcement Learning* (2007)
24. Strehl, A.L., Diuk, C., Littman, M.L.: Efficient structure learning in factored state mdps. *Proceedings of the 22nd AAAI Conference on Artificial Intelligence* pp. 645–650 (2007)
25. Strehl, A.L., Littman, M.L.: Online linear regression and its application to model-based reinforcement learning. *Proceedings of the 21st Annual Conference on Neural Information Processing Systems* (2008)
26. Sutton, R.S., Barto, A.G.: *Reinforcement Learning: An Introduction*. MIT Press (2018)
27. Taylor, M.E., Stone, P.: Transfer learning for reinforcement learning domains: A survey. *Journal of Machine Learning Research* **10**(1), 1633–1685 (2009)
28. Tesauro, G.J.: Td-gammon, a self-teaching backgammon program, achieves master-level play. *Neural Computation* **6**(2), 215–219 (2004)

1 MLJ Contribution Information Sheet

Our previous work introduced the Deictic OO-MDP [1] framework. Deictic OO-MDPs are more expressive than the previously introduced Propositional OO-MDP framework [2]. This additional expressive power permits both a compact and transferable representation of the transition dynamics for certain classes of domains, that is not possible under a purely propositional approach.

In conjunction, the Propositional OO-MDP framework introduces a provably efficient learning algorithm (under the KWIK protocol [3]) to learn the transition dynamics under the formalism. In our previous work we extend this algorithm to the Deictic OO-MDP setting so that provably efficient learning is possible under our framework as well [1].

Both the Propositional OO-MDP and Deictic OO-MDP frameworks focus on learning transition dynamics, *while assuming known reward dynamics*. In the most general case of model-based reinforcement learning, the agent must learn both these components.

1.1 What is the main claim of the paper? Why is this an important contribution to the machine learning literature?

The main contributions of this paper are as follows:

- We extend the Propositional OO-MDP formalism to incorporate reward dynamics. As in the case of transition dynamics, the reward dynamics under our formalism transfer across all tasks of a domain.
- We extend the previously introduced learning algorithms to incorporate efficient learning of reward dynamics under our proposed formalism.
- We integrate the frameworks as follows: we use the deictic formalism to efficiently learn transition dynamics and the propositional formalism to efficiently learn reward dynamics. To the best of our knowledge, this is the first paper to introduce a *complete model-based framework under object-oriented representations with a provably efficient KWIK-bound*.
- We provide more details on the Deictic OO-MDP framework that were not included in the original paper. This includes: more thorough examples, full proofs and a complete suite of learning algorithms.

These contributions are important because Reinforcement Learning, while appealing, struggles to scale to larger, real-world tasks. Transfer learning is an

important area of research in Reinforcement Learning as it provides a mechanism to accelerate learning in new tasks by learning from previously solved source tasks. In this paper, we show that under an object-oriented approach it is possible to efficiently learn both the transferable transition and reward dynamics models in some simple source tasks of a domain and then zero-shot transfer them to more complex tasks of the domain to greatly improve sample efficiency.

1.2 What is the evidence you provide to support your claim?

Our claim is supported as follows:

- In section 7 we provide a formal proof of our algorithm’s KWIK-bound. This supports the claim that our approach is provably efficient.
- In section 8 we run experiments on an extended version of the Taxi and Sokoban domains. In both these domains we show empirically that we can efficiently learn the transition and reward dynamics models on some simple source tasks of the domain and then zero-shot transfer them to more complex tasks of the domain to improve learning performance. This supports the claim that we can learn and transfer both transition and reward dynamics under our proposed integrated framework.

1.3 What papers by other authors make the most closely related contributions, and how is your paper related to them?

Our work builds on the foundations of the research that introduces the Propositional OO-MDP framework [2]. Our previous research introduced the Deictic OO-MDP framework [1] that extended the Propositional OO-MDP framework to be more expressive while still allowing for provably efficient learning.

Both the Propositional OO-MDP and Deictic OO-MDP frameworks focus on learning transition dynamics, while assuming known reward dynamics. In this paper we extend the Propositional OO-MDP framework to incorporate the compact representation and efficient learning of reward dynamics. We then integrate the two frameworks by using a deictic formalism to efficiently learn transition dynamics and a propositional formalism to efficiently learn reward dynamics.

1.4 Have you published parts of your paper before, for instance in a conference?

Yes. Deictic OO-MDPs have previously been introduced [1]. The research that introduces Deictic OO-MDPs focuses on learning transition dynamics while assuming known reward dynamics.

In this paper we introduce a new formalism to efficiently learn reward dynamics under a propositional approach. We then integrate the two frameworks to efficiently learn transition dynamics under a deictic formalism and learn reward dynamics under a propositional formalism.

In addition, we provide a considerably more detailed coverage of Deictic OO-MDPs than was included in the first paper. These include: more thorough treatment of examples, tables detailing complete dynamics under our formalism, formal proofs and a complete suite of learning algorithms.

References

- [1] O. Marom and B. S. Rosman. Zero-shot transfer with deictic object-oriented representation in reinforcement learning. *Proceedings of the 32nd Conference on Neural Information Processing Systems*, 2018.
- [2] C. Diuk, A. Cohen, and M. L. Littman. An object-oriented representation for efficient reinforcement learning. *Proceedings of the 25th International Conference on Machine learning*, pages 240–247, 2008.
- [3] L. Li, M. L. Littman, and T. J. Walsh. Knows what it knows: A framework for self-aware learning. *Proceedings of the 25th International Conference on Machine Learning*, pages 568–575, 2008.