

Lifted Model Checking for Relational MDPs

Wen-Chi Yang¹ wenchi.yang@cs.kuleuven.be
Jean-François Raskin² jean-francois.raskin@ulb.be
Luc De Raedt^{1,3} luc.deraedt@cs.kuleuven.be

Received: date / Accepted: date

Abstract Probabilistic model checking has been developed for verifying systems that have stochastic and nondeterministic behavior. Given a probabilistic system, a probabilistic model checker takes a property and checks whether or not the property holds in that system. For this reason, probabilistic model checking provide rigorous guarantees. So far, however, probabilistic model checking has focused on propositional models where a state is represented by a symbol. On the other hand, it is commonly required to make relational abstractions in planning and reinforcement learning. Various frameworks handle relational domains, for instance, STRIPS planning and relational Markov Decision Processes. Using propositional model checking in relational settings requires one to ground the model, which leads to the well known state explosion problem and intractability. We present pCTL-REBEL, a lifted model checking approach for verifying pCTL properties of relational MDPs. It extends REBEL, a relational model-based reinforcement learning technique, toward relational pCTL model checking. PCTL-REBEL is lifted, which means that rather than grounding, the model exploits symmetries to reason about a group of objects as a whole at the relational level. Theoretically, we show that pCTL model checking is decidable for relational MDPs that have a possibly infinite domain, provided that the states have a bounded size. Practically, we contribute algorithms and an implementation of lifted relational model checking, and we show that the lifted approach improves the scalability of the model checking approach.

Keywords model checking · probabilistic computation tree logic (pCTL) · first-order logic · lifted inference · relational MDPs

¹Department of Computer Science, KU Leuven, Celestijnenlaan 200a - box 2402, 3001 Leuven, Belgium

²Université libre de Bruxelles - Campus de la Plaine, CP212 - 1050 Bruxelles - Belgium

³Centre for Applied Autonomous Sensor Systems, Örebro University, Sweden.

1 Introduction

Probabilistic model checking aims at deciding whether a stochastic model satisfies a given probabilistic property (Baier and Katoen, 2008; Forejt et al., 2011). By doing so, it provides rigorous guarantees about the model. Markov Decision Processes (MDPs) and Probabilistic Computational Tree Logic (pCTL) (Forejt et al., 2011; Kwiatkowska et al., 2011; Dehnert et al., 2017) are standard formalisms for specifying the model and the properties, respectively. MDPs are commonly used for modeling sequential decision making problems where the actions have stochastic effects. PCTL is a temporal logic that expresses model properties over time and allows for probabilistic quantification. A property can be *the machine gives a warning before shutting down with a probability higher than 0.95*, or *the probability is higher than 0.9 that the emergency power supply, after giving a warning, continues to function for at least 10 more minutes*.

It is common in planning and reinforcement learning to make abstraction of the domain elements in order to compactly define models and speed up the computation. However, model checking methods most often operate on explicit-state MDPs (Baier and Katoen, 2008) thus do not allow for such abstractions. This is undesirable since the number of states increases exponentially with the domain size, making it infeasible to explicitly traverse the state space (Slaney and Thiébaux, 2001). In such large domains, it is impractical to apply model checking techniques as they lead to a state explosion (Otterlo, 2004). For instance, the well known blocks world has 501 states for 5 blocks but over 58 million states for only 10 blocks (Slaney and Thiébaux, 2001). In this paper, we aim at mitigating such state explosions in probabilistic model checking by making relational abstractions and using lifted inference.

Lifting is the key to scalability in relational domains (Kersting, 2012; Van den Broeck et al., 2011; de Salvo Braz et al., 2007). It is also central to statistical relational AI (StarAI) (De Raedt et al., 2016). Lifting implies reasoning about a group of objects as a whole at the first-order level, and exploiting the shared relational structures and symmetries in the model. This is done by making abstraction of irrelevant details of the objects. As an illustration, an object’s full identity (e.g. a block’s ID number) can be left out as long as it satisfies the property description (e.g. blue). There has been a significant interest in such relational representations in reinforcement learning and planning. For instance, Džeroski et al. (2001) introduced *Relational Markov Decision Processes (RMDP)*, a first-order generalization of MDPs that succinctly formulates relational models by implicitly defining states in terms of objects and relations (Džeroski et al., 2001; Otterlo, 2004). RMDPs have been used in reinforcement learning and planning to compute first-order policies without explicitly constructing the underlying state space (Džeroski et al., 2001; Kersting et al., 2004; Wang et al., 2008; Kersting and De Raedt, 2004; Driessens and Džeroski, 2004; Boutilier et al., 2001; Sanner and Boutilier, 2009; Yoon et al., 2012). One especially interesting example is *REBEL* (Kersting et al., 2004), the RElational BELlman operator, which we will extend in this paper. REBEL is a model-based reinforcement learning technique for constructing an optimal policy in a given RMDP. It is also a lifted inference technique that alleviates state explosions.

Motivated by the success of temporal logics and MDPs in probabilistic model checking and in planning, we investigate whether it is possible to lift these ap-

proaches to RMDPs. We show that the answer is positive by introducing *pCTL-REBEL*, a new framework that augments REBEL (Kersting et al., 2004) with pCTL. More specifically, pCTL-REBEL is a relational model checking approach that checks relational pCTL formulae in RMDPs. In addition to mitigating state explosions, it is important to take one step further to investigate lifted probabilistic model checking for infinite models. Although model checking for infinite models is generally undecidable (Gabbay, 2003), a rich body of research has discussed the *state-boundedness* assumption that yields decidable verification of infinite systems (Bagheri Hariri et al., 2013; Belardinelli et al., 2011, 2012; De Giacomo et al., 2012; Calvanese et al., 2018). These studies almost exclusively focus on non-probabilistic actions. Nevertheless, they provide great insight into relational model checking with pCTL properties. We extend the work of Belardinelli et al. (2012) to the probabilistic setting to obtain the decidability of the model checking problem for a subclass of infinite MDPs.

The key contribution of this paper is twofold. First, we introduce a lifted model checking algorithm, pCTL-REBEL, to mitigate the state explosion problem of checking relational MDPs. PCTL-REBEL is fully automated and provides a complete, lifted solution for relational model checking. In order to adapt to the model checking framework, pCTL-REBEL introduces an alternative interpretation of REBEL (Kersting et al., 2004) in which a state value corresponds to the *probability* that a given formula is satisfied by executions starting from that state. Second, while model checking is generally undecidable for infinite MDPs, we provide decidability results for a class of infinite MDPs under the state-boundedness condition. In particular, we prove that a finite relational abstraction exists for RMDPs that have an infinitely large domain, and that checking the relational abstraction is equivalent to checking the infinite MDP. This means that strong guarantees for such infinite MDPs can be provided via lifted model checking, as implemented in pCTL-REBEL.

This paper is structured as follows. Section 2 provides an overview on basic notions of relational MDPs. Section 3 reviews basic notions of model checking and introduces relational pCTL. Section 3.3 defines the problem statement of this paper, that is, probabilistic model checking for relational MDPs. Section 4 defines a relational Bellman update operator (a generalization of REBEL (Kersting et al., 2004)) for relational value iteration. Based on Section 4, Section 5 introduces the main algorithm, pCTL-REBEL, for relational model checking. Section 6 provides theoretic results about obtaining decidability for a subclass of infinite MDPs. Section 7 reports on experimental evaluation. Section 8 discusses related work and Section 9 highlights the link between relational model checking and safe reinforcement learning. Finally, Section 10 concludes the work.

2 Relational Notions and Relational MDPs

This section defines the basic notions of relational MDPs. These notions will be used in the remainder of this paper with the *blocks world* running example. We closely follow the notions of the standard first-order logic (Nienhuys-Cheng and Wolf, 1997) and relational MDPs (Kersting et al., 2004).

2.1 Relational Logic

Relational logic generalizes propositional logic with *variables* such that a variable represents a set of constants. This section provides an overview of relational logic, following the notions of Nienhuys-Cheng and Wolf (1997).

An *alphabet* is a tuple $\Sigma = \langle R, D \rangle$ where R is a finite set of relation symbols and D is a possibly infinite set of constants. Each relation symbol $p \in R$ has an arity $m \geq 0$. An *atom* $p(\tau_1, \dots, \tau_m)$ is a relation symbol p followed by an m -tuple of terms τ_i . A *term* is a variable W or a constant c . A variable (resp. constant) is expressed by a string that starts with an upper (resp. lower) case letter. A conjunction is a set of atoms, which is implicitly assumed to be *existentially quantified*. A definite clause $H \leftarrow B$ consists of an atom H and a conjunction B , stating that *H is true if B is true*. Given an expression E , $vars(E)$ (resp. $consts(E)$, $terms(E)$) denotes the set of all variables (resp. constants, terms) in E . An expression is called *ground* if it contains no variables. We shall make the **unique name assumption**, that states all constants are unequal, that is, $c_1 \neq c_2$ holds for different constants c_1 and c_2 . A substitution θ is a set of bindings $\{W_1/\tau_1, \dots, W_n/\tau_n\}$ that assigns terms τ_i to variables W_i . A grounding substitution θ assigns constants to variables in an expression E such that $E\theta$ contains no variables.

We shall use the **Object Identity subsumption** framework (OI-subsumption) (Ferilli et al., 2002), which means that any two terms in an expression are unequal, and the pairwise inequalities should be added. For instance, under OI-subsumption, the conjunction $\{c1(\mathbf{b}), on(\mathbf{X}, \mathbf{Y})\}$ denotes the expression $\{c1(\mathbf{b}), on(\mathbf{X}, \mathbf{Y}), \mathbf{X} \neq \mathbf{Y}, \mathbf{X} \neq \mathbf{b}, \mathbf{Y} \neq \mathbf{b}\}$. For ease of writing, when the context is clear, we shall not write these inequalities explicitly. A conjunction A is *OI-subsumed* by conjunction B , denoted by $A \preceq_{\theta} B$, if there exists a substitution θ such that $B\theta \subseteq A$. Only substitutions that satisfy the inequality constraints are allowed.

A *unifier* θ of two conjunctions A and B under OI-subsumption is a substitution such that $A\theta = B\theta$. For example, the conjunctions $\{c1(\mathbf{X}), on(\mathbf{y}, \mathbf{Z})\}$ and $\{c1(\mathbf{x}), on(\mathbf{Y}, \mathbf{Z})\}$ have a unifier $\{\mathbf{X}/\mathbf{x}, \mathbf{Y}/\mathbf{y}\}$. A *maximally general specialization (mgs)* of two conjunctions A and B under OI-subsumption is a conjunction that is OI-subsumed by A and B , and is not OI-subsumed by any other specializations. The mgs operation is not always unique under OI-subsumption. For example, the conjunctions $\{c1(\mathbf{X})\}$ and $\{on(\mathbf{Y}, \mathbf{Z})\}$ have maximally general specializations $\{c1(\mathbf{X}), on(\mathbf{X}, \mathbf{Y})\}$ and $\{c1(\mathbf{X}), on(\mathbf{Y}, \mathbf{Z})\}$ that do not OI-subsume one another.

The *Herbrand base* of an alphabet $\Sigma = \langle R, D \rangle$, denoted by HB^{Σ} , is the set of all ground atoms that can be constructed from Σ . A *Herbrand interpretation* s is a subset of HB^{Σ} where all atoms in s are *true* and all others are *false*. The set of all Herbrand interpretations determined by Σ is denoted by S^{Σ} . We shall write S instead of S^{Σ} when the context is clear. When the domain D is infinite, the set of all Herbrand interpretations S^{Σ} must be infinite.

2.2 Relational MDP

Relational MDPs (RMDPs) generalize explicit-state MDPs in a twofold manner. First, RMDPs have structured states. More specifically, an RMDP state is represented by a conjunction of ground atoms whereas an explicit state is represented by a single constant. Second, RMDPs allows variables in the state description. In

consequence, a set of RMDP states can be represented by one single *abstract state*, which enables reasoning about a set of states as a whole. In this paper, an RMDP is a variant of the standard RMDP (Kersting et al., 2004; Boutilier et al., 2001) that varies by allowing the domain to be infinite. This section formally defines RMDPs, following the notions of Kersting et al. (2004).

An RMDP is a tuple $K = \langle \Sigma, \Delta \rangle$ where the alphabet $\Sigma = \langle R, D \rangle$ contains a set of relations and a domain, and Δ is a finite set of abstract transitions. The alphabet Σ determines the state space. A state $s \in S^\Sigma$ is a Herbrand interpretation. An **abstract state** s' is then a conjunction of atoms, representing a set of states, denoted by $s'\Theta = \{s \in S^\Sigma \mid s \preceq_\theta s'\}$.

Example 1 Consider a blocks world with the alphabet $\Sigma = \langle R, D \rangle$ where the relations are $R = \{\text{cl}/1, \text{on}/2\}$ and the domain is $D = \{\mathbf{a}, \mathbf{b}, \mathbf{c}\}$, the abstract state $s = \{\text{cl}(\mathbf{A}), \text{cl}(\mathbf{C}), \text{on}(\mathbf{A}, \mathbf{B})\}$ represents the following six ground states.

$$\begin{aligned} s\Theta = \{ & \{\text{cl}(\mathbf{a}), \text{cl}(\mathbf{c}), \text{on}(\mathbf{a}, \mathbf{b})\}, \{\text{cl}(\mathbf{a}), \text{cl}(\mathbf{b}), \text{on}(\mathbf{a}, \mathbf{c})\}, \\ & \{\text{cl}(\mathbf{b}), \text{cl}(\mathbf{c}), \text{on}(\mathbf{b}, \mathbf{a})\}, \{\text{cl}(\mathbf{b}), \text{cl}(\mathbf{a}), \text{on}(\mathbf{b}, \mathbf{c})\}, \\ & \{\text{cl}(\mathbf{c}), \text{cl}(\mathbf{b}), \text{on}(\mathbf{c}, \mathbf{a})\}, \{\text{cl}(\mathbf{c}), \text{cl}(\mathbf{a}), \text{on}(\mathbf{c}, \mathbf{b})\} \} \end{aligned}$$

An **abstract action** $\alpha \notin R$ is an atom for an action relation. An **abstract transition** $\delta \in \Delta$ (based on an abstract action α) is a finite set of probabilistic transition rules $\delta = \{H_1 \xleftarrow{p_1:\alpha} B, \dots, H_n \xleftarrow{p_n:\alpha} B\}$ where B (resp. H_i) is an abstract state, representing the precondition (resp. postcondition), and $p_i \in [0, 1]$ is the transition probability. These transition rules δ denote a proper probability distribution over H_i , that is, $\sum_{i=1}^n p_i = 1$. To ensure that all abstract transitions rely only on information in the current state, we assume that all variables in H_i are also in B , that is, $\text{vars}(H_i) \subseteq \text{vars}(B)$.

Example 2 A blocks world is defined by an RMDP $K = \langle \Sigma, \Delta \rangle$ where $\Sigma = \langle R, D \rangle$, $R = \{\text{cl}/1, \text{on}/2\}$ and Δ contains the following abstract transition δ_{move} .

$$\delta_{\text{move}} \begin{cases} \text{cl}(\mathbf{A}), \text{cl}(\mathbf{C}), \text{on}(\mathbf{A}, \mathbf{B}) \xleftarrow{0.9:\text{move}(\mathbf{A}, \mathbf{B}, \mathbf{C})} \text{cl}(\mathbf{A}), \text{cl}(\mathbf{B}), \text{on}(\mathbf{A}, \mathbf{C}) \\ \text{cl}(\mathbf{A}), \text{cl}(\mathbf{B}), \text{on}(\mathbf{A}, \mathbf{C}) \xleftarrow{0.1:\text{move}(\mathbf{A}, \mathbf{B}, \mathbf{C})} \text{cl}(\mathbf{A}), \text{cl}(\mathbf{B}), \text{on}(\mathbf{A}, \mathbf{C}) \end{cases}$$

The abstract action $\text{move}(\mathbf{A}, \mathbf{B}, \mathbf{C})$ expresses *moving block A to block B from block C*. The action succeeds with probability 0.9 and fails with probability 0.1. When the action fails, the state stays the same. A graphical illustration is in Figure 1.

2.3 Grounding an RMDP

The semantics of an RMDP is defined at the ground level such that any RMDP (including the infinite ones) implicitly defines an underlying ground MDP. This section formally defines the construction of the underlying ground MDP. In the end of this section, we will briefly discuss the decidability issue of infinite RMDPs.

Given an RMDP $K = \langle \Sigma, \Delta \rangle$, the underlying ground MDP is a tuple $M = \langle S, A, T \rangle$ where S is a set of ground states, A is a set of ground actions and $T : S \times A \times S \rightarrow [0, 1]$ is a ground transition function. Every **ground state** $s \in S$ is a Herbrand interpretation of Σ . A ground state s has a set of available

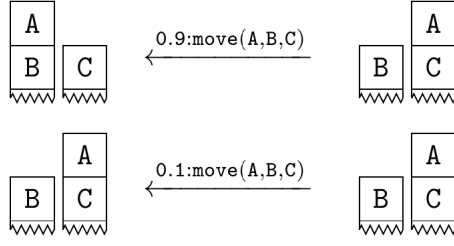


Fig. 1: The abstract transition δ_{move} moves block A to block B from block C with probability 0.9. It fails to move the block with probability 0.1.

ground actions, denoted by $A(s) \subseteq A$. These ground actions $A(s)$ are defined by grounding the abstract actions. Formally,

$$A(s) := \{\alpha\theta \mid H_i \xleftarrow{p_i:\alpha} B \in \Delta, s \preceq_\theta B\}$$

Similarly, given a ground state s and a ground action $\alpha\theta \in A(s)$, the set of **ground transitions** $T(s, \alpha\theta)$ are defined by grounding the abstract transitions. Formally,

$$T(s, \alpha\theta) := \{h_i \xleftarrow{p_i:\alpha\theta} s \mid H_i \xleftarrow{p_i:\alpha} B \in \Delta, s \preceq_\theta B, h_i = (s \setminus B\theta) \cup H_i\theta\}$$

After taking action $\alpha\theta$ in state s , the transition probability distribution over all ground states $s' \in S$ is then

$$T(s, \alpha\theta)(s') := p_i, \text{ where } s' = (s \setminus B\theta) \cup H_i\theta$$

Since Δ is a proper abstract transition function, T must be a proper probability distribution, i.e. $\sum_{s' \in S} T(s, \alpha\theta)(s') = 1$.

Example 3 Consider a blocks world defined by an RMDP $K = \langle \Sigma, \Delta \rangle$ where $\Sigma = \langle R, D \rangle$, $R = \{\text{c1}/1, \text{on}/2\}$, $D = \{\mathbf{a}, \mathbf{b}, \mathbf{c}, \mathbf{d}, \mathbf{e}\}$ and $\Delta = \{\delta_{\text{move}}\}$ (see Figure 1). The RMDP K defines the underlying MDP $M = \langle S, A, T \rangle$. One of the resulting ground transitions is as follows. Let a ground state be $s = \{\text{c1}(\mathbf{a}), \text{c1}(\mathbf{b}), \text{c1}(\mathbf{d}), \text{on}(\mathbf{a}, \mathbf{c}), \text{on}(\mathbf{d}, \mathbf{e})\} \in S$, and let an action be $\text{move}(\mathbf{a}, \mathbf{b}, \mathbf{c}) \in A(s)$, the resulting next state must be $s' = \{\text{c1}(\mathbf{a}), \text{c1}(\mathbf{c}), \text{c1}(\mathbf{d}), \text{on}(\mathbf{a}, \mathbf{b}), \text{on}(\mathbf{d}, \mathbf{e})\} \in S$. By taking the $\text{move}(\mathbf{a}, \mathbf{b}, \mathbf{c})$ action in state s , the probability of reaching s' is 0.9, and the probability of staying in s is 0.1, as illustrated in Figure 2.

In general, for an RMDP that has an infinite domain, the underlying ground MDP has an infinitely large state space and action space. Hence, it is infeasible to explicitly traverse the state space. Moreover, such RMDPs have an unbounded branching behavior such that a state has infinitely many available actions, leading to infinitely many other states. For example, one can move a clear block to any of the (infinitely many) other clear blocks. Therefore, the model checking problem for infinite RMDPs is generally undecidable. To obtain decidability, we will later identify a special class of infinite RMDPs such that the branching behavior is bounded. More details are in Section 6.

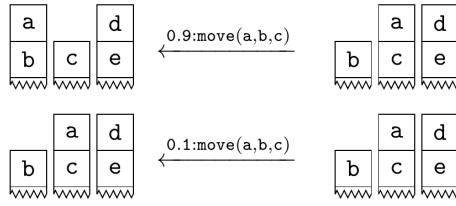


Fig. 2: The ground action $\text{move}(a, b, c)$ moves block a to block b from block c with probability 0.9. The action fails with probability 0.1.

3 Model Checking for Relational MDPs

This section defines the problem statement of this paper, namely the *model checking problem of relational MDPs*. More specifically, Section 3.1 reviews the fundamentals of model checking, Section 3.2 defines the relational pCTL language that will be used to specify properties throughout this paper, and Section 3.3 defines the main problem statement. Later, Section 4 and Section 5 will provide a solution for the problem in Section 3.3.

3.1 Probabilistic Reachability

This section defines the *probabilistic reachability property*, the most fundamental property in model checking. The probabilistic reachability refers to the maximum probability of reaching a set of goal states from a given initial state. This property is required to define the relational pCTL language in Section 3.2. The content of this section is standard in model checking (Baier and Katoen, 2008; Kwiatkowska et al., 2011) and is included to make this paper self-contained.

We first define *paths* in an MDP and their probability assignment. Given an MDP $M = \langle S, A, T \rangle$, a path of length n is denoted by $\rho_n = s_1 \xrightarrow{a_1} s_2 \xrightarrow{a_2} \dots \xrightarrow{a_{n-1}} s_n$ where $s_i \in S$, $a_i \in A(s_i)$ and $T(s_i, a_i)(s_{i+1}) > 0$. Similarly, an infinite path is denoted by $\rho = s_1 \xrightarrow{a_1} s_2 \xrightarrow{a_2} \dots$. The set of all finite (resp. infinite) paths is denoted by $FPath_M$ (resp. $IPath_M$), and the set of all finite (resp. infinite) paths starting from state s is denoted by $FPath_{M,s}$ (resp. $IPath_{M,s}$). We denote all paths starting from s in M as $Path_{M,s} = FPath_{M,s} \cup IPath_{M,s}$. The i -th state of a path ρ is denoted by $\rho(i)$. The last state of a path ρ is denoted by $last(\rho)$. To project a set of paths $Path_{M,s}$ to a probability space, it is required to remove nondeterministic actions by a policy. A policy $\pi : FPath_M \times A \rightarrow [0, 1]$ takes a finite path ρ_n and specifies a probability distribution over all available actions $A(last(\rho_n))$. We write $\pi(\rho_n)$ to denote all possible actions selected by the policy, and we write $\pi(\rho_n, a)$ to denote the probability of action a being selected.

Definition 1 (cf. Forejt et al. (2011)) *Given an MDP M and a policy π , the probability $P_M^\pi(\rho_n)$ of a finite path $\rho_n = s_1 \xrightarrow{a_1} s_2 \xrightarrow{a_2} \dots \xrightarrow{a_{n-1}} s_n \in FPath_{M,s_1}$ is*

inductively defined by

$$\begin{aligned} P_M^\pi(\rho_1) &= 1 \\ P_M^\pi(\rho_k) &= P_M^\pi(\rho_{k-1}) \sum_{a \in \pi(\rho_{k-1})} \pi(\rho_{k-1}, a) T(s_{k-1}, a)(s_k) \end{aligned} \quad (1)$$

where ρ_k denotes ρ_n 's prefix of length k . Let $C_{\rho_n} \in \mathcal{IPath}_{M, s_1}$ be the set of all infinite paths that have a prefix ρ_n (also known as the basic cylinder), the probability $P_M^\pi(C_{\rho_n})$ is then defined as the probability of ρ_n , i.e. $P_M^\pi(C_{\rho_n}) = P_M^\pi(\rho_n)$.

After defining the probability for an MDP path, we can now formally define probabilistic reachability.

Definition 2 (cf. Forejt et al. (2011)) Given an MDP $M = \langle S, A, T \rangle$, an initial state $s \in S$ and a set of goal states $G \subseteq S$, the set of all paths from s to G , denoted by $\text{Path}_{M, s}(G)$, is formally defined as follows.

$$\text{Path}_{M, s}(G) := \{\rho \in \text{Path}_{M, s} \mid \exists i \in \mathbb{N}, \rho(i) \in G\}$$

By following an optimal policy π , one can generate a path from s to G with a maximum probability. Any other policies in the policy space Π does not achieve a probability larger than the optimal policy π does. The maximum reachability is defined as follows.

$$P_M^{\max}(\text{Path}_{M, s}(G)) := \sup_{\pi \in \Pi} P_M^\pi(\text{Path}_{M, s}(G)) \quad (2)$$

In general, $P_M^{\max}(\cdot)$ denotes the maximum reachability in M .

It is known that a deterministic, stationary ϵ -optimal policy suffices to achieve the maximum probabilistic reachability in a finite MDP (Baier and Katoen, 2008). This means the policy will simplify to a function $\pi : S \rightarrow A$ that maps a path's last state to a single action. For an infinite RMDP that has a finite abstraction (cf. Section 6), we assume there exists an optimal deterministic, stationary policy for the maximum probabilistic reachability.

3.2 Relational pCTL

This section introduces *Relational Probabilistic Computational Tree Logic* (relational pCTL), a temporal logic that describes system behavior over time and allows for probabilistic quantification. We will use relational pCTL to specify properties of an RMDP. Relational pCTL is a variant of the standard pCTL (cf. Forejt et al. (2011) and Baier and Katoen (2008)) that varies by allowing variables in atoms and only allowing negations in front of atoms.

The syntax of the relational pCTL is as follows. A property is always specified by a state formula ϕ .

$$\begin{aligned} \text{state formula } \phi &::= \mathbf{true} \mid l \mid \neg l \mid \phi \wedge \phi \mid \phi \vee \phi \mid \mathbb{P}_{\geq p}[\psi] \\ \text{path formula } \psi &::= \mathbf{X} \phi \mid \phi \mathbf{U}^{\leq k} \phi \mid \phi \mathbf{U} \phi \end{aligned}$$

where l is an atom (that can contain variables), \mathbf{p} is a probability such that $0 \leq \mathbf{p} \leq 1$, $\mathbf{k} \in \mathbb{N}$ is a step bound and $\bowtie \in \{\leq, <, \geq, >\}$. Here, relational pCTL generalizes the standard pCTL by letting l be a relational atom instead of a constant.

The semantics of the relational pCTL resembles the standard pCTL (Baier and Katoen, 2008). A state either satisfies or violates a state formula ϕ , resulting in a boolean evaluation for each state. The \mathbf{X} operator stands for *next*, and the \mathbf{U} stands for *until*. A path formula $\mathbf{X}\phi$ is satisfied if ϕ is satisfied in the next state; $\phi_1 \mathbf{U}^{\leq \mathbf{k}} \phi_2$ is satisfied if ϕ_2 is satisfied within \mathbf{k} steps and ϕ_1 holds before then; $\phi_1 \mathbf{U} \phi_2$ is satisfied if ϕ_2 is eventually satisfied and ϕ_1 holds before then.

The semantics of the relational pCTL is defined at the ground level. Given an RMDP $K = \langle \Sigma, \Delta \rangle$ that defines a ground MDP $M = \langle S, A, T \rangle$, we say a ground state $s \in S$ satisfies a state formula ϕ , denoted by $s \models \phi$, if and only if there exists a grounding substitution θ for all free variables in ϕ such that s satisfies ϕ under θ , i.e. $s \models \phi \Leftrightarrow \exists \theta. s \models^\theta \phi$. All substitutions must respect OI-subsumption, i.e. any two terms $\mathbf{t}_1, \mathbf{t}_2$ in a conjunction must be unequal. Formally, the pCTL satisfiability relation \models^θ is inductively defined as follows.

$$\begin{aligned}
s \models^\theta \mathbf{true} & & & \\
s \models^\theta l & & \Leftrightarrow s \leq_\theta l & \\
s \models^\theta \neg l & & \Leftrightarrow s \not\leq_\theta l & \\
s \models^\theta \phi_1 \wedge \phi_2 & & \Leftrightarrow s \models^\theta \phi_1 \wedge s \models^\theta \phi_2 & \\
s \models^\theta \phi_1 \vee \phi_2 & & \Leftrightarrow s \models^\theta \phi_1 \vee s \models^\theta \phi_2 & \\
s \models^\theta \mathbf{P}_{\bowtie \mathbf{p}}[\psi] & & \Leftrightarrow P_M^{max}(\{\rho \in Path_{M,s} \mid \rho \models \psi\}) \bowtie \mathbf{p} &
\end{aligned}$$

where

$$\begin{aligned}
\rho \models \mathbf{X} \phi & & \Leftrightarrow \rho(2) \models \phi & \\
\rho \models \phi_1 \mathbf{U}^{\leq \mathbf{k}} \phi_2 & & \Leftrightarrow \exists i \leq \mathbf{k} + 1. [\rho(i) \models \phi_2 \wedge \forall j < i. \rho(j) \models \phi_1] & \\
\rho \models \phi_1 \mathbf{U} \phi_2 & & \Leftrightarrow \exists i \in \mathbb{N}. [\rho(i) \models \phi_2 \wedge \forall j < i. \rho(j) \models \phi_1] &
\end{aligned}$$

With the aforementioned operators, additional operators can be defined as follows where \mathbf{F} stands for *eventually*.

$$\mathbf{false} \equiv \neg \mathbf{true} \qquad \mathbf{F}^{\leq \mathbf{k}} \phi \equiv \mathbf{true} \mathbf{U}^{\leq \mathbf{k}} \phi \qquad \mathbf{F} \phi \equiv \mathbf{true} \mathbf{U} \phi$$

Example 4 Consider a formula $\phi = \neg \mathbf{cl}(\mathbf{A})$ that states "there exists an unclear block \mathbf{A} in the state". The satisfiability relation $s \models \phi$ can be rewritten as follows.

$$\begin{aligned}
s \models \neg \mathbf{cl}(\mathbf{A}) & \\
\Leftrightarrow \exists \theta. s \models^\theta \neg \mathbf{cl}(\mathbf{A}) & \\
\Leftrightarrow \exists \theta. s \not\leq_\theta \mathbf{cl}(\mathbf{A}) &
\end{aligned}$$

For the ground state $s_1 = \{\mathbf{cl}(\mathbf{a})\}$, the above evaluates to *false* as no θ exists for the satisfiability relation. For another ground state $s_2 = \{\mathbf{cl}(\mathbf{a}), \mathbf{on}(\mathbf{a}, \mathbf{b})\}$, the above evaluates to *true* with $\theta = \{\mathbf{A}/\mathbf{b}\}$.

Example 5 Consider a formula $\phi = P_{\geq 0.9}[\mathbf{X} \text{cl}(\mathbf{A})]$ that states "there exists a block \mathbf{A} that can become clear in the next state with a probability ≥ 0.9 ". The satisfiability relation $s \models \phi$ can be rewritten as follows.

$$\begin{aligned} s &\models P_{\geq 0.9}[\mathbf{X} \text{cl}(\mathbf{A})] \\ \Leftrightarrow \exists \theta. s &\models^\theta P_{\geq 0.9}[\mathbf{X} \text{cl}(\mathbf{A})] \\ \Leftrightarrow \exists \theta. s &\models^\theta P_M^{max}(\{\rho \in \text{Path}_{M,s} \mid \rho(2) \models \text{cl}(\mathbf{A})\}) \geq 0.9 \end{aligned}$$

For the ground state $s = \{\text{cl}(\mathbf{a}), \text{on}(\mathbf{a}, \mathbf{b})\}$, the above evaluates to *true* with $\theta = \{\mathbf{A}/\mathbf{b}\}$ as block \mathbf{b} can be clear after taking the action $\text{move}(\mathbf{a}, \text{fl}, \mathbf{b})$ where fl stands for floor.

It is assumed that the scope of OI-subsumption is within the conjunction. That is, no term inequalities are assumed across different conjunctions. For example, the formula $P_{\geq 0.7}[\mathbf{X} \text{cl}(\mathbf{A})] \wedge P_{\geq 0.95}[\mathbf{F} \text{on}(\mathbf{B}, \mathbf{C})]$ has two conjunctions, and term inequalities such as $\mathbf{A} = \mathbf{B}$ or $\mathbf{A} \neq \mathbf{C}$ do not exist, but $\mathbf{B} \neq \mathbf{C}$ holds under OI-subsumption.

3.3 The Relational Model Checking Problem

This section defines the model checking problem for RMDPs, using the definition of RMDP (Section 2.2) and relational pCTL (Section 3.2). Later, Section 4 and Section 5 will illustrate techniques for solving this model checking problem.

Relational model checking resembles the standard model checking problem (Baier and Katoen, 2008). Given a model M and a pCTL formula ϕ , relational model checking computes *all states in M that satisfy ϕ* , denoted by $\text{Sat}_M(\phi)$. The significance of relational model checking is that it computes $\text{Sat}_M(\phi)$ at a lifted level by using relational states to represent groups of underlying ground states. Hence, relational model checking finds a set of *abstract states* that represents $\text{Sat}_M(\phi)$.

Definition 3 *Given an RMDP $K = \langle \Sigma, \Delta \rangle$ that defines the underlying MDP $M = \langle S, A, T \rangle$ and a relational pCTL formula ϕ , the relational model checking problem is to determine all ground states $\text{Sat}_M(\phi) \subseteq S$ that satisfy ϕ , i.e. $\text{Sat}_M(\phi) = \{s \in S \mid s \models \phi\}$. It does so by finding a set of abstract states $\text{Sat}_K(\phi)$ in K that represents $\text{Sat}_M(\phi)$. Formally,*

$$s \in \text{Sat}_M(\phi) \Leftrightarrow \exists s' \in \text{Sat}_K(\phi). s \preceq_\theta s'$$

In this work, we will solve the relational model checking problem of two types of RMDPs. The first type is the RMDPs that have a finite domain (i.e. finite RMDPs). The second type is a special class of RMDPs that have an infinite domain (i.e. infinite RMDPs).

We discuss the decidability of these two types of RMDPs. The model checking problem for a finite RMDP is decidable. This is because when the domain is finite, the state space is also finite, i.e. the underlying ground MDP is finite. One can thus enumerate all states and collect the states that satisfy the given property. In contrast, an infinite RMDP contains infinitely many states, which makes enumerating all states infeasible. In this case, we focus on a class of infinite RMDPs that have a finite abstraction. More details are given in Section 6.

4 PCTL Relational Bellman Operator

This section defines the **pCTL relational Bellman operator** (pCTL-REBEL), the essential building block for solving the relational model checking problem. Given a pCTL formula $P_{\triangleright p}[\psi]$ and an RMDP, pCTL-REBEL evaluates a function $V^p : S \rightarrow [0, 1]$ that assigns a probability to each RMDP state. A probability $V^p(s)$ represents the probability that state s satisfies the path formula ψ . If the probability is within the bound, i.e. $V^p(s) \bowtie p$, then state s satisfies $P_{\triangleright p}[\psi]$ and belongs to the solution set, i.e. $s \in \text{Sat}_K(P_{\triangleright p}[\psi])$.

At this point, it is important to remark that PCTL-REBEL is a variant of REBEL, but does not consider a reward structure. Indeed, REBEL (Kersting et al., 2004) is a model-based relational reinforcement learning technique that operates on a *reward structure* and computes an optimal policy for reaching a set of goal states, which can be seen as a *reward-based* reachability property. However, since we do not consider a reward structure and are interested in *probabilistic* properties, an alternative interpretation of REBEL is required. Section 4.1 introduces an alternative interpretation of the relational Bellman operator. Based on which, Section 4.2, 4.3 and 4.4 respectively describe in detail the three components of pCTL-REBEL. Section 4.5 then gives an illustration of pCTL-REBEL with an example.

4.1 PCTL Relational Bellman Operator

Given an RMDP and a pCTL formula $P_{\triangleright p}[\psi]$, the task of pCTL-REBEL is to compute a *state probability function* $V^p : S \rightarrow [0, 1]$ that assigns a probability to each state. Similar to the original REBEL, pCTL-REBEL takes an initial state probability function V_0^p and iteratively computes V_1^p, V_2^p , etc for a number of steps, depending on the given formula. When the formula has a *step bound* k , then pCTL-REBEL is applied for k times¹. When the formula is unbounded, pCTL-REBEL is applied for arbitrarily many times until the probabilities converge.

The state probability function V^p is similar to REBEL’s state value function V but interprets state values as probabilities rather than as expected rewards. In a similar way, the state-action probability function Q^p is related to REBEL’s Q function but interprets state-action values as probabilities. More details will be given later. In order to maintain the connection to the original REBEL and to leave room for extending the present model checking approach to incorporate rewards, we use the same notation V and Q for these functions as in REBEL. For clarification, we add a superscript p to denote that V^p and Q^p interpret values as probabilities. We now formally define these functions and the pCTL relational Bellman operator, following the notations of Kersting et al. (2004).

Definition 4 (cf. Kersting et al. (2004)) *A state probability function $V^p : S \rightarrow [0, 1]$ is an ordered set of V^p -rules of the form of $c \leftarrow B$ where B is an abstract state and $c \in [0, 1]$, representing the probability of reaching a goal state from B . The value $V^p(s)$ of a ground state s is assigned by the first rule that subsumes s , i.e. $s \preceq_{\theta} B$.*

¹ A step bound is commonly called a finite horizon in AI.

Given an abstract goal state G , the initial state probability function V_0^p is defined as

$$\begin{aligned} 1.0 &\leftarrow G \\ 0 &\leftarrow \emptyset \end{aligned}$$

The first rule expresses that any ground state subsumed by the goal state G , by definition, satisfies G with probability 1. The second rule expresses that any other states that are not captured by the first rule satisfy G with probability 0. The rule of $0 \leftarrow \emptyset$ ensures that all states are assigned a value. Hence, it is often the last V^p -rule to capture the states that are not captured by any previous rules.

Definition 5 (cf. Kersting et al. (2004)) *A state-action probability function $Q^p : S \times A \rightarrow [0, 1]$ is an ordered set of Q^p -rules of the form of $c : A \leftarrow B$ where A is an abstract action and B is an abstract state, representing the probability of reaching a goal state when A is taken in B . The value $Q^p(s, a)$ of a ground state s and an action a is assigned by the first rule that subsumes s and a , i.e. $s \preceq_\theta B$ and $a \preceq_\theta A$.*

The pCTL-REBEL operator is listed in Equation 3. By iteratively applying pCTL-REBEL, we compute the state probability functions V_1^p, V_2^p , etc. Notice that pCTL-REBEL is a special case of the original REBEL that sets the discount factor to 1, has no reward structure, and connects a single reward 1 to the target condition². As a result, all V^p values are interpreted as probabilities in $[0, 1]$. This alternative interpretation allows to capture the probability that a formula is satisfied, which is essential for adapting the original REBEL framework into a model checking setting.

$$V_{t+1}^p \underbrace{(s)}_{\textcircled{1}} = \max_{a \in A(s)} \underbrace{\sum_{s'} T(s, a, s') V_t^p(s')}_{\textcircled{2}} = \max_{a \in A(s)} Q_{t+1}^p(s, a) \quad \textcircled{3} \quad (3)$$

PCTL-REBEL (Equation 3) is implemented by `OneIteration` (Algorithm 1). This algorithm makes use of the following three components.

- ① **Regression** (Algorithm 2): Deriving the abstract states s in V_{t+1}^p
- ② **Q^p Rules** (Algorithm 3): Computing $Q_{t+1}^p(s, a)$ for all actions $a \in A(s)$
- ③ **V^p Rules** (Algorithm 4): Updating V_{t+1}^p by maximizing over $Q_{t+1}^p(s, a)$.

These three components modify the algorithms of Kersting et al. (2004) to provide support to pCTL operators. The modified parts in the algorithms will be marked **blue**. We now describe these components in detail in Section 4.2, 4.3 and 4.4, respectively.

² Our work addresses a special case of relational model-based reinforcement learning. More details will be given in Section 9.

Algorithm 1: OneIteration. This algorithm implements Equation 3.

Require: V_t^p : state probability function
 ψ : pCTL formula $[S_1 \cup S_2]$ or $[X S_2]$
Return : V_{t+1}^p : the next probability function

- 1 $Q_{t+1}^p := Q^p\text{Rules}(V_t^p, S_1)$
 - 2 $V_{t+1}^p := V^p\text{Rules}(Q_{t+1}^p, \psi)$
-

4.2 Logical Regression

Logical regression is a standard technique that reasons about abstract transitions at the relational level (Boutillier et al., 2001; Kersting et al., 2004; Sanner and Boutillier, 2009). This technique is essential to scalability as it mitigates state explosions by operating on the relational state space instead of the underlying ground state space. Concretely, logical regression searches *backwards* the possible *pre-states* that can reach a given state after taking a number of transitions.

This section extends the standard logical regression to implement Equation 3 ①. The task is to identify the pre-states that can reach a given state by taking *one transition*, which is related to computing pCTL formulae with a step bound 1. These formulae are of the form of $[X S_2]$, $[F^{\leq 1} S_2]$ or $[S_1 \cup^{\leq 1} S_2]$. These three formulae are different in terms of constraints. First, $[F^{\leq 1} S_2]$ is a reachability property that does not impose any constraints. This formula is the basis for the other two formulae. Second, $[X S_2]$ imposes the constraint that " S_2 must be reached after taking *exactly one transition*". Finally, $[S_1 \cup^{\leq 1} S_2]$ is a *constrained reachability property* that imposes the constraint that " S_2 must be reached by going through *only the states where S_1 holds*".

Example 6 To illustrate logical regression, consider an abstract state $\{\text{on}(\mathbf{a}, \mathbf{b})\}$ and the following δ_{move_1} transition rule (also shown in Figure 1).

$$\delta_{\text{move}_1} : \text{cl}(\mathbf{A}), \text{cl}(\mathbf{C}), \text{on}(\mathbf{A}, \mathbf{B}) \xleftarrow{0.9:\text{move}(\mathbf{A}, \mathbf{B}, \mathbf{C})} \text{cl}(\mathbf{A}), \text{cl}(\mathbf{B}), \text{on}(\mathbf{A}, \mathbf{C})$$

By applying the substitution $\theta = \{\mathbf{A}/\mathbf{a}, \mathbf{B}/\mathbf{b}\}$ to δ_{move_1} , we obtain the following rule. This rule describes that $\{\text{on}(\mathbf{a}, \mathbf{b})\}$ can be reached from any ground state that is subsumed by $\{\text{cl}(\mathbf{a}), \text{cl}(\mathbf{b}), \text{on}(\mathbf{a}, \mathbf{C})\}$ after taking the action $\text{move}(\mathbf{a}, \mathbf{b}, \mathbf{C})$.

$$\text{on}(\mathbf{a}, \mathbf{b}) \xleftarrow{0.9:\text{move}(\mathbf{a}, \mathbf{b}, \mathbf{C})} \text{cl}(\mathbf{a}), \text{cl}(\mathbf{b}), \text{on}(\mathbf{a}, \mathbf{C})$$

Logical regression can be applied multiple times. For example, by applying twice the δ_{move_1} transition rule, we obtain all abstract states that can reach $\{\text{on}(\mathbf{a}, \mathbf{b})\}$ within 2 steps, as shown in Figure 3.

Logical regression is implemented in **Regression** (Algorithm 2) that identifies all states that can reach a goal state after taking one transition. **Regression** generalizes the **WEAKESTPRE** algorithm in REBEL (Kersting et al., 2004) in two ways (marked **blue**). First, it provides support to *until* formulae $[S_1 \cup S_2]$ by ensuring that all pre-states satisfy S_1 . Second, it provides support to infinite RMDPs by filtering out the pre-states that exceed a given state bound b . More details will be given in Section 6.3.

³ The state bound $b \in \mathbb{N}$ is for infinite RMDPs. More details in Section 6.

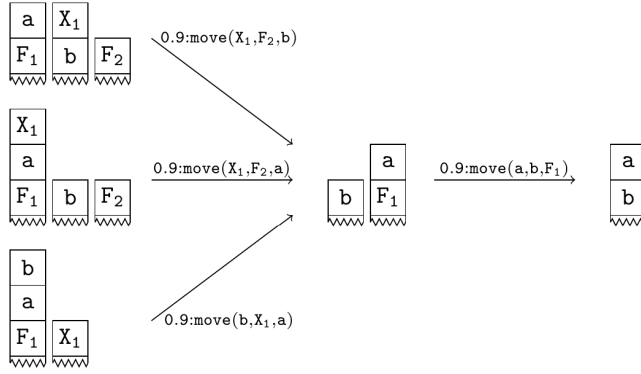


Fig. 3: The figure shows the derived paths by applying two logical regressions on $\{\text{on}(a, b)\}$ with the δ_{move_1} transition rule. The abstract states from left to right reach $\{\text{on}(a, b)\}$ after 2, 1, and 0 steps respectively.

Algorithm 2: Regression. This algorithm implements Equation 3 ①. All returned states must be subsumed by S_1 , and can reach S_2 by following the δ transition rule. All returned states must be legal and respect OI-subsumption.

Require: δ : transition rule $\delta = H_i \xleftarrow{p_i \cdot \alpha} B$
 S_1 : all returned states must be subsumed by S_1
 S_2 : goal state
Return : PreS : states that can reach S_2 by following δ

```

1 PreS :=  $\emptyset$ 
2 foreach  $S'_2 \subseteq S_2$  and  $H'_1 \subseteq H_i$  s.t.  $\theta = \text{unifier}(S'_2, H'_1)$  exists do
3    $S := (S_2 \theta \setminus H'_1 \theta) \cup B \theta$  // identify a pre-state S
4   foreach  $l \in (S_2 \theta \setminus H'_1 \theta)$  and  $l' \in H_i \theta \cup B \theta$  s.t.  $\text{unifier}(l, l')$  exists do
5      $\lfloor$  add  $l \neq l'$  to S // add object inequalities to S
6    $S_{\text{pre}} := \text{mgs}(S, S_1)$  // apply the  $S_1$  constraint
7   foreach  $s_{\text{pre}} \in S_{\text{pre}}$  do
8      $\lfloor$  if  $|\text{terms}(s_{\text{pre}})| \leq b$  then add  $s_{\text{pre}}$  to PreS // apply the state bound3
9    $\rfloor$ 

```

4.3 Q^p -Rules Generation

Given a state probability function V^p , this section computes Q^p probabilities for state-action pairs (cf. Equation 3 ②). $Q^p(s, a)$ denotes the probability that a formula is satisfied by a path that starts from s and has the first action a . Similar to the state probability function V^p , the Q^p function is an ordered set of Q^p -rules. These Q^p -rules are an intermediate representation, from which we can compute the next state probability function (cf. Section 4.4).

Deriving a Q^p function is complex for two reasons. First, we must consider all transition rules that indicate different action postconditions. For example, action `move/3` has two possible outcomes, namely succeeding and failing. As could be

expected, it is not sufficient to consider only the case where the action always succeeds as in Figure 3. We must consider all transition rules in order to correctly calculate the Q^p -rules. Second, an abstract action can *diverge* and produce multiple preconditions. For example, Figure 3 shows three different preconditions are derived by applying different substitutions to the δ_{move_1} transition. It is required to consider all these preconditions. Fortunately, to resolve these issues, we could follow the procedure provided by the original REBEL, with limited adaptations.

$Q^p\text{Rules}$ (Algorithm 3) implements the procedure of computing a set of Q^p -rules (cf. Equation 3 (2)). The process is as follows. First, for each transition rule and current states in V_t^p , $Q^p\text{Rules}$ computes a set of *partial rules* (cf. line 2-5). This procedure results in multiple sets of partial rules, and each set considers one single postcondition of the action. Then, we must combine these partial rule sets to get a complete Q^p -rule set. This is done by unifying the partial rules (cf. line 6-14). Finally, $Q^p\text{Rules}$ returns a set of complete Q^p -rules.

$Q^p\text{Rules}$ modifies the QRULES algorithm in REBEL (Kersting et al., 2004) in two ways, marked **blue**. First, it provides support to the *until* formulae $[S_1 \text{ U } S_2]$ by passing an extra parameter S_1 to Algorithm 2 (cf. Section 4.2). Second, it achieves the probabilistic interpretation of the Bellman operator by discarding the reward component and setting the discount factor to 1.

Algorithm 3: $Q^p\text{Rules}$. This algorithm implements Equation 3 (2). All returned Q^p -rules must be legal and respect OI-subsumption. All states in the returned Q^p -rules must be subsumed by S_1 .

```

Require:  $V_t^p$  : state probability function
            $S_1$  : all states in the returned  $Q^p$ -rules must be subsumed by  $S_1$ 
Return :  $Q^p_{t+1}$  : the next state-action probability function

```

```

1  $Q^p_{t+1} := \emptyset$ 
  /* fix an action  $\alpha$  */
2 foreach  $H_i \stackrel{p_i:\alpha}{\leftarrow} B$  for  $\alpha$  do
3   foreach  $v^p \leftarrow V \in V_t^p$  do
4      $\text{PreS} := \text{Regression}(H_i \stackrel{p_i:\alpha}{\leftarrow} B, S_1, V)$ 
      /* derive a subset of partial rules */
5      $\text{partial}Q^p := \{q^p : \alpha\theta \leftarrow S \mid S \in \text{PreS} \text{ and } q^p = p_i \times v^p\}$ 
6     if  $Q^p_{t+1} = \emptyset$  then  $Q^p_{t+1} := \text{partial}Q^p$ 
7     else
8        $\text{new}Q^p := \emptyset$ 
9       for all pairs  $q_1^p : A_1 \leftarrow S_1 \in Q^p_{t+1}$  and
10         $q_2^p : A_2 \leftarrow S_2 \in \text{partial}Q^p$  do
11         /* unify the partial rules */
12         foreach  $A \leftarrow S = \text{mgs}(A_1 \leftarrow S_1, A_2 \leftarrow S_2)$  do
13            $q^p := q_1^p + q_2^p$ 
14           add  $q^p : A \leftarrow S$  to  $\text{new}Q^p$ 
15        $Q^p_{t+1} := \text{new}Q^p$ 

```

4.4 V^P -Rules Generation

This section calculates the new state probability function V_{t+1}^P , given Q_{t+1}^P , by maximizing over the actions (cf. Equation 3 (3)). Recall that Q_{t+1}^P is an ordered set of Q^P -rules of the form of $q^P : A \leftarrow S$. The task is to derive V_{t+1}^P , an ordered set of V^P -rules of the form of $v^P \leftarrow S$.

Trivially, turning Q_{t+1}^P into V_{t+1}^P takes three steps. First, the Q^P -rules must be sorted such that a rule connected to a high probability has a high priority as we are interested in the *maximum* probability (as defined by Equation 2). Second, the redundant Q^P -rules must be removed. A Q^P -rule is redundant if it is subsumed by another Q^P -rule that has a higher priority. Third, the remaining Q^P -rules are turned into V^P -rules by removing the action in the rule.

V^P Rules (Algorithm 4) implements Equation 3 (3). The process is as follows. First, the Q^P -rules are ordered decreasingly so that a state is always assigned a *maximum* probability (line 2). Second, to remove redundant Q^P -rules, an *absorbing rule* is required when the formula has an absorbing goal in that no more transitions occur after the goal is reached. Hence, given an absorbing goal, any rule concerning transitions that start from the goal is redundant and should be removed. Concretely, an *until* formula $[S_1 \text{ U } S_2]$ has an absorbing goal S_2 such that an execution stops once S_2 is reached. Hence, an absorbing rule $1.0 \leftarrow S_2$ must be inserted to the beginning of the Q^P -rules. On the other hand, a *next* formula $[X S_2]$ does not have an absorbing goal as it is possible to follow exactly one transition from S_2 to another state where S_2 may or may not hold. Therefore, no absorbing rules are inserted. Finally, redundant rules in the ordered set are removed (line 6-10).

V^P Rules generalizes the VRULES algorithm in REBEL (Kersting et al., 2004) to handle non-absorbing *next* formulae $[X S_2]$ as the original REBEL considers only absorbing goals. The generalized part is marked [blue](#).

Algorithm 4: V^P Rules. This function implements Equation 3 (3). All returned V^P -rules must be legal and respect OI-subsumption.

Require: Q_{t+1}^P : state-action probability function
 ψ : pCTL formula $[S_1 \text{ U } S_2]$ or $[X S_2]$
Return : V_{t+1}^P : the next state probability function

```

1  $V_{t+1}^P := \emptyset$ 
2 sort  $Q_{t+1}^P$  in decreasing order of  $Q^P$  probability
   /* Add absorbing  $Q^P$ -rules */
3 if  $\psi = S_1 \text{ U}^{\leq k} S_2$  then  $Q_{abs}^P := \{1.0 : \emptyset \leftarrow S | S \in S_2\}$ 
4 else  $Q_{abs}^P := \emptyset$ 
5  $Q_{t+1}^P := Q_{abs}^P + Q_{t+1}^P$  // attach absorbing  $Q^P$ -rules to the top of  $Q_{t+1}^P$ 
6 while  $Q_{t+1}^P \neq \emptyset$  do
7   remove the top  $Q^P$ -rule  $d : A \leftarrow B$  from their  $Q_{t+1}^P$ 
8   if no other rule  $d : A' \leftarrow B'$  in  $Q_{t+1}^P$  exists s.t.  $B'$  subsumes  $B$  then
9     add  $d \leftarrow B$  to  $V_{t+1}^P$ 
   /* remove redundant  $Q^P$ -rules */
10  remove all rules  $d'' \leftarrow B''$  from  $Q_{t+1}^P$  s.t.  $B$  subsumes  $B''$ 

```

4.5 PCTL-REBEL Illustration

This section illustrates pCTL-REBEL (cf. Equation 3), namely, taking a state probability function V_t^p to compute the next state probability function V_{t+1}^p . Clearly, different pCTL formulae require different numbers of iterations. That is, $[X S_2]$ requires one iteration, $[S_1 U^{\leq k} S_2]$ requires k iterations, and $[S_1 U S_2]$ requires an arbitrary number iterations to obtain an accurate enough approximation. For simplicity, we illustrate with path formulae (without probabilities) that require one iteration. Section 5 will cover the full pCTL language.

Formula 1 $\psi_1 = [X \text{ on}(\mathbf{a}, \mathbf{b})]$: Find all states that reach $\{\text{on}(\mathbf{a}, \mathbf{b})\}$ after 1 step.

Formula 2 $\psi_2 = [\text{on}(\mathbf{c}, \mathbf{d}) U^{\leq 1} \text{on}(\mathbf{a}, \mathbf{b})]$: Find all states that reach $\{\text{on}(\mathbf{a}, \mathbf{b})\}$ within 1 step by going through only the states where $\{\text{on}(\mathbf{c}, \mathbf{d})\}$ holds.

We consider the δ_{move} transition in the blocks world (also in Figure 1).

$$\begin{aligned} \delta_{\text{move}_1} &: \text{cl}(\mathbf{A}), \text{cl}(\mathbf{C}), \text{on}(\mathbf{A}, \mathbf{B}) \xleftarrow{0.9:\text{move}(\mathbf{A}, \mathbf{B}, \mathbf{C})} \text{cl}(\mathbf{A}), \text{cl}(\mathbf{B}), \text{on}(\mathbf{A}, \mathbf{C}) \\ \delta_{\text{move}_2} &: \text{cl}(\mathbf{A}), \text{cl}(\mathbf{B}), \text{on}(\mathbf{A}, \mathbf{C}) \xleftarrow{0.1:\text{move}(\mathbf{A}, \mathbf{B}, \mathbf{C})} \text{cl}(\mathbf{A}), \text{cl}(\mathbf{B}), \text{on}(\mathbf{A}, \mathbf{C}) \end{aligned}$$

PCTL-REBEL on Formula 1: $[X \text{ on}(\mathbf{a}, \mathbf{b})]$

For the formula $\psi_1 = [X \text{ on}(\mathbf{a}, \mathbf{b})]$, the initial function V_0^p is

$$\begin{aligned} 1.0 &\leftarrow \text{on}(\mathbf{a}, \mathbf{b}) \\ 0.0 &\leftarrow \emptyset \end{aligned}$$

① **Regression**. Given the initial state probability function V_0^p , to obtain all possible pre-states, **Regression** is called with all combinations of transition rules and V^p -rules, e.g. **Regression**($\delta_{\text{move}_1}, \emptyset, \text{on}(\mathbf{a}, \mathbf{b})$)⁴. This results in two sets of partial Q^p -rules such that each set considers one outcome of the **move** action. The resulting partial Q^p -rules are listed below. The $\langle 1 \cdot \rangle$ rules correspond to the successful outcome (i.e. δ_{move_1}) and the $\langle 2 \cdot \rangle$ rules correspond to the unsuccessful outcome (i.e. δ_{move_2}). Table 1 shows the the corresponding transitions and states. All partial Q^p -rules respect OI-subsumption.

$$\begin{aligned} \langle 1a \rangle 0.9 &: \text{move}(\mathbf{a}, \mathbf{b}, \mathbf{Z}) \leftarrow \text{cl}(\mathbf{a}), \text{cl}(\mathbf{b}), \text{on}(\mathbf{a}, \mathbf{Z}) \\ \langle 1b \rangle 0.9 &: \text{move}(\mathbf{X}, \mathbf{a}, \mathbf{Z}) \leftarrow \text{cl}(\mathbf{X}), \text{cl}(\mathbf{a}), \text{on}(\mathbf{X}, \mathbf{Z}), \text{on}(\mathbf{a}, \mathbf{b}) \\ \langle 1c \rangle 0.9 &: \text{move}(\mathbf{X}, \mathbf{Y}, \mathbf{a}) \leftarrow \text{cl}(\mathbf{X}), \text{cl}(\mathbf{Y}), \text{on}(\mathbf{X}, \mathbf{a}), \text{on}(\mathbf{a}, \mathbf{b}) \\ \langle 1d \rangle 0.9 &: \text{move}(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) \leftarrow \text{cl}(\mathbf{X}), \text{cl}(\mathbf{Y}), \text{on}(\mathbf{X}, \mathbf{Z}), \text{on}(\mathbf{a}, \mathbf{b}) \\ \langle 1e \rangle 0.0 &: \text{move}(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) \leftarrow \text{cl}(\mathbf{X}), \text{cl}(\mathbf{Y}), \text{on}(\mathbf{X}, \mathbf{Z}) \\ \langle 2a \rangle 0.1 &: \text{move}(\mathbf{a}, \mathbf{Y}, \mathbf{b}) \leftarrow \text{cl}(\mathbf{a}), \text{cl}(\mathbf{Y}), \text{on}(\mathbf{a}, \mathbf{b}) \\ \langle 2b \rangle 0.1 &: \text{move}(\mathbf{X}, \mathbf{a}, \mathbf{Z}) \leftarrow \text{cl}(\mathbf{X}), \text{cl}(\mathbf{a}), \text{on}(\mathbf{X}, \mathbf{Z}), \text{on}(\mathbf{a}, \mathbf{b}) \\ \langle 2c \rangle 0.1 &: \text{move}(\mathbf{X}, \mathbf{Y}, \mathbf{a}) \leftarrow \text{cl}(\mathbf{X}), \text{cl}(\mathbf{Y}), \text{on}(\mathbf{X}, \mathbf{a}), \text{on}(\mathbf{a}, \mathbf{b}) \\ \langle 2d \rangle 0.1 &: \text{move}(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) \leftarrow \text{cl}(\mathbf{X}), \text{cl}(\mathbf{Y}), \text{on}(\mathbf{X}, \mathbf{Z}), \text{on}(\mathbf{a}, \mathbf{b}) \\ \langle 2e \rangle 0.0 &: \text{move}(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) \leftarrow \text{cl}(\mathbf{X}), \text{cl}(\mathbf{Y}), \text{on}(\mathbf{X}, \mathbf{Z}) \end{aligned}$$

⁴ The second argument is \emptyset as the X operator is not absorbing.

$\text{on}(\mathbf{a}, \mathbf{b})$	δ_{move_1}				δ_{move_2}			
	$\langle 1a \rangle$	$\langle 1b \rangle$	$\langle 1c \rangle$	$\langle 1d \rangle$	$\langle 2a \rangle$	$\langle 2b \rangle$	$\langle 2c \rangle$	$\langle 2d \rangle$
\emptyset			$\langle 1e \rangle$				$\langle 2e \rangle$	

Table 1: Different combinations of transition rules and states result in different partial Q^p -rules $\langle 1a \rangle - \langle 2e \rangle$. The states on the leftmost column comes from V_0^p .

② Q^p Rules. To compute the Q^p -rules, the two sets of partial Q^p -rules must be combined. To do so, we consider all possible combinations of $\langle 1 \cdot \rangle$ and $\langle 2 \cdot \rangle$. The resulting Q^p -rules are listed as follows, along with how they are created. All rules respect OI-subsumption.

$\langle 1 \rangle$ 1.0 : $\text{move}(\mathbf{X}, \mathbf{a}, \mathbf{Z}) \leftarrow \text{cl}(\mathbf{X}), \text{cl}(\mathbf{a}), \text{on}(\mathbf{X}, \mathbf{Z}), \text{on}(\mathbf{a}, \mathbf{b})$	$\langle 1b \rangle + \langle 2b \rangle$
$\langle 2 \rangle$ 1.0 : $\text{move}(\mathbf{X}, \mathbf{Y}, \mathbf{a}) \leftarrow \text{cl}(\mathbf{X}), \text{cl}(\mathbf{Y}), \text{on}(\mathbf{X}, \mathbf{a}), \text{on}(\mathbf{a}, \mathbf{b})$	$\langle 1c \rangle + \langle 2c \rangle$
$\langle 3 \rangle$ 1.0 : $\text{move}(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) \leftarrow \text{cl}(\mathbf{X}), \text{cl}(\mathbf{Y}), \text{on}(\mathbf{X}, \mathbf{Z}), \text{on}(\mathbf{a}, \mathbf{b})$	$\langle 1d \rangle + \langle 2d \rangle$
$\langle 4 \rangle$ 0.9 : $\text{move}(\mathbf{a}, \mathbf{b}, \mathbf{Z}) \leftarrow \text{cl}(\mathbf{a}), \text{cl}(\mathbf{b}), \text{on}(\mathbf{a}, \mathbf{Z})$	$\langle 1a \rangle + \langle 2e \rangle$
$\langle 5 \rangle$ 0.9 : $\text{move}(\mathbf{X}, \mathbf{a}, \mathbf{Z}) \leftarrow \text{cl}(\mathbf{X}), \text{cl}(\mathbf{a}), \text{on}(\mathbf{X}, \mathbf{Z}), \text{on}(\mathbf{a}, \mathbf{b})$	$\langle 1b \rangle + \langle 2e \rangle$
$\langle 6 \rangle$ 0.9 : $\text{move}(\mathbf{X}, \mathbf{Y}, \mathbf{a}) \leftarrow \text{cl}(\mathbf{X}), \text{cl}(\mathbf{Y}), \text{on}(\mathbf{X}, \mathbf{a}), \text{on}(\mathbf{a}, \mathbf{b})$	$\langle 1c \rangle + \langle 2e \rangle$
$\langle 7 \rangle$ 0.9 : $\text{move}(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) \leftarrow \text{cl}(\mathbf{X}), \text{cl}(\mathbf{Y}), \text{on}(\mathbf{X}, \mathbf{Z}), \text{on}(\mathbf{a}, \mathbf{b})$	$\langle 1d \rangle + \langle 2e \rangle$
$\langle 8 \rangle$ 0.1 : $\text{move}(\mathbf{a}, \mathbf{Y}, \mathbf{b}) \leftarrow \text{cl}(\mathbf{a}), \text{cl}(\mathbf{Y}), \text{on}(\mathbf{a}, \mathbf{b})$	$\langle 1e \rangle + \langle 2a \rangle$
$\langle 9 \rangle$ 0.1 : $\text{move}(\mathbf{X}, \mathbf{a}, \mathbf{Z}) \leftarrow \text{cl}(\mathbf{X}), \text{cl}(\mathbf{a}), \text{on}(\mathbf{X}, \mathbf{Z}), \text{on}(\mathbf{a}, \mathbf{b})$	$\langle 1e \rangle + \langle 2b \rangle$
$\langle 10 \rangle$ 0.1 : $\text{move}(\mathbf{X}, \mathbf{Y}, \mathbf{a}) \leftarrow \text{cl}(\mathbf{X}), \text{cl}(\mathbf{Y}), \text{on}(\mathbf{X}, \mathbf{a}), \text{on}(\mathbf{a}, \mathbf{b})$	$\langle 1e \rangle + \langle 2c \rangle$
$\langle 11 \rangle$ 0.1 : $\text{move}(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) \leftarrow \text{cl}(\mathbf{X}), \text{cl}(\mathbf{Y}), \text{on}(\mathbf{X}, \mathbf{Z}), \text{on}(\mathbf{a}, \mathbf{b})$	$\langle 1e \rangle + \langle 2d \rangle$
$\langle 12 \rangle$ 0.0 : $\text{move}(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) \leftarrow \text{cl}(\mathbf{X}), \text{cl}(\mathbf{Y}), \text{on}(\mathbf{X}, \mathbf{Z})$	$\langle 1e \rangle + \langle 2e \rangle$

The Q^p -rules must be ordered by their probabilities as above. The rules $\langle 1 \rangle$ - $\langle 3 \rangle$ are interchangeable as they have the same probability. Similarly, $\langle 4 \rangle$ - $\langle 7 \rangle$ and $\langle 8 \rangle$ - $\langle 11 \rangle$ are interchangeable, respectively.

③ V^p Rules. Now we can derive the new V^p -rules by removing redundant Q^p -rules and dropping the action components. The resulting V^p -rules are shown below where the numbering inherits the one of the Q^p -rules. Rules $\langle 5 \rangle$ - $\langle 7 \rangle$ are redundant because they are subsumed by $\langle 1 \rangle$ - $\langle 3 \rangle$, respectively. Similarly, rules $\langle 9 \rangle$ - $\langle 11 \rangle$ are redundant as they are subsumed by $\langle 1 \rangle$ - $\langle 3 \rangle$, respectively.

$\langle 1 \rangle$ 1.0 $\leftarrow \text{cl}(\mathbf{X}), \text{cl}(\mathbf{a}), \text{on}(\mathbf{X}, \mathbf{Z}), \text{on}(\mathbf{a}, \mathbf{b})$
$\langle 2 \rangle$ 1.0 $\leftarrow \text{cl}(\mathbf{X}), \text{cl}(\mathbf{Y}), \text{on}(\mathbf{X}, \mathbf{a}), \text{on}(\mathbf{a}, \mathbf{b})$
$\langle 3 \rangle$ 1.0 $\leftarrow \text{cl}(\mathbf{X}), \text{cl}(\mathbf{Y}), \text{on}(\mathbf{X}, \mathbf{Z}), \text{on}(\mathbf{a}, \mathbf{b})$
$\langle 4 \rangle$ 0.9 $\leftarrow \text{cl}(\mathbf{a}), \text{cl}(\mathbf{b}), \text{on}(\mathbf{a}, \mathbf{Z})$
$\langle 8 \rangle$ 0.1 $\leftarrow \text{cl}(\mathbf{a}), \text{cl}(\mathbf{Y}), \text{on}(\mathbf{a}, \mathbf{b})$
$\langle 12 \rangle$ 0.0 $\leftarrow \text{cl}(\mathbf{X}), \text{cl}(\mathbf{Y}), \text{on}(\mathbf{X}, \mathbf{Z})$

Given $\psi_1 = [\mathbf{X} \text{ on}(\mathbf{a}, \mathbf{b})]$, and the initial state probability function V_0^p , we have applied pCTL-REBEL and obtained V_1^p , which assigns to all states a maximum probability of reaching $\text{on}(\mathbf{a}, \mathbf{b})$ after exactly one step.

PCTL-REBEL on Formula 2: $[\text{on}(c, d) \text{U}^{\leq 1} \text{on}(a, b)]$

For the formula $\psi_2 = [\text{on}(c, d) \text{U}^{\leq 1} \text{on}(a, b)]$, the initial probability function V_0^P is

$$\begin{aligned} 1.0 &\leftarrow \text{on}(a, b) \\ 0.0 &\leftarrow \emptyset \end{aligned}$$

① **Regression.** Given the initial state probability function V_0^P , to obtain all possible pre-states, **Regression** is called with all combinations of transition rules and V^P -rules, e.g. **Regression**(δ_{move_1} , $\text{on}(c, d)$, $\text{on}(a, b)$). This results in two sets of partial Q^P -rules such that each set considers one outcome of the **move** action. We list some of the resulting partial Q^P -rules below. The $\langle 1 \cdot \rangle$ rules correspond to the successful outcome (i.e. δ_{move_1}) and the $\langle 2 \cdot \rangle$ rules correspond to the unsuccessful outcome (i.e. δ_{move_2}). Table 2 shows the the corresponding transitions and states. All partial Q^P -rules respect OI-subsumption.

$$\begin{aligned} \langle 1a \rangle & 0.9 : \text{move}(a, b, c) \leftarrow \text{cl}(a), \text{cl}(b), \text{on}(a, c), \text{on}(c, d) \\ \langle 1b \rangle & 0.9 : \text{move}(a, b, Z) \leftarrow \text{cl}(a), \text{cl}(b), \text{on}(a, Z), \text{on}(c, d) \\ \langle 1c \rangle & 0.9 : \text{move}(c, a, d) \leftarrow \text{cl}(a), \text{cl}(c), \text{on}(a, b), \text{on}(c, d) \\ \langle 1d \rangle & 0.9 : \text{move}(c, Y, d) \leftarrow \text{cl}(c), \text{cl}(Y), \text{on}(a, b), \text{on}(c, d) \\ & \dots \\ \langle 1l \rangle & 0.0 : \text{move}(X, c, Z) \leftarrow \text{cl}(X), \text{cl}(c), \text{on}(X, Z), \text{on}(c, d) \\ \langle 1m \rangle & 0.0 : \text{move}(c, Y, d) \leftarrow \text{cl}(c), \text{cl}(Y), \text{on}(c, d) \\ \langle 1n \rangle & 0.0 : \text{move}(X, Y, c) \leftarrow \text{cl}(X), \text{cl}(Y), \text{on}(X, c), \text{on}(c, d) \\ \langle 1o \rangle & 0.0 : \text{move}(X, Y, Z) \leftarrow \text{cl}(X), \text{cl}(Y), \text{on}(X, Z), \text{on}(c, d) \\ \langle 2a \rangle & 0.1 : \text{move}(c, a, d) \leftarrow \text{cl}(a), \text{cl}(c), \text{on}(a, b), \text{on}(c, d) \\ \langle 2b \rangle & 0.1 : \text{move}(c, Y, d) \leftarrow \text{cl}(c), \text{cl}(Y), \text{on}(a, b), \text{on}(c, d) \\ & \dots \\ \langle 2m \rangle & 0.0 : \text{move}(c, Y, d) \leftarrow \text{cl}(c), \text{cl}(Y), \text{on}(c, d) \\ \langle 2n \rangle & 0.0 : \text{move}(X, Y, c) \leftarrow \text{cl}(X), \text{cl}(Y), \text{on}(X, c), \text{on}(c, d) \\ \langle 2o \rangle & 0.0 : \text{move}(X, Y, Z) \leftarrow \text{cl}(X), \text{cl}(Y), \text{on}(X, Z), \text{on}(c, d) \end{aligned}$$

	δ_{move_1}	δ_{move_2}
$\text{on}(a, b)$	$\langle 1a \rangle - \langle 1k \rangle$	$\langle 2a \rangle - \langle 2k \rangle$
\emptyset	$\langle 1l \rangle - \langle 1o \rangle$	$\langle 2l \rangle - \langle 2o \rangle$

Table 2: Different combinations of transition rules and states result in different partial Q^P -rules $\langle 1a \rangle - \langle 2o \rangle$. The states on the leftmost column comes from V_0^P .

② **Q^P Rules.** To compute the Q^P -rules, the two sets of partial Q^P -rules must be combined. To do so, we consider all possible combinations of $\langle 1 \cdot \rangle$ and $\langle 2 \cdot \rangle$. Some

of the resulting Q^p -rules are listed below, along with how they are created. All rules respect OI-subsumption.

$\langle 1 \rangle$ 1.0 : $\text{move}(c, a, d) \leftarrow \text{cl}(a), \text{cl}(c), \text{on}(a, b), \text{on}(c, d)$	$\langle 1c \rangle + \langle 2a \rangle$
$\langle 2 \rangle$ 1.0 : $\text{move}(c, Y, d) \leftarrow \text{cl}(c), \text{cl}(Y), \text{on}(a, b), \text{on}(c, d)$	$\langle 1d \rangle + \langle 2b \rangle$
...	
$\langle 10 \rangle$ 0.9 : $\text{move}(a, b, c) \leftarrow \text{cl}(a), \text{cl}(b), \text{on}(a, c), \text{on}(c, d)$	$\langle 1a \rangle + \langle 2n \rangle$
$\langle 11 \rangle$ 0.9 : $\text{move}(a, b, Z) \leftarrow \text{cl}(a), \text{cl}(b), \text{on}(a, Z), \text{on}(c, d)$	$\langle 1b \rangle + \langle 2o \rangle$
...	
$\langle 33 \rangle$ 0.0 : $\text{move}(c, Y, d) \leftarrow \text{cl}(c), \text{cl}(Y), \text{on}(c, d)$	$\langle 1m \rangle + \langle 2m \rangle$
$\langle 34 \rangle$ 0.0 : $\text{move}(X, Y, c) \leftarrow \text{cl}(X), \text{cl}(Y), \text{on}(X, c), \text{on}(c, d)$	$\langle 1n \rangle + \langle 2n \rangle$
$\langle 35 \rangle$ 0.0 : $\text{move}(X, Y, Z) \leftarrow \text{cl}(X), \text{cl}(Y), \text{on}(X, Z), \text{on}(c, d)$	$\langle 1o \rangle + \langle 2o \rangle$

The Q^p -rules must be ordered by their probabilities as above.

③ V^p Rules. Now we can derive the new V^p -rules. Since $[\text{on}(a, b)]$ is absorbing, the following absorbing rule must be inserted to the beginning of the Q^p -rule set.

$$\langle 0 \rangle 1.0 : \emptyset \leftarrow \text{on}(a, b)$$

After removing redundant Q^p -rules and dropping the action components in the rules, we obtain the resulting V^p -rules below. The numbering of the V^p -rules inherits the one of the Q^p -rules. Most redundant rules are subsumed by the absorbing rule $\langle 0 \rangle$.

$\langle 0 \rangle$ 1.0 $\leftarrow \text{on}(a, b)$
$\langle 10 \rangle$ 0.9 $\leftarrow \text{cl}(a), \text{cl}(b), \text{on}(a, c), \text{on}(c, d)$
$\langle 11 \rangle$ 0.9 $\leftarrow \text{cl}(a), \text{cl}(b), \text{on}(a, Z), \text{on}(c, d)$
$\langle 33 \rangle$ 0.0 $\leftarrow \text{cl}(c), \text{cl}(Y), \text{on}(c, d)$
$\langle 34 \rangle$ 0.0 $\leftarrow \text{cl}(X), \text{cl}(Y), \text{on}(X, c), \text{on}(c, d)$
$\langle 35 \rangle$ 0.0 $\leftarrow \text{cl}(X), \text{cl}(Y), \text{on}(X, Z), \text{on}(c, d)$

5 Main Contribution – a Relational Model Checker

This section introduces the main algorithm in this paper, the relational model checking algorithm, that solves the relational model checking problem (Section 3.3). Formally, given a pCTL formula ϕ and an RMDP K , the relational model checker identifies the set $Sat_K(\phi)$ of abstract states that represent all ground states that satisfy ϕ . Since the relational model checker is based on PCTL-REBEL (Section 4), it operates at the relational level.

Different from Section 4 that computes one single iteration, this section allows for full relational pCTL formulae that requires multiple iterations and can be *nested*. It is standard practice to represent a nested pCTL formula as a parse tree (Baier and Katoen, 2008). In a parse tree, each leaf node is an abstract state and each inner node contains exactly one operator.

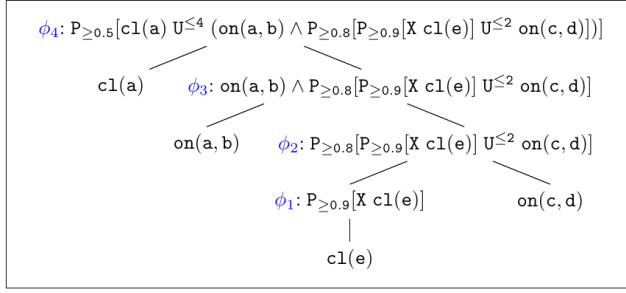


Fig. 4: The parse tree of ϕ_{nested} . Each inner node of the parse tree is annotated with a subformula ϕ_i and handles one operator. ϕ_1 handles X , ϕ_2 handles $U^{\leq 2}$, ϕ_3 handles \wedge and ϕ_4 handles $U^{\leq 4}$.

Example 7 A nested relational pCTL formula example is as follows.

$$\phi_{\text{nested}} = P_{\geq 0.5}[c1(a) U^{\leq 4}(\text{on}(a, b) \wedge P_{\geq 0.8}[P_{\geq 0.9}[X c1(e)] U^{\leq 2} \text{on}(c, d)])]$$

The nested formula ϕ_{nested} is complex at first sight, however, it can be represented as a *parse tree* in Figure 4. A state s satisfies ϕ_{nested} if and only if s satisfies all the following three conditions. (1) A path starting from s must reach $\{\text{on}(a, b)\}$ within 4 steps with a probability greater than or equal to 0.5 by going through the states where $\{c1(a)\}$ holds. (2) Then, the path must reach $\{\text{on}(c, d)\}$ within 2 steps with a probability greater than or equal to 0.8 by going through the states that satisfy $P_{\geq 0.9}[X c1(e)]$. (3) The states that satisfy $P_{\geq 0.9}[X c1(e)]$ are the ones that can transition to $\{c1(e)\}$ with a probability greater than or equal to 0.9 after exactly one step.

PCTL-REBEL always use a parse tree to evaluate a given formula. A parse tree is recursively evaluated *upwards*. That is, an inner node considers its child leaf nodes to evaluate a subformula ϕ_i , resulting in a set of states $Sat(\phi_i)$. Then, the inner node *folds* the sub-tree so that its parent node can be activated. The final set of states $Sat(\phi)$ is produced by the root. For example, in Figure 4, an intermediate state set $Sat(\phi_i)$ is used to evaluate $Sat(\phi_{i+1})$. The state set $Sat(\phi_4)$ is returned as the solution.

Section 5.1 defines the pCTL-REBEL model checking algorithm, which is composed of 3 mutually recursive algorithms. Then, Section 5.2 gives an overview of the properties of the model checking algorithm.

5.1 The pCTL-REBEL Model Checking Algorithm

The pCTL-REBEL model checking algorithm reforms the satisfiability relation \models^θ (cf. Section 3.2). Given an RMDP K and a pCTL formula ϕ , the pCTL-REBEL model checker $Sat_K(\phi)$, i.e. all states in K that satisfy ϕ . The model checker consists of 3 algorithms **Check**, **CheckUntil** and **CheckNext** that call one another recursively. In particular, **Check** is the main algorithm and is mutually recursive

with the other two algorithms. **CheckUntil** handles the *until* formulae of the form of $P_{\bowtie p}[\phi_1 U^{\leq k} \phi_2]$. **CheckNext** handles the *next* formulae of the form of $P_{\bowtie p}[X \phi_2]$.

Given a formula ϕ , **Check** (Algorithm 5) evaluates ϕ in a recursive fashion as follows. If ϕ is a relation atom l (resp. $\neg l$), it returns a single abstract state $\{l\}$ (resp. $\{\neg l\}$) (line 1-2). If ϕ is a conjunction of the form of $\phi_1 \wedge \phi_2$, it first computes ϕ_1 and ϕ_2 separately to get two sets of abstract states Sat_1 and Sat_2 (line 4-5). Then, for each possible pair of abstract states $s_1 \in Sat_1$ and $s_2 \in Sat_2$, it collects all maximally general specializations $mgs(s_1, s_2)$ (line 6). Since all $mgs(s_1, s_2)$ are OI-subsumed by s_1 and s_2 by definition, they automatically satisfy $\phi_1 \wedge \phi_2$. If ϕ is a disjunction of the form of $\phi_1 \vee \phi_2$, it returns the union of the solutions to ϕ_1 and ϕ_2 (line 7-8). If ϕ is an *until* formula, it calls **CheckUntil** (Algorithm 6) (line 9-10). If ϕ is a *next* formula, it calls **CheckNext** (Algorithm 7) (line 11-12).

Algorithm 5: Check. The main algorithm to solve the relational model checking problem. This algorithm implements the satisfiability relation (Section 3.2). **Check** is mutually recursive with **CheckUntil** and **CheckNext**.

Require: ϕ : pCTL formula
Return : $Sat_K(\phi)$: abstract states that satisfy ϕ

```

1 if  $\phi = l$  then  $Sat_K(\phi) = \{l\}$ 
2 else if  $\phi = \neg l$  then  $Sat_K(\phi) = \{\neg l\}$ 
3 else if  $\phi = \phi_1 \wedge \phi_2$  then
4   |  $Sat_1 = \mathbf{Check}(\phi_1)$ 
5   |  $Sat_2 = \mathbf{Check}(\phi_2)$ 
6   |  $Sat_K(\phi) = \{mgs(s_1, s_2) \mid s_1 \in Sat_1, s_2 \in Sat_2\}$ 
7 else if  $\phi = \phi_1 \vee \phi_2$  then
8   |  $Sat_K(\phi) = \mathbf{Check}(\phi_1) \cup \mathbf{Check}(\phi_2)$ 
9 else if  $\phi = P_{\bowtie p}[\phi_1 U^{\leq k} \phi_2]$  then
10  |  $Sat_K(\phi) = \mathbf{CheckUntil}(\phi)$ 
11 else if  $\phi = P_{\bowtie p}[X \phi_2]$  then
12  |  $Sat_K(\phi) = \mathbf{CheckNext}(\phi)$ 

```

CheckUntil (Algorithm 6) is mutually recursive with **Check** (Algorithm 5). It computes a formula of the form of $P_{\bowtie p}[\phi_1 U^{\leq k} \phi_2]$ or $P_{\bowtie p}[\phi_1 U \phi_2]$ where ϕ_1 and ϕ_2 are pCTL formulae. To begin the process, **CheckUntil** computes ϕ_1 and ϕ_2 separately to get two sets of abstract states (line 1-2). Then, it sets the step bound for pCTL-REBEL (line 3). When the given formula has a bounded $U^{\leq k}$ operator, the step bound is k . Otherwise, the step bound is set to infinity and an arbitrary number of pCTL-REBEL iterations are applied until convergence (line 5-11). The convergence condition is twofold. First, the abstract states in the probability function do not change, i.e. states in V_t^p and V_{t-1}^p are the same. Second, the state probabilities have converged with respect to a given threshold ϵ , i.e. $\max_{s \in V_t^p} |V_t^p(s) - V_{t-1}^p(s)| < \epsilon$. Then, **CheckUntil** collects and returns the abstract states in V_t^p that satisfy the given probability threshold $\bowtie p$ (line 12-14).

CheckNext (Algorithm 7) is mutually recursive with **Check** (Algorithm 5). It computes a formula of the form of $\phi = P_{\bowtie p}[X \phi_2]$ where ϕ_2 is another pCTL formula. **CheckNext** is a similar but simpler than **CheckUntil**. It performs one single value iteration to obtain the probability function V_1^p (line 2-3). Then it

Algorithm 6: CheckUntil. CheckUntil is mutually recursive with Check.

Require: ϕ : pCTL formula $P_{\bowtie p}[\phi_1 U^{\leq k} \phi_2]$ or $P_{\bowtie p}[\phi_1 U \phi_2]$
Return : $Sat_K(\phi)$: abstract states that satisfy ϕ

```

1  $S_1 := \text{Check}(\phi_1)$ 
2  $S_2 := \text{Check}(\phi_2)$ 
3 if  $k$  is not given then  $k := \infty$  // set  $k$  as infinity for unbounded U
4
5  $V_0^p := \{1.0 \leftarrow s | s \in S_2\} \cup \{0 \leftarrow \emptyset\}$  // Initialize  $V_0^p$ 
6  $t := 0$  //  $t$  is the number of iterations so far
7  $\epsilon' := \infty$  //  $\epsilon'$  is the distance between  $V_t^p$  and  $V_{t+1}^p$ 
8 while  $t < k$  or  $\epsilon' \geq \epsilon$  do
9    $V_{t+1}^p := \text{OneIteration}(V_t^p, S_1 U^{\leq k} S_2)$ 
10   $\epsilon' := V_{t+1}^p - V_t^p$ 
11   $t := t + 1$ 
12  $Sat_K(\phi) := \emptyset$ 
13 foreach  $p_i \leftarrow s_i \in V_t^p$  do
14   if  $p_i \bowtie p$  then add  $s_i$  to  $Sat_K(\phi)$ 

```

collects and returns the abstract states in V_1^p that satisfy the given probability threshold $\bowtie p$ (line 4-6). CheckNext needs only one iteration as it considers exactly one transition.

Algorithm 7: CheckNext. CheckNext is mutually recursive with Check.

Require: ϕ : pCTL formula $P_{\bowtie p}[X \phi_2]$
Return : $Sat_K(\phi)$: abstract states that satisfy ϕ

```

1  $S_2 = \text{Check}(\phi_2)$ 
2  $V_0^p := \{1.0 \leftarrow s | s \in S_2\} \cup \{0 \leftarrow \emptyset\}$ 
3  $V_1^p := \text{OneIteration}(V_0^p, X S_2)$ 
4  $Sat_K(\phi) := \emptyset$ 
5 foreach  $p_i \leftarrow s_i \in V_1^p$  do
6   if  $p_i \bowtie p$  then add  $s_i$  to  $Sat_K(\phi)$ 

```

5.2 Properties of PCTL-REBEL

PCTL-REBEL is a relational model checking algorithm that finds all states that satisfy a given pCTL formula. We discuss the properties of pCTL-REBEL.

Lifted instead of operating at the ground level, pCTL-REBEL performs lifted inference as both the formula and the states are specified at an abstract level using relational representations. Using lifted inference allows pCTL-REBEL to exploit relational symmetries in the model and make abstraction of the domain, hence mitigate the state explosion problem. The lifted inference was discussed in detail in Section 4.

Sound for step-bounded pCTL formulae PCTL-REBEL is sound for finite RMDPs (that have a finite domain) and any step-bounded pCTL formulae. PCTL-REBEL is not sound for indefinite-horizon formulae (e.g. \mathcal{U}), just like many other value iteration algorithms. This is because pCTL-REBEL uses a naive termination criterion with some arbitrary convergence threshold ϵ (see `CheckUntil`, Algorithm 6). Although this naive termination criterion is not sound, it achieves precise approximation in practice. For further details, please refer to Haddad and Monmege (2014).

Complete PCTL-REBEL is complete for finite RMDPs such that the state probability function V^p assigns a probability to all states. PCTL-REBEL captures the entire state space by using an ordered set of relational V^p -rules. Those states that are not captured by any other V^p -rules are guaranteed to be covered by the last V^p -rule $0 \leftarrow \emptyset$ (see Section 4.1).

6 Relational Model Checking for Infinite MDPs

It is clear that when the domain is finite, the model checking problem is decidable as one can enumerate all states in the model, as discussed in Section 3.3. However, when the domain is infinite, the state space is typically infinite, making enumerating all states infeasible. In this section, we obtain decidability for a special class of infinite RMDPs. We will prove that under the *state-boundedness* condition (Belardinelli et al., 2011), a *finite abstraction* of such RMDPs can be constructed checked by pCTL-REBEL. The main idea is to generate a finite abstraction that captures all relevant information of the underlying infinite RMDP with respect to a pCTL formula. Accordingly, checking the abstraction is equivalent to checking the infinite RMDP. By checking the finite abstraction, the model checking problem becomes decidable. This section adopts the approach of Belardinelli et al. (2011) to construct finite abstractions of infinite RMDPs. Furthermore, we show that pCTL-REBEL can naturally handle such abstractions as they are structurally similar to RMDPs.

Section 6.1 defines the state-boundedness condition and under which, the finite abstraction of an infinite RMDP. Section 6.2 proves that checking a pCTL formula ϕ against the finite abstraction is equivalent to checking ϕ against the corresponding infinite RMDP. Section 6.3 discusses properties of pCTL-REBEL when handling such infinite RMDPs.

6.1 Abstract MDP: a Finite Abstraction of an Infinite RMDP

Given an infinite RMDP, the relational model checking problem is generally undecidable due to the possibly infinite domain. To obtain decidability, this section constructs a finite abstraction of a given infinite RMDP, called an *abstract RMDP* (ARMDP). The purpose of an ARMDP is to use a finite model to capture all relevant information about a pCTL formula. An ARMDP must be constructed under the *state-boundedness* condition.

The state-boundedness condition states that any state concerns only a finite number of objects. For example, consider a blocks world that has infinitely many blocks and a table with a capacity b . An agent can take a block away or put a

new block to the table, but no more than b blocks can be on the table at any moment. Hence, the state bound is b . Since any two states can describe totally different blocks, the model still contains *infinitely many* states. We say the blocks on the table are in the *active domain*. We now formally define active domain and state-boundedness.

Definition 6 A bounded state s_b is a finite subset of a ground state s in some RMDP, i.e. $s_b \subseteq s$. The active domain of s_b , denoted $\text{adom}(s_b)$, is the set of all domain objects in s_b .

An active domain $\text{adom}(s_b)$ is by definition finite as a bounded state s_b is finite. A bounded state is similar to a ground state but concerns only a finite number of objects. That is, all atoms in s_b are *true* and all the others are *false*.

Definition 7 Given an RMDP $K = \langle \Sigma, \Delta \rangle$, its underlying MDP $M = \langle S, A, T \rangle$, a state bound b , and a starting state s_0 such that $|\text{adom}(s_0)| \leq b$, an MDP $M_b = \langle S_b, A_b, T_b \rangle$ can be defined by including all states that are reachable from s_0 and contain at most b constants. Formally, $M_b = \langle S_b, A_b, T_b \rangle$ is defined as

$$\begin{aligned} S_b &:= \{s \mid s \in S, |\text{adom}(s)| \leq b\} \\ T_b &:= \{h \xrightarrow{p:a} b \in T \mid h \in S_b, b \in S_b\} \\ A_b &:= \{a \mid h \xrightarrow{p:a} b \in T_b\} \end{aligned}$$

If $b \in \mathbb{N}$, then M_b is called *state-bounded* or *b-bounded*.

A b -bounded MDP M_b is uniquely defined by an RMDP and a state bound. Roughly speaking, M_b a *sub-MDP* of M that concerns at most b objects in any state.

Example 8 Consider a blocks world $K = \langle \Sigma, \Delta \rangle$ with $\Sigma = \langle R, D \rangle$ where $R = \{\text{cl}/1, \text{on}/2\}$, $D = \{\text{bl}_i \mid i \in \mathbb{N}\}$ is infinite, and Δ contains the following rules (as in Figure 1).

$$\begin{aligned} \delta_{\text{move}_1} &: \text{cl}(\mathbf{A}), \text{cl}(\mathbf{C}), \text{on}(\mathbf{A}, \mathbf{B}) \xrightarrow{0.9:\text{move}(\mathbf{A}, \mathbf{B}, \mathbf{C})} \text{cl}(\mathbf{A}), \text{cl}(\mathbf{B}), \text{on}(\mathbf{A}, \mathbf{C}) \\ \delta_{\text{move}_2} &: \text{cl}(\mathbf{A}), \text{cl}(\mathbf{B}), \text{on}(\mathbf{A}, \mathbf{C}) \xrightarrow{0.1:\text{move}(\mathbf{A}, \mathbf{B}, \mathbf{C})} \text{cl}(\mathbf{A}), \text{cl}(\mathbf{B}), \text{on}(\mathbf{A}, \mathbf{C}) \end{aligned}$$

Given a state bound b , the RMDP K defines a b -bounded MDP $M_b = \langle S_b, A_b, T_b \rangle$ such that each bounded state $s_b \in S_b$ contains at most b blocks. We will use this blocks world example throughout this section.

Having defined state-boundedness, let us now move on to constructing an abstract MDP (ARMDP). With respect to a pCTL formula, some domain objects in an MDP are irrelevant. For example, $\text{P}_{\geq 0.9}[\mathbf{X} \text{cl}(\mathbf{a})]$ does not concern any particular blocks other than block \mathbf{a} . By abstracting away irrelevant objects of a pCTL formula ϕ , we can construct an ARMDP M_ϕ that captures all necessary information to check ϕ . Concretely, to construct an ARMDP, all ϕ -*relevant* domain objects are preserved, and all other objects are abstracted using variables. A domain object is ϕ -*relevant* if and only if it is in ϕ or in a transition rule. Furthermore, if a formula ϕ concerns a finite number of objects, the corresponding ARMDP M_ϕ is finite. In this paper, we assume all formulae contain at most b objects where b is the state bound.

Definition 8 For a b -bounded MDP $M_b = \langle S_b, A_b, T_b \rangle$ of an RMDP $K = \langle \Sigma, \Delta \rangle$ and a relational pCTL formula ϕ , a b -bounded abstract RMDP $M_\phi = \langle S_\phi, \Sigma, \Delta, W \rangle$ is defined where

$$\begin{aligned} S_\phi &= \{s_\phi \mid \exists s \in S_b. s \preceq_\theta s_\phi, \\ &\quad \text{consts}(s_\phi) \subseteq \text{consts}(\phi) \cup \text{consts}(\Delta), \\ &\quad \text{vars}(s_\phi) \subseteq W\} \end{aligned}$$

and W is a set of $b - |\text{consts}(\phi) \cup \text{consts}(\Delta)|$ distinct variables. All terms in $s_\phi \in S_\phi$ are from a finite set of terms

$$\text{terms}(s_\phi) := W \cup \text{consts}(\phi) \cup \text{consts}(\Delta)$$

It is assumed that $\text{vars}(\phi) \in W$. It is assumed that states in S_ϕ are not syntactic variants, that is, they are not a variable renaming of one another.

An ARMDP is finite as only a finite number of terms are allowed in the state description. An ARMDP is b -bounded as each state $s_\phi \in S_\phi$ contains at most b terms, i.e. $|\text{terms}(s_\phi)| \leq b$. Since an ARMDP state is b -bounded, it has bounded branching behavior such that the number of available actions in any state is finite. The finite state space of ARMDPs is crucial to obtaining decidability.

An ARMDP has a similar structure as an RMDP but has a more abstract state space. That is, unlike RMDPs, the state space of an ARMDP is not connected to an explicit domain as any ARMDP state concerns only ϕ -relevant constants and variables. By replacing the variables by domain constants, an ARMDP state can enumerate infinitely many ground states in the underlying RMDP.

Example 9 (Cont. Example 8) A b -bounded MDP $M_b = \langle S_b, A_b, T_b \rangle$ and the formula $\phi = \text{c1}(\mathbf{a})$ defines an ARMDP $M_\phi = \langle S_\phi, \Sigma, \Delta, W \rangle$. A state $s_b \in S_b$ concerns at most b terms, namely $\{\mathbf{B1}_1, \mathbf{B1}_2, \dots, \mathbf{B1}_{b-1}, \mathbf{a}\}$. The abstract states $s_1 = \{\text{c1}(\mathbf{a})\} \in S_\phi$ and $s_2 = \{\text{c1}(\mathbf{B1}_1)\} \in S_\phi$ represent the set of all infinitely many ground states that have at least one clear block. These ground states can be enumerated by assigning domain objects to $\mathbf{B1}_1$.

This section has shown that given a state-bounded MDP and a pCTL formula, an finite abstraction can be constructed. Such abstraction is called an ARMDP. An ARMDP is similar to an RMDP but is state-bounded, finite, and captures only ϕ -relevant information. The next section will show that checking an ARMDP is equivalent to checking the underlying model.

6.2 Decidable Model Checking for ARMDPs

This section proves the decidability of the model checking problem for a special class of infinite RMDPs, namely, the ARMDP, defined in Section 6.1. Decidability is obtained by proving that checking the ARMDP yields equivalent results as checking the underlying infinite, state-bounded MDP, namely, Theorem 1.

Theorem 1 For a b -bounded MDP M_b and its corresponding ARMDP M_ϕ based on a pCTL sentence ϕ , checking ϕ against M_ϕ is equivalent to checking ϕ against M_b , formally,

$$M_b \models \phi \Leftrightarrow M_\phi \models \phi$$

Theorem 1 is proven by using probabilistic bisimulation. Probabilistic bisimulation is a standard model checking technique that compares two probabilistic transition systems (Baier and Katoen, 2008). Given a pCTL formula ϕ , if two states from different transition systems are *probabilistic bisimilar*, then their behaviors are indistinguishable. Hence, checking ϕ in either system yields identical results. Theorem 1 extends the theorem by Belardinelli et al. (2011)[Theorem 2] for the probabilistic setting⁵.

We prove Theorem 1 in two steps. First, we define indistinguishable states (Definition 9) and probabilistic bisimulation (Definition 10). Second, we show that an ARMDP and its underlying state-bounded MDP actually define a probabilistic bisimulation (Proposition 1 and 2).

Definition 9 *Given the b -bounded MDP $M_b = \langle S_b, A_b, T_b \rangle$ of a RMDP $K = \langle \Sigma, \Delta \rangle$ and a pCTL sentence ϕ , consider two ground states s_1 and s_2 in S_b such that $C_1 = \text{consts}(s_1) \subset D$ and $C_2 = \text{consts}(s_2) \subset D$. Let C be the set of all ϕ -relevant domain objects $C := \text{consts}(\phi) \cup \text{consts}(\Delta) \subseteq C_1 \cap C_2$. The two states s_1 and s_2 are called **indistinguishable**⁶ under C , if and only if s_1 and s_2 are renamings of one another under a bijection $f : C_1 \setminus C \mapsto C_2 \setminus C$. The bijection f renames ϕ -irrelevant domain objects. Formally,*

$$s_1 \sim_C s_2 \iff f(s_1) = s_2$$

We abuse the notation and let $f(s_1)$ be the state obtained by renaming constants in s_1 .

Definition 9 has defined indistinguishable states under a pCTL formula ϕ . Indistinguishable states are essentially renamings of one another and share the same properties with respect to a formula ϕ , hence, they can be represented by one abstract state in an ARMDP. That is, an ARMDP state represents infinitely many indistinguishable states in the underlying MDP.

Example 10 (Cont. Example 9) Consider a state-bounded MDP $M_b = \langle S_b, A_b, T_b \rangle$ with a state bound $b = 2$, a pCTL formula $\phi = \text{c1}(\mathbf{a})$ and two ground states $s_3 = \{\text{c1}(\mathbf{b}), \text{on}(\mathbf{b}, \mathbf{a})\} \in S_b$ and $s_4 = \{\text{c1}(\mathbf{c}), \text{on}(\mathbf{c}, \mathbf{a})\} \in S_b$. States s_3 and s_4 are variable renamings of each other thus indistinguishable under ϕ . They both express *some block is on block a* and can be represented by the abstract state $s_5 = \{\text{c1}(\mathbf{B1}_1), \text{on}(\mathbf{B1}_1, \mathbf{a})\}$.

Now we must define *probabilistic bisimulation*.

Definition 10 (cf. Baier and Katoen (2008)) *Consider a b -bounded MDP $M_b = \langle S_b, A_b, T_b \rangle$, a pCTL formula ϕ and the finite set C of ϕ -relevant constants, $\mathcal{R} \subseteq S_b \times S_b$ is a **probabilistic bisimulation** if for any pair of states $\langle s_1, s_2 \rangle \in \mathcal{R}$:*

1. s_1 and s_2 are indistinguishable under C , i.e. $s_1 \sim_C s_2$
2. s_1 and s_2 have the identical probabilities of going to any other distinguishable states, i.e. $\forall a \in A : T(s_1, a)(s) = T(s_2, a)(s)$ for each $s \in S/\mathcal{R}$

⁵ More specifically, as compared to their work, this paper uses RMDPs instead of artifact systems, MDPs instead of Kripke structures, ground states instead of databases, and probabilistic bisimulation instead of bisimulation.

⁶ Indistinguishable states are also commonly called isomorphic states.

States s_1 and s_2 are called *probabilistic bisimilar*.

To prove Theorem 1, we must prove that M_ϕ and M_b indeed form a probabilistic bisimulation. Propositions 1 and 2 will together prove Theorem 1, i.e. checking an ARMDP is equivalent to checking its underlying state-bounded MDP.

Proposition 1 *For an ARMDP $M_\phi = \langle S_\phi, \Sigma, \Delta, W \rangle$ based on a b -bounded MDP $M_b = \langle S_b, A_b, T_b \rangle$ and a pCTL sentence ϕ , if two states are indistinguishable, then they either both can fire a transition rule or both cannot fire a transition rule. Formally, for any transition rule of the form*

$$H_i \xrightarrow{p_i:\alpha} B$$

if $s_1 \sim_C s_2$, then

$$s_1 \preceq_{\theta_1} B \iff s_2 \preceq_{\theta_2} B$$

for two indistinguishable substitutions θ_1 and θ_2 .

Proposition 2 *For an ARMDP $M_\phi = \langle S_\phi, \Sigma, \Delta, W \rangle$ based on a b -bounded MDP $M_b = \langle S_b, A_b, T_b \rangle$ and a pCTL sentence ϕ , if two states are indistinguishable, then they have the identical probabilities of going to any other distinguishable states. Formally, for any transition rule $\delta \in \Delta$ of the form*

$$\{H_1 \xrightarrow{p_1:\alpha} B, \dots, H_n \xrightarrow{p_n:\alpha} B\}$$

if $s_1 \sim_C s_2$, then they define two sets of ground transition rules $T(s_1, \alpha\theta_1)$ and $T(s_2, \alpha\theta_2)$ that share the same transition probabilities.

$$T(s_1, \alpha\theta_1) := \{h_{1,i} \xrightarrow{p_i:\alpha\theta_1} s_1 \mid H_i \xrightarrow{p_i:\alpha} B \in \delta, s_1 \preceq_{\theta_1} B, \\ h_{1,i} = (s_1 \setminus B\theta_1) \cup H_i\theta_1\}$$

$$T(s_2, \alpha\theta_2) := \{h_{2,i} \xrightarrow{p_i:\alpha\theta_2} s_2 \mid H_i \xrightarrow{p_i:\alpha} B \in \delta, s_2 \preceq_{\theta_2} B, \\ h_{2,i} = (s_2 \setminus B\theta_2) \cup H_i\theta_2\}$$

Propositions 1 and 2 have shown that the indistinguishable relation \sim_C between an ARMDP and its underlying state-bounded MDP is a probabilistic bisimulation. \sim_C is an equivalence relation such that two probabilistic bisimilar states exhibit identical behavior. In other words, the ARMDP and the underlying MDP denote a mutual, step-wise simulation of indistinguishable states. With a similar reasoning as in (Belardinelli et al., 2011), we conclude that the model checking problem for infinite, state-bounded MDPs is decidable as it can be done on its corresponding finite ARMDP.

6.3 PCTL-REBEL for Infinite MDPs

PCTL-REBEL can handle an infinite, state-bounded RMDP by reasoning about its finite abstraction. More formally, given an infinite RMDP with a step bound and a pCTL formula, a finite ARMDP can be constructed and naturally checked by pCTL-REBEL. In fact, as an ARMDP is just like an RMDP that concerns at most b objects in a state (see Section 6.1), pCTL-REBEL requires only one adaption to restrict the state size. That is, all states that have more than b objects must be eliminated. The number of objects in a state bound is simply obtained by counting all variables and constants as OI-subsumption is imposed. The adaption of maintaining the state bound is implemented in **Regression** (Algorithm 2).

PCTL-REBEL for an infinite, state-bounded MDP is lifted and complete, just as like pCTL-REBEL for a finite RMDP (see Section 5.2). We discuss properties of pCTL-REBEL for state-bounded MDPs in detail.

Lifted PCTL-REBEL checks a state-bounded, infinite RMDP that has a finite abstraction (namely, an ARMDP) at a lifted level. Specifically, pCTL-REBEL operates on the ARMDP that is bisimilar to the underlying ground MDP. Checking an ARMDP is the same as checking an RMDP, except that the state bound must be maintained.

Sound for step-bounded pCTL formulae PCTL-REBEL is sound for step-bounded pCTL formulae in an ARMDP. This is a consequence of Theorem 1 that shows checking an ARMDP yields exactly the same results as checking its underlying MDP. PCTL-REBEL is not sound for indefinite-horizon formulae but it achieves precise approximation in practice.

Complete PCTL-REBEL is complete for checking ARMDPs as all states are assigned a probability. PCTL-REBEL captures the entire state space by using an ordered set of relational V^p -rules.

7 Experiments

We aim to answer the following questions in this section. Q1 and Q2 focus on the benefits of pCTL-REBEL, and Q3-Q5 focus on the limitations of pCTL-REBEL.

- Q1** What formulae can pCTL-REBEL check in practice?
- Q2** How does pCTL-REBEL compare with state-of-the-art model checkers?
- Q3** How well does pCTL-REBEL handle indefinite-horizon formulae?
- Q4** How well does pCTL-REBEL handle a complex relational transition function?
- Q5** What are the computational costs of different pCTL operators?

We implemented and validated an unoptimized pCTL-REBEL research prototype using SWI-Prolog 8.0.2, with the *constraint handling rules* library. Experiments were run on a 2.4 GHz Intel i5 processor. We use the blocks world dataset ⁷ and the box world dataset ⁸. We compare pCTL-REBEL with the state-of-the-art model checkers PRISM (Kwiatkowska et al., 2011) and STORM (Dehnert et al.,

⁷ <https://qcomp.org/benchmarks/#blocksworld>

⁸ <https://qcomp.org/benchmarks/#boxworld>

2017)⁹. We set a time-out of 1800 seconds for all model checkers. PCTL-REBEL’s state bound is sometimes referred as *number of blocks/cities* so as to provide a direct comparison with PRISM and STORM that operate on ground models.

The blocks world and the box world datasets originate from IPPC-2008¹⁰ and are originally specified in the ppddl language. STORM and PRISM operate on the converted prism-format files¹¹, and pCTL-REBEL operates on the converted prolog-format files¹². We discard the reward structure in the model. The blocks world dataset is simplified to have the relations `on/2`, `clear/1` and the action `move/3` that has a success probability of 0.9. The box world dataset has the relations `bin/2`, `on/2`, `tin/2`, `can-drive/2` and the actions `drive/3`, `load/2`, `unload/2`. In the box world, `bin(B, C)` expresses a box is in a city, `on(B, T)` expresses a box is on a truck, and `tin(T, C)` expresses a truck is in a city. An atom `can-drive(C1, C2)` expresses a road that directly connects C1 and C2. A box can be loaded on (`load(B, T)`) or unloaded from a truck (`unload(B, T)`). A truck can travel from a city to another (`drive(T, C1, C2)`). The `load/2` and `unload/2` actions succeed with probability 0.9. The `drive/2` action succeeds with probability 0.8.

Q1: What formulae can pCTL-REBEL check in practice?

Although pCTL-REBEL can check any formulae of the pCTL language as discussed in Section 3.2, in practice, some formulae are more costly than others. This experiment aims at evaluating the computational costs of different pCTL formulae. Table 3 includes formulae ranging from classic planning properties, i.e. reachability properties (Kersting et al., 2004; Boutilier et al., 2001; Sanner and Boutilier, 2009), to more complex, nested pCTL properties, e.g. the following $\phi_{\text{nested}}(i, j)$ formula. Recall that $\phi_{\text{nested}}(4, 2)$ was discussed in Example 7, and its parse tree is in Figure 4. Table 3 also includes a formula that has an indefinite horizon, namely Property 5.

$$\phi_{\text{nested}}(i, j) = P_{\geq 0.5}[\text{cl}(\mathbf{a})U^{\leq i}(\text{on}(\mathbf{a}, \mathbf{b}) \wedge P_{\geq 0.8}[P_{\geq 0.9}[\mathbf{X} \text{cl}(\mathbf{e})]U^{\leq j} \text{on}(\mathbf{c}, \mathbf{d})])]$$

We now discuss Table 3 in detail. First, all formulae in Table 3 are checked against models that have *infinitely many* objects, which is possible only because pCTL-REBEL uses lifted inference. Second, pCTL-REBEL can handle classic planning tasks by formulating them as reachability formulae (Property 1,2,5, Table 3) or constrained reachability formulae (Property 3,6, Table 3). For example, Figure 5 visualizes the solution to Property 2 in Table 3. Third, pCTL-REBEL can check, in an infinite model, formulae that have an infinite horizon (Property 5, Table 3). This property cannot be evaluated by an explicit-state model checker as state space is infinite. Finally, pCTL-REBEL can handle nested formulae such as $\phi_{\text{nested}}(i, j)$ (Property 4, Table 3). Under a time-out of 1800 seconds, pCTL-REBEL can compute $\phi_{\text{nested}}(3, 1)$, which is cheaper than computing $\phi_{\text{nested}}(4, 2)$.

⁹ The experiments will be made public once the paper is accepted

¹⁰ http://ippc-2008.loria.fr/wiki/index.php/Results.html#Fully_Observable_Non-Deterministic_28FOND_29_track_2

¹¹ The ppddl models are first converted to the jani format by the ppddl2jani tool (attached in the datasets), then converted to the prism format by ePMC (Hahn et al., 2014).

¹² The ppddl models are translated to prolog models as PCTL-REBEL is implemented in prolog.

domain	property	formula	runtime(sec)
blocks	bounded reachability with det. actions	$P_{\geq 0.5}[F^{\leq 10}\text{on}(\mathbf{a}, \mathbf{b})]$	11 ^a
blocks	bounded reachability	$P_{\geq 0.5}[F^{\leq 10}\text{on}(\mathbf{a}, \mathbf{b})]$	226
blocks	bounded constrained reachability	$P_{\geq 0.5}[\text{on}(\mathbf{c}, \mathbf{d})U^{\leq 5}\text{on}(\mathbf{a}, \mathbf{b})]$	7.1
blocks	nested	$\phi_{\text{nested}}(3, 1)$	1554
box	unbounded reachability	$P_{\geq 0.5}[F \text{bin}(\mathbf{b1}, \text{paris})]$	0.7
box	bounded constrained reachability	$P_{\geq 0.5}[\text{tin}(\mathbf{t1}, \text{paris})U^{\leq 5}\text{bin}(\mathbf{b1}, \text{paris})]$	0.365

Table 3: PCTL-REBEL can perform a range of pCTL properties in the *infinite* blocks and box worlds under a time-out of 1800 seconds.

^a Kersting et al. (2004) reported a runtime of roughly 10 minutes for this task

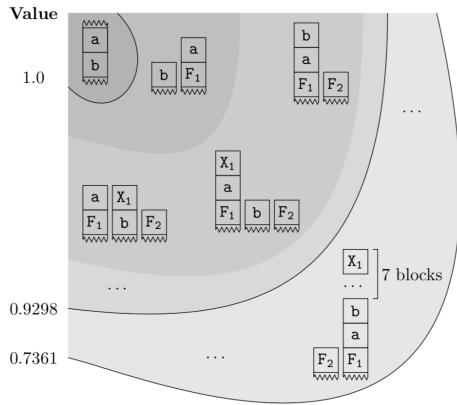


Fig. 5: The abstract states that satisfy $P_{\geq 0.5}[F^{\leq 10}\text{on}(\mathbf{a}, \mathbf{b})]$. Any state subsumed by $\{\text{on}(\mathbf{a}, \mathbf{b})\}$ is a goal state thus has a value 1. The states that are further away from the upper-left corner require more successful actions to reach $\{\text{on}(\mathbf{a}, \mathbf{b})\}$ thus have a smaller probability.

Nonetheless, Figure 6 visualizes the results of $\phi_{\text{nested}}(4, 2)$ by relaxing the time-out. More discussion about nested formulae will be in Q5.

Q2: How does pCTL-REBEL compare with state-of-the-art model checkers?

State-of-the-art model checkers STORM and PRISM do not operate on an RMDP as in Q1. They require the RMDP to be grounded, which leads to a state explosion. This experiment evaluates how well STORM and PRISM handle such state explosions in practice. In particular, we show that pCTL-REBEL, by lifting, performs better in mitigating state explosions in relational domains. While Q1 assumes that the given RMDP has an infinitely large domain, in order to compare pCTL-REBEL with STORM and PRISM, this section considers finite domains such that explicit-state models can be generated. Furthermore, we increment the domain size by 1 to track the trends in computation time for all model checkers.

The state explosion problem is illustrated in Figure 7, showing the number of ground states (Slaney and Thiébaux, 2001) and the corresponding ground transi-

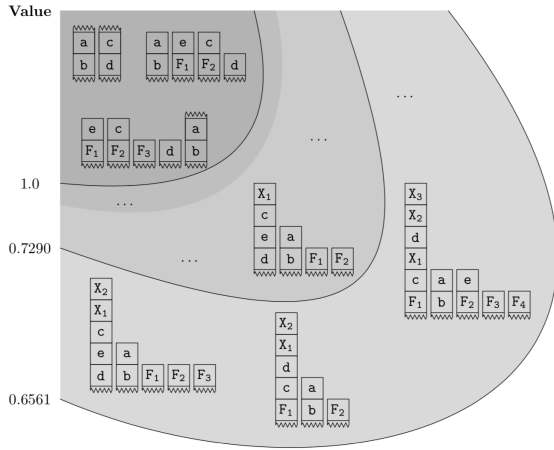


Fig. 6: Similar to Figure 5, this illustrates the abstract states that satisfy $\phi_{\text{nested}}(4, 2)$ in the blocks world.

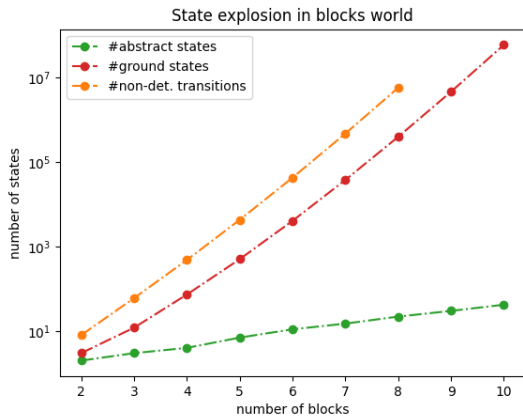


Fig. 7: The blocks world stats. When the number of blocks grows, the number of ground states and the number of non-deterministic transitions explode more than the minimum number of abstract states.

tions¹³ grow much faster than the minimum number of abstract states¹⁴. Even though all numbers grow exponentially, the minimum number of abstract states is the most resistant to the growth of the domain size.

We compare the scalability of the three model checkers (including 4 PRISM engines and 3 STORM engines) in the blocks world. We could essentially select any pCTL property for the comparison. Nonetheless, for simplicity, we use a simple reachability property $P_{\geq 0.5}[\mathbf{F}^{\leq 10} \text{on}(\mathbf{a}, \mathbf{b})]$ (Property 2, Table 3). We measure how the runtime grows with respect to the domain size, shown in Figure 8. The results show that pCTL-REBEL is more scalable than any PRISM or STORM engine in relational domains. Under a time-out of 1800 seconds, pCTL-REBEL can handle 15 blocks (6.6e13 ground states) whereas PRISM can handle at most 8 blocks

¹³ The number of ground transitions is obtained from the STORM model checker.

¹⁴ The minimum number of abstract states is the number of integer partitions that capture all relational structures. For example, for three blocks, there are three abstract states $\{\text{c1}(\mathbf{A}), \text{c1}(\mathbf{B}), \text{c1}(\mathbf{C})\}$, $\{\text{c1}(\mathbf{A}), \text{c1}(\mathbf{B}), \text{on}(\mathbf{B}, \mathbf{C})\}$ and $\{\text{c1}(\mathbf{A}), \text{on}(\mathbf{A}, \mathbf{B}), \text{on}(\mathbf{B}, \mathbf{C})\}$.

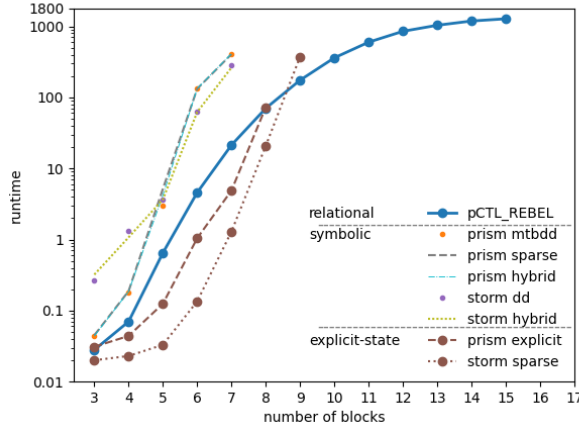


Fig. 8: The property $P_{\geq 0.5}[\mathbb{F}^{\leq 10} \text{on}(\mathbf{a}, \mathbf{b})]$ is checked by three model checkers. With a time-out of 1800 seconds, PRISM can handle at most 8 blocks with the explicit engine, STORM can handle at most 9 blocks with the sparse engine, and pCTL-REBEL can handle 15 blocks.

($4.0e5$ ground states) and STORM can handle at most 9 blocks ($4.6e6$ ground states).

We analyze Figure 8 to compare lifting with other optimization techniques. Observe that Figure 8 categorizes all engines in three groups: (1) *relational* model checking: pCTL-REBEL, (2) *symbolic* model checking: PRISM’s mtbdd, sparse, hybrid engines and STORM’s dd, hybrid engines, and (3) *explicit-state* model checking: PRISM’s explicit engine and STORM’s sparse engine. Compared to pCTL-REBEL, not only the explicit-state approaches but also the symbolic approaches are much less resilient to the growth of the domain size. Clearly, although symbolic model checking methods are known to be able to handle large domains, their ability of handling state explosions in relational domains is limited. Furthermore, unlike other engines, the runtime of pCTL-REBEL saturates when the domain grows to a certain extent. This phenomenon is a consequence of lifted inference, reflecting that the number of relational structures will eventually saturate.

Q3: How well does pCTL-REBEL handle indefinite-horizon formulae?

PCTL-REBEL suffers from runtime explosions when checking a formula that has an indefinite horizon. In specific, the runtime grows exponentially with respect to the horizon. This experiment aims at analyzing the cause of the runtime explosion. We use the reachability property $P_{\geq 0.5}[\mathbb{F} \text{ on}(\mathbf{a}, \mathbf{b})]$ that has an indefinite horizon.

To analyze the cause of the runtime explosion, we split the *iterations* of checking an indefinite-horizon formula into two stages: (1) the *state recognition stage* and (2) the *value convergence stage*. In the state recognition stage, pCTL-REBEL

identifies new satisfying states in each iteration. This stage ends when no more new satisfying states are found. In the value convergence stage, pCTL-REBEL does not discover new states, and only updates probabilities of all identified states until convergence. Notice that these two stages are just for analyzing purposes as pCTL-REBEL conducts exactly the same computation in all iterations in both stages.

We show the runtime growth of $P_{\geq 0.5}[\mathbf{F} \text{ on}(\mathbf{a}, \mathbf{b})]$ in Figure 9 with respect to the domain size. Figure 9(a) shows the runtime explodes with respect to *the number of required iterations*, and illustrates that most of the runtime is taken by the value convergence stage. Moreover, Figure 9(b) shows that within one checking process, later iterations are more expensive than earlier iterations. This means that although the job of the value convergence state is seemingly easier than the state recognition stage, in practice, it takes more time. This extra computation time comes from pCTL-REBEL’s attempts at identifying new states in the value convergence stage.

In summary, pCTL-REBEL suffers from a runtime explosion when handling indefinite formulae. The main cause is the redundant computation conducted in the value convergence stage. That is, pCTL-REBEL performs the same computation without being aware of the two-stage nature of the model checking process. This problem can be mitigated by optimizing the algorithm so as to reduce the time taken by the second stage. The optimization should benefit all indefinite-horizon formulae.

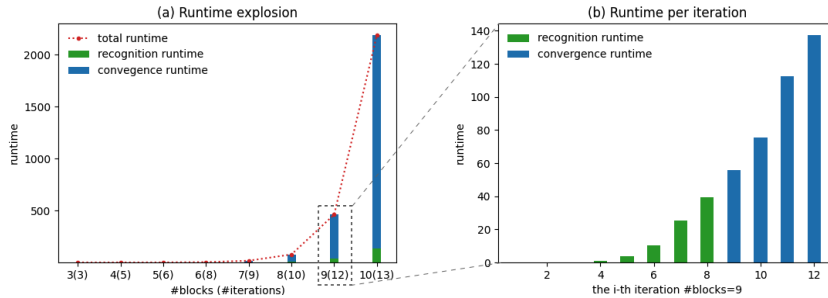


Fig. 9: The results obtained from checking the reachability property $P_{\geq 0.5}[\mathbf{F} \text{ on}(\mathbf{a}, \mathbf{b})]$ in the blocks world. (a) The runtime increases exponentially to the number of required iterations. Most of the runtime is taken by the value convergence stage (b) The later iterations are more expensive than the earlier iterations.

Q4: How well does pCTL-REBEL handle a complex relational transition function?

The use of abstract actions is key to the scalability of RMDPs as it allows for reasoning about a large group of ground actions as a whole. Roughly speaking, the more ground actions an abstract action captures, the more scalable pCTL-REBEL can be. For example, the blocks world domain has a simple move abstract

action (see Figure 1) that captures countless ground actions. Therefore, pCTL-REBEL is very scalable to the growth of the domain size. However, it is common to have constraints on the action, which significantly harms scalability. An action constraint harms scalability by prohibiting certain substitutions for an abstract action. In other words, it forces the abstract action to depend on the true identity of some objects. Action constraints break relational structures and result in non-symmetric transition functions. The resulting transition function consists of more rules and cannot be represented as compactly anymore. We call this constrained transition functions *relationally* complex.

The following example illustrates action constraints. Consider a box world domain where the action $\text{drive}(T, C1, C2)$ is constrained by the underlying road map below ¹⁵. A truck T can drive from $C1$ to $C2$ only if $\text{can-drive}(C1, C2)$ exists.

```

can-drive(city0,city2). can-drive(city2,city1). can-drive(city5,city4).
can-drive(city0,city1). can-drive(city2,city6). can-drive(city5,city1).
can-drive(city0,city6). can-drive(city3,city1). can-drive(city5,city3).
can-drive(city1,city0). can-drive(city3,city6). can-drive(city6,city0).
can-drive(city1,city6). can-drive(city3,city4). can-drive(city6,city1).
can-drive(city1,city3). can-drive(city3,city5). can-drive(city6,city2).
can-drive(city1,city4). can-drive(city4,city1). can-drive(city6,city3).
can-drive(city1,city2). can-drive(city4,city3). can-drive(city6,city4).
can-drive(city1,city5). can-drive(city4,city6).
can-drive(city2,city0). can-drive(city4,city5).

```

This experiment examines how well pCTL-REBEL handles complex relational transition functions. We use the reachability property $\phi = P_{\geq 0.5}[F \text{ bin}(b1, \text{city0})]$ and automatically generated road maps. Figure 10 shows the results that pCTL-REBEL can handle a road map containing at most 8 cities under a time-out of 1800 seconds. The runtime increases exponentially with respect to the size of the road network. Recall that without a road network, it takes only 0.7 seconds to check ϕ (see Property 5, Table 3). Since the road network makes the model dynamics more complex, more transition rules are required, resulting in runtime explosions. In short, pCTL-REBEL works the best when the dynamics of the objects is simple, e.g. in the blocks world or in the box world without a road map.

Q5: What are the computational costs of different pCTL operators?

It is crucial to know what kind of formulae pCTL-REBEL can handle efficiently when designing properties for a model. This experiment aims at giving insight into how different factors affect pCTL-REBEL's performance.

Some path operators are more costly than others. e.g. the X operator is the cheapest since it requires only one iteration. The reachability operator F is generally more expensive than the constrained reachability operator U since U is a stricter operator that prunes out more states. For example, checking $P_{\leq 0.9}[F^{\leq 7} \text{ on}(a, b)]$ (18.494 seconds) is more costly than checking $P_{\leq 0.9}[\text{on}(c, d) U^{\leq 7} \text{ on}(a, b)]$ (0.733 seconds). A nested formula is more costly than a flat formula since it requires several recursive checking processes.

¹⁵ The road networks are automatically generated by the box world dataset.

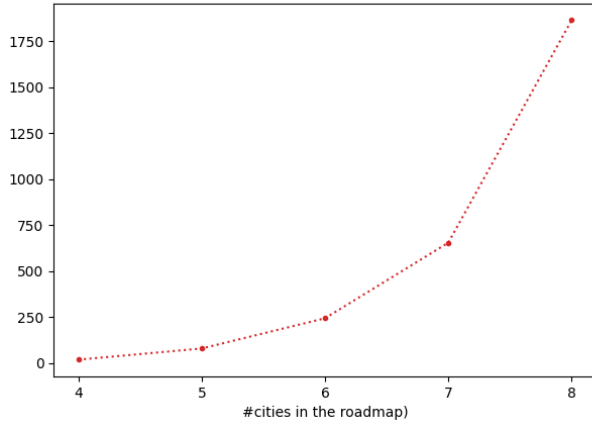


Fig. 10: The results obtained from checking the reachability property $P_{\geq 0.5}[F \text{ bin}(\mathbf{b1}, \text{city0})]$ in the box world. The runtime increases exponentially when the transition structure becomes more complex.

The formula structure also influences runtime. Specifically, changing an inner step bound has a larger impact than changing an outer step bound. Taking the formula $\phi_{\text{nested}}(i, j)$ (defined in Q1) for example, the parameters i and j are both step bounds for U formulae, but i is attached to an outer subformula and j is attached to the inner subformula. Hence, the total runtime is more sensitive to the change in j than to the change in i , as shown in Figure 11. The cause of this difference is that the inner formulae are computed before the outer formulae. By changing an inner step bound, the number of satisfying states for the inner subformula changes, which directly influence the overhead of computing the outer subformulae.

8 Related work

We have introduced pCTL-REBEL, a relational model checking technique that reasons at the relational level and mitigates the state explosion problem in relational MDPs. PCTL-REBEL extends existing frameworks. First, pCTL-REBEL extends pCTL model checking with relational representations. As far as the authors know, pCTL-REBEL is the first relational technique for pCTL model checking. Second, pCTL-REBEL extends the decidability results of infinite, state-bounded MDPs to the probabilistic setting.

To mitigate the state explosion problem in model checking, several *grouping* techniques have been proposed. Different from our work, the following methods exploit various types of symmetries other than relational structures. *Symbolic model checking* (McMillan, 1993) is a form of model checking that symbolically represents the state space and the set of states in which a formula holds. Symbolic model checking is often implemented using BDDs that are Boolean formulae in canonical

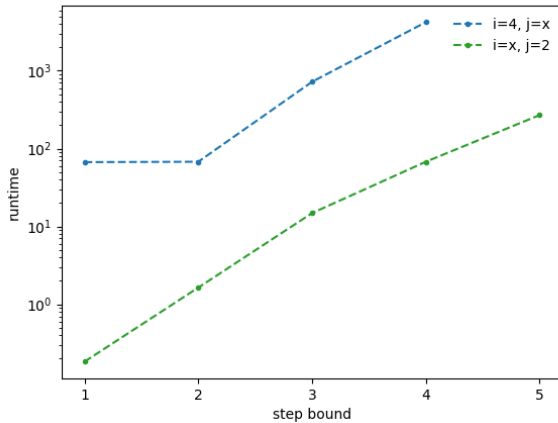


Fig. 11: The runtime of $\phi_{\text{nested}}(i, j)$ for different values of i and j . The argument i is the step bound of the outer U operator, and j is the step bound of the inner U operator. The runtime grows exponentially when either step bound increases. However, the runtime is more sensitive to the inner step bound j .

form, in which isomorphic subformulae are shared, to allow efficient operations on the sets of states that they represent. Symbolic model checking has been first used for model checking of plain non-probabilistic finite state systems (McMillan, 1993), then later adapted to the probabilistic case in Forejt et al. (2011). Contrary to our algorithm, these symbolic methods only apply to finite state MDPs. Another method to mitigate the state explosion problem is *abstraction-refinement*, which successively generates abstractions by partitioning the state space into *regions*. Therefore, the model checking algorithms handle regions instead of individual states (de Alfaro and Roy, 2007; Roy et al., 2008). These techniques focus on grouping states into regions which have similar local properties but such regions are usually not exact as the states in a region can have different dynamics and future evaluations. Similarly, Marthi (2007) uses *abstract MDP* to aggregate not only states but also actions by exploiting the temporal dependency of actions. In addition, *game-based abstraction* techniques construct abstractions of MDPs by merging concrete transitions into abstract transitions (Kattenbelt et al., 2008). Different from our work, for the abstractions, they calculate approximate upper and lower bounds instead of exact values.

Many studies integrate model checking techniques into other AI fields. In planning, authors compute the optimal policy at the first-order level using value iteration (Boutilier et al., 2001; Sanner and Boutilier, 2009; Kersting et al., 2004), which can be seen as a probabilistic reachability task, namely a special case of the general probabilistic model checking problem. PCTL-REBEL generalizes this special case with a temporal logic thus fits in the Planning as Model Checking paradigm (Giunchiglia and Traverso, 2000). In robotics, model checking techniques

are used for online motion planning (Lahijanian et al., 2012; Maly et al., 2013; He et al., 2015) where the states are labeled with propositions. These propositions result in a multidimensional state space that grows exponentially with the domain size (He et al., 2015), which is the exact problem that this paper tackles. Generally, our work can be applied a wide range of frameworks that have a large model, particularly in relational domains.

There is a fruitful literature on theoretical results of infinite systems. Particularly, the state-boundedness condition for decidable verification has been defined in first-order mu-calculus (Bagheri Hariri et al., 2013; Calvanese et al., 2018; De Giacomo et al., 2012; De Giacomo et al., 2015) and first-order CTL (Belardinelli et al., 2011, 2012). However, these studies focus on the non-probabilistic setting. Our work contributes to extending the theoretical results to stochastic models by introducing the abstract relational MDP. A limitation of the relational pCTL language is that it does not support quantification across conjunctions, which could be investigated as described in the work of Belardinelli et al. (2013).

9 Towards Safe Reinforcement Learning

Recently, a converging interest has emerged about the pursuit of general dynamic systems that can autonomously *reason* and *learn* by interacting with the environment (Giacomo, 2019). Since such systems require the ability of reasoning about first-order state representations and safety (Amodei et al., 2016), pCTL-REBEL fits in as a first-order model checking technique that is also particularly relevant in *safe reinforcement learning*.

Safe reinforcement learning is the process of learning a policy that maximizes the expected reward in domains where safety constraints must be respected. In particular, *safe exploration*, i.e. guaranteeing an agent’s safety in the exploration phase, is a nontrivial problem. Commonly used exploration strategies such as ϵ -greedy and softmax sometimes select a *random* action, which can result in catastrophic situations (Amodei et al., 2016). Existing research recognizes the importance of providing safety guarantees in reinforcement learning (Garcia and Fernández, 2015; Pecka and Svoboda, 2014; Lahijanian et al., 2012; Teichteil-Königsbuch, 2012; Sprauel et al., 2014; Giunchiglia and Traverso, 2000; Mason et al., 2018; Hasanbeig et al., 2019; Alshiekh et al., 2018; Jansen et al., 2020). However, most studies are based on techniques that require explicit state exploration whereas only few studies touch on safe exploration in relational domains (Driessens and Džeroski, 2004; Martínez et al., 2017) Furthermore, several studies have investigated augmenting reinforcement learning with model checking techniques, including preventing the agent from taking risky actions (Alshiekh et al., 2018; Jansen et al., 2020; Fulton and Platzer, 2018), synthesizing an initial safe *partial* policy for learning (Leonetti et al., 2012), learning a policy that maximizes the probability of satisfying a temporal formula (Hasanbeig et al., 2019), and shaping the reward function based on a temporal formula (De Giacomo et al., 2019). These approaches can be augmented with our work for scalability in large relational domains.

PCTL-REBEL can be used as a safe model-based reinforcement learning algorithm as it performs value iteration in a setting where goal states and safety constraints are formulated as pCTL formulae. More precisely, it derives the states

in an RMDP that satisfy a given (relational) pCTL formula that encodes goal states and safety constraints. Indeed, pCTL-REBEL tackles a special safe reinforcement learning task where the reward is 1 for goal states and 0 for all other states, and the discount factor is 1 (Kersting et al., 2004; Yoon et al., 2012; Otterlo, 2004). Formally, pCTL-REBEL in this paper

$$V_0^p(s) = \begin{cases} 0 & , \forall s \notin G \\ 1 & , \forall s \in G \end{cases}$$

$$V_{t+1}^p(s) = \max_a \sum_{s'} T(s, a, s') V_t^p(s')$$

is a special case of the standard Bellman operator

$$V_0^p(s) = \begin{cases} 0 & , \forall s \notin G \\ r & , \forall s \in G \end{cases}$$

$$V_{t+1}^p(s) = \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_t^p(s')]$$

where r is a reward value in \mathbb{R} .

It is an interesting avenue for further research to investigate representing general reward functions within this framework. Furthermore, pCTL-REBEL can be naturally combined with relational reinforcement learning (Dzeroski et al., 2001) to achieve safe relational reinforcement learning as both frameworks use relational representations and work at an abstract level rather than at the ground level.

10 Conclusions

We have introduced a framework for lifted model checking in relational domains. To this aim, relational Markov Decision Processes have been integrated with model checking principles for pCTL. The result is a very expressive framework for model checking in probabilistic planning domains that are relational, that is, involve objects as well as the relations among them. The framework is lifted in that it does not require to first ground the relational MDP and then exhaustively check all possible paths, but rather works at a more abstract relational level where variables are only instantiated whenever needed. The resulting algorithm is based on the relational Bellman operator REBEL. It is quite complex but manages to rather compactly compress enormous spaces of ground states in a couple of 10s of rules. The algorithm has not yet been optimized, and one route for further work is to combine pCTL-REBEL with expressive first order decision diagrams (Wang et al., 2008) to gain efficiency and to further compress the abstract states. Another route for further work is to explore the use of pCTL-REBEL for reasoning and reinforcement learning in safety-critical contexts.

Acknowledgements This work was supported by the FNRS-FWO joint programme under EOS No. 30992574. It has also received funding from the Flemish Government under the ‘‘Onderzoeksprogramma Artificiële Intelligentie (AI) Vlaanderen’’ programme, the EU H2020 ICT48 project ‘‘TAILOR’’ under contract #952215, the Wallenberg AI, Autonomous Systems and Software Program (WASP) funded by the Knut and Alice Wallenberg Foundation, and the KU Leuven Research fund.

References

- de Alfaro L, Roy P (2007) Magnifying-lens abstraction for markov decision processes. In: Damm W, Hermanns H (eds) *Computer Aided Verification*, Springer Berlin Heidelberg, Berlin, Heidelberg, pp 325–338
- Alshiekh M, Bloem R, Ehlers R, Könighofer B, Niekum S, Topcu U (2018) Safe reinforcement learning via shielding. In: *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), the 30th innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18)*, New Orleans, Louisiana, USA, February 2-7, 2018, pp 2669–2678
- Amodei D, Olah C, Steinhardt J, Christiano P, Schulman J, Mané D (2016) Concrete problems in ai safety. 1606.06565
- Bagheri Hariri B, Calvanese D, De Giacomo G, Deutsch A, Montali M (2013) Verification of relational data-centric dynamic systems with external services. In: *Proceedings of the 32Nd ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems*, ACM, New York, NY, USA, PODS '13, pp 163–174, DOI: 10.1145/2463664.2465221, URL <http://doi.acm.org/10.1145/2463664.2465221>
- Baier C, Katoen JP (2008) *Principles of Model Checking (Representation and Mind Series)*. The MIT Press
- Belardinelli F, Lomuscio A, Patrizi F (2011) Verification of deployed artifact systems via data abstraction. In: Kappel G, Maamar Z, Motahari-Nezhad HR (eds) *Service-Oriented Computing*, Springer Berlin Heidelberg, Berlin, Heidelberg, pp 142–156
- Belardinelli F, Lomuscio A, Patrizi F (2012) An abstraction technique for the verification of artifact-centric systems. 13th International Conference on the Principles of Knowledge Representation and Reasoning, KR 2012
- Belardinelli F, Lomuscio A, Patrizi F (2013) Verification of agent-based artifact systems. CoRR abs/1301.2678, URL <http://arxiv.org/abs/1301.2678>, 1301.2678
- Boutilier C, Reiter R, Price B (2001) Symbolic dynamic programming for first-order mdps. In: *Proceedings of the 17th International Joint Conference on Artificial Intelligence - Volume 1*, Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, IJCAI'01, pp 690–697, URL <http://dl.acm.org/citation.cfm?id=1642090.1642184>
- Van den Broeck G, Taghipour N, Meert W, Davis J, De Raedt L (2011) Lifted probabilistic inference by first-order knowledge compilation. In: *Proceedings of the Twenty-Second international joint conference on Artificial Intelligence, AAAI Press/International Joint Conferences on Artificial Intelligence; Menlo ...*, pp 2178–2185
- Calvanese D, Giacomo GD, Montali M, Patrizi F (2018) First-order μ -calculus over generic transition systems and applications to the situation calculus. *Information and Computation* 259:328 – 347, DOI: <https://doi.org/10.1016/j.ic.2017.08.007>, URL <http://www.sciencedirect.com/science/article/pii/S0890540117301487>, 22nd International Symposium on Temporal Representation and Reasoning
- De Giacomo G, Lespérance Y, Patrizi F (2012) Bounded situation calculus action theories and decidable verification. *Proc of KR 12*

- De Giacomo G, Lespérance Y, Patrizi F (2015) Bounded situation calculus action theories. CoRR abs/1509.02012, URL <http://arxiv.org/abs/1509.02012>, 1509.02012
- De Giacomo G, Iocchi L, Favorito M, Patrizi F (2019) Foundations for restraining bolts: Reinforcement learning with ltl/ldf restraining specifications. Proceedings of the International Conference on Automated Planning and Scheduling 29(1):128–136, URL <https://ojs.aaai.org/index.php/ICAPS/article/view/3549>
- De Raedt L, Kersting K, Natarajan S, Poole D (2016) Statistical relational artificial intelligence: Logic, probability, and computation. Synthesis Lectures on Artificial Intelligence and Machine Learning 10(2):1–189, DOI: 10.2200/S00692ED1V01Y201601AIM032, URL <https://doi.org/10.2200/S00692ED1V01Y201601AIM032>, <https://doi.org/10.2200/S00692ED1V01Y201601AIM032>
- Dehnert C, Junges S, Katoen JP, Volk M (2017) A storm is coming: A modern probabilistic model checker. In: Majumdar R, Kunčák V (eds) Computer Aided Verification, Springer International Publishing, Cham, pp 592–600
- Driessens K, Džeroski S (2004) Integrating guidance into relational reinforcement learning. Machine Learning 57:271–304, DOI: 10.1023/B:MACH.0000039779.47329.3a
- Džeroski S, De Raedt L, Driessens K (2001) Relational reinforcement learning. Machine learning 43(1-2):7–52
- Ferilli S, Fanizzi N, Mauro ND, Basile TMA (2002) Efficient theta-subsumption under object identity. In: In Atti del Workshop AI*IA su Apprendimento Automatico
- Forejt V, Kwiatkowska M, Norman G, Parker D (2011) Automated Verification Techniques for Probabilistic Systems, Springer Berlin Heidelberg, Berlin, Heidelberg, pp 53–113. DOI: 10.1007/978-3-642-21455-4_3, URL https://doi.org/10.1007/978-3-642-21455-4_3
- Fulton N, Platzer A (2018) Safe reinforcement learning via formal methods: Toward safe control through proof and learning. In: AAAI, pp 6485–6492, URL <https://www.aaai.org/ocs/index.php/AAAI/AAAI18/paper/view/17376>
- Gabbay DM (2003) Many-Dimensional Modal Logics: Theory and Applications. Elsevier North Holland
- García J, Fernández F (2015) A comprehensive survey on safe reinforcement learning. J Mach Learn Res 16:1437–1480
- Giacomo GD (2019) Queryable self-deliberating dynamic systems, URL <https://www.cse.ust.hk/pg/seminars/S19/giacomo.html>, iJCAI
- Giunchiglia F, Traverso P (2000) Planning as model checking. In: Biundo S, Fox M (eds) Recent Advances in AI Planning, Springer Berlin Heidelberg, Berlin, Heidelberg, pp 1–20
- Haddad S, Monmege B (2014) Reachability in MDPs: Refining Convergence of Value Iteration. In: 8th International Workshop on Reachability Problems (RP’14), Springer-Verlag, Oxford, United Kingdom, Proceedings of the 8th International Workshop on Reachability Problems (RP’14), vol 8762, pp 125–137, DOI: 10.1007/978-3-319-11439-2_10, URL <https://hal.archives-ouvertes.fr/hal-01091122>
- Hahn EM, Li Y, Schewe S, Turrini A, Zhang L (2014) IscasMC: A web-based probabilistic model checker. In: Nineteenth international symposium of the For-

- mal Methods Europe association (FM), Springer, Lecture Notes in Computer Science, vol 8442, pp 312–317
- Hasanbeig M, Kantaros Y, Abate A, Kroening D, Pappas GJ, Lee I (2019) Reinforcement learning for temporal logic control synthesis with probabilistic satisfaction guarantees. In: 2019 IEEE 58th Conference on Decision and Control (CDC), pp 5338–5343
- He K, Lahijanian M, Kavraki LE, Vardi MY (2015) Towards manipulation planning with temporal logic specifications. In: 2015 IEEE International Conference on Robotics and Automation (ICRA), pp 346–352, DOI: 10.1109/ICRA.2015.7139022
- Jansen N, Könighofer B, Junges S, Serban A, Bloem R (2020) Safe reinforcement learning using probabilistic shields. In: Konnov I, Kovacs L (eds) 31st International Conference on Concurrency Theory, CONCUR 2020, Schloss Dagstuhl-Leibniz-Zentrum für Informatik GmbH, Dagstuhl Publishing, pp 31–316, DOI: 10.4230/LIPIcs.CONCUR.2020.3
- Kattenbelt M, Kwiatkowska M, Norman G, Parker D (2008) Game-based probabilistic predicate abstraction in prism. *Electronic Notes in Theoretical Computer Science* 220(3):5–21, DOI: <https://doi.org/10.1016/j.entcs.2008.11.016>, URL <https://www.sciencedirect.com/science/article/pii/S1571066108004532>, proceedings of the Sixth Workshop on Quantitative Aspects of Programming Languages (QAPL 2008)
- Kersting K (2012) Lifted probabilistic inference. In: ECAI, pp 33–38
- Kersting K, De Raedt L (2004) Logical markov decision programs and the convergence of logical $td(\lambda)$. In: Camacho R, King R, Srinivasan A (eds) *Inductive Logic Programming*, Springer Berlin Heidelberg, Berlin, Heidelberg, pp 180–197
- Kersting K, Otterlo MV, De Raedt L (2004) Bellman goes relational. In: *Proceedings of the Twenty-first International Conference on Machine Learning*, ACM, New York, NY, USA, ICML '04, pp 59–, DOI: 10.1145/1015330.1015401, URL <http://doi.acm.org/10.1145/1015330.1015401>
- Kwiatkowska M, Norman G, Parker D (2011) PRISM 4.0: Verification of probabilistic real-time systems. In: Gopalakrishnan G, Qadeer S (eds) *Proc. 23rd International Conference on Computer Aided Verification (CAV'11)*, Springer, LNCS, vol 6806, pp 585–591
- Lahijanian M, Andersson SB, Belta C (2012) Temporal logic motion planning and control with probabilistic satisfaction guarantees. *IEEE Transactions on Robotics* 28(2):396–409, DOI: 10.1109/TR0.2011.2172150
- Leonetti M, Iocchi L, Patrizi F (2012) Automatic generation and learning of finite-state controllers. In: Ramsay A, Agre G (eds) *Artificial Intelligence: Methodology, Systems, and Applications*, Springer Berlin Heidelberg, Berlin, Heidelberg, pp 135–144
- Maly MR, Lahijanian M, Kavraki LE, Kress-Gazit H, Vardi MY (2013) Iterative temporal motion planning for hybrid systems in partially unknown environments. In: *Proceedings of the 16th International Conference on Hybrid Systems: Computation and Control*, Association for Computing Machinery, New York, NY, USA, HSCC '13, p 353–362, DOI: 10.1145/2461328.2461380, URL <https://doi.org/10.1145/2461328.2461380>
- Marthi B (2007) Automatic shaping and decomposition of reward functions. In: *Proceedings of the 24th International Conference on Machine Learning*, Association for Computing Machinery, New York, NY, USA, ICML '07, p 601–608, DOI:

- 10.1145/1273496.1273572, URL <https://doi.org/10.1145/1273496.1273572>
- Martínez D, Alenyà G, Torras C (2017) Relational reinforcement learning with guided demonstrations. *Artificial Intelligence* 247:295 – 312, DOI: <https://doi.org/10.1016/j.artint.2015.02.006>, URL <http://www.sciencedirect.com/science/article/pii/S0004370215000284>, special Issue on AI and Robotics
- Mason G, Calinescu R, Kudenko D, Banks A (2018) Assurance in Reinforcement Learning Using Quantitative Verification, Springer International Publishing, Cham, pp 71–96. DOI: 10.1007/978-3-319-66790-4_5, URL https://doi.org/10.1007/978-3-319-66790-4_5
- McMillan KL (1993) *Symbolic Model Checking*, Springer US, Boston, MA, pp 25–60. DOI: 10.1007/978-1-4615-3190-6_3, URL https://doi.org/10.1007/978-1-4615-3190-6_3
- Nienhuys-Cheng SH, Wolf Rd (1997) *Foundations of Inductive Logic Programming*. Springer-Verlag, Berlin, Heidelberg
- Otterlo MV (2004) Reinforcement learning for relational mdps. In: In Proceedings of the Machine Learning Conference of Belgium and the Netherlands
- Pecka M, Svoboda T (2014) Safe exploration techniques for reinforcement learning – an overview. In: Hodicky J (ed) *Modelling and Simulation for Autonomous Systems*, Springer International Publishing, Cham, pp 357–375
- Roy P, Parker D, Norman G, De Alfaro L (2008) Symbolic magnifying lens abstraction in markov decision processes. pp 103–112, DOI: 10.1109/QEST.2008.41
- de Salvo Braz R, Amir E, Roth D (2007) Lifted first-order probabilistic inference. *Statistical relational learning* p 433
- Sanner S, Boutelier C (2009) Practical solution techniques for first-order mdps. *Artificial Intelligence* 173(5):748 – 788, DOI: <https://doi.org/10.1016/j.artint.2008.11.003>, URL <http://www.sciencedirect.com/science/article/pii/S0004370208001884>, advances in Automated Plan Generation
- Slaney J, Thiébaux S (2001) Blocks world revisited. *Artificial Intelligence* 125(1):119 – 153, DOI: [https://doi.org/10.1016/S0004-3702\(00\)00079-5](https://doi.org/10.1016/S0004-3702(00)00079-5), URL <http://www.sciencedirect.com/science/article/pii/S0004370200000795>
- Sprauel J, Kolobov A, Teichteil-Königsbuch F (2014) Saturated path-constrained mdp: Planning under uncertainty and deterministic model-checking constraints. In: *Twenty-Eighth AAAI Conference on Artificial Intelligence*, AAAI Press, URL <https://www.microsoft.com/en-us/research/publication/saturated-path-constrained-mdp-planning-uncertainty-deterministic-model-checking-constraints/>
- Teichteil-Königsbuch F (2012) Path-Constrained Markov Decision Processes: bridging the gap between probabilistic model-checking and decision-theoretic planning. In: *20th European Conference on Artificial Intelligence (ECAI 2012)*, MONTPELLIER, France, URL <https://hal-onera.archives-ouvertes.fr/hal-01060349>
- Wang C, Joshi S, Khaldon R (2008) First order decision diagrams for relational mdps. *Journal of Artificial Intelligence Research* 31:431–472
- Yoon SW, Fern A, Givan R (2012) Inductive policy selection for first-order mdps. 1301.0614