

POWER

Technical Report 2018-01

Title: **Discounted Markov Automata**

Author: Yuliya Butkova

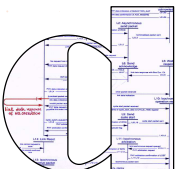
Report Number: 2018-01

ERC Project: Power to the People. Verified.

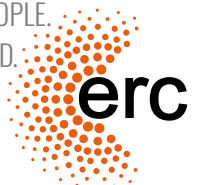
ERC Project ID: 695614

Funded Under: H2020-EU.1.1. – EXCELLENT SCIENCE

Host Institution: Universität des Saarlandes, Dependable Systems and Software
Saarland Informatics Campus



POWER TO THE PEOPLE.
VERIFIED.



Discounted Markov Automata^{*}

Yuliya Butkova

Saarland University, Saarbrücken, Germany
butkova@cs.uni-saarland.de

Abstract. Markov automata (MA) are a rich modelling formalism for complex systems combining compositionality with probabilistic choices and continuous stochastic timing. Model checking algorithms for different classes of properties involving probabilities and rewards have been devised for MA, opening up a spectrum of applications in dependability engineering and artificial intelligence, reaching out into economy and finance. In the latter more general contexts, several quantities of considerable importance are based on the idea of discounting reward expectations, so that the near future is more important than the far future. This paper introduces the expected discounted reward value for MA and develops effective iterative algorithms to quantify it, based on value- as well as policy-iteration. To arrive there, we reduce the problem to the computation of expected discounted rewards and expected total rewards in Markov decision processes. This allows us to adapt well-known algorithms to the MA setting. Experimental results clearly show that our algorithms are efficient and scale to MA with hundred thousands of states.

1 Introduction

The design and analysis of complex systems operating in uncertain environments requires a powerful modelling language. It is desirable to support compositionality for constructing large models from individual components; nondeterminism for abstraction and representing unknown behaviour of the environment; continuous stochastic timing and probabilistic choices. *Markov automata* (MA) [8] combine all these aspects in one formalism. They have been extended to express costs and rewards, yielding Markov reward automata [11]. Efficient algorithms for the automatic analysis of Markov (reward) automata are available for a broad range of properties like time- and cost-bounded reachability probabilities [10,13], expected rewards [11], long-run averages [10,5] and properties expressed in the temporal logic CSL [12]. Tool support is available: Both IMCA [9] and Storm [6] support Markov automata model checking. This makes Markov automata well suited not only as a modelling formalism by itself, but also as the semantical foundation of higher-level formalisms like dynamic fault trees [4] and generalized stochastic Petri nets [7].

^{*} This work was partly supported by the ERC Advanced Grant POWVER (695614) and by the Sino-German project CAP (GZ 1023).

All the measures considered so far for Markov reward automata do not make a difference between a unit of reward being accumulated early or late along the evolution of the system. In economics, in artificial intelligence, and in the theory of optimal control [3], it is a well-understood practice to discount the future, that is, to give more weight to the near future than to the far away future. This view is natural in model checking quantitative temporal logic properties of systems [2,1], including continuous-time Markov decision processes (CTMDPs) [16]. It appears equally natural for Markov automata, but as yet, there is neither a theory nor an algorithmic approach for discounting in Markov reward automata.

The present paper investigates discounting on MRA. We first settle the foundational basis of what discounting actually means for Markov automata. Due to the continuous nature of time in MRA we define discounting analogously to the way it is defined for CTMDPs.

Our findings are rooted in the observation that we can view any MRA as a representation encoding a possibly exponentially larger CTMDP, preserving discounted reward values. This enables one to quantify the discounted reward in MRA by computing the respective value on its value-preserving CTMDP, however at the price of possibly exponential time and space requirements.

Overall our approach has similarities in spirit to the one introduced to quantify long-run average rewards on MRA [5], but the constructions needed have to be entirely different due to the dependency of the discounted reward on time. Instead of the naïve approach, we show that the exponential blow-up can be avoided by recognising that the value requiring exponentially many computational steps as the expected total reward in a specific linear-sized discrete-time MDP. Using classic dynamic programming for the latter then turns the exponential naïve approach into an effective polynomial characterisation. In this way we derive the Bellman equation characterising the expected total reward in the presence of discounting in MRA. The Bellman equation in turn is the basis for value- and policy-iteration algorithms quantifying the discounted reward on MRA. The efficiency of the approach is demonstrated with examples of MRAs with hundreds of thousands of states.

2 Foundations

Given a finite set S , a *probability distribution* over S is a function $\mu : S \rightarrow [0, 1]$ with $\sum_{s \in S} \mu(s) = 1$. We denote the set of all probability distributions over S by $\text{Dist}(S)$. ξ_s is the *Dirac distribution* on s , i.e. $\xi_s(s) = 1$ and $\xi_s(s') = 0$ for $s' \neq s$.

Definition 1. A Markov reward automaton (MRA) \mathcal{M} is a tuple $\mathcal{M} = (S, s_{\text{init}}, \text{Act}, \hookrightarrow, \rightsquigarrow, \mathbf{r}, \rho)$ s.t. S is a finite set of states; $s_{\text{init}} \in S$ is the initial state; Act is a finite set of actions; $\hookrightarrow \subseteq S \times \text{Act} \times \text{Dist}(S)$ is a finite probabilistic transition relation; $\rightsquigarrow \subseteq S \times \mathbb{R}_{>0} \times S$ is a finite Markovian transition relation; $\mathbf{r} : \hookrightarrow \rightarrow \mathbb{R}_{\geq 0}$ is a transition reward function; and $\rho : S \rightarrow \mathbb{R}_{\geq 0}$ is a state reward function.

We abbreviate $(s, \alpha, \mu) \in \hookrightarrow$ by $s \xrightarrow{\alpha} \mu$ and write $s \xrightarrow{\lambda} s'$ instead of $(s, \lambda, s') \in \rightsquigarrow$. $Act(s) = \{\alpha \in Act \mid \exists \mu \in \text{Dist}(S) : s \xrightarrow{\alpha} \mu\}$ denotes the set of actions that are enabled in state $s \in S$. A state s is *probabilistic (Markovian)*, if it has at least one probabilistic (Markovian) transition $s \xrightarrow{\alpha} \mu$ ($s \xrightarrow{\lambda} s'$, resp.). States can be both probabilistic and Markovian. We denote the set of probabilistic states by $PS_{\mathcal{M}}$ and the Markovian states by $MS_{\mathcal{M}}$. We assume w.l.o.g. that actions of probabilistic transitions of a state are pairwise different¹. Therefore we will write $r(s, \alpha)$ instead of $r(s, \alpha, \mu)$. The successors of a state $s \in S$ are given by $\text{succ}(s) = \{s' \in S \mid \exists \alpha \in Act \exists \mu \in \text{Dist}(S) : s \xrightarrow{\alpha} \mu \wedge \mu(s') > 0 \vee \exists \lambda \in \mathbb{R}_{>0} : s \xrightarrow{\lambda} s'\}$ and its predecessors by $\text{pred}(s) = \{s' \in S \mid s \in \text{succ}(s')\}$.

For a Markovian state $s \in MS_{\mathcal{M}}$, the value $R(s, s') := \sum_{(s, \lambda, s') \in \rightsquigarrow} \lambda$ is called the *transition rate* from s to s' . The *exit rate* of a Markovian state s is $E(s) := \sum_{s' \in S} R(s, s')$. We require $E(s) < \infty$ for all $s \in MS_{\mathcal{M}}$.

For $s \in PS_{\mathcal{M}}$ with $s \xrightarrow{\alpha} \mu$ for some α , we set $\mathbb{P}[s, \alpha, s'] := \mu(s')$. For $s \in MS_{\mathcal{M}}$ with $E(s) > 0$, the branching probability distribution when leaving the state through a Markovian transition is denoted by $\mathbb{P}[s, \cdot] \in \text{Dist}(S)$ and defined by $\mathbb{P}[s, s'] := R(s, s')/E(s)$.

The evolution of an MRA starts in its initial state. Whenever the system encounters a Markovian state $s \in MS_{\mathcal{M}}$, its sojourn time in s is governed by an exponential distribution, i.e. the probability of leaving s within $t \geq 0$ time units is given by $1 - e^{-E(s) \cdot t}$, after which the next state is chosen according to $\mathbb{P}[s, \cdot]$.

The behaviour of the system in probabilistic states is different. In this paper we consider *closed* MRA, which are not subject to further composition operations that could delay the execution of probabilistic transitions. Therefore we can make the usual *urgency assumption: Probabilistic transitions happen instantaneously*. The residence time in probabilistic states is therefore always 0. Whenever the system is in state s with $Act(s) \neq \emptyset$ and an action $\alpha \in Act(s)$ is chosen, the successor s' is selected according to the distribution $\mathbb{P}[s, \alpha, \cdot]$ and the system moves instantaneously from s to s' . As the execution of a probabilistic transition is instantaneous and because the probability that a Markovian transition is triggered immediately is 0, we can assume that the probabilistic transitions take precedence over Markovian transitions. We therefore assume $PS_{\mathcal{M}} \cap MS_{\mathcal{M}} = \emptyset$. The way Markov automata choose actions will be covered shortly.

During its evolution a Markov reward automaton collects rewards. The transition reward $r(s, \alpha)$ is granted immediately for taking the (probabilistic)

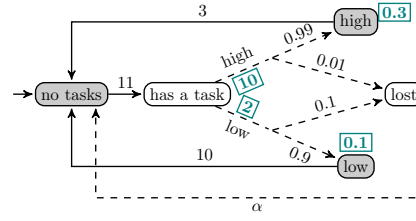


Fig. 1: An example MRA

¹ This can be achieved by renaming the actions and does not affect compositionality properties of MRA due to the fact that only *closed* MRA are considered in this work.

transition $s \xrightarrow{\alpha} \mu$, while the state rewards are accumulated over time, i.e. for staying in a state s for $t > 0$ time units, a reward of $\rho(s) \cdot t$ is granted.

Example 1. Figure 1 shows an example MRA. Grey and white colouring of states indicates the sets $MS_{\mathcal{M}}$ and $PS_{\mathcal{M}}$, resp.; they are disjoint here. Labels *high*, *low*, and α denote actions. Dashed transitions are probabilistic; solid transitions are Markovian. We omitted Dirac distributions from probabilistic transitions. Rewards associated with states and transitions are depicted as numbers in green rectangular frames.

The MRA models a task processing system aiming at maximising the revenue on the long run. Tasks arrive at rate 11; this is modelled by a Markovian transition with rate 11. Whenever there is a task to process, the system decides whether to handle it with high or low reliability. In the former case, the system receives an immediate reward of 10, modelled by $r(\text{has a task, high}) = 10$. A low reliability task produces a reward of 2 only. The tasks are sent for processing to a remote server over lossy channels. The high reliability channel loses tasks with probability 0.01, while low reliability tasks are lost ten times more often. Whenever a task is lost, no further reward for it is paid. Processing high reliability tasks takes more time, which is modelled with exit rate 3 of state *high*, however, it generates reward proportional to the processing time (modelled by state reward of 0.3). Low reliability tasks are faster to process but produce less reward.

An MRA is *non-Zeno* iff no *maximal end component* (see [10]) of only probabilistic states is reachable with probability > 0 . This excludes models in which there is a chance to get trapped in an infinite number of transitions occurring in finite time. In this work, similarly to [10], we restrict ourselves to non-Zeno models; Zenoness is typically considered a modelling error.

For ease of representation, we additionally assume that all states have at least one outgoing transition. This can be easily achieved by adding a Markovian self-loop to the state, with reward 0 and arbitrary non-zero rate.

Paths and Schedulers. A (*timed*) *path* in \mathcal{M} is a finite or infinite sequence $\pi = s_0 \xrightarrow{\alpha_0, t_0} s_1 \xrightarrow{\alpha_1, t_1} \dots \xrightarrow{\alpha_k, t_k} s_{k+1} \xrightarrow{\alpha_{k+1}, t_{k+1}} \dots$ with $s_0 = s_{\text{init}}$, and for all $i \geq 0$: $s_i \in S$, $t_i \in \mathbb{R}_{\geq 0}$, and $\alpha_i \in \text{Act} \cup \{\perp\}$. Here $s_i \xrightarrow{\alpha_i, 0} s_{i+1}$ s.t. $\alpha_i \in \text{Act}(s_i)$ is a probabilistic transition via action α_i , and $s_i \xrightarrow{\perp, t_i} s_{i+1}$ s.t. $t_i > 0$ and there exists a transition $s_i \xrightarrow{\lambda} s_{i+1}$, denotes a Markovian transition with sojourn time t_i in state s_i . We define $\pi[i] := s_i$, $\alpha[\pi, i] := \alpha_i$ and for an infinite path π and $k \geq 0$, the elapsed time $\tau[\pi, k]$ until entering $\pi[k]$ is defined by $\tau[\pi, 0] := 0$ and $\tau[\pi, k] := t_0 + \dots + t_{k-1}$. Whenever it is clear from the context, we omit π and just use $\alpha[i]$ and $\tau[k]$ instead. The set of all finite (infinite) paths of \mathcal{M} is denoted by $\text{Paths}_{\mathcal{M}}^*$ ($\text{Paths}_{\mathcal{M}}$). The length $|\pi|$ of a finite path π is the number of its transitions; its last state is denoted by $\pi\downarrow$.

In order to resolve the nondeterminism in probabilistic states of an MRA, we need the notion of a scheduler. A (*measurable*) *scheduler* (or *policy*) $\sigma : \text{Paths}_{\mathcal{M}}^* \rightarrow \text{Dist}(\hookrightarrow)$ is a measurable function, s.t. $\sigma(\pi)$ assigns positive probability only to

transitions $(\pi \downarrow, \alpha, \mu) \in \hookrightarrow$, for some α, μ . The set of all measurable schedulers is denoted by $GM_{\mathcal{M}}$. A *(deterministic) stationary scheduler* is a function $\sigma : PS_{\mathcal{M}} \rightarrow \hookrightarrow$, s. t. $\sigma(s)$ chooses only from transitions $(s, \alpha, \mu) \in \hookrightarrow$, for some α, μ .

For the definition of the probability measure on MRA we refer to [14, Sect. 3.2].

Markov Decision Processes.

Definition 2. A continuous-time Markov decision process (CTMDP) is a tuple $\mathcal{C} = (S, s_{\text{init}}, Act, R)$, where S is a finite set of states, $s_{\text{init}} \in S$ is an initial state, Act is a finite set of actions, and $R : S \times Act \times S \rightarrow \mathbb{R}_{\geq 0}$ is a rate function.

The set $Act(s) = \{\alpha \in Act \mid \exists s' \in S : R(s, \alpha, s') > 0\}$ is the set of *enabled actions* in state s . A path in a CTMDP is a finite or infinite sequence $\pi = s_0 \xrightarrow{\alpha_0, t_0} s_1 \xrightarrow{\alpha_1, t_1} \dots \xrightarrow{\alpha_{k-1}, t_{k-1}} s_k \dots$, where $s_0 = s_{\text{init}}$, $\alpha_i \in Act(s_i)$ and t_i denotes the residence time of the system in state s_i . $E(s, \alpha) := \sum_{s' \in S} R(s, \alpha, s')$ and $\mathbb{P}_{\mathcal{C}}[s, \alpha, s'] := \frac{R(s, \alpha, s')}{E(s, \alpha)}$. The notions of $Paths_{\mathcal{C}}^*$, $Paths_{\mathcal{C}}$, $\pi \downarrow$ and schedulers are defined analogously to corresponding definitions for an MRA. A *reward structure* on a CTMDP \mathcal{C} is a tuple $(\rho_{\mathcal{C}}, r_{\mathcal{C}})$, where $\rho_{\mathcal{C}} : S \rightarrow \mathbb{R}_{\geq 0}$ and $r_{\mathcal{C}} : S \times Act \rightarrow \mathbb{R}_{\geq 0}$.

The counterpart of CTMDPs and MRA in discrete time are (discrete-time) Markov decision processes:

Definition 3. A Markov decision process (MDP) is a tuple $\mathcal{D} = (S_{\mathcal{D}}, s_{\text{init}}, Act_{\mathcal{D}}, \mathbb{P}_{\mathcal{D}})$ where $S_{\mathcal{D}}$ is a finite set of states, s_{init} is the initial state, $Act_{\mathcal{D}}$ is a finite set of actions and $\mathbb{P}_{\mathcal{D}} : S_{\mathcal{D}} \times Act_{\mathcal{D}} \rightarrow \text{Dist}(S_{\mathcal{D}})$ is a probabilistic transition function.

The definitions of paths, schedulers, etc. are discrete analogues of those definitions for CTMDPs. A reward structure on an MDP is a function $r_{\mathcal{D}} : S_{\mathcal{D}} \times Act_{\mathcal{D}} \rightarrow \mathbb{R}_{\geq 0}$.

In the following a special subclass of MDPs – acyclic MDPs – will be of particular importance. A state of an MDP is called *terminal* if all its outgoing transitions are self-loops with probability 1 and reward 0. We call an MDP *acyclic* if the self-loops of terminal states are the only loops appearing in the MDP.

3 Discounted Reward for Markov Automata

In this section, we define the discounted reward value for Markov reward automata and consider its relation to discounted rewards on CTMDPs.

3.1 Continuous Discounting

Markov automata are a continuous-time model and we therefore define discounting in a classical way via the continuous exponential decay over time. Yet special care has to be taken when dealing with probabilistic states due to the fact that time in those states does not pass. Essentially the definition is lifted to the MRA setting from CTMDPs [18].

Let $\mathcal{M} = (S, s_{\text{init}}, Act, \hookrightarrow, \rightsquigarrow, r, \rho)$ be a Markov reward automaton, $\beta > 0$, and $\pi = s_0 \xrightarrow{\alpha_0, t_0} s_1 \xrightarrow{\alpha_1, t_1} \dots \xrightarrow{\alpha_k, t_k} s_{k+1} \xrightarrow{\alpha_{k+1}, t_{k+1}} \dots$ an infinite path in \mathcal{M} .

Then $\text{rew}_{\mathcal{M},\beta}^N(\pi)$ is the *discounted reward with rate β* of the path π within $N \in \mathbb{N}$ steps, where $\text{rew}_{\mathcal{M},\beta}^0(\pi) := 0$ and

$$\text{rew}_{\mathcal{M},\beta}^N(\pi) := \sum_{k=0}^{N-1} \left[e^{-\beta \cdot \tau[k]} \cdot r(s_k, \alpha_k) + \int_{\tau[k]}^{\tau[k]+t_k} e^{-\beta \cdot t} \cdot \rho(s_k) dt \right].$$

Example 2. Consider the MRA from Fig. 1 and its path $\pi = (nt) \xrightarrow{\perp, t_1} (ht) \xrightarrow{high, 0} (lost) \xrightarrow{\alpha, 0} (nt) \xrightarrow{\perp, t_2} (ht) \longrightarrow \dots$, where (nt) stands for *(no tasks)* and (ht) for *(has a task)*. Then the discounted reward collected over this path for $N = 4$ is:

$$\text{rew}_{\mathcal{M},\beta}^4(\pi) = \int_0^{t_1} \rho(nt) \cdot e^{-\beta \cdot \tau} d\tau + e^{-\beta \cdot t_1} (r(ht, high) + r(lost, \alpha)) + \int_{t_1}^{t_1+t_2} \rho(nt) \cdot e^{-\beta \cdot \tau} d\tau.$$

The *optimal expected cumulative discounted reward* (or just *discounted reward*) in \mathcal{M} with *discount rate β* is:

$$\text{dR}_{\mathcal{M},\beta}^{\text{opt}} := \text{opt}_{\sigma \in GM} \left\{ \lim_{N \rightarrow \infty} \mathbf{E}_{\sigma}[\text{rew}_{\mathcal{M},\beta}^N] \right\} = \text{opt}_{\sigma \in GM} \left\{ \lim_{N \rightarrow \infty} \int_{\text{Paths}_{\mathcal{M}}} \text{rew}_{\mathcal{M},\beta}^N(\pi) \cdot \Pr_{\mathcal{M},\sigma}[d\pi] \right\},$$

where $\text{opt} \in \{\sup, \inf\}$. A scheduler σ is called *optimal for $\text{dR}_{\mathcal{M},\beta}^{\text{opt}}$* , if it attains optimum in this equation. In the following, $\text{dR}_{\mathcal{M},\beta}^{\text{opt}}(s)$ denotes the discounted reward collected in \mathcal{M} assuming that the initial state is s . Whenever it is clear from the context, we omit the subscripts and use notation dR^{opt} .

Lemma 1. *The value $\text{dR}_{\mathcal{M},\beta}^{\text{opt}}$ exists.*

3.2 Relation to Discounted Rewards on CTMDP

We now show that for each MRA there exists a (possibly exponentially larger) CTMDP that preserves the discounted reward property. We first need to introduce *uniform* and *normalised* MRA:

Definition 4. *An MRA \mathcal{M} is uniform if $\exists \eta \in \mathbb{R}_{>0}$, s. t. $\forall s \in MS_{\mathcal{M}} : E(s) = \eta$.*

Definition 5. *An MRA \mathcal{M} is called normalised if*

1. *the initial state of \mathcal{M} is probabilistic;*
2. *every Markovian state s has only probabilistic predecessors: $\text{pred}(s) \subseteq PS_{\mathcal{M}}$;*
3. *probabilistic states of \mathcal{M} have either only probabilistic or only Markovian predecessors: $\forall s \in PS_{\mathcal{M}} : \text{pred}(s) \subseteq PS_{\mathcal{M}} \vee \text{pred}(s) \subseteq MS_{\mathcal{M}}$.*

Lemma 2. *For any MRA \mathcal{M} , $\eta \geq \max_{s \in S} E(s)$ there exists a uniform normalised MRA $\overline{\mathcal{M}}_{\eta}$, s. t. $\text{dR}_{\overline{\mathcal{M}}_{\eta},\beta}^{\text{opt}} = \text{dR}_{\mathcal{M},\beta}^{\text{opt}}$ and its size is linear in the size of \mathcal{M} .*

Informally, $\overline{\mathcal{M}}_\eta$ is obtained by first uniformising the Markovian states. This is performed via the well-known approach from [17] by adding self-loop transitions to them. Then this uniform MRA is normalised by introducing probabilistic states of zero reward (i) in between each pair of states that violate Properties 2 or 3 of the definition above, or (ii) as a new initial state as detailed in Appendix A.3.

In the following, we assume that the MRA at hand is uniform and normalised and show how to construct a value-preserving CTMDP for it. Before proceeding we need to introduce some notation:

- $\Pi_{\setminus B}(s, s')$ is the set of all *untimed* paths $\pi = s \xrightarrow{\alpha} s_1 \xrightarrow{\alpha_1} \dots s_k \xrightarrow{\alpha_k} s'$ (paths of \mathcal{M} with abstracted timing information), such that $\forall i = 1..k, s_i \notin B$;
- $PS_{\setminus B}(s)$ is the set of states containing s and all states $s' \in PS \setminus B$ that are related to s via the transitive closure of relation \hookrightarrow ;
- $\mathbb{P}[\pi] := \prod_{i=1}^{|\pi|-1} \mathbb{P}[\pi[i], \alpha[i], \pi[i+1]]$, $r(\pi) := \sum_{i=0}^{|\pi|-1} r(\pi[i], \alpha[i])$, $\rho(\pi) := \sum_{i=0}^{|\pi|-1} \rho(\pi[i])$.

Value-Preserving CTMDP. Let $\mathcal{M} = (S, s_{\text{init}}, Act, \hookrightarrow, \rightsquigarrow, r, \rho)$ be a uniform normalised MRA with exit rate η . We define the CTMDP $\mathcal{C}(\mathcal{M}) := (S_{\mathcal{C}}, s_{\text{init}}, Act_{\mathcal{C}}, R_{\mathcal{C}})$ and reward structure $(\rho_{\mathcal{C}}, r_{\mathcal{C}})$ as follows:

- $S_{\mathcal{C}}$:** The state space of $\mathcal{C}(\mathcal{M})$ is the set $S_{\mathcal{C}} \subseteq PS$ that contains the initial state s_{init} and all probabilistic states of \mathcal{M} that are successors of a Markovian state in \mathcal{M} : $S_{\mathcal{C}} = \{s \in PS \mid s = s_{\text{init}} \text{ or } \exists s' \in MS : s' \rightsquigarrow s\}$. We define the set of *marked* states as $S_{\text{mrk}} := S_{\mathcal{C}}$.
- $Act_{\mathcal{C}}$:** An action of a state s in this CTMDP is a mapping $A : PS_{\setminus S_{\text{mrk}}}(s) \rightarrow Act$, such that $A(s') \in Act(s')$. Then the set of all enabled actions $Act_{\mathcal{C}}(s)$ is the set of all possible functions A , and $Act_{\mathcal{C}} = \bigcup_{s \in S_{\text{mrk}}} Act_{\mathcal{C}}(s)$.
- $R_{\mathcal{C}}$:** Let $s, s' \in S_{\text{mrk}}$ and $\Pi_{\setminus S_{\text{mrk}}}(s, A, s') \subseteq \Pi_{\setminus S_{\text{mrk}}}(s, s')$ be the set of all paths from $\Pi_{\setminus S_{\text{mrk}}}(s, s')$ that select actions according to A , i. e. for each path $\pi : \pi[i] \in PS \Rightarrow \alpha[i] = A(\pi[i])$. Then $R_{\mathcal{C}}(s, A, s') := \eta \cdot \sum_{\pi \in \Pi_{\setminus S_{\text{mrk}}}(s, A, s')} \mathbb{P}[\pi]$.
- $\rho_{\mathcal{C}}$:** The state reward of a state s in $\mathcal{C}(\mathcal{M})$ is the expected state reward gathered in \mathcal{M} on paths between s and any other $s' \in S_{\text{mrk}}$ that is a successor of s in $\mathcal{C}(\mathcal{M})$: $\rho_{\mathcal{C}}(s) := \sum_{s' \in S_{\text{mrk}}} \sum_{\pi \in \Pi_{\setminus S_{\text{mrk}}}(s, s')} \rho(\pi) \cdot \mathbb{P}[\pi]$.
- $r_{\mathcal{C}}$:** The transition reward of a state s and action A in $\mathcal{C}(\mathcal{M})$ is the expected transition reward gathered in \mathcal{M} on paths between s and any other $s' \in S_{\text{mrk}}$ that is a successor of s in $\mathcal{C}(\mathcal{M})$: $r_{\mathcal{C}}(s, A) := \sum_{s' \in S_{\text{mrk}}} \sum_{\pi \in \Pi_{\setminus S_{\text{mrk}}}(s, A, s')} r(\pi) \cdot \mathbb{P}[\pi]$.

The main idea of this construction is to lump together states that are all entered at the same time point. For example, in a sequence of probabilistic states followed by a Markovian state all of the states of the sequence will be entered at the same time due to the fact that probabilistic states are left instantaneously upon entry. The construction is similar in spirit to the construction of a value-preserving CTMDP for the *long-run average reward* problem from [5]. It differs, however, in the treatment of Markovian and probabilistic states due to the fact that timing information effects collected rewards and thus has to be preserved.

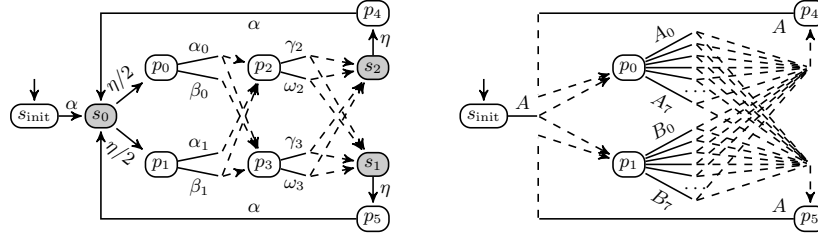


Fig. 2: An example of $\mathcal{C}(\mathcal{M})$ (on the right) for an MRA \mathcal{M} (on the left). We omitted the probabilities of the probabilistic transitions of \mathcal{M} . Here $A_0 = [p_0 \rightarrow \alpha_0, p_2 \rightarrow \gamma_2, p_3 \rightarrow \gamma_3]$ and other actions of $\mathcal{C}(\mathcal{M})$ are constructed analogously. If all the probabilistic distributions are uniform, then $R_{\mathcal{C}}(p_0, A_0, p_4) = \eta \cdot [0.5 \cdot 0.5 \cdot 1 + 0.5 \cdot 0.5 \cdot 1] = 0.5 \cdot \eta$.

An example of this transformation is depicted in Fig. 2. One can easily see that even in small examples the amount of transitions of $\mathcal{C}(\mathcal{M})$ can grow extremely fast. Let $s \in S_{\text{mrk}}$ and $|PS_{\setminus S_{\text{mrk}}}(s)| = n$. If every probabilistic state of \mathcal{M} has two enabled actions, then the set of enabled actions of s in $\mathcal{C}(\mathcal{M})$ is 2^n . This growth is therefore exponential in the worst case.

Theorem 1. *For any uniform normalised MRA \mathcal{M} we have $\text{dR}_{\mathcal{M},\beta}^{\text{opt}} = \text{dR}_{\mathcal{C}(\mathcal{M}),\beta}^{\text{opt}}^2$, and there is an optimal scheduler for \mathcal{M} that is stationary.*

4 Bellman Equation

In this section, we introduce the Bellman equation for the discounted reward problem on MRA.

First of all, due to the results obtained in Sect. 3.2, one could obtain the Bellman equation for an MRA by constructing the value-preserving CTMDP and using the Bellman equation for this CTMDP³ [18]. However, since the construction of $\mathcal{C}(\mathcal{M})$ is exponential in the size of \mathcal{M} , using the thus obtained Bellman equation for quantifying dR^{opt} is not efficient for general MRA.

First we informally discuss the reasons why using this method would be inefficient for \mathcal{M} . First of all, in order to use this approach one would need to construct $\mathcal{C}(\mathcal{M})$, which may require exponentially many computations. Additionally, having constructed $\mathcal{C}(\mathcal{M})$, the solution of its Bellman equation requires computing an extremum of an operator $F_{\mathcal{C}(\mathcal{M})}$ over all enabled actions: $V^* := \text{opt}_{A \in \text{Act}_{\mathcal{C}}} F_{\mathcal{C}(\mathcal{M})}(A)$. The definition of $F_{\mathcal{C}(\mathcal{M})}$ is irrelevant for the current discussion. Since the number of enabled actions in $\mathcal{C}(\mathcal{M})$ is in the worst case exponential in the size of \mathcal{M} , this operation is essentially a brute-force check over exponentially many options. However, we can show that this optimisation problem on $\mathcal{C}(\mathcal{M})$ when mirrored back to \mathcal{M} itself reduces to the computation of the *expected total reward* $\text{tR}_{\mathcal{D}(\mathcal{M})}^{\text{opt}}$

² Here $\text{dR}_{\mathcal{C},\beta}^{\text{opt}}$ denotes discounted reward on a CTMDP \mathcal{C} [18].

³ For details we refer to Appendix A.4.

on a discrete-time Markov decision process $\mathcal{D}(\mathcal{M})$, whose size is linear in the size of \mathcal{M} . Computing $\text{tR}_{\mathcal{D}(\mathcal{M})}^{\text{opt}}$ is a well-studied problem on MDPs that admits an efficient solution via dynamic programming. Thus instead of naïvely brute-forcing $\sup_{A \in \text{Act}_{\mathcal{C}}} F_{\mathcal{C}(\mathcal{M})}(A)$, the value V^* can be efficiently computed by well-known dynamic programming techniques for $\text{tR}_{\mathcal{D}(\mathcal{M})}^{\text{opt}}$ [3,18]. To formalise this result we need to introduce the expected total reward tR^{opt} on MDPs, as defined in [18].

Expected total reward. Let \mathcal{D} be a (discrete-time) MDP and X_i^s, Y_i^s be random variables denoting the state occupied by \mathcal{D} and the action chosen at step i starting from state s . Then the value

$$\text{tR}_{\mathcal{D}, r_{\mathcal{D}}}^{\text{opt}}(s) := \underset{\sigma \in GM_{\mathcal{D}}}{\text{opt}} \mathbf{E}_{s, \sigma} \left[\lim_{N \rightarrow \infty} \sum_{i=0}^{N-1} r_{\mathcal{D}}(X_i^s, Y_i^s) \right],$$

where $\text{opt} \in \{\text{sup}, \text{inf}\}$, denotes the *optimal expected total reward* on \mathcal{D} with reward structure $r_{\mathcal{D}}$, starting from state s .

Terminal MDP. We now construct the discrete MDP and the reward structure on it that enables us to substitute the naïve brute-force approach with the efficient computation of the expected total reward.

Let \mathcal{M} be a uniform normalised MRA. Informally, we keep the structure of \mathcal{M} , but for each marked state $s \in S_{\text{mrk}}$, we introduce a copy state s_{cp} and redirect all the transitions leading to s to the new copy state s_{cp} . These copy states have only transitions with probability 1 to a new terminal state t . Formally, the *terminal MDP* of \mathcal{M} is $\mathcal{D}(\mathcal{M}) := (S_{\mathcal{D}}, s_{\text{init}}, \text{Act} \dot{\cup} \{\perp\}, \mathbb{P}_{\mathcal{D}})$, where $S_{cp} = \{s_{cp} \mid s \in S_{\text{mrk}}\}$, $S_{\mathcal{D}} = S \dot{\cup} S_{cp} \dot{\cup} \{t\}$ and

$$\mathbb{P}_{\mathcal{D}}[s, \alpha, s'] = \begin{cases} \mathbb{P}[s, \alpha, s'] & \text{for } s \in PS, s' \notin S_{cp}, \\ \mathbb{P}[s, p] & \text{for } s \in MS, s' = p_{cp} \in S_{cp}, \alpha = \perp, \\ 1 & \text{for } (s \in S_{cp} \text{ or } s = t), s' = t, \alpha = \perp, \\ 0 & \text{otherwise.} \end{cases}$$

Figure 3(b) depicts the terminal MDP for the MRA from Fig. 1.

We can now present an efficient characterisation of the discounted reward on MRA. Let $\mathcal{D}(\mathcal{M}) = (S_{\mathcal{D}}, s_{\text{init}}, \text{Act}_{\mathcal{D}}, \mathbb{P}_{\mathcal{D}})$ be the terminal MDP of \mathcal{M} and $h : S_{\text{mrk}} \rightarrow \mathbb{R}$. We define a reward structure $\text{rew}_{\mathcal{D}(\mathcal{M}), h}$ for $\mathcal{D}(\mathcal{M})$ as follows:

$$\text{rew}_{\mathcal{D}(\mathcal{M}), h}(s, \alpha) := \begin{cases} r(s, \alpha) & \text{for } s \in PS, \alpha \in \text{Act}_{\mathcal{D}}(s), \\ \frac{\rho(s)}{\beta + \eta} & \text{for } s \in MS, \alpha = \perp, \\ \frac{\eta}{\beta + \eta} h(s) & \text{for } s \in S_{cp}, \alpha = \perp, \\ 0 & \text{for } s = t, \alpha = \perp. \end{cases}$$

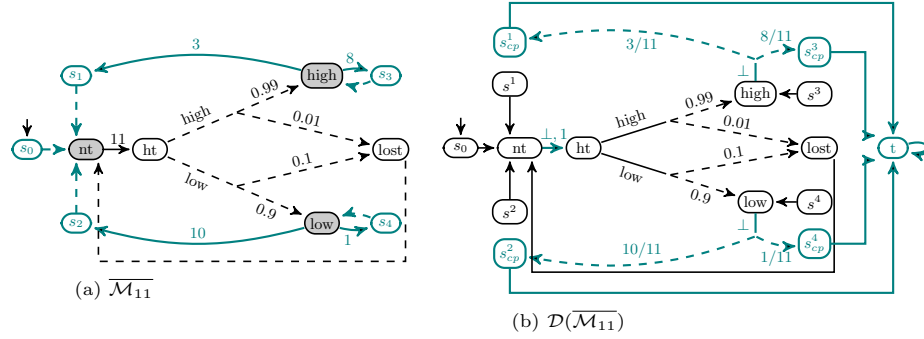


Fig. 3: Figure (a) depicts the uniform normalised MRA from Fig. 1 and (b) shows the corresponding terminal MDP. All the added/updated transitions and states are highlighted in green color. The figure omits Dirac distributions and rewards.

Theorem 2. Let \mathcal{M} be a uniform normalised MRA with exit rate η and $\mathcal{D}(\mathcal{M})$ the corresponding terminal MDP. Then the vector $\mathbf{dR}_{\mathcal{M},\beta}^{\text{opt}} := (\text{dR}_{\mathcal{M},\beta}^{\text{opt}}(s)), \forall s \in S_{\text{mrk}}$, is the unique solution to the Bellman equation:

$$\forall s \in S_{\text{mrk}} : \quad v(s) = \text{tR}_{\mathcal{D}(\mathcal{M}), \text{rew}_{\mathcal{D}(\mathcal{M}),v}}^{\text{opt}}(s). \quad (1)$$

The reason for this characterisation being efficient is the right-hand side of Equation (1). Quantification of the expected total reward on MDPs is a well-established problem that admits such algorithms as policy-iteration and linear programming [18]. Moreover, for a subclass of models it can be solved in time linear in the size of the MRA. Those are models that have no cycles consisting of only probabilistic states. Such cycles (even non-Zeno ones) almost never happen in real-world applications and are usually considered a modelling mistake. In fact, we are not aware of any practical example where that case occurs.

5 Numerical Solution

Having a Bellman equation at hand, one can normally derive three types of algorithms: based on value-iteration, policy-iteration, and linear programming. Due to scalability issues of the latter we do not consider it and present here only value- and policy-iteration algorithms.

In the following let $sp(v) := |\max_{s \in S_{\text{mrk}}} v(s) - \min_{s \in S_{\text{mrk}}} v(s)|$. Additionally, we denote with $\text{ExpectedTotalReward}(\mathcal{D}, \text{rew}, \sigma)$ the function that computes the expected total reward on an MDP \mathcal{D} for reward structure rew . The last parameter σ can be either a stationary deterministic scheduler or one of $\{\text{sup}, \text{inf}\}$. In the latter case, the optimal expected total reward is computed and in the former the expected total reward for scheduler σ .

Algorithms 1 and 2 are value- and modified policy-iteration algorithms that compute the value $\text{dR}_{\mathcal{M},\beta}^{\text{opt}}$ for an arbitrary MRA \mathcal{M} and discounting rate $\beta > 0$.

Algorithm 1: ValueIteration

input : MRA $\mathcal{M} = (S, s_{\text{init}}, Act, \hookrightarrow, \rightsquigarrow, r, \rho)$, $\text{opt} \in \{\text{sup}, \text{inf}\}$, $\beta > 0$,
approximation error $\varepsilon > 0$
output : v such that $\|v - \text{dR}_{\mathcal{M}, \beta}^{\text{opt}}\| < \varepsilon$, and the ε -optimal scheduler σ

- 1 $\overline{\mathcal{M}}_\eta := \text{normalise}(\text{uniformise}(\mathcal{M}, \text{rate } \eta := \max_{s \in S} E(s)))$;
- 2 $\mathcal{D}(\overline{\mathcal{M}}_\eta) := \text{terminal MDP for } \overline{\mathcal{M}}_\eta$;
- 3 $v_0 := \overline{0}$; /* vector of zeros */
- 4 **for** ($n := 0$; $sp(v_{n+1} - v_n) < \frac{\varepsilon \cdot \beta}{\eta}$; $n++$) **do**
- 5 $(v_{n+1}, \sigma) := \text{ExpectedTotalReward}(\mathcal{D}(\overline{\mathcal{M}}_\eta), \text{rew}_{\mathcal{D}(\overline{\mathcal{M}}_\eta), v_n}, \text{opt})$;
- 6 **return** $v_{n+1}(s_{\text{init}}), \sigma$;

The standard policy-iteration algorithm in which the policy evaluation step is performed exactly is also possible in the setting of MRA. However, this requires the exact solution of a linear equation system with one variable per state of \mathcal{M} . Since this is an expensive operation for hundreds of thousands of states, we choose the modified policy-iteration instead. The latter bypasses this issue by performing the policy evaluation step numerically. The algorithm depends on a sequence of natural numbers called *order sequence* $(m_n)_{n \in \mathbb{N}_{\geq 0}}$; it converges however for an arbitrary sequence.

Theorem 3. *Algorithms 1 and 2 are sound and complete.*

Algorithms 1 and 2 are essentially the respective algorithms on CTMDPs [18], in which the extremum value over enabled actions is searched through the solution of the expected total reward problem $\text{tR}_{\mathcal{D}(\mathcal{M})}^{\text{opt}}$ on the terminal MDP $\mathcal{D}(\mathcal{M})$. Therefore in both algorithms the complexity of an iteration equals the complexity of computing the value $\text{tR}_{\mathcal{D}(\mathcal{M})}^{\text{opt}}$ (which is polynomial), and the convergence rate is the same as the convergence rate of the corresponding CTMDP algorithms.

Computation of the expected total reward. Notice that the presented algorithms are guaranteed to converge whenever the expected total reward of the terminal MDP is computed precisely. Exact quantification of this value can be achieved with policy-iteration or linear programming algorithms [18]. Moreover, for models that have no cycles and consist of only probabilistic states, the expected total reward can be solved efficiently in time $O(|\rightsquigarrow| + |\hookrightarrow|)$.

6 Experiments

Here we present the empirical evaluation of the discussed algorithms. Both algorithms were implemented as part of the IMCA/MAMA toolset [9]. All experiments were run on a single core of Intel Core i7-4790 with 8 GB of RAM.

Algorithm 2: ModifiedPolicyIteration

input : MRA $\mathcal{M} = (S, s_{\text{init}}, Act, \hookrightarrow, \rightsquigarrow, r, \rho), \text{opt} \in \{\text{sup}, \text{inf}\}, \beta > 0, \varepsilon > 0$,
order sequence $(m_n)_{n \in \mathbb{N}_{\geq 0}}$

output : v such that $\|v - \text{dR}_{\mathcal{M}, \beta}^{\text{opt}}\| < \varepsilon$, and the ε -optimal scheduler σ

- 1 $\overline{\mathcal{M}}_\eta := \text{normalise}(\text{uniformise}(\mathcal{M}, \text{rate } \eta \leftarrow \max_{s \in S} E(s)))$;
- 2 $\mathcal{D}(\overline{\mathcal{M}}_\eta) := \text{terminal MDP for } \overline{\mathcal{M}}_\eta$;
- 3 $v_0 := \vec{0}$; /* vector of zeros */
- 4 **stop** := *false*; $n := 0$;
- 5 **while** ($\neg \text{stop}$) **do**
- 6 /* Policy improvement */
 $(u_n^0, \sigma_{n+1}) := \text{ExpectedTotalReward}(\mathcal{D}(\overline{\mathcal{M}}_\eta), \text{rew}_{\mathcal{D}(\overline{\mathcal{M}}_\eta), v_n}, \text{opt})$;
- 7 /* Partial policy evaluation */
- 8 **if** $sp(u_n^0 - v_n) < \frac{\varepsilon \cdot \beta}{\eta}$ **then**
- 9 **stop** := *true*; **break**;
- 10 **for** ($k := 0$; $k < m_n$; $k++$) **do**
- 11 $u_n^{k+1} := \text{ExpectedTotalReward}(\mathcal{D}(\overline{\mathcal{M}}_\eta), \text{rew}_{\mathcal{D}(\overline{\mathcal{M}}_\eta), v_n}, \sigma_{n+1})$;
- 12 $v_{n+1} := u_n^{m_n}$;
- 13 $n := n + 1$
- 14 **return** $u_n^0(s_{\text{init}}, \sigma_{n+1})$;

Table 1: Parameters of some of the benchmarks.

	$ S $	$ PS $	$ MS $	$ \hookrightarrow $	$ \rightsquigarrow $	$\mathbb{A}^{\mathcal{M}}$	$E(\mathcal{M})$
FTWC-resp-50-40	92,819	20,806	72,013	72,007	305,613	5	6.35
PS-256-3-4	131,529	87,605	43,924	189,129	72,965	3	14
QS-256-256	465,177	398,096	67,081	530,966	200,208	2	26

Benchmarks. We have evaluated our approach on a collection of published benchmark models: the *Polling System* [9,20], *Queuing System* [12], and the *Fault Tolerant Workstation Cluster* [15]. Discounting for the selected benchmarks naturally models the decrease of the value of costs over time. In order to address a case study with a specific set of parameters we use the same notation as in [5]. We used the tool SCOOP [19] to generate those models and for this reason the degree of variation of some parameters is restricted by its runtime/space requirements.

Table 1 shows the parameters of some of the used models. We use the symbols $\mathbb{A}^{\mathcal{M}}$ to denote the maximal number of enabled actions in probabilistic states of \mathcal{M} , and $E(\mathcal{M})$ shows the maximal exit rate of Markovian states of \mathcal{M} .

Empirical Evaluation. The space complexity of both algorithms is polynomial. Therefore, we have evaluated the effect of varying model size, precision, and the discounting rate on their runtime only. In plots, whenever the experiment covers several benchmarks, we use the symbol “X” to denote respective part of

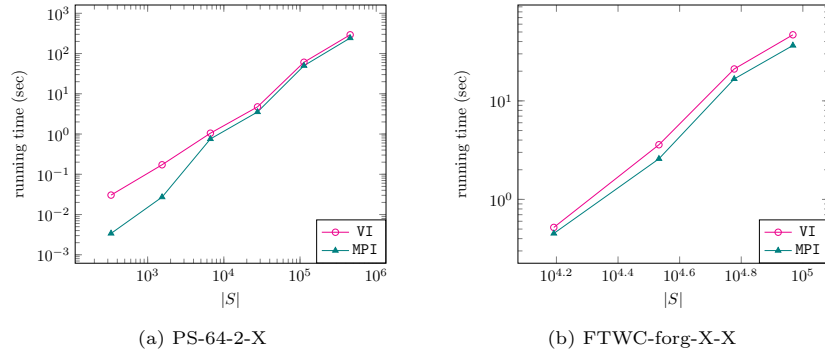


Fig. 4: Runtime complexity of VI and MPI w.r.t. the increase of the model size in log-log scale. For these experiments the discount rate β was set to 0.05 and $\varepsilon = 10^{-8}$.

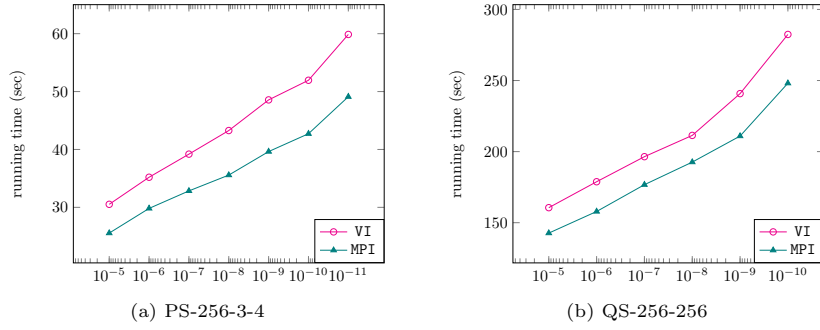


Fig. 5: Observed dependency of VI and MPI on the precision parameter ε in reversed logarithmic x -axis. In these experiments $\beta = 0.1$.

the model name, e. g. PS-2-X. In this section we will refer to the value-iteration algorithm 1 as VI and to the modified policy-iteration algorithm 2 as MPI.

As we have mentioned in Sect. 5, the modified policy-iteration algorithm depends on a parameter called order sequence. The optimal choice of the order sequence is an open question [18]. In this section, we present the best results we could achieve with different order sequences. Let us notice that if every element of the order sequence is 0 then MPI is the same as VI. When the values grow infinitely large, the algorithm turns into standard policy-iteration. Almost always in our experiments we could find an order sequence that led to lower running times than those of VI. Moreover, relatively small order sequences achieved the best value, e. g. $m_n = 100, \forall n > 0$ was sufficient for most of the models. Selecting a sequence with larger values, however, quite often led to running times significantly worse than those of VI.

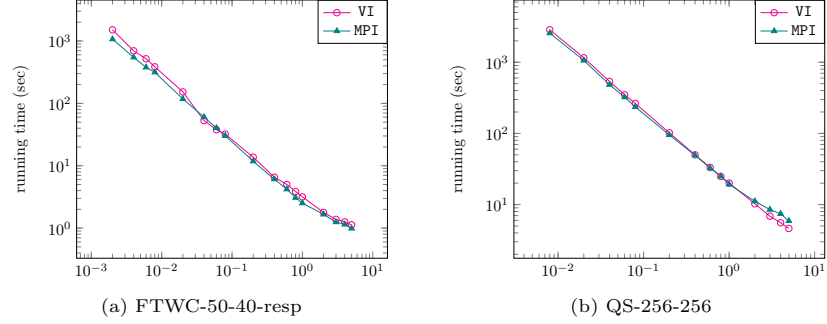


Fig. 6: Observed dependency on discounting rate β in log-log scale. Here $\varepsilon = 10^{-8}$.

Model size. Figure 4 shows the dependency of the running time of both algorithms on the size of the state space. Both algorithms exhibit polynomial dependency (linear in log-log scale) on the depicted size range. This agrees with theoretical expectations since the computation of expected total rewards is polynomial in the size of the state space and the convergence rate does not depend on the model size.

Precision. Figure 5 shows the dependency of the computation time on the precision parameter ε . The theoretical convergence rate of both algorithms resembles that of respective algorithms on CTMDPs. The expected complexity of VI is logarithmic in ε , which is supported by the observed results. Regarding MPI, we observed that the function of the running time repeats that of VI, possibly due to the relatively small values used for the order sequence.

Discounting rate. Figure 6 depicts the dependency of the running time of the algorithms on the discounting rate β . The observed dependency of VI follows the theoretical bound of $O(\frac{1}{1-\beta})$. Similarly to the previous case, the function of the running time of policy-iteration repeats that of value iteration.

7 Conclusion

While discounting is a standard concept on Markov chains and Markov decision processes, this is the first paper to consider discounting for the more general model of Markov reward automata. We have discussed that computing discounted rewards on MRA can be reduced to the same task on a possibly exponentially larger CTMDP. Constructing and optimizing over this large CTMDP can be avoided by recognising the essential computation as determining the expected total reward in a linear-sized discrete-time MDP. This in turn is a well-understood problem enabling an efficient solution. Experiments clearly demonstrate the efficiency of our approach, being able to handle MRAs with hundred thousands of states.

References

1. de Alfaro, L., Faella, M., Henzinger, T.A., Majumdar, R., Stoelinga, M.: Model checking discounted temporal properties. *Theor. Comp. Sci.* 345(1), 139–170 (2005)
2. de Alfaro, L., Henzinger, T.A., Majumdar, R.: Discounting the future in systems theory. In: *ICALP 2003*. LNCS, vol. 2719, pp. 1022–1037. Springer (2003)
3. Bertsekas, D.P.: *Dynamic Programming and Optimal Control*. Athena Scientific, 2nd edn. (2000)
4. Boudali, H., Crouzen, P., Stoelinga, M.: A rigorous, compositional, and extensible framework for dynamic fault tree analysis. *IEEE Trans. Dependable Sec. Comput.* 7(2), 128–143 (2010)
5. Butkova, Y., Wimmer, R., Hermanns, H.: Long-run rewards for Markov automata. In: *TACAS 2017, Part II*. LNCS, vol. 10206, pp. 188–203. Springer (2017)
6. Dehnert, C., Junges, S., Katoen, J., Volk, M.: A Storm is coming: A modern probabilistic model checker. In: *CAV 2017, Part II*. LNCS, vol. 10427, pp. 592–600. Springer (2017)
7. Eisentraut, C., Hermanns, H., Katoen, J., Zhang, L.: A semantics for every GSPN. In: *Petri Nets 2013*. LNCS, vol. 7927, pp. 90–109. Springer (2013)
8. Eisentraut, C., Hermanns, H., Zhang, L.: On probabilistic automata in continuous time. In: *LICS 2010*. pp. 342–351. IEEE CS (2010)
9. Guck, D., Hatefi, H., Hermanns, H., Katoen, J., Timmer, M.: Modelling, reduction and analysis of Markov automata. In: *QEST 2013*. LNCS, vol. 8054, pp. 55–71. Springer (2013)
10. Guck, D., Hatefi, H., Hermanns, H., Katoen, J., Timmer, M.: Analysis of timed and long-run objectives for Markov automata. *Log. Meth. Comput. Sci.* 10(3) (2014)
11. Guck, D., Timmer, M., Hatefi, H., Ruijters, E., Stoelinga, M.: Modelling and analysis of Markov reward automata. In: *ATVA 2014*. LNCS, vol. 8837, pp. 168–184. Springer (2014)
12. Hatefi, H., Hermanns, H.: Model checking algorithms for Markov automata. *Electronic Communication of the EASST* 53 (2012)
13. Hatefi, H., Wimmer, R., Braitling, B., Fioriti, L.M.F., Becker, B., Hermanns, H.: Cost vs. time in stochastic games and markov automata. *Formal Aspects of Computing* 29(4), 629–649 (2017)
14. Hatefi Ardakani, H.: *Finite Horizon Analysis of Markov Automata*. Ph.D. thesis, Universität des Saarlandes, Saarbrücken, Germany (2017)
15. Haverkort, B.R., Hermanns, H., Katoen, J.: On the use of model checking techniques for dependability evaluation. In: *SRDS 2000*. pp. 228–237. IEEE CS (2000)
16. Jansen, D.N.: More or less true DCTL for continuous-time MDPs. In: *FORMATS 2013*. LNCS, vol. 8053, pp. 137–151. Springer (2013)
17. Jensen, A.: Markoff chains as an aid in the study of Markoff processes. *Scandinavian Actuarial Journal* 1953, 87–91 (1953)
18. Puterman, M.L.: *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley & Sons, Inc., New York, NY, USA, 1st edn. (1994)
19. Timmer, M.: SCOOP: A tool for symbolic optimisations of probabilistic processes. In: *QEST 2011*. pp. 149–150. IEEE CS (2011)
20. Timmer, M., van de Pol, J., Stoelinga, M.: Confluence reduction for Markov automata. In: *FORMATS 2013*. LNCS, vol. 8053, pp. 243–257. Springer (2013)

A Appendix

A.1 Proof of Lemma 1

Let $\mathcal{M} = (S, s_{\text{init}}, Act, \hookrightarrow, \rightsquigarrow, r, \rho)$ be a uniform normalised MRA. In the following we will use the following abbreviations:

- $1_S(s) := 1$ if $s \in S$ and 0 otherwise.
- For a finite path $\pi = s_0 \xrightarrow{\alpha_0, t_0} s_1 \xrightarrow{\alpha_1, t_1} \dots s_n$, $t[\pi, i] := t_i$ (π is omitted when it is clear from the context) and $T(\pi) := \sum_{i=0}^{|\pi|-1} t[i]$.

Lemma 1. *The value $\text{dR}_{\mathcal{M}, \beta}^{\text{opt}}$ exists.*

Proof. We first show that for every scheduler $\sigma \in GM$ the value $\text{dR}_{\mathcal{M}, \beta}^{\sigma}$ exists. By definition:

$$\begin{aligned} \text{dR}_{\mathcal{M}, \beta}^{\sigma} &= \lim_{N \rightarrow \infty} \mathbf{E}_{\sigma}[\text{rew}_{\mathcal{M}, \beta}^N(\pi)] \\ &= \lim_{N \rightarrow \infty} \mathbf{E}_{\sigma} \left[\sum_{k=0}^{N-1} [e^{-\beta \cdot \tau[k]} \cdot r(\pi[k], \alpha[k]) + \int_{\tau[k]}^{\tau[k]+t[k]} e^{-\beta \cdot t} \cdot \rho(\pi[k]) dt] \right]. \end{aligned}$$

Let $\pi' \in \text{Paths}_{\mathcal{M}}$, such that $\pi' \downarrow \in PS_{\mathcal{M}}$. We denote with $\mathbf{E}_{\sigma, MS_{\mathcal{M}}}[\text{rew}_{\mathcal{M}, \beta} | \pi']$ the expected discounted reward collected starting from state $\pi' \downarrow$ until encountering a Markovian state s and following the scheduler σ that assumes that path π' has been taken until reaching s . Then

$$\begin{aligned} \mathbf{E}_{\sigma, MS_{\mathcal{M}}}[\text{rew}_{\mathcal{M}, \beta} | \pi'] &= \mathbf{E}_{\sigma, MS_{\mathcal{M}}} \left[\sum_{k=0}^{\infty} e^{-\beta \cdot \tau[k]} \cdot r(\pi[k], \alpha[k]) \mid \pi' \right] \\ &= \mathbf{E}_{\sigma, MS_{\mathcal{M}}} \left[\sum_{k=0}^{\infty} e^{-\beta \cdot \tau[\pi' \downarrow]} \cdot r(\pi[k], \alpha[k]) \mid \pi' \right] \\ &= e^{-\beta \cdot \tau[\pi' \downarrow]} \cdot \mathbf{E}_{\sigma, MS_{\mathcal{M}}} \left[\sum_{k=0}^{\infty} r(\pi[k], \alpha[k]) \mid \pi' \right]. \end{aligned}$$

The value $\mathbf{E}_{\sigma, MS_{\mathcal{M}}} \left[\sum_{k=0}^{\infty} r(\pi[k], \alpha[k]) \mid \pi' \right]$ is the total expected reward gathered in the MDP \mathcal{D}_s formed by probabilistic states reachable from s via transition relation \hookrightarrow . Where the reward value of states of this MDP $r_{\mathcal{D}}(s, \alpha) := r(s, \alpha)$.

Since \mathcal{M} is non-Zeno, all the probabilistic states are transient. This is a sufficient condition [18] for the value $\mathbf{E}_{\sigma, MS_{\mathcal{M}}} \left[\sum_{k=0}^{\infty} r(\pi[k], \alpha[k]) \mid \pi' \right]$ to exist and be finite. Therefore there exists a value M , such that

$$\mathbf{E}_{\sigma, MS_{\mathcal{M}}}[\text{rew}_{\mathcal{M}, \beta} | \pi'] \leq e^{-\beta \cdot \tau[\pi' \downarrow]} \cdot M.$$

Therefore

$$\begin{aligned}
dR_{\mathcal{M},\beta}^{\sigma} &= \lim_{N \rightarrow \infty} \mathbf{E}_{\sigma} \left[\sum_{k=0}^{N-1} \left[e^{-\beta \cdot \tau[k]} \cdot \mathbf{r}(\pi[k], \alpha[k]) + \int_{\tau[k]}^{\tau[k]+t[k]} e^{-\beta \cdot t} \cdot \rho(\pi[k]) \, dt \right] \right] \\
&= \lim_{N \rightarrow \infty} \sum_{k=0}^{N-1} \left(\sum_{s \in MS_{\mathcal{M}}} \mathbf{E}_{s,\sigma} \left[\int_{\tau[k]}^{\tau[k]+t[k]} e^{-\beta \cdot t} \cdot \rho(s) \, dt \right] \right. \\
&\quad \left. + \sum_{s \in PS_{\mathcal{M}}} \mathbf{E}_{s,\sigma} \left[e^{-\beta \cdot \tau[k]} \cdot \mathbf{r}(s, \alpha[k]) \right] \right) \\
&\leq \lim_{N \rightarrow \infty} \sum_{k=0}^{N-1} \left(\sum_{s \in MS_{\mathcal{M}}} \mathbf{E}_{s,\sigma} \left[\int_{\tau[k]}^{\tau[k]+t[k]} e^{-\beta \cdot t} \cdot \rho(s) \, dt \right] \right. \\
&\quad \left. + \sum_{s \in S_{\text{mrk}}} \mathbf{E}_{s,\sigma} \left[\mathbf{E}_{\sigma, MS_{\mathcal{M}}}[\text{rew}_{\mathcal{M},\beta} | \pi'] \right] \right) \\
&\leq \lim_{N \rightarrow \infty} \sum_{k=0}^{N-1} \left(\sum_{s \in MS_{\mathcal{M}}} \mathbf{E}_{s,\sigma} \left[\int_{\tau[k]}^{\tau[k]+t[k]} e^{-\beta \cdot t} \cdot \rho(s) \, dt \right] \right. \\
&\quad \left. + \sum_{s \in S_{\text{mrk}}} \mathbf{E}_{s,\sigma} \left[e^{-\beta \cdot \tau[k]} \cdot M \right] \right) \\
&= \lim_{N \rightarrow \infty} \sum_{k=0}^{N-1} \sum_{s \in MS_{\mathcal{M}}} \mathbf{E}_{s,\sigma} \left[\int_{\tau[k]}^{\tau[k]+t[k]} e^{-\beta \cdot t} \cdot \rho(s) \, dt + e^{-\beta \cdot \tau[k]} \cdot M \right] \\
&= \lim_{N \rightarrow \infty} \mathbf{E}_{\sigma} \left[\sum_{k=0}^{N-1} 1_{MS_{\mathcal{M}}}(\pi[k]) \left(\int_{\tau[k]}^{\tau[k]+t[k]} e^{-\beta \cdot t} \cdot \rho(\pi[k]) \, dt + e^{-\beta \cdot \tau[k]} \cdot M \right) \right] \\
&= \lim_{N \rightarrow \infty} \mathbf{E}_{\sigma} \left[\sum_{k=0}^{N-1} 1_{MS_{\mathcal{M}}}(\pi[k]) \cdot e^{-\beta \cdot \tau[k]} \cdot \left(\int_0^{t[k]} e^{-\beta \cdot t} \cdot \rho(\pi[k]) \, dt + M \right) \right] \\
&= \lim_{N \rightarrow \infty} \mathbf{E}_{\sigma} \left[\sum_{k=0}^{N-1} 1_{MS_{\mathcal{M}}}(\pi[k]) \cdot e^{-\beta \cdot \tau[k]} \cdot \left(\frac{\rho(\pi[k])}{\beta + E(\pi[k])} + M \right) \right] \\
&\leq \lim_{N \rightarrow \infty} \mathbf{E}_{\sigma} \left[\sum_{k=0}^{N-1} 1_{MS_{\mathcal{M}}}(\pi[k]) \cdot e^{-\beta \cdot \tau[k]} \cdot C \right] \\
&\leq \lim_{N \rightarrow \infty} C \cdot \mathbf{E}_{\sigma} \left[\int_0^{\infty} e^{-\beta \cdot t} \, dt \right] \\
&\leq \frac{C}{\beta}.
\end{aligned}$$

Therefore

$$\text{dR}_{\mathcal{M},\beta}^{\text{opt}} = \sup_{\sigma \in GM} \{\text{dR}_{\mathcal{M},\beta}^{\sigma}\} \leq \frac{C}{\beta}$$

and the value $\text{dR}_{\mathcal{M},\beta}^{\text{opt}}$ exists and is finite. \square

A.2 Proof of Theorem 1

Let $\mathcal{M} = (S, s_{\text{init}}, \text{Act}, \hookrightarrow, \rightsquigarrow, r, \rho)$ be a uniform normalised MRA and $\mathcal{C}(\mathcal{M}) = (S_{\text{mrk}}, s_{\text{init}}, \text{Act}_{\mathcal{C}}, R_{\mathcal{C}})$ the CTMDP obtained as shown in Section 3.2. In the following we will denote paths of \mathcal{M} with symbol π and paths of \mathcal{C} with symbol φ . The following is the list of abbreviations that we will use for this proof:

- Let $\pi \in \text{Paths}_{\mathcal{M}}^* \cup \text{Paths}_{\mathcal{M}}$. Then $\pi_{\mathcal{C}}[i] = s$, s. t. $s \in S_{\text{mrk}}$ and s is the i -th state of π that belongs to S_{mrk} .
- $\text{Paths}_{\mathcal{M}}^{N(\mathcal{C})}(s)$ denotes the set of finite paths in \mathcal{M} that have exactly N positions i such that $\pi[i] \in S_{\text{mrk}}$.
- $\text{rew}_{\mathcal{M},\beta}^{T,N(\mathcal{C})}(\pi) := \text{rew}_{\mathcal{M},\beta}^{T,K}(\pi)$, where K is the length of the prefix of π that contains exactly N states from S_{mrk} .
- Let $\varphi \in \text{Paths}_{\mathcal{C}}^*$, then $\text{Paths}_{\mathcal{M}}(\varphi)$ denotes all finite paths π in \mathcal{M} , such that $\forall i : \pi_{\mathcal{C}}[i] = \varphi[i]$.
- Let π, π' be two finite paths, such that $\pi \downarrow = \pi'[0]$, then $\pi \cdot \pi'$ denotes the concatenation of the two paths.
- Let Π be a set of paths (finite or infinite), then $\text{ta}(\Pi) := \{\pi = s_0 \xrightarrow{\alpha_0} s_1 \xrightarrow{\alpha_1} \dots \mid s_0 \xrightarrow{\alpha_0, t_0} s_1 \xrightarrow{\alpha_1, t_1} \dots \in \Pi\}$.

We will now redefine the reward collected over a path to take into account a time shift. Let $T > 0$, then $\text{rew}_{\mathcal{M},\beta}^{T,0}(\pi) := 0$, and:

$$\text{rew}_{\mathcal{M},\beta}^{T,N}(\pi) := \sum_{k=0}^{N-1} \left[e^{-\beta \cdot (\tau[k]+T)} \cdot r_{\mathcal{M}}(\pi[k], \alpha[k]) + \int_{\tau[k]+T}^{\tau[k]+\tau[k]+T} e^{-\beta \cdot t} \cdot \rho_{\mathcal{M}}(\pi[k]) dt \right].$$

Obviously $\text{rew}_{\mathcal{M},\beta}^N(\pi) = \text{rew}_{\mathcal{M},\beta}^{0,N}(\pi)$. We additionally redefine the discounted reward to take into account a history $\pi \in \text{Paths}_{\mathcal{M}}^*$:

$$\text{dR}_{\mathcal{M},\beta}^{\sigma}(\pi, N) := \int_{\text{Paths}_{\mathcal{M}}} \text{rew}_{\mathcal{M}}^{T(\pi),N}(\pi') \cdot \Pr_{\sigma}[\text{d}\pi'],$$

Obviously $\text{dR}_{\mathcal{M},\beta}^{\sigma}(N) = \text{dR}_{\mathcal{M},\beta}^{\sigma}(\pi, N)$. And analogously we redefine respective values for CTMDPs.

We define $\text{dR}_{\mathcal{M}}^{\sigma}(\pi, N(\mathcal{C})) := \int_{\text{Paths}_{\mathcal{M}}} \text{rew}_{\mathcal{M}}^{T(\pi),N}(\pi') \cdot \Pr_{\sigma}[\text{d}\pi']$. Let σ be a scheduler in \mathcal{M} . We define scheduler δ in \mathcal{C} as follows:

$$\delta(\varphi, A) := \int_{\text{Paths}_{\mathcal{M}}(\varphi)} \prod_{s \in PS \setminus S_{\text{mrk}}(\varphi \downarrow)} \sum_{\pi' \in \Pi \setminus S_{\text{mrk}}(\varphi \downarrow, s)} \sigma(\pi \cdot \pi', A(s)) \, d\pi.$$

Let $\pi, \pi' \in \text{Paths}_{\mathcal{M}}^*(\varphi)$, then $\sigma_{\pi}(\pi', \alpha) := \sigma(\pi \cdot \pi', \alpha)$. And analogously for δ_{φ} .

Lemma 3. Let $\varphi \in \text{Paths}_{\mathcal{C}}^*$ and $\pi \in \text{Paths}_{\mathcal{M}}^*(\varphi)$, then

$$\begin{aligned} \int_{\text{Paths}_{\mathcal{M}}(\varphi)} \Pr_{\sigma}[\text{d}\pi] \int_{\text{Paths}_{\mathcal{M}}^{1(\mathcal{C})}(\pi \downarrow)} \text{rew}_{\mathcal{M}}^{T(\pi), 1(\mathcal{C})}(\pi') \cdot \Pr_{\sigma_{\pi}}[\text{d}\pi'] = \\ \int_{\text{Paths}_{\mathcal{C}}^1(\varphi \downarrow)} \text{rew}_{\mathcal{C}}^{T(\varphi), 1}(\pi') \Pr_{\delta_{\varphi}}[\text{d}\pi']. \end{aligned}$$

Proof. Let $\pi \in \text{Paths}_{\mathcal{M}}^*(\varphi)$ and $Y_{\varphi, \pi} : \text{Act}_{\mathcal{C}}(\varphi \downarrow) \rightarrow [0, 1]$ – a random variable, such that $\Pr[Y_{\varphi, \pi} = A] = \prod_{s \in PS \setminus S_{\text{mrk}}(\varphi \downarrow)} \sum_{\pi' \in \Pi \setminus S_{\text{mrk}}(\varphi \downarrow, s)} \sigma(\pi \cdot \pi', A(s))$.

Then

$$\begin{aligned} \int_{\text{Paths}_{\mathcal{M}}^{1(\mathcal{C})}(\pi \downarrow)} \text{rew}_{\mathcal{M}}^{T(\pi), 1(\mathcal{C})}(\pi') \cdot \Pr_{\sigma_{\pi}}[\text{d}\pi'] = \\ \int_0^{\infty} dt \sum_{\pi' \in \text{ta}(\text{Paths}_{\mathcal{M}}^{1(\mathcal{C})}(\pi \downarrow))} \text{rew}_{\mathcal{M}}^{T(\pi), 1(\mathcal{C})}(\pi'(t)) \cdot \Pr_{\sigma_{\pi}}[\pi'], \end{aligned}$$

where $\pi'(t)$ is the timed path, obtained from the untimed path $\pi' \in \text{ta}(\text{Paths}_{\mathcal{M}}^{1(\mathcal{C})}(\pi \downarrow))$ by assigning sojourn time t to the only Markovian state of π' , and sojourn time 0 to all the rest states of the path. Consider the inner summation of the last expression:

$$\begin{aligned} \sum_{\pi' \in \text{ta}(\text{Paths}_{\mathcal{M}}^{1(\mathcal{C})}(\pi \downarrow))} \text{rew}_{\mathcal{M}}^{T(\pi), 1(\mathcal{C})}(\pi'(t)) \cdot \underbrace{\Pr_{\sigma_{\pi}}[\pi']}_{= \mathbb{P}[\pi'] \cdot \prod_{i=0}^{|\pi'|-1} \sigma_{\pi}(\pi'[i], \alpha[\pi', i])} &= \\ \sum_{\pi' \in \text{ta}(\text{Paths}_{\mathcal{M}}^{1(\mathcal{C})}(\pi \downarrow))} \text{rew}_{\mathcal{M}}^{T(\pi), 1(\mathcal{C})}(\pi'(t)) \cdot \mathbb{P}[\pi'] \cdot \underbrace{\sigma_{\pi}(\pi')}_{:= \prod_{i=0}^{|\pi'|-1} \sigma_{\pi}(\pi'[i], \alpha[\pi', i])} &= \\ \sum_{\pi' \in \text{ta}(\text{Paths}_{\mathcal{M}}^{1(\mathcal{C})}(\pi \downarrow))} \text{rew}_{\mathcal{M}}^{T(\pi), 1(\mathcal{C})}(\pi'(t)) \cdot \mathbb{P}[\pi'] \cdot \sum_{\substack{A \\ \forall i: \alpha[\pi', i] = A(\pi'[i])}} \mathbb{P}[Y_{\varphi, \pi} = A] &= \\ \sum_{\pi' \in \text{ta}(\text{Paths}_{\mathcal{M}}^{1(\mathcal{C})}(\pi \downarrow))} \text{rew}_{\mathcal{M}}^{T(\pi), 1(\mathcal{C})}(\pi'(t)) \cdot \mathbb{P}[\pi'] \cdot \sum_{A \in \text{Act}_{\mathcal{C}}(s)} 1_{\pi'}(A) \cdot \mathbb{P}[Y_{\varphi, \pi} = A], \end{aligned}$$

where $1_{\pi'}(A) = 1$ if $\forall i : \alpha[\pi', i] = A(\pi'[i])$ and 0 otherwise. The set of paths time abstract paths π' on which $1_{\pi'}(A) = 1$ we denote as $\text{ta}(\text{Paths}_{\mathcal{M}}^{1(\mathcal{C})}(\pi \downarrow, A))$. Then

$$\begin{aligned}
& \sum_{\pi' \in \text{ta}(\text{Paths}_{\mathcal{M}}^{1(\mathcal{C})}(\pi \downarrow))} \text{rew}_{\mathcal{M}}^{T(\pi), 1(\mathcal{C})}(\pi'(t)) \cdot \Pr_{\sigma_{\pi}}[\pi'] = \\
& \sum_{A \in \text{Act}_{\mathcal{C}}(s)} \mathbb{P}[Y_{\varphi, \pi} = A] \sum_{\pi' \in \text{ta}(\text{Paths}_{\mathcal{M}}^{1(\mathcal{C})}(\pi \downarrow))} \text{rew}_{\mathcal{M}}^{T(\pi), 1(\mathcal{C})}(\pi'(t)) \cdot \mathbb{P}[\pi'] \cdot 1_{\pi'}(A) = \\
& \sum_{A \in \text{Act}_{\mathcal{C}}(s)} \mathbb{P}[Y_{\varphi, \pi} = A] \underbrace{\sum_{\pi' \in \text{ta}(\text{Paths}_{\mathcal{M}}^{1(\mathcal{C})}(\pi \downarrow, A))} \text{rew}_{\mathcal{M}}^{T(\pi), 1(\mathcal{C})}(\pi'(t)) \cdot \mathbb{P}[\pi']}_{= \text{rew}_{\mathcal{C}}^{T(\pi), 1}(s, A, t)} = \\
& \sum_{A \in \text{Act}_{\mathcal{C}}(s)} \mathbb{P}[Y_{\varphi, \pi} = A] \cdot \text{rew}_{\mathcal{C}}^{T(\pi), 1}(s, A, t).
\end{aligned}$$

Therefore

$$\begin{aligned}
& \int_{\text{Paths}_{\mathcal{M}}(\varphi)} \Pr_{\sigma}[\text{d}\pi] \int_{\text{Paths}_{\mathcal{M}}^{1(\mathcal{C})}(\pi \downarrow)} \text{rew}_{\mathcal{M}}^{T(\pi), 1(\mathcal{C})}(\pi') \cdot \Pr_{\sigma_{\pi}}[\text{d}\pi'] = \\
& \int_{\text{Paths}_{\mathcal{M}}(\varphi)} \Pr_{\sigma}[\text{d}\pi] \int_0^{\infty} \text{d}t \sum_{A \in \text{Act}_{\mathcal{C}}(s)} \mathbb{P}[Y_{\varphi, \pi} = A] \cdot \text{rew}_{\mathcal{C}}^{T(\pi), 1}(s, A, t) = \\
& \int_0^{\infty} \text{d}t \sum_{A \in \text{Act}_{\mathcal{C}}(s)} \delta_{\varphi}(A) \cdot \text{rew}_{\mathcal{C}}^{T(\varphi), 1}(s, A, t) = \\
& \int_{\text{Paths}_{\mathcal{C}}^1(\varphi \downarrow)} \text{rew}_{\mathcal{C}}^{T(\varphi), 1}(s, A, t) \cdot \Pr_{\delta_{\varphi}}[\text{d}\varphi'].
\end{aligned}$$

□

Theorem 1. For any uniform normalised MRA \mathcal{M} we have $\text{dR}_{\mathcal{M}, \beta}^{\text{opt}} = \text{dR}_{\mathcal{C}(\mathcal{M}), \beta}^{\text{opt}}$ ⁴, and there is an optimal scheduler for \mathcal{M} that is stationary.

Proof. For convenience in the following we denote the states of the MRA with the symbol s , and states of the CTMDP with q . The initial states of both models coincide and we therefore denote it as s_{init} .

We now prove by induction over N , that $\forall \varphi \in \text{Paths}_{\mathcal{C}}^*$, $\forall N \in \mathbb{N}_{\geq 0}$:

$$\int_{\pi \in \text{Paths}_{\mathcal{M}}(\varphi)} \text{dR}_{\mathcal{M}}^{\varphi}(\pi, N(\mathcal{C})) \cdot \Pr[\text{d}\pi] = \text{dR}_{\mathcal{C}}^{\delta}(\varphi, N).$$

⁴ Here $\text{dR}_{\mathcal{C}, \beta}^{\text{opt}}$ denotes discounted reward on a CTMDP \mathcal{C} [18].

- *Induction Base:* Obvious, since $\forall \pi \in Paths_{\mathcal{M}}^* : dR_{\mathcal{M}}^{\sigma}(\pi, 0(\mathcal{C})) = dR_{\mathcal{M}}^{\sigma}(\pi, 0) = 0$ and $\forall \varphi : dR_{\mathcal{C}}^{\delta}(\varphi, 0) = 0$;
- *Induction Step:* Let $\varphi \in Paths_{\mathcal{C}}^*$ and $\pi \in Paths_{\mathcal{M}}^*(\varphi)$. Then

$$\begin{aligned}
& dR_{\mathcal{M}}^{\sigma}(\pi, N(\mathcal{C})) \\
&= \int_{Paths_{\mathcal{M}}^{1(\mathcal{C})}(\pi \downarrow)} \Pr_{\sigma_{\pi}}[d\pi'] \times \\
&\quad \int_{Paths_{\mathcal{M}(\pi' \downarrow)}} \Pr_{\sigma_{\pi \cdot \pi'}}[d\pi''] \underbrace{\text{rew}_{\mathcal{M}}^{T(\pi), N(\mathcal{C})}(\pi' \cdot \pi'')}_{= \text{rew}_{\mathcal{M}}^{T(\pi), 1(\mathcal{C})}(\pi') + \text{rew}_{\mathcal{M}}^{T(\pi) + T(\pi'), (N-1)(\mathcal{C})}(\pi'')} \\
&= \int_{Paths_{\mathcal{M}}^{1(\mathcal{C})}(\pi \downarrow)} \text{rew}_{\mathcal{M}}^{T(\pi), 1(\mathcal{C})}(\pi') \cdot \Pr_{\sigma_{\pi}}[d\pi'] + \\
&\quad \int_{Paths_{\mathcal{M}}^{1(\mathcal{C})}(\pi \downarrow)} \Pr_{\sigma_{\pi}}[d\pi'] \\
&\quad \int_{Paths_{\mathcal{M}(\pi' \downarrow)}} \text{rew}_{\mathcal{M}}^{T(\pi) + T(\pi'), (N-1)(\mathcal{C})}(\pi'') \cdot \Pr_{\sigma_{\pi \cdot \pi'}}[d\pi''] \\
&= \int_{Paths_{\mathcal{M}}^{1(\mathcal{C})}(\pi \downarrow)} \text{rew}_{\mathcal{M}}^{T(\pi), 1(\mathcal{C})}(\pi') \Pr_{\sigma_{\pi}}[d\pi'] + \\
&\quad \int_{Paths_{\mathcal{M}}^{1(\mathcal{C})}(\pi \downarrow)} dR_{\mathcal{M}}^{\sigma}(\pi \cdot \pi', (N-1)(\mathcal{C})) \cdot \Pr_{\sigma_{\pi}}[d\pi'].
\end{aligned}$$

Therefore

$$\begin{aligned}
& \int_{Paths_{\mathcal{M}}(\varphi)} dR_{\mathcal{M}}^{\sigma}(\pi, N(\mathcal{C})) \cdot \Pr_{\sigma}[d\pi] \\
&= \int_{Paths_{\mathcal{M}}(\varphi)} \int_{Paths_{\mathcal{M}}^{1(\mathcal{C})}(\pi \downarrow)} \text{rew}_{\mathcal{M}}^{T(\pi), 1(\mathcal{C})}(\pi') \cdot \Pr_{\sigma_{\pi}}[d\pi'] \cdot \Pr_{\sigma}[d\pi] + \\
&\quad \int_{Paths_{\mathcal{M}}(\varphi)} \int_{Paths_{\mathcal{M}}^{1(\mathcal{C})}(\pi \downarrow)} dR_{\mathcal{M}}^{\sigma}(\pi \cdot \pi', (N-1)(\mathcal{C})) \cdot \Pr_{\sigma_{\pi}}[d\pi'] \cdot \Pr_{\sigma}[d\pi].
\end{aligned}$$

Let $q \in S_{\text{mrk}}, A_q \in \text{Act}_{\mathcal{C}}(q), t > 0$ and $\varphi \xrightarrow{A_q, t} q := \varphi \cdot (\pi \downarrow \xrightarrow{A_q, t} q)$. Then

$$\begin{aligned}
& \int_{\text{Paths}_{\mathcal{M}}(\varphi)} dR_{\mathcal{M}}^{\sigma}(\pi, N(\mathcal{C})) \cdot \Pr_{\sigma_{\pi}}[d\pi] \\
&= \int_{\text{Paths}_{\mathcal{M}}(\varphi)} \int_{\text{Paths}_{\mathcal{M}}^{1(\mathcal{C})}(\pi \downarrow)} \text{rew}_{\mathcal{M}}^{T(\pi), 1(\mathcal{C})}(\pi') \cdot \Pr_{\sigma_{\pi}}[d\pi'] \cdot \Pr_{\sigma}[d\pi] + \\
& \quad \sum_{q \in S_{\text{mrk}}} \sum_{A_q \in \text{Act}_{\mathcal{C}}(q)} \int_0^{\infty} dt \int_{\text{Paths}_{\mathcal{M}}(\varphi \xrightarrow{A_q, t} q)} dR_{\mathcal{M}}^{\sigma}(\pi, (N-1)(\mathcal{C})) \cdot \Pr_{\sigma}[d\pi] \\
&\stackrel{IH}{=} \int_{\text{Paths}_{\mathcal{M}}(\varphi)} \int_{\text{Paths}_{\mathcal{M}}^{1(\mathcal{C})}(\varphi \downarrow)} \text{rew}_{\mathcal{M}}^{T(\pi), 1(\mathcal{C})}(\pi') \cdot \Pr_{\sigma_{\pi}}[d\pi'] \cdot \Pr_{\sigma}[d\pi] + \\
& \quad \sum_{q \in S_{\text{mrk}}} \sum_{A_q \in \text{Act}_{\mathcal{C}}(q)} \int_0^{\infty} dR_{\mathcal{C}}^{\delta}(\varphi \xrightarrow{A_q, t} q, N-1) \cdot dt \\
&\stackrel{\text{Lemma 3}}{=} \int_{\text{Paths}_{\mathcal{C}}^1(\varphi \downarrow)} \text{rew}_{\mathcal{C}}^{T(\varphi), 1}(\pi') \Pr_{\delta_{\varphi}}[d\pi'] + \\
& \quad \sum_{q \in S_{\text{mrk}}} \sum_{A_q \in \text{Act}_{\mathcal{C}}(q)} \int_0^{\infty} dR_{\mathcal{C}}^{\delta}(\varphi \xrightarrow{A_q, t} q, N-1) \cdot dt \\
&= \int_{\text{Paths}_{\mathcal{C}}^1(\varphi \downarrow)} \text{rew}_{\mathcal{C}}^{T(\varphi), 1}(\pi') \Pr_{\delta_{\varphi}}[d\pi'] + \\
& \quad \int_{\text{Paths}_{\mathcal{C}}^1(\varphi \downarrow)} \int_{\text{Paths}_{\mathcal{C}}(\pi' \downarrow)} \text{rew}_{\mathcal{C}}^{T(\varphi) + T(\pi'), (N-1)}(\pi) \Pr_{\delta_{\pi'}}[d\pi] \cdot \Pr_{\delta_{\varphi}}[d\pi'] \\
&= \int_{\text{Paths}_{\mathcal{C}}(\varphi \downarrow)} \text{rew}_{\mathcal{C}}^{T(\varphi), N}(\pi') \Pr_{\delta_{\varphi}}[d\pi'] \\
&= dR_{\mathcal{C}}^{\delta}(\varphi, N).
\end{aligned}$$

Thus, for each scheduler $\sigma \in GM$ on the MRA \mathcal{M} we built a scheduler δ on $\mathcal{C}(\mathcal{M})$, such that $dR_{\mathcal{M}, \beta}^{\sigma} = dR_{\mathcal{C}, \beta}^{\delta}$. Therefore, since the optimal scheduler for discounted reward on CTMDPs is stationary deterministic [18], the same holds for Markov reward automata and its discounted reward. \square

A.3 Uniformisation and Normalisation

Here we will show how to uniformise and normalise an arbitrary MRA.

Uniformisation. Markovian states of \mathcal{M} may have different exit rates. In order to achieve uniform exit rates we will use a well-known approach [17] by adding self-loop transitions to them. Let $\eta \geq \max_{s \in S} E(s)$. The *uniformisation* of \mathcal{M} to the rate η is the MRA $\mathcal{M}_\eta := (S_\eta, s_{\text{init}}, Act, \hookrightarrow, \rightsquigarrow_\eta, r, \rho)$, such that

$$\begin{aligned} s \xrightarrow{\mu} s', s \neq s' &\iff s \xrightarrow{\mu}_\eta s', s \neq s' \\ \eta \geq E(s) &\iff s \xrightarrow{\mu'}_\eta s, \text{ where } \mu' = \eta - (E(s) - R(s, s)). \end{aligned}$$

Lemma 4. \mathcal{M}_η is uniform, its size is linear in the size of \mathcal{M} and $\text{dR}_{\mathcal{M}_\eta, \beta}^{\text{opt}} = \text{dR}_{\mathcal{M}, \beta}^{\text{opt}}$.

Proof. It is obvious that thus obtained MRA is uniform. Since the transformation adds at most as many transitions as there are Markovian states, the transformation is linear. We will show now that the discounted reward is preserved. Let \mathcal{M} be a Markov reward automaton and \mathcal{M}_η - the uniformised version of \mathcal{M} with the uniformisation rate η .

We show now, that for every scheduler σ on \mathcal{M} there exists a scheduler σ_η on \mathcal{M}_η that achieves the same discounted reward value.

We denote with $[\pi]_\eta \subseteq \text{Paths}_{\mathcal{M}_\eta}$ the set of paths π' of \mathcal{M}_η that imitate π fully for probabilistic transitions, and partially for Markovian transitions. Namely, for all transitions leading from a state $s \in MS$ to state $s' \in S$, such that $s' \neq s$ π' takes exactly the same transition, however it may additionally have several self-loop transitions on s before moving to s' .

For each scheduler $\sigma \in GM$ on \mathcal{M} we define a scheduler $\sigma_\eta \in GM$ on \mathcal{M}_η , such that $\sigma_\eta(\pi', \alpha) = \sigma(\pi, \alpha)$, where $\pi' \in [\pi]_\eta$.

Let $\mathbf{E}[\text{rew}_{\mathcal{M}, \beta} | v, s' \neq v]$ denote the expected discounted reward, received by \mathcal{M} conditional on starting in a Markovian state v and until taking a transition to state $s' \neq v$.

As shown in [18] (Proposition 11.5.1), $\mathbf{E}[\text{rew}_{\mathcal{M}, \beta} | v, s' \neq v] = \mathbf{E}[\text{rew}_{\mathcal{M}_\eta, \beta} | v, s' \neq v]$. Therefore, the expected total discounted rewards during the residence in v agree for the original and the uniformised processes.

Additionally coincide the distributions of sojourn times for Markovian states in \mathcal{M} and \mathcal{M}_η , as well as the discrete distributions $\mathbb{P}_{\mathcal{M}}[v, s']$ and $\mathbb{P}_{\mathcal{M}_\eta}[v, s']$, for all $s' \neq v$.

Since the scheduler σ_η mimics σ , the expected reward gathered from probabilistic states in \mathcal{M}_η , as well as the discrete probability distributions over successor states, induced by σ_η coincide with those of σ on \mathcal{M} .

Thus, $\text{dR}_{\mathcal{M}_\eta, \beta}^{\sigma_\eta} = \text{dR}_{\mathcal{M}, \beta}^{\sigma}$. □

Normalisation. We will now *normalise* \mathcal{M} . Informally, in order to achieve the desired effect we simply introduce probabilistic states of zero reward. There are surely many ways of normalising an MRA. Here we present one of the possibilities to do it.

The *normalisation* of \mathcal{M} is an MRA $\overline{\mathcal{M}} = (\overline{S}, \overline{s_{\text{init}}}, \text{Act} \uplus \{\alpha^*\}, \overline{\hookrightarrow}, \overline{\rightsquigarrow}, \overline{r}, \overline{\rho})$. We start with \mathcal{M} and in the process of transformation will add new states and transitions to it, as well as modify existing ones:

1. Repeat until \mathcal{M} satisfies normalisation properties 2 and 3:
 - (a) If $s \in MS_{\mathcal{M}} \wedge \text{pred}(s) \cap MS_{\mathcal{M}} \neq \emptyset$, or $s \in PS_{\mathcal{M}} \wedge \text{pred}(s) \cap PS_{\mathcal{M}} \neq \emptyset \wedge \text{pred}(s) \cap MS_{\mathcal{M}} \neq \emptyset$, then create a new probabilistic state q and place it in between s and its Markovian predecessors. Or, formally:
 - i. add transition $q \xrightarrow{\alpha^*} \xi_s$;
 - ii. for all transitions $s' \xrightarrow{\lambda} s$, add a new transition $s' \xrightarrow{\lambda} q$ and remove the transition $s' \xrightarrow{\lambda} s$;
 - (b) If $s \in MS_{\mathcal{M}}$ and $\text{succ}(s) \cap MS_{\mathcal{M}} \neq \emptyset$, then then create a new probabilistic state q and
 - i. set $\lambda' := \sum_{s \xrightarrow{\lambda} s'} \lambda$ and add transition $s \xrightarrow{\lambda'} q$;
 - ii. for all transitions $s \xrightarrow{\lambda} s'$, add a new transition $q \xrightarrow{\alpha^*} \mu$, where $\mu(s') = \frac{\lambda}{\lambda'}$ and remove the transition $s \xrightarrow{\lambda} s'$;
2. if $s_{\text{init}} \in PS_{\mathcal{M}}$ and $\text{pred}(s_{\text{init}}) \subseteq MS_{\mathcal{M}}$, then $\overline{s_{\text{init}}} = s_{\text{init}}$. Otherwise create a new state s_{init}' , add it to S , add transition $(s_{\text{init}}', \alpha^*, \xi_{s_{\text{init}}})$ to \hookrightarrow , and set $\overline{s_{\text{init}}} = s_{\text{init}}'$;
3. set the rewards of all newly added states to 0.

Lemma 5. $\overline{\mathcal{M}}$ satisfies the normalisation properties, its size is linear in the size of \mathcal{M} and $\text{dR}_{\overline{\mathcal{M}}, \beta}^{\text{opt}} = \text{dR}_{\mathcal{M}, \beta}^{\text{opt}}$.

Proof. Obviously thus constructed MRA is normalised. The amount of added states is not higher than the amount of Markovian transitions of the original MRA +1 and the transformation is therefore linear. We will show now that it also preserves the discounted reward value.

We first show that there exists a one-to-one correspondence between the paths of the original MRA \mathcal{M} and its normalisation \mathcal{M}' . Let $\pi = s_0 \xrightarrow{\alpha_0, t_0} s_1 \xrightarrow{\alpha_1, t_1} \dots \in \text{Paths}_{\mathcal{M}}$ be a path in \mathcal{M} . Then there exists a path in the normalised MRA $\pi' \in \text{Paths}_{\mathcal{M}'}$ such that

- let $s_0 \xrightarrow{\alpha, t} s'$ be the very first transition from s_0 on π . If a new initial state s'_0 was added to \mathcal{M} during the normalisation, then the corresponding subpath of π' is $s'_0 \xrightarrow{\alpha', t'} s_0 \xrightarrow{\alpha, t} s'$, else $\pi'[0] = \pi[0]$;
- for each subpath $s \xrightarrow{\perp, t} s' \xrightarrow{\perp, t'} s''$ of π , the corresponding subpath of π' is $s \xrightarrow{\perp, t} q \xrightarrow{\alpha^*, 0} s' \xrightarrow{\perp, t'} s''$;
- if for some $k : \pi[k] = s' \in PS_{\mathcal{M}}$ and s' has both Markovian and probabilistic predecessors, then for each subpath $s \xrightarrow{\perp, t} s'$ of π , the corresponding subpath of π' is $s \xrightarrow{\perp, t} q \xrightarrow{\alpha^*, 0} s'$;
- π' repeats π in all other cases.

Analogously, for each path π' of \mathcal{M}' there exists a path π in \mathcal{M} . Since all the introduced states are probabilistic states, that have only 1 action leading to only 1 successor, the probability measure of a set of paths E in \mathcal{M} equals the probability measure of the corresponding set of paths E' in \mathcal{M}' . Since the introduced states have reward 0, then $\text{rew}(\pi) = \text{rew}(\pi')$. Thus the total expected reward value is preserved under normalisation. \square

Lemma 2. *For any MRA $\mathcal{M}, \eta \geq \max_{s \in S} E(s)$ there exists a uniform normalised MRA $\overline{\mathcal{M}}_\eta$, s. t. $\text{dR}_{\overline{\mathcal{M}}_\eta, \beta}^{\text{opt}} = \text{dR}_{\mathcal{M}, \beta}^{\text{opt}}$ and its size is linear in the size of \mathcal{M} .*

Proof. The proof follows from Lemmas 5 and 4.

A.4 Bellman Equations.

Bellman Equation derived from $\mathcal{C}(\mathcal{M})$

Theorem 4 (Bellman equation via $\mathcal{C}(\mathcal{M})$). *Let \mathcal{M} be a normalized uniformized MRA with exit rate η , $\mathcal{C}(\mathcal{M}) = (S_{\text{mrk}}, s_{\text{init}\mathcal{C}}, \text{Act}_{\mathcal{C}}, \text{R}_{\mathcal{C}})$ and $(\rho_{\mathcal{C}}, \text{r}_{\mathcal{C}})$ – the value preserving CTMDP and the corresponding reward structure. Then the vector $\text{dR}_{\mathcal{M}, \beta}^{\text{opt}} := (\text{dR}_{\mathcal{M}, \beta}^{\text{opt}}(s)), \forall s \in S_{\text{mrk}}$, is the unique solution to the Bellman equation:*

$$\forall s \in S_{\text{mrk}} : \quad v(s) = \underset{A_s \in \text{Act}(s)}{\text{opt}} \left\{ \text{r}_{\mathcal{C}}(s, A_s) + \frac{\rho_{\mathcal{C}}(s)}{\beta + \eta} + \frac{\eta}{\beta + \eta} \sum_{s' \in S_{\text{mrk}}} \mathbb{P}_{\mathcal{C}}[s, \alpha, s'] \cdot v(s') \right\}. \quad (2)$$

Proof. The proof follows directly from Theorem 1 and the corresponding proof for CTMDPs from [18].

Bellman Equation. Efficient Way.

Lemma 6. *Let \mathcal{D} and $\text{r}_{\mathcal{D}}$ be a MDP and a reward structure, such that the value $\text{tR}_{\mathcal{D}, \text{r}_{\mathcal{D}}}^{\text{opt}}(s)$ exists. Let $N > 0$, then*

$$\mathbf{E}_{s, D} \left[\sum_{i=0}^{N-1} \text{r}_{\mathcal{D}}(X_i^s, Y_i^s) \right] = \sum_{\substack{\pi \in \text{Paths}_{\mathcal{D}} \\ |\pi| = N}} \text{Pr}_D[\pi] \cdot \text{r}_{\mathcal{D}, D}(\pi),$$

where $\text{Pr}_D[\pi] = \prod_{i=0}^{|\pi|-1} \mathbb{P}_{\mathcal{D}}[\pi[i], D(\pi[i]), \pi[i+1]]$ and $\text{r}_{\mathcal{D}, D}(\pi) = \sum_{i=0}^{|\pi|-1} \text{r}_{\mathcal{D}}(\pi[i], D(\pi[i]))$.

Proof.

$$\begin{aligned}
\mathbf{E}_{s,D} \left[\sum_{i=0}^{N-1} r_{\mathcal{D}}(X_i^s, Y_i^s) \right] &= \sum_{i=0}^{N-1} \mathbf{E}_{s,D} [r_{\mathcal{D}}(X_i^s, Y_i^s)] \\
&= \sum_{i=0}^{N-1} \sum_{s' \in S_{\mathcal{D}}} \Pr_D[X_i^s = s' | X_0^s = s] \cdot r_{\mathcal{D}}(X_i^s, Y_i^s) \\
&= \sum_{i=0}^{N-1} \sum_{s' \in S_{\mathcal{D}}} \sum_{\pi, \pi[i]=s'} \Pr_D[\pi] \cdot r_{\mathcal{D}}(s', D(s')) \\
&= \sum_{i=0}^{N-1} \sum_{\substack{\pi \in \text{Paths}_{\mathcal{D}} \\ |\pi| = N}} \Pr_D[\pi] \cdot r_{\mathcal{D}}(\pi[i], D(\pi[i])) \\
&= \sum_{\substack{\pi \in \text{Paths}_{\mathcal{D}} \\ |\pi| = N}} \Pr_D[\pi] \cdot \sum_{i=0}^{N-1} r_{\mathcal{D}}(\pi[i], D(\pi[i])) \\
&= \sum_{\substack{\pi \in \text{Paths}_{\mathcal{D}} \\ |\pi| = N}} \Pr_D[\pi] \cdot r_{\mathcal{D},D}(\pi).
\end{aligned}$$

□

Theorem 2. Let \mathcal{M} be a uniform normalised MRA with exit rate η and $\mathcal{D}(\mathcal{M})$ the corresponding terminal MDP. Then the vector $\mathbf{dR}_{\mathcal{M},\beta}^{\text{opt}} := (\mathbf{dR}_{\mathcal{M},\beta}^{\text{opt}}(s)), \forall s \in S_{\text{mrk}}$, is the unique solution to the Bellman equation:

$$\forall s \in S_{\text{mrk}} : \quad v(s) = \mathbf{tR}_{\mathcal{D}(\mathcal{M}), \text{rew}_{\mathcal{D}(\mathcal{M}),v}}^{\text{opt}}(s). \quad (1)$$

Proof. Consider the right-hand side of equation 2:

$$\begin{aligned}
&r_{\mathcal{C}}(s, A_s) + \frac{\rho_{\mathcal{C}}(s)}{\beta + \eta} + \frac{\eta}{\beta + \eta} \sum_{s' \in S_{\text{mrk}}} \mathbb{P}_{\mathcal{C}}[s, \alpha, s'] \cdot v(s') \\
&= \sum_{s' \in S_{\text{mrk}}} \sum_{\pi \in \Pi_{S_{\text{mrk}}}(s, s')} r(\pi) \cdot \mathbb{P}[\pi] \\
&\quad + \frac{1}{\beta + \eta} \cdot \sum_{s' \in S_{\text{mrk}}} \sum_{\pi \in \Pi_{S_{\text{mrk}}}(s, A, s')} \rho(\pi) \cdot \mathbb{P}[\pi] \\
&\quad + \frac{\eta}{\beta + \eta} \sum_{s' \in S_{\text{mrk}}} \sum_{\pi \in \Pi_{S_{\text{mrk}}}(s, A, s')} \mathbb{P}[\pi] \cdot v(s') \\
&= \sum_{s' \in S_{\text{mrk}}} \sum_{\pi \in \Pi_{S_{\text{mrk}}}(s, s')} \mathbb{P}[\pi] \cdot \left[r(\pi) + \frac{\rho(\pi)}{\beta + \eta} + \frac{\eta}{\beta + \eta} \cdot v(s') \right] \\
&\stackrel{\text{Lemma 4}}{=} \mathbf{tR}_{\mathcal{D}(\mathcal{M}), \text{rew}_{\mathcal{D}(\mathcal{M}),v}}^{\text{opt}}(s).
\end{aligned}$$

□

A.5 Proof of Theorem 3

Theorem 3. *Algorithms 1 and 2 are sound and complete.*

Proof. The proof of Theorem 2 shows that the right-hand side of equation (2) equals the right-hand side of equation (1). Thus, a value- or modified policy-iteration algorithm for the Bellman equation from Theorem 4 in which the computation of the right-hand side of (2) is substituted with the right-hand side of (1) is as sound and complete as the original algorithm, given that the expected total reward is quantified precisely. The proof follows from results for respective algorithms for CTMDPs, see [18].