

# On reduction criteria for probabilistic reward models

Marcus Größer<sup>a</sup>, Gethin Norman<sup>b</sup>, Christel Baier<sup>a</sup>,  
Frank Ciesinski<sup>a</sup>, Marta Kwiatkowska<sup>b</sup>, David Parker<sup>b</sup>

<sup>a</sup> *Universität Bonn, Institut für Informatik I, Germany*  
{baier | groesser | ciesinsk}@cs.uni-bonn.de

<sup>b</sup> *University of Birmingham, School of Computer Science, Edgbaston, United Kingdom*  
{kwiatkowska | norman | parker}@cs.bham.ac.uk

**Abstract.** In recent papers, the partial order reduction approach has been adapted to reason about the probabilities for temporal properties in concurrent systems with probabilistic behaviours. This paper extends these results by presenting reduction criteria for a probabilistic branching time logic that allows specification of constraints on quantitative measures given by a reward or cost function for the actions of the system.

## 1 Introduction

Partial order reduction [14, 27, 33] is one of the most prominent techniques for tackling the state explosion problem for concurrent software systems. It has been implemented in many tools and successfully applied to a large number of case studies, see e.g. [15, 20]. Recently, the ample-set method [26] has been extended for concurrent probabilistic systems, both in the setting of quantitative linear time [4, 9] and branching time [3] properties. The underlying models used in this work are Markov decision processes (MDPs), an extension of transition systems where nondeterminism can be used e.g. to model the interleaving of concurrent activities, to represent the interface with an unknown system environment or for abstraction purposes, and where probability serves e.g. to model coin tossing actions or to specify the frequency of exceptional (faulty) behaviour (such as losing messages from a buffer). Thus, MDPs arise as natural operational models for randomized distributed algorithms and communication or security protocols. Equipped with reward or cost functions, MDPs are also standard models in many other areas, such as reinforcement learning, robot path planning, or operations research.

The contribution of this paper is reduction criteria which are shown to be sound for an extension of probabilistic computation tree logic (PCTL) [7] that serves to reason about rewards or costs. Our logic, called  $PCTL_r$ , essentially agrees with the logic suggested by de Alfaro [11, 10]. ( $PCTL_r$  is also similar to the logic PRCTL [1, 25] which relies on a Markov chain semantics, while  $PCTL_r$ -formulae are interpreted over MDPs.)  $PCTL_r$  allows specifications regarding e.g. the packet loss characteristics of a queueing system, the energy consumption, or the average number of unsuccessful attempts to find a leader in a distributed system. We first explain how the ample-set conditions suggested in [3] for PCTL can be modified to treat reward-based properties specified in  $PCTL_r$  and then identify a fragment of  $PCTL_r$  (which still contains a wide range of

non-trivial reward properties) where the weaker criteria of [3] are sufficient. We also present results on a new logic  $\text{PCTL}_c$ , that treats the rewards with a discounting semantics. As in the case of previous publications on partial order reduction for probabilistic systems, the major difficulty was to provide the proof of correctness. The general proof technique follows the line of [13, 3] by establishing a bisimulation between the full and the reduced system. However, we depart here from these approaches by introducing a new variant of bisimulation equivalence for MDPs which borrows ideas from [24, 16, 5] and relies on the concept of norm functions. This new type of bisimulation equivalence preserves  $\text{PCTL}_r$ -properties and might be useful also for other purposes.

**Organization of the paper.** Section 2 summarizes the basic definitions concerning Markov decision processes, reward structures and  $\text{PCTL}_r$ . Section 2 also recalls the partial order reduction approach for MDPs without reward structure and  $\text{PCTL}$  of [3] which we then extend to reason about rewards in Section 3. Section 4 identifies a class of reward-based properties that are preserved when using the weaker conditions of [3]. In Section 5 we discuss our approach in the setting of discounted rewards and Section 6 concludes the paper.

## 2 Preliminaries

**Markov decision processes (MDPs), see e.g. [29].** An MDP is a tuple  $\mathcal{M} = (S, \text{Act}, \mathbf{P}, s_{\text{init}}, \text{AP}, L, \text{rew})$  where  $S$  is a finite state space,  $s_{\text{init}} \in S$  is the initial state,  $\text{Act}$  a finite set of actions,  $\text{AP}$  a set of atomic propositions,  $L : S \rightarrow 2^{\text{AP}}$  a labelling function,  $\mathbf{P} : S \times \text{Act} \times S \rightarrow [0, 1]$  the three-dimensional transition probability matrix such that  $\sum_{u \in S} \mathbf{P}(s, \alpha, u) \in \{0, 1\}$  for all states  $s$  and actions  $\alpha$ , and a function  $\text{rew}$  that assigns to each action  $\alpha \in \text{Act}$  a reward  $\text{rew}(\alpha) \in \mathbb{R}$ .

Action  $\alpha$  is called enabled in state  $s$  if  $\sum_{u \in S} \mathbf{P}(s, \alpha, u) = 1$ . We write  $\text{Act}(s)$  for the set of actions that are enabled in  $s$ . The states  $t$  with  $\mathbf{P}(s, \alpha, t) > 0$  are called  $\alpha$ -successors of  $s$ . For technical reasons, we require that  $\text{Act}(s) \neq \emptyset$  for all states  $s$ . Action  $\alpha$  is called a *stutter action* iff for all  $s \in S$  where  $\alpha$  is enabled in  $s$ ,  $L(s) = L(u)$  for all  $\alpha$ -successors  $u$  of  $s$ . That is, stutter actions do not change the state labelling. Action  $\alpha$  is called *non-probabilistic* iff for all states  $s$ , there is at most one  $\alpha$ -successor. That is, if  $\alpha$  is enabled in  $s$  then there is a state  $s_\alpha$  with  $\mathbf{P}(s, \alpha, s_\alpha) = 1$ , while  $\mathbf{P}(s, \alpha, u) = 0$  for all other states  $u$ . In particular, if  $\alpha \in \text{Act}(s)$  is a non-probabilistic stutter action then  $L(s) = L(s_\alpha)$ .

An infinite *path* in an MDP is a sequence  $\zeta = s_0 \xrightarrow{\alpha_1} s_1 \xrightarrow{\alpha_2} s_2 \xrightarrow{\alpha_3} \dots$  such that  $\alpha_i \in \text{Act}(s_{i-1})$  and  $\mathbf{P}(s_{i-1}, \alpha_i, s_i) > 0$  for all  $i \geq 1$ . We denote by  $\text{first}(\zeta) = s_0$  the starting state of  $\zeta$  and write  $\text{state}(\zeta, i)$  for the  $(i+1)$ th state in  $\zeta$  and  $\rho(\zeta, i)$  for the cumulative reward obtained through the first  $i$  actions. That is, if  $\zeta$  is as above then  $\text{state}(\zeta, i) = s_i$  and  $\rho(\zeta, i) = \text{rew}(\alpha_1 \dots \alpha_i)$  where  $\text{rew}(\alpha_1 \dots \alpha_i) = \text{rew}(\alpha_1) + \dots + \text{rew}(\alpha_i)$ . If  $T \subseteq S$  is a set of states then  $\text{Rew}(\zeta, T)$  denotes the reward that is earned until a  $T$ -state is visited the first time. Formally, if  $\text{state}(\zeta, i) \in T$  and  $\text{state}(\zeta, j) \notin T$  for all  $j < i$  then  $\text{Rew}(\zeta, T) = \rho(\zeta, i)$ . If  $\text{state}(\zeta, i) \notin T$  for all  $i \geq 0$  we set  $\text{Rew}(\zeta, T) = \infty$ . Finite paths (denoted by  $\sigma$ ) are finite prefixes of infinite paths that end in a state. We use the notations  $\text{first}(\sigma)$ ,  $\text{state}(\sigma, i)$  and  $\rho(\sigma, i)$  as for infinite paths and  $|\sigma|$  for the length (number of

actions).  $Paths_{\text{fin}}(s)$  (resp.  $Paths_{\omega}(s)$ ) denotes the set of all finite (resp. infinite) paths of  $\mathcal{M}$  with  $first(\cdot) = s$ .

A *scheduler*, also often called policy, strategy or adversary, denotes an instance that resolves the nondeterminism in the states, and thus yields a Markov chain and a probability measure on the paths. We shall use here *history-dependent randomized schedulers* in the classification of [29]. They are defined as functions  $A$  that take as input a finite path  $\sigma$  and return a distribution over the actions  $\alpha \in Act(last(\sigma))$ .<sup>1</sup> A scheduler  $A$  is called deterministic if it chooses a unique action (with probability 1) for all finite paths. An  $A$ -path denotes an infinite or finite path  $\sigma$  that can be generated by  $A$ . Given a state  $s$  and a scheduler  $A$ , the behaviour of  $\mathcal{M}$  under  $A$  can be formalised by a (possibly infinite-state) Markov chain.  $\text{Pr}^{A,s}$  denotes the standard probability measure on the Borel field of the infinite  $A$ -paths  $\zeta$  with  $first(\zeta) = s$ . If  $T \subseteq S$  then  $\mathbb{E}^{A,s}(\diamond T)$  denotes the expected value under  $A$  with starting state  $s$  for the random function  $\zeta \mapsto \text{Rew}(\zeta, T)$ . Recall that  $\text{Rew}(\zeta, T)$  denotes the reward that is earned by the prefix of  $\zeta$  that leads from the starting state  $s$  to a state in  $T$  and that  $\text{Rew}(\zeta, T)$  equals  $\infty$  if  $\zeta$  does not reach  $T$ . Thus, if there is a positive probability of not reaching  $T$  under scheduler  $A$  (from state  $s$ ), then  $\mathbb{E}^{A,s}(\diamond T) = \infty$ . If  $s = s_{\text{init}}$  we simply write  $\text{Pr}^A$  and  $\mathbb{E}^A$ .

**Probabilistic computation tree logic.** PCTL is a probabilistic variant of CTL [8] which has been introduced first for Markov chains [18] and then for Markovian models with non-determinism [17, 7, 32]. We follow here the approach of de Alfaro [11, 10] and extend PCTL with an operator  $\mathcal{R}$  to reason about expected rewards. As partial order reduction relies on identifying stutter equivalent paths which might be distinguishable by the next step operator, we do not include the next step operator in the logic. PCTL $_r$ -state formulae are therefore given by the grammar:

$$\Phi ::= \text{true} \mid a \mid \Phi \wedge \Phi \mid \neg\Phi \mid \mathcal{P}_J(\Phi_1 U_I \Phi_2) \mid \mathcal{R}_I(\Phi)$$

Here,  $a \in AP$  is an atomic proposition,  $J \subseteq [0, 1]$  is a probability interval and  $I \subseteq \mathbb{R} \cup \{-\infty, \infty\}$  a reward interval. We refer to the terms  $\Phi_1 U_I \Phi_2$  as PCTL $_r$ -path formulae.  $U_I$  denotes the standard until operator with a reward bound. The meaning of the path formula  $\varphi = \Phi_1 U_I \Phi_2$  is that a  $\Phi_2$ -state will be reached via a finite path  $\sigma$  where the cumulative reward is in  $I$ , while all states in  $\sigma$ , possibly except the last one, fulfil  $\Phi_1$ . The state formula  $\mathcal{P}_J(\varphi)$  holds for state  $s$  if for each scheduler  $A$  the probability measure of all infinite paths starting in  $s$  and fulfilling the path formula  $\varphi$  meets the probability bound given by  $J$ . On the other hand,  $\mathcal{R}_I(\Phi)$  asserts that for any scheduler  $A$  the expected reward that is earned until a  $\Phi$ -state has been reached meets the reward bound given by  $I$ . For instance,  $\mathcal{R}_{[0,17]}(\text{goal})$  asserts that independent of the scheduling policy the average costs to reach a goal state do not exceed 17. The formula  $\mathcal{P}_{(0.9,1]}(\text{true } U_{[0,4]} \text{ delivered})$  requires that the probability of a message being delivered with at most 4 retransmissions is greater than 0.9.

If  $\mathcal{M}$  is an MDP and  $s$  a state in  $\mathcal{M}$  then we write  $s \models \Phi$  to denote that state-formula  $\Phi$  holds in state  $s$ , and similarly,  $\zeta \models \varphi$  to denote that path formula  $\varphi$  holds for

<sup>1</sup> By a distribution on a finite set  $X$  we mean a function  $\mu : X \rightarrow [0, 1]$  such that  $\sum_{x \in X} \mu(x) = 1$  and refer to  $\mu(x)$  as the probability for  $x$ .

the infinite path  $\zeta$ . The formal semantics is formalised by:

$$\begin{aligned}
s &\models \text{true} \\
s &\models a && \Leftrightarrow a \in L(s) \\
s &\models \Phi_1 \wedge \Phi_2 && \Leftrightarrow s \models \Phi_1 \text{ and } s \models \Phi_2 \\
s &\models \neg\Phi && \Leftrightarrow s \not\models \Phi \\
s &\models \mathcal{P}_J(\Phi_1 U_I \Phi_2) && \Leftrightarrow \text{for all schedulers } A: \Pr^{A,s}\{\zeta \in \text{Paths}_\omega(s) : \zeta \models \Phi_1 U_I \Phi_2\} \in J \\
s &\models \mathcal{R}_I(\Phi) && \Leftrightarrow \text{for all schedulers } A: E^{A,s}(\diamond \text{Sat}(\Phi)) \in I
\end{aligned}$$

If  $\zeta = s_0 \xrightarrow{\alpha_1} s_1 \xrightarrow{\alpha_2} s_2 \xrightarrow{\alpha_3} \dots$  then  $\zeta \models \Phi_1 U_I \Phi_2$  iff  $\exists i \geq 0$  s.t.  $s_i \models \Phi_2 \wedge \rho(\zeta, i) \in I \wedge \forall j < i. s_j \models \Phi_1$ . The satisfaction set of  $\Phi$  in  $\mathcal{M}$  is  $\text{Sat}(\Phi) = \{s \in S : s \models \Phi\}$ . State formula  $\Phi$  is said to hold for an MDP if the initial state satisfies  $\Phi$ .

Note that one could also give the  $\mathcal{R}_I$  operator a different semantics as follows.  $s \models \mathcal{R}_I(\Phi)$  if and only if for all schedulers  $A$ , such that the probability to reach  $\text{Sat}(\Phi)$  from  $s$  equals 1, it holds that  $E^{A,s}(\diamond \text{Sat}(\Phi)) \in I$ . But this is irrelevant for our purposes.

*Derived operators.* Other Boolean connectives, such as disjunction  $\vee$ , implication  $\rightarrow$ , can be derived as usual. The temporal operator eventually  $\diamond$  is obtained in the standard way by  $\diamond_I \Phi = \text{true } U_I \Phi$ . The always-operator can be derived as in PCTL by the duality of lower and upper probability bounds. E.g., we may define  $\mathcal{P}_{[0,p]}(\square_I \Phi) = \mathcal{P}_{(1-p,1]}(\diamond_I \neg \Phi)$ . The formula  $\diamond_I \Phi$  holds for a path  $\zeta$  if one can reach a state which satisfies  $\Phi$  and the reward cumulated up until this point is in the interval  $I$ . On the other hand, the dual always-operator  $\square_I \Phi$  holds for a path  $\zeta$  if whenever the cumulated reward is in the interval  $I$  the formula  $\Phi$  holds. Thus, for the trivial reward-interval  $I = (-\infty, \infty)$ , we obtain the standard eventually and always operators. The same holds for the until operator. We simply write  $U$ ,  $\diamond$  and  $\square$  rather than  $U_{(-\infty, \infty)}$ ,  $\diamond_{(-\infty, \infty)}$  and  $\square_{(-\infty, \infty)}$ , respectively.

PCTL denotes the sublogic of  $\text{PCTL}_r$ , that does not use the  $\mathcal{R}$ -operator and where the path-formulae have the trivial reward interval. Since the reward structure is irrelevant for PCTL-formulae, they can be interpreted over MDPs without reward structure.

**The ample set method for PCTL [3].** Before presenting the partial order reduction criteria for  $\text{PCTL}_r$  in Section 3, we briefly summarize the results of [3] for applying the ample-set method to PCTL model checking. The starting point is an MDP  $\mathcal{M} = (S, \text{Act}, \mathbf{P}, s_{\text{init}}, \text{AP}, L)$ , without reward structure, to be verified against a PCTL-formula. Following Peled's ample-set method [26], the idea is to assign to any reachable state  $s$  a nonempty action-set  $\text{ample}(s) \subseteq \text{Act}(s)$  and to construct a reduced MDP  $\hat{\mathcal{M}}$  by using the action-sets  $\text{ample}(s)$  instead of  $\text{Act}(s)$ . Formally, given a function  $\text{ample} : S \rightarrow 2^{\text{Act}}$  with  $\emptyset \neq \text{ample}(s) \subseteq \text{Act}(s)$  for all states  $s$ , the state space of the reduced MDP  $\hat{\mathcal{M}} = (\hat{S}, \text{Act}, \hat{\mathbf{P}}, s_{\text{init}}, \text{AP}, \hat{L})$  induced by  $\text{ample}$  is the smallest set  $\hat{S} \subseteq S$  that contains  $s_{\text{init}}$  and any state  $u$  where  $\mathbf{P}(s, \alpha, u) > 0$  for some  $s \in \hat{S}$  and  $\alpha \in \text{ample}(s)$ . The labelling function  $\hat{L} : \hat{S} \rightarrow 2^{\text{AP}}$  is the restriction of the original labelling function  $L$  to the state-set  $\hat{S}$ . The transition probability matrix of  $\hat{\mathcal{M}}$  is given by  $\hat{\mathbf{P}}(s, \alpha, t) = \mathbf{P}(s, \alpha, t)$  if  $\alpha \in \text{ample}(s)$  and 0 otherwise. State  $s$  is called fully expanded if  $\text{ample}(s) = \text{Act}(s)$ .

The main ingredient of any partial order reduction technique in the non-probabilistic or probabilistic setting is an adequate notion for the independence of actions. The definition for the independence of actions  $\alpha$  and  $\beta$  in the composed transition system (which

- A1 (Stutter-condition)** If  $ample(s) \neq Act(s)$  then all actions  $\alpha \in ample(s)$  are stutter actions.
- A2 (Dependence-condition)** For each path  $\sigma = s \xrightarrow{\alpha_1} \dots \xrightarrow{\alpha_n} s_n \xrightarrow{\gamma} \dots$  in  $\mathcal{M}$  where  $\gamma$  is dependent on  $ample(s)$  there exists an index  $i \in \{1, \dots, n\}$  such that  $\alpha_i \in ample(s)$ .
- A3 (Cycle-condition)** On each cycle  $s \xrightarrow{\alpha_1} s_1 \xrightarrow{\alpha_2} \dots \xrightarrow{\alpha_n} s_n = s$  in  $\hat{\mathcal{M}}$  there exists a state  $s_i$  which is fully expanded, i.e.,  $ample(s_i) = Act(s_i)$ .
- A4 (Branching condition)** If  $ample(s) \neq Act(s)$  then  $ample(s)$  is a singleton consisting of a non-probabilistic action.

**Fig. 1. Conditions for the ample-set method for PCTL [3]**

captures the semantics of the parallel composition of all processes that run in parallel) relies on recovering the interleaving ‘diamonds’. Formally, two distinct actions  $\alpha$  and  $\beta$  are called independent (in  $\mathcal{M}$ ) iff for all states  $s \in S$  with  $\{\alpha, \beta\} \subseteq Act(s)$ ,

- (I1)  $\mathbf{P}(s, \alpha, u) > 0$  implies  $\beta \in Act(u)$ ,  
 (I2)  $\mathbf{P}(s, \beta, u) > 0$  implies  $\alpha \in Act(u)$   
 (I3)  $\mathbf{P}(s, \alpha\beta, w) = \mathbf{P}(s, \beta\alpha, w)$  for all  $w \in S$ ,

where  $\mathbf{P}(s, \gamma\delta, w) = \sum_{u \in S} \mathbf{P}(s, \gamma, u) \cdot \mathbf{P}(u, \delta, w)$  for  $\gamma, \delta \in Act$ .

Two different actions  $\alpha$  and  $\beta$  are called dependent iff  $\alpha$  and  $\beta$  are not independent. If  $D \subseteq Act$  and  $\alpha \in Act \setminus D$  then  $\alpha$  is called independent of  $D$  iff for all actions  $\beta \in D$ ,  $\alpha$  and  $\beta$  are independent. Otherwise,  $\alpha$  is called dependent on  $D$ .

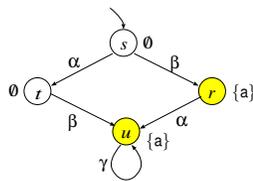
To preserve PCTL properties, [3] use the four conditions in Fig. 1. These rely on a slight modification of the conditions by Gerth et al [13] for preserving CTL-properties.

**Theorem 1 ([3]).** *If (A1)-(A4) hold then  $\mathcal{M}$  and  $\hat{\mathcal{M}}$  fulfil the same PCTL-formulae.*

### 3 Reduction criteria for rewards

In the sequel, we assume that we are given an MDP  $\mathcal{M}$  and discuss the partial order reduction approach for properties specified in  $PCTL_r$ . We first show that the conditions (A1)-(A4) are not sufficient to preserve  $PCTL_r$  properties with nontrivial reward bounds. To treat full  $PCTL_r$ , we shall need a modification of the branching condition (A4).

*Example 1.* We begin with a simple example illustrating that (A1)-(A4) cannot ensure that all  $PCTL_r$ -formulae are preserved. Consider the following MDP with the actions  $\alpha, \beta, \gamma$  that are all non-probabilistic and where  $\text{rew}(\alpha) = \text{rew}(\beta) = \text{rew}(\gamma) = 1$ .



Since  $\alpha$  and  $\beta$  are independent and  $\alpha$  is a stutter action, (A1)-(A4) allow for a reduction obtained through  $ample(s) = \{\alpha\}$ . Thus,  $\hat{S} = \{s, t, u\}$ . Consider the  $PCTL_r$  formula  $\Phi = \mathcal{R}_{[2, \infty)}(a)$ . Then, the reduced system  $\hat{\mathcal{M}}$  satisfies  $\Phi$ , while the original system  $\mathcal{M}$  does not, because  $\mathcal{M}$  might choose action  $\beta$  in  $s$  which yields the expected reward 1 to reach an a-state.  $\square$

We now discuss how to strengthen conditions (A1)-(A4) such that reward-based properties are preserved. We start with some simple observations. First, as  $\hat{\mathcal{M}}$  is a sub-MDP of the original system  $\mathcal{M}$ , any scheduler  $A$  for  $\hat{\mathcal{M}}$  is also a scheduler for  $\mathcal{M}$ . Thus:

**Lemma 1.** *Let  $\Phi_1, \Phi_2$  be PCTL<sub>r</sub>-formulae with  $Sat_{\mathcal{M}}(\Phi_i) \cap \hat{S} = Sat_{\hat{\mathcal{M}}}(\Phi_i)$ ,  $i = 1, 2$ . Then:*

- (i)  $\mathcal{M} \models \mathcal{R}_I(\Phi_1) \Rightarrow \hat{\mathcal{M}} \models \mathcal{R}_I(\Phi_1)$ ,
- (ii)  $\mathcal{M} \models \mathcal{P}_J(\Phi_1 U_I \Phi_2) \Rightarrow \hat{\mathcal{M}} \models \mathcal{P}_J(\Phi_1 U_I \Phi_2)$ .

The converse directions in Lemma 1 do not hold in general as  $\mathcal{M}$  might have “more” schedulers than  $\hat{\mathcal{M}}$ . To get a feeling of how to modify the reduction criteria for PCTL<sub>r</sub>, let us first give some informal explanations. In [3], the soundness proof of (A1)-(A4) for PCTL establishes a kind of bisimulation between the full MDP  $\mathcal{M}$  and the reduced MDP  $\hat{\mathcal{M}}$  which allows one to transform any scheduler  $A$  for  $\mathcal{M}$  into a scheduler  $B$  for  $\hat{\mathcal{M}}$  such that  $A$  and  $B$  yield the same probabilities for PCTL-path formulae. As in the case of the ample-set method for verifying linear time properties (where (A1)-(A3) and a weaker form of (A4) are sufficient [4, 9]) this scheduler-transformation yields a transformation of the  $A$ -paths into “corresponding”  $B$ -paths. Let us look at this path-transformation “path  $\zeta$  in  $\mathcal{M} \rightsquigarrow$  path  $\hat{\zeta}$  in  $\hat{\mathcal{M}}$ ” which, in fact, is already known from the non-probabilistic case [26]. The path  $\hat{\zeta}$  in  $\hat{\mathcal{M}}$  is obtained through a sequence of paths  $\zeta_0, \zeta_1, \zeta_2, \dots$  in  $\mathcal{M}$  such that the first  $i$ -steps in  $\zeta_i$  and  $\zeta_{i+1}$  agree and are composed of transitions in  $\hat{\mathcal{M}}$ . The switch from  $\zeta_i$  to  $\zeta_{i+1}$  is performed as follows.

Let  $\pi = s_1 \xrightarrow{\alpha_1} s_2 \xrightarrow{\alpha_2} \dots$  be the suffix of  $\zeta_i$  starting with the  $(i+1)$ th step (by the above,  $s_1$  is a state in  $\hat{\mathcal{M}}$ ). Our goal is to construct a stutter equivalent path  $\hat{\pi}$  from  $s_1$  that starts with an action in  $ample(s_1)$ . We then may compose the prefix of  $\zeta$  from  $s$  to  $s_1$  with  $\hat{\pi}$  to obtain the path  $\zeta_{i+1}$ . If  $\alpha_1 \in ample(s_1)$  then we may put  $\pi = \hat{\pi}$ . Let us now assume that  $\alpha_1 \notin ample(s_1)$ . Then, by (A4),  $ample(s_1)$  consists of a single non-probabilistic action.

- (T1) If there is some index  $j \geq 2$  such that  $\alpha_j \in ample(s_1)$  then choose the smallest such index  $j$  and replace the action sequence  $\alpha_1 \dots \alpha_{j-1} \alpha_j \alpha_{j+1} \dots$  with  $\alpha_j \alpha_1 \dots \alpha_{j-1} \alpha_{j+1} \dots$ . This is possible since by (A2) the actions  $\alpha_1, \dots, \alpha_{j-1}$  are independent of  $\alpha_j$ . The resulting path  $\hat{\pi}$  is stutter-equivalent to  $\pi$  by condition (A1).
- (T2) If  $\alpha_j \notin ample(s_1)$  for all  $j \geq 1$  and  $ample(s) = \{\beta\}$  then replace the action sequence  $\alpha_1 \alpha_2 \dots$  with  $\beta \alpha_1 \alpha_2 \dots$ . Again, (A2) ensures that each  $\alpha_j$  is independent of  $\beta$ . (A1) yields the stutter-equivalence of  $\pi$  and the resulting path  $\hat{\pi}$ .

Note, that the insertion of the additional action in transformation (T2) possibly changes the cumulative reward. Since we are interested in the cumulative reward that is gained until a certain state labelling is reached, the action permutation in transformation (T1) possibly changes this reward, as can be seen in Example 1 (note that a stutter action is permuted to the front of the action sequence).

To establish the equivalence of  $\mathcal{M}$  and  $\hat{\mathcal{M}}$  for PCTL<sub>r</sub> it seems to be sufficient to ensure that, in transformation (T2), the additional action  $\beta$  has zero reward, and in transformation (T1), the stutter action  $\alpha_j$ , that is permuted to the front of the action sequence, has zero reward. This motivates the following stronger branching condition:

**A4' (New branching condition)** If  $\text{ample}(s) \neq \text{Act}(s)$  then  $\text{ample}(s) = \{\beta\}$  for some non-probabilistic action with  $\text{rew}(\beta) = 0$ .

**Theorem 2.** *If (A1)-(A3), (A4') hold then  $\mathcal{M}$  and  $\hat{\mathcal{M}}$  satisfy the same  $\text{PCTL}_r$  formulae.*

The remainder of this section is concerned with the proof of Theorem 2. The above argument only applies to the path-level (and linear time properties), while the correctness proof of (A1)-(A3), (A4') for  $\text{PCTL}_r$  is not obvious since it requires reasoning about the probabilities for path-sets induced by schedulers rather than single paths.

As is the case for many other types of (bi)simulation relations for probabilistic systems, our notion of bisimulation equivalence will use the concept of *weight functions* [21, 22]. Let  $S, S'$  be finite sets and  $R \subseteq S \times S'$ . If  $\mu$  and  $\mu'$  are distributions on  $S$  and  $S'$  respectively then a weight function for  $(\mu, \mu')$  with respect to  $R$  denotes a function  $w : S \times S' \rightarrow [0, 1]$  such that  $\{(s, s') : w(s, s') > 0\} \subseteq R$ ,  $\sum_{u' \in S'} w(s, u') = \mu(s)$  and  $\sum_{u \in S} w(u, s') = \mu'(s')$  for all  $s \in S, s' \in S'$ . We write  $\mu \sqsubseteq_R \mu'$  iff there exists a weight function for  $(\mu, \mu')$  with respect to  $R$  and refer to  $\sqsubseteq_R$  as the lifting of  $R$  to distributions.

**Definition 1 (Normed reward (bi)simulation).** Let  $\mathcal{M} = (S_{\mathcal{M}}, \text{Act}, \mathbf{P}_{\mathcal{M}}, s_{\text{init}}^{\mathcal{M}}, \text{AP}, L_{\mathcal{M}}, \text{rew})$  and  $\mathcal{N} = (S_{\mathcal{N}}, \text{Act}, \mathbf{P}_{\mathcal{N}}, s_{\text{init}}^{\mathcal{N}}, \text{AP}, L_{\mathcal{N}}, \text{rew})$  be two MDPs with the same set of atomic propositions, the same action set  $\text{Act}$  and the same reward structure  $\text{rew} : \text{Act} \rightarrow \mathbb{R}_{\geq 0}$ . A normed reward simulation for  $(\mathcal{M}, \mathcal{N})$  with respect to  $\text{rew}$  is a triple  $(R, \eta_1, \eta_2)$  consisting of a binary relation  $R \subseteq S_{\mathcal{M}} \times S_{\mathcal{N}}$  and functions  $\eta_1, \eta_2 : R \rightarrow \mathbb{N}$  such that  $(s_{\text{init}}^{\mathcal{M}}, s_{\text{init}}^{\mathcal{N}}) \in R$  and for each pair  $(s, s') \in R$  the following conditions hold.

(N1)  $L_{\mathcal{M}}(s) = L_{\mathcal{N}}(s')$

(N2) If  $\alpha \in \text{Act}_{\mathcal{M}}(s)$  then at least one of the following conditions holds:

(N2.1)  $\alpha$  is enabled in  $s'$  (i.e.,  $\alpha \in \text{Act}_{\mathcal{N}}(s')$ ) and  $\mathbf{P}_{\mathcal{M}}(s, \alpha, \cdot) \sqsubseteq_R \mathbf{P}_{\mathcal{N}}(s', \alpha, \cdot)$ ,

(N2.2)  $\alpha$  is a non-probabilistic stutter action with  $\text{rew}(\alpha) = 0$  such that  $(s_{\alpha}, s') \in R$  and  $\eta_1(s_{\alpha}, s') < \eta_1(s, s')$ .

(N2.3) There is a non-probabilistic stutter action  $\beta \in \text{Act}_{\mathcal{N}}(s')$  with  $\text{rew}(\beta) = 0$ ,  $(s, s'_{\beta}) \in R$  and  $\eta_2(s, s'_{\beta}) < \eta_2(s, s')$ .

A normed reward bisimulation for  $(\mathcal{M}, \mathcal{N})$  is a tuple  $(R, \eta_1, \eta_2, \eta_1^-, \eta_2^-)$  such that  $(R, \eta_1, \eta_2)$  is a normed reward simulation for  $(\mathcal{M}, \mathcal{N})$  and  $(R^{-1}, \eta_1^-, \eta_2^-)$  is a normed reward simulation for  $(\mathcal{N}, \mathcal{M})$ .  $\square$

We write  $\mathcal{M} \approx_{nrB} \mathcal{N}$  iff there exists a normed reward bisimulation for  $\mathcal{M}$  and  $\mathcal{N}$ . To speak about normed reward bisimulation equivalence classes we consider the MDP  $\mathcal{M} \uplus \mathcal{N}$  which arises through the disjoint union of  $\mathcal{M}$  and  $\mathcal{N}$  and then consider the equivalence (also denoted  $\approx_{nrB}$ ) which identifies all states that are contained in some normed reward bisimulation for  $(\mathcal{M} \uplus \mathcal{N}, \mathcal{M} \uplus \mathcal{N})$ .

Our goal is to show that (i) normed reward bisimulation equivalent MDPs fulfil the same  $\text{PCTL}_r$  formulae and (ii) if (A1)-(A3), (A4') hold then  $\mathcal{M} \approx_{nrB} \hat{\mathcal{M}}$ . First, we establish the preservation property stated in (i). If  $\mathcal{M} \approx_{nrB} \mathcal{N}$  then there is a transformation “scheduler  $A$  for  $\mathcal{M}$   $\mapsto$  scheduler  $B$  for  $\mathcal{N}$ ” such that  $A$  and  $B$  are equivalent for  $\text{PCTL}_r$ -properties, i.e., they have the same expected reward for all  $\text{PCTL}_r$ -formulae  $\Phi$  and the same probabilities for the path formulae  $\Phi_1 U_I \Phi_2$ . The formal arguments are rather technical and very similar to the techniques worked out by Segala [31] for other (weak) bisimulation-relations. We obtain:

**Lemma 2.** *If  $\mathcal{M} \approx_{nrB} \mathcal{N}$  then  $\mathcal{M}$  and  $\mathcal{N}$  fulfil the same PCTL<sub>r</sub>-formulae.*

(See appendix for a rough proof sketch.) The goal is now to show that  $\mathcal{M} \approx_{nrB} \hat{\mathcal{M}}$  where,  $\mathcal{M} = (S, Act, \mathbf{P}, s_{init}, AP, L, rew)$  is the given “full” MDP and  $\hat{\mathcal{M}} = (\hat{S}, Act, \hat{\mathbf{P}}, s_{init}, AP, \hat{L}, rew)$  the sub-MDP that arises from  $\mathcal{M}$  through the definition of ample-sets. The following argument is similar to those in [13, 3], and uses the concept of *forming paths*. Let  $s, s' \in S$ . A forming path from  $s$  to  $s'$  is a finite path  $s = s_0 \xrightarrow{\beta_0} s_1 \xrightarrow{\beta_1} \dots \xrightarrow{\beta_{n-1}} s_n = s'$  where  $\beta_0, \dots, \beta_{n-1}$  are non-probabilistic stutter actions with reward 0, and, for  $0 \leq i < n$ , the singleton action-set  $\{\beta_i\}$  fulfils the dependence condition (A2) for state  $s_i$ . That is, for any finite path  $s_i \xrightarrow{\alpha_0} t_1 \xrightarrow{\alpha_1} \dots \xrightarrow{\alpha_{m-1}} t_m \xrightarrow{\gamma} \dots$  where  $\gamma$  is dependent on  $\beta_i$  there exists  $j \in \{0, 1, \dots, m-1\}$  such that  $\alpha_j = \beta_i$ . We write  $s \rightsquigarrow s'$  iff there exists a forming path from  $s$  to  $s'$  and  $\sqsubseteq$  for the lifting of  $\rightsquigarrow$  to distributions on  $S$  via weight functions, i.e.,  $\sqsubseteq$  agrees with  $\sqsubseteq_{\rightsquigarrow}$ . The length of a shortest forming path from  $s$  to  $\hat{s}$  in  $\mathcal{M}$  is denoted by  $|s \rightsquigarrow \hat{s}|$ . A forming path in  $\hat{\mathcal{M}}$  means a forming path as above where  $s_0, s_1, \dots, s_n \in \hat{S}$  and  $\beta_i \in ample(s_i)$ ,  $i = 0, 1, \dots, n-1$ . If (A1), (A2), (A3) and (A4') hold then forming paths enjoy the following properties:

- (i) If there is a forming path  $\sigma$  from  $s$  to  $s'$  where  $\alpha \in Act(s)$  does not occur then  $\alpha$  is independent on  $\sigma$ 's actions which yields  $\alpha \in Act(s')$  and  $\mathbf{P}(s, \alpha, \cdot) \sqsubseteq \mathbf{P}(s', \alpha, \cdot)$ .
- (ii) For all states  $\hat{s}$  in  $\hat{\mathcal{M}}$  there is a forming path from  $\hat{s}$  in  $\hat{\mathcal{M}}$  to some fully expanded state.

Let  $d(\hat{s})$  denote the length of a shortest forming path in  $\hat{\mathcal{M}}$  from  $\hat{s}$  to a fully expanded state in  $\hat{S}$ . By (ii),  $d(\hat{s}) < \infty$ .

**Lemma 3.** *The tuple  $(R, \eta_1, \eta_2, \eta_1^-, \eta_2^-)$ , where  $R = \{(s, \hat{s}) \in S \times \hat{S} : s \rightsquigarrow \hat{s}\}$ ,  $\eta_1(s, \hat{s}) = |s \rightsquigarrow \hat{s}|$ ,  $\eta_2(s, \hat{s}) = d(\hat{s})$ ,  $\eta_2^-(\hat{s}, s) = \eta_1(s, \hat{s}) = |s \rightsquigarrow \hat{s}|$  and  $\eta_1^-(\hat{s}, s)$  is arbitrary, is a normed reward bisimulation for  $(\mathcal{M}, \hat{\mathcal{M}})$ .*

*Proof.* Clearly,  $(s_{init}, s_{init}) \in R$  since we may consider a forming path of length 0. We show that for any pair  $(s, \hat{s}) \in R$  conditions (N1) and (N2) hold for  $(s, \hat{s}) \in R, \eta_1$  and  $\eta_2$  and for  $(\hat{s}, s) \in R^{-1}, \eta_1^-$  and  $\eta_2^-$ . The labelling condition (N1) is obvious as all actions on a forming path are stutter actions. Thus, all states on a forming path have the same labelling. Let  $(s, \hat{s}) \in R$  and  $\alpha \in Act(s)$ .

*Case 1:*  $\alpha$  does not occur on some forming path from  $s$  to  $\hat{s}$ . Then,  $\alpha \in Act(\hat{s})$  and  $\mathbf{P}(s, \alpha, \cdot) \sqsubseteq \mathbf{P}(\hat{s}, \alpha, \cdot)$  by (i). Hence, if  $\alpha \in ample(\hat{s})$  then case (N2.1) applies. If  $\alpha \notin ample(\hat{s})$  then we choose the first action  $\beta$  of a shortest forming path in  $\hat{\mathcal{M}}$  from  $\hat{s}$  to some fully expanded state. Then,  $(s, \hat{s}_\beta) \in R$  and  $\eta_2(s, \hat{s}_\beta) = d(\hat{s}_\beta) = d(\hat{s}) - 1 < \eta_2(s, \hat{s})$ . Hence, we are in case (N2.3).

*Case 2:*  $\alpha$  occurs in some shortest forming path from  $s$  to  $\hat{s}$ . Let  $\sigma$  be such a shortest forming path from  $s$  to  $\hat{s}$  and let  $\gamma_1 \dots \gamma_{i-1} \alpha \gamma_{i+1} \dots \gamma_n$  be the underlying action sequence. Then,  $\alpha$  is independent from  $\gamma_j$  for  $1 \leq j < i$ . Hence, there is a shortest forming path from  $s$  to  $\hat{s}$  that uses the action sequence  $\alpha \gamma_1 \dots \gamma_{i-1} \gamma_{i+1} \dots \gamma_n$ . This yields  $(s_\alpha, \hat{s}) \in R$ ,  $\eta_1(s_\alpha, \hat{s}) = \eta_1(s, \hat{s}) - 1$ . Hence, we are in case (N2.2).

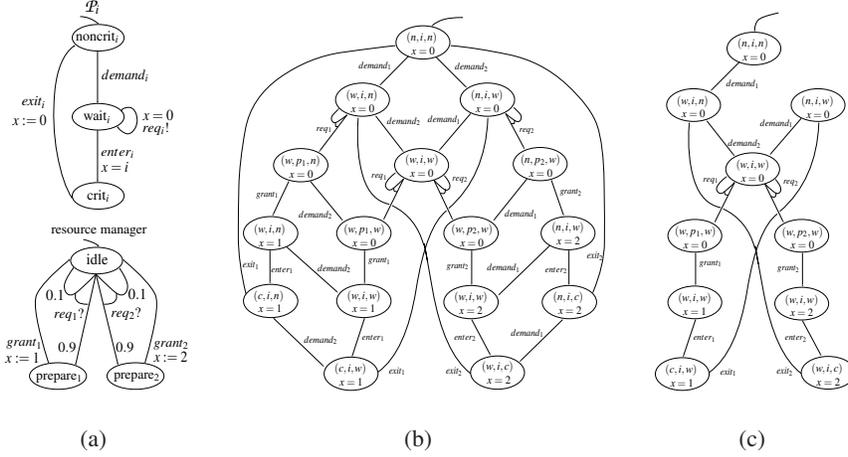


Fig. 2. Mutual exclusion example: (a) components, (b) full system and (c) reduced MDP

It remains to show that (N2) holds for  $(\hat{s}, s) \in R^{-1}$ ,  $\eta_1^-$  and  $\eta_2^-$ . Let  $\alpha \in \text{ample}(\hat{s})$ . If  $s = \hat{s}$  then case (N2.1) applies. If  $s \neq \hat{s}$  then we regard the first action  $\beta$  of a shortest forming path from  $s$  to  $\hat{s}$ . Then,  $(\hat{s}, s_\beta) \in R^{-1}$ ,  $\eta_2^-(\hat{s}, s_\beta) = \eta_1(s_\beta, \hat{s}) = \eta_1(s, \hat{s}) - 1 = \eta_2^-(\hat{s}, s) - 1$ . Thus, case (N2.3) applies.  $\square$

Lemma 2 and Lemma 3 complete the proof of Theorem 2.

*Example 2.* To illustrate our approach we consider a simple mutual exclusion protocol in which the processes  $P_1$  and  $P_2$  attempt to access a common resource controlled by a resource manager. A shared variable  $x$  is used to guarantee mutual exclusion and we assume that the communication is unreliable (requests to the resource manager are corrupted/lost with probability 0.1). Fig. 2(a) presents the different components of the system. Associating a reward of 1 with the actions  $req_1$  and  $req_2$  and 0 with all other actions, using PCTL<sub>r</sub> one can, for example, specify:

- $R_{\leq 1.4}(\text{crit}_1 \vee \text{crit}_2)$  : the expected number of requests before a process enters the critical section is at most 1.4;
- $P_{> 0.7}(\text{true } U_{[0,6]} \text{crit}_1 \vee \text{crit}_2)$ : the probability that a process enters its critical section after at most 6 requests have been issued is strictly greater than 0.7.

Fig. 2(b) gives the full MDP for the system and (assuming  $\text{AP} = \{\text{crit}_1, \text{crit}_2\}$ ), by applying (A1)-(A4') one can construct the reduced system given in Fig. 2(c).  $\square$

#### 4 A preservation result for (A1)-(A4) and reward-based properties

We now turn to the question of which properties with nontrivial reward bounds are preserved by (A1)-(A3) and the original branching condition (A4) in Fig. 1. Let us

again look at the path transformation described in (T1) and (T2) where, given a path  $\pi$  in  $\mathcal{M}$  a path  $\hat{\pi}$  is generated, where either the action sequence of  $\hat{\pi}$  is a permutation of the action sequence of  $\pi$  (T1) or  $\hat{\pi}$  starts with a non-probabilistic stutter action and then performs the same action sequence as the original path  $\pi$  (T2). As the rewards are in  $\mathbb{R}$  we do not know, how the cumulative reward of  $\hat{\pi}$  has changed compared to that of  $\pi$ . If we however require that the rewards of all actions are *non-negative*, along the modified path  $\hat{\pi}$  a reward equal or greater will be earned than that along  $\pi$ . This yields an informal explanation why the additional power of  $\mathcal{M}$  can lead to smaller minimal expected rewards, but the maximal expected rewards agree in  $\mathcal{M}$  and  $\hat{\mathcal{M}}$ . Similarly, we might expect that the minimal probabilities for events of the form  $a_1 U_{[0,r]} a_2$  agree under  $\mathcal{M}$  and  $\hat{\mathcal{M}}$ . The same holds for maximal probabilities for events of the form  $\square_{[0,r]} a$ . This motivates the definition of the following sublogic of  $\text{PCTL}_r$ .

Let  $\text{PCTL}_r^-$  be the sublogic of  $\text{PCTL}_r$  which only uses the  $\mathcal{R}$ -operator with upper reward bounds, i.e., formulae of the form  $\mathcal{R}_{[0,r]}(\Phi)$ , and where the probabilistic operator is only used in combination with PCTL-path formulae  $\Phi_1 U \Phi_2$  or with the until-operator in combination with upper reward and lower probability bounds or in combination with lower reward and upper probability bounds or with the always-operator in combination with upper reward and upper probability bounds or in combination with lower reward and lower probability bounds, e.g.  $\mathcal{P}_{[0,p]}(\square_{[0,r]} \Phi)$  or  $\mathcal{P}_{(p,1]}(\Phi_1 U_{[0,r]} \Phi_2)$ . Note that PCTL is contained in  $\text{PCTL}_r^-$ . (The result stated in Theorem 3 would still hold when dealing with a release- or weak until operator rather than the always-operator.)

**Theorem 3.** *If (A1)-(A4) hold and  $\text{rew}(\alpha) \geq 0$  for all  $\alpha \in \text{Act}$  then  $\mathcal{M}$  and  $\hat{\mathcal{M}}$  satisfy the same  $\text{PCTL}_r^-$  formulae.*

*Proof.* We choose the same proof technique as for Theorem 2 and deal with modified notions of normed reward (bi)simulation and forming paths that are obtained by simply ignoring the reward-constraint. That is, the formal definition of normed (bi)simulation that we will use now relies on conditions (N1) and (N2) as in Def. 1, but skipping the constraint “with reward 0” in (N2.2) and (N2.3). Let  $\approx_{nb}$  denote the resulting equivalence. Similarly, a forming path from  $s$  to  $\hat{s}$  means a path  $s = s_0 \xrightarrow{\beta_0} s_1 \xrightarrow{\beta_1} \dots \xrightarrow{\beta_{n-1}} s_n = \hat{s}$  where  $\beta_0, \dots, \beta_{n-1}$  are non-probabilistic stutter actions, and, for  $0 \leq i < n$ , the singleton action-set  $\{\beta_i\}$  fulfils the dependence condition (A2) for state  $s_i$ . We also modify the definition of shortest forming paths from  $s$  to  $\hat{s}$  by which we now mean a forming path from  $s$  to  $\hat{s}$  where the cumulative reward is minimal under all forming paths from  $s$  to  $\hat{s}$  and where the length (number of actions) is minimal under all forming paths with minimal cumulative reward. We will write  $\mu(s, \hat{s})$  for the cumulative reward of all/some shortest forming path from  $s$  to  $\hat{s}$ . As before,  $s \rightsquigarrow s'$  denotes the existence of a forming path from  $s$  to  $s'$  and we put  $R = \{(s, \hat{s}) \in S \times \hat{S} : s \rightsquigarrow \hat{s}\}$ .

Inspecting the above proof of Lemma 3 yields the following. Let  $(s, \hat{s}) \in R$  and  $\alpha \in \text{Act}(s)$ . If  $\alpha$  does not occur in some shortest forming path from  $s$  to  $\hat{s}$  then at least one of the following conditions holds:

- (i)  $\alpha \in \text{ample}(\hat{s}) \wedge \mathbf{P}(s, \alpha, \cdot) \sqsubseteq \mathbf{P}(\hat{s}, \alpha, \cdot)$  and there is a weight function  $w$  for  $(\mathbf{P}(s, \alpha, \cdot), \mathbf{P}(\hat{s}, \alpha, \cdot))$  with respect to  $R$  such that  $w(u, \hat{u}) > 0$  implies  $\mu(u, \hat{u}) \leq \mu(s, \hat{s})$ . This latter constraint holds since any forming path from  $s$  to  $\hat{s}$  can be “lifted” to forming paths connecting “corresponding”  $\alpha$ -successors.

- (ii) There exists a non-probabilistic stutter action  $\beta$  with  $(s, \hat{s}_\beta) \in R$  and  $\mu(s, \hat{s}_\beta) \leq \mu(s, \hat{s}) + \text{rew}(\beta)$ .

If  $\alpha$  occurs in some shortest forming path from  $s$  to  $\hat{s}$  then  $\alpha$  is the first action of some shortest forming path from  $s$  to  $\hat{s}$  and we have  $(s_\alpha, \hat{s}) \in R$  and  $\mu(s_\alpha, \hat{s}) = \mu(s, \hat{s}) - \text{rew}(\alpha)$ .

We now sketch a scheduler transformation “scheduler  $A$  for  $\mathcal{M} \mapsto$  scheduler  $B$  for  $\mathcal{N} = \hat{\mathcal{M}}$ ” such that  $\mathcal{M}$ ’s behaviour under  $A$  is mimicked by  $\mathcal{N}$ ’s behaviour under  $B$ , while  $B$  earns equal or more reward while simulating  $A$ ’s paths. For simplicity, we assume  $A$  to be deterministic. Let  $s$  be the current state of  $A$  and  $\hat{s}$  the current state of  $B$  where  $(s, \hat{s}) \in R$ . If  $A$  chooses an action  $\alpha$  for  $s$  that does not occur on a shortest forming path from  $s$  to  $\hat{s}$  then  $B$  selects  $\alpha$  for  $\hat{s}$  if case (i) applies and  $B$  selects  $\beta$  as in (ii) if (i) does not apply. If the action  $\alpha$  chosen by  $A$  does appear on a shortest forming path from  $s$  to  $\hat{s}$  then  $B$  waits until  $A$  reaches a state where it chooses an action that does not appear on a shortest forming path from  $s$  to  $\hat{s}$ . Note that in the latter step where  $B$  “waits”  $A$  performs a prefix  $\lambda_0$  of a shortest forming path from  $s$  to  $\hat{s}$  and reaches a state  $s_0$  with  $(s_0, \hat{s}) \in R$  and  $\mu(s_0, \hat{s}) = \mu(s, \hat{s}) - \text{rew}(\lambda_0)$ . Using the concept of weight functions on the path-level, we may continue in the same way to define a (randomized) scheduler  $B$  for  $\hat{\mathcal{M}}$  that mimics  $A$ ’s behaviour and enjoys the following property. If  $\sigma$  is a finite  $A$ -path starting in  $s$  then  $B$  has generated a finite path  $\hat{\sigma}$  starting in  $\hat{s}$  such that  $\rho(\sigma) \leq \rho(\hat{\sigma}) + \mu(s, \hat{s})$ . In fact, this also holds in a quantitative setting in the following sense. If  $(s, \hat{s}) \in R$  then

$$\Pr^{A,s}(\Pi(s, r + \mu(s, \hat{s}), C_1, \dots, C_n)) \geq \Pr^{B,\hat{s}}(\Pi(\hat{s}, r, C_1, \dots, C_n)) \quad (*)$$

and  $\Pr^{A,s}(\Pi(s, C_1, \dots, C_n)) = \Pr^{B,\hat{s}}(\Pi(\hat{s}, C_1, \dots, C_n))$ . Here, we used the following notation. Let  $u \in \mathcal{S}$ ,  $C_1, C_2, \dots, C_n$  be a sequence of  $\approx_{nb}$ -equivalence classes with  $C_i \neq C_{i+1}$  for  $1 \leq i < n$  and  $r \geq 0$ . Then,  $\Pi(u, r, C_1, \dots, C_n)$  denotes the set of all infinite paths that have a finite prefix of the form

$$u_0 \xrightarrow{*}_{C_1} \tilde{u}_1 \xrightarrow{\gamma_1} u_2 \xrightarrow{*}_{C_2} \tilde{u}_2 \xrightarrow{\gamma_2} \dots \xrightarrow{\gamma_{n-2}} u_{n-1} \xrightarrow{*}_{C_{n-1}} \tilde{u}_{n-1} \xrightarrow{\gamma_{n-1}} u_n$$

where  $u_0 = u$  and the total reward is  $\leq r$  and  $u_n \in C_n$ . The actions  $\gamma_i$  are arbitrary. In this context,  $v \xrightarrow{*}_C \tilde{v}$  means a finite path built out of non-probabilistic stutter actions such that  $v, \tilde{v}$  and all intermediate states of that path belong to  $C$ .  $\Pi(u, C_1, \dots, C_n)$  stands for the union of the path-sets  $\Pi(u, r, C_1, \dots, C_n)$  for arbitrary  $r \geq 0$ . For  $s = s_{\text{init}} = \hat{s}$  we have  $\mu(s, \hat{s}) = \mu(s_{\text{init}}, s_{\text{init}}) = 0$ .

The above yields that for each scheduler  $A$  for  $\mathcal{M}$  there exists a scheduler  $B$  for  $\hat{\mathcal{M}}$  such that  $\Pr^A(\Pi(s_{\text{init}}, r, C_1, \dots, C_n)) \geq \Pr^B(\Pi(s_{\text{init}}, r, C_1, \dots, C_n))$  for all  $r \geq 0$  and all  $\approx_{nb}$ -equivalence classes  $C_1, \dots, C_n$ . From this we may derive that:

- (a)  $\mathbb{E}^A(\diamond a) \leq \mathbb{E}^B(\diamond a)$ ,
- (b)  $\Pr^A\{\zeta \in \text{Paths}_\omega(s_{\text{init}}) : \zeta \models a_1 U_{[0,r]} a_2\} \geq \Pr^B\{\zeta \in \text{Paths}_\omega(s_{\text{init}}) : \zeta \models a_1 U_{[0,r]} a_2\}$ ,
- (c)  $\Pr^A\{\zeta \in \text{Paths}_\omega(s_{\text{init}}) : \zeta \models \square_{[0,r]} a\} \leq \Pr^B\{\zeta \in \text{Paths}_\omega(s_{\text{init}}) : \zeta \models \square_{[0,r]} a\}$ ,

where  $a, a_1, a_2 \in \text{AP}$ . Thus:

- (a)  $\hat{\mathcal{M}} \models \mathcal{R}_{[0,r]}(a)$  implies  $\mathcal{M} \models \mathcal{R}_{[0,r]}(a)$ ,

- (b)  $\hat{\mathcal{M}} \models \mathcal{P}_{[p,1]}(a_1 U_{[0,r]} a_2)$  implies  $\mathcal{M} \models \mathcal{P}_{[p,1]}(a_1 U_{[0,r]} a_2)$ ,
- (c)  $\hat{\mathcal{M}} \models \mathcal{P}_{[0,p]}(\Box_{[0,r]} a)$  implies  $\mathcal{M} \models \mathcal{P}_{[0,p]}(\Box_{[0,r]} a)$ ,

Similarly it can be shown that

- (d)  $\hat{\mathcal{M}} \models \mathcal{P}_{[0,p]}(a_1 U_{[r,\infty]} a_2)$  implies  $\mathcal{M} \models \mathcal{P}_{[0,p]}(a_1 U_{[r,\infty]} a_2)$  and
- (e)  $\hat{\mathcal{M}} \models \mathcal{P}_{[p,1]}(\Box_{[r,\infty]} a)$  implies  $\mathcal{M} \models \mathcal{P}_{[p,1]}(\Box_{[r,\infty]} a)$ .

From Lemma 1, we can then derive that  $\mathcal{M}$  and  $\hat{\mathcal{M}}$  fulfil the same  $\text{PCTL}_r^-$  formulae.  $\square$

*Example 3.* Let us return to Example 2 and redefine the rewards such that the only nonzero rewards are for actions  $demand_1$  and  $demand_2$  which have reward 1. Now, in this situation the reduced MDP in Fig. 2(c) can no longer be constructed using (A1)-(A4'). However, this construction is still possible under (A1)-(A4).

This is demonstrated by the fact that both the reduced and full MDP satisfy the  $\text{PCTL}_r^-$  property  $\mathcal{R}_{[0,2]}(\text{crit}_1 \vee \text{crit}_2)$  (the maximum expected number of processes that can attempt to enter the critical section before one of them does so is at most 2), while only the reduced model satisfies the  $\text{PCTL}_r$  property  $\mathcal{R}_{[2,\infty]}(\text{crit}_1 \vee \text{crit}_2)$  (the minimum expected number is at least 2).  $\square$

## 5 Reward properties w.r.t discounted rewards

In many research areas (e.g. economics, operations research, control theory) rewards are treated with a different semantics, namely as so-called *discounted rewards* [29, 6], where given a discount factor  $0 < c < 1$ , the reward of the  $i$ -th action of a path is multiplied with  $c^{i-1}$ . This interpretation of rewards reflects the fact that a reward (e.g. a payment) in the future is not worth quite as much as it is now (e.g. due to inflation). In this Section we investigate our partial order approach for discounted rewards.

Given a path  $\zeta = s_0 \xrightarrow{\alpha_1} s_1 \xrightarrow{\alpha_2} s_2 \xrightarrow{\alpha_3} \dots$  and a discount factor  $c \in (0, 1)$ , we denote by  $\rho_c(\zeta, i) = \text{rew}_c(\alpha_1 \dots \alpha_i) = c^0 \cdot \text{rew}(\alpha_1) + c^1 \cdot \text{rew}(\alpha_2) + \dots + c^{i-1} \cdot \text{rew}(\alpha_i)$  the cumulative discounted reward obtained through the first  $i$  actions.

With this on hand we can define the logic  $\text{PCTL}_c$ , which is a variant of  $\text{PCTL}_r$ . In  $\text{PCTL}_c$ , we use the new operators  $U_I^c$  and  $\mathcal{R}_I^c$  instead of  $U_I$  and  $\mathcal{R}_I$ , where instead of the cumulative reward  $\rho(\zeta, i)$  the cumulative discounted reward  $\rho_c(\zeta, i)$  is used in the semantics of those new operators. The semantics of the  $U_I^c$  operator is as follows. Given a path  $\zeta = s_0 \xrightarrow{\alpha_1} s_1 \xrightarrow{\alpha_2} s_2 \xrightarrow{\alpha_3} \dots$ , we say that  $\zeta \models \Phi_1 U_I^c \Phi_2$  iff  $\exists i \geq 0$  s.t.  $s_i \models \Phi_2 \wedge \forall j < i : s_j \models \Phi_1 \wedge \rho_c(\zeta, i) \in I$ . Similarly, given a set of states  $T \subseteq S$  we denote by  $\text{Rew}_c(\zeta, T)$  the discounted reward that is earned until a  $T$ -state is visited the first time. Formally, if  $\text{state}(\zeta, i) \notin T$  for all  $i \geq 0$  then  $\text{Rew}_c(\zeta, T) = \infty$ . If  $\text{state}(\zeta, i) \in T$  and  $\text{state}(\zeta, j) \notin T$  for all  $j < i$  then  $\text{Rew}_c(\zeta, T) = \rho_c(\zeta, i)$ . For  $T \subseteq S$  and a scheduler  $A$ ,  $\mathbb{E}_c^{A,s}(\diamond T)$  denotes the expected value under  $A$  with starting state  $s$  for the random function  $\zeta \mapsto \text{Rew}_c(\zeta, T)$ . Then  $s \models \mathcal{R}_I^c(\Phi)$  iff  $\forall$  schedulers  $A : \mathbb{E}_c^{A,s}(\diamond \text{Sat}(\Phi)) \in I$ .

A simple example shows that theorem 2 does not hold for  $\text{PCTL}_c$  (even if all rewards are nonnegative). Consider the MDP  $\mathcal{M}$  in example 1 on page 5. We assign the following rewards :  $\text{rew}(\alpha) = 0, \text{rew}(\beta) = \text{rew}(\gamma) = 1$ . Choosing  $\text{ample}(s) = \{\alpha\}$ , conditions

(A1)-(A3) and (A4') are satisfied. However, if we consider the formula  $\Phi = \mathcal{R}_{[0,c]}^c(a)$ , we gain that the reduced system  $\hat{\mathcal{M}}$  satisfies  $\Phi$  while the original system  $\mathcal{M}$  does not, because  $\mathcal{M}$  might choose action  $\beta$  in state  $s$  which yields the expected discounted reward to reach an  $a$ -state to be  $c^0 \cdot \text{rew}(\beta) = 1 > c$ .

The reader should notice that due to the discounting, the transformations (T1) and (T2) described in Section 3 on page 6 change the reward of a given path, even under condition (A4') which requires the ample set of a non-fully expanded state to be a singleton consisting of a non-probabilistic action with zero reward. Nevertheless, the following holds: given an MDP  $M$  with only *non-negative* rewards, ample-sets that satisfy (A1)-(A3) and (A4') and a path  $\zeta$  in  $\mathcal{M}$ , let  $\hat{\zeta}$  be a path that emanates from  $\zeta$  by applying transformation (T1) or (T2). Then  $\rho_c(\hat{\zeta}, i) \leq \rho_c(\zeta, i)$ . Similarly as in Section 4 this informally explains that the additional power of  $\hat{\mathcal{M}}$  can lead to greater maximal expected rewards, but the minimal expected rewards agree in  $\mathcal{M}$  and  $\hat{\mathcal{M}}$ . Also, the maximal probabilities for events of the form  $a_1 U_{[0,r]}^c a_2$  agree under  $\mathcal{M}$  and  $\hat{\mathcal{M}}$ . This motivates the definition of the following sublogic of  $PCTL_c$ .

Let  $PCTL_c^-$  be the sublogic of  $PCTL_c$  which uses the  $\mathcal{R}^c$  operator only with lower reward bounds (i.e.  $\mathcal{R}_{[r,\infty)}^c(\Phi)$ ) and where the probabilistic operator is only used in combination with PCTL-path formulae

- $\Phi_1 U \Phi_2$  or
- with the until-operator in combination with lower reward and lower probability bounds or in combination with upper reward and upper probability bounds or
- with the always-operator in combination with upper reward and lower probability bounds or in combination with lower reward and upper probability bounds,

e.g.  $\mathcal{P}_{[0,p]}(\Box_{[r,\infty)}\Phi)$  or  $\mathcal{P}_{[0,p]}(\Phi_1 U_{[0,r]}\Phi_2)$ . Note that PCTL is contained in  $PCTL_c^-$ .

**Theorem 4.** *If (A1)-(A3) and (A4') hold and  $\text{rew}(\alpha) \geq 0$  for all  $\alpha \in \text{Act}$  then  $\mathcal{M}$  and  $\hat{\mathcal{M}}$  satisfy the same  $PCTL_c^-$  formulae.*

The proof of Theorem 4 follows the same lines as the proof of Theorem 3, but instead of inequality (\*) one can establish the following inequality:

$$\Pr^{A,s}(\Pi(s, c^{\mu(s,\hat{s})} \cdot r, C_1, \dots, C_n)) \leq \Pr^{B,\hat{s}}(\Pi(\hat{s}, r, C_1, \dots, C_n))$$

## 6 Conclusion

The goal of this paper was to study the theoretical foundations of the ample-set approach for the logic  $PCTL_r$ , a variant of PCTL with reward-bounded temporal modalities and an expectation operator. The main results of this paper are that the ample-set conditions presented in [3] for PCTL preserve a class of non-trivial reward-based properties (Theorem 3) and that a slight modification of the conditions of [3] are sufficient to treat full  $PCTL_r$  (Theorem 2). The proofs of these results have been established by means of a new notion of weak bisimulation for rMDPs which preserves  $PCTL_r$  and – since it is simpler than other notions of weak bisimulation equivalence for MDPs – might also be useful for other purposes. Moreover we investigated the logic  $PCTL_c$ , a variant of

$PCTL_r$ , where the rewards are given a discounting semantics. We presented ample-set conditions that preserve a non-trivial subset of  $PCTL_c$  properties if all given rewards are non-negative (Theorem 4).

We concentrated here on the probabilities and expectations of cumulative rewards. However, we claim that the criteria (A1)-(A4) are also sufficient to treat long run average properties formalized by P-experiments [10, 12].

Besides being of theoretical interest, the results of this paper also have a practical impact. First experimental results on the ample set approach for MDPs (without reward structure) with the forthcoming model checker LiQuor [2] show that although the criteria needed for probabilistic systems are stronger than in the non-probabilistic case, good reductions (up to 80%) can be obtained. Furthermore, the bottleneck in analysis of probabilistic systems modelled by MDPs are the required techniques for solving linear programs. Since the amount of time required for the construction of the reduced MDP is negligible compared to the running time of linear program solvers, even small reductions can increase the efficiency of the quantitative analysis.

In future work, we plan to integrate the partial order reduction techniques suggested here in the symbolic MTBDD-based model checker PRISM [19] by constructing a syntactic representation of the reduced MDP at compile time, in the style of static partial order reduction [23] which permits a combination of partial order reduction with symbolic BDD-based model checking.

## References

1. S. Andova, H. Hermanns, and J.-P. Katoen. Discrete-time rewards model-checked. In *Proc. FORMATS*, volume 2791 of *Lecture Notes in Computer Science*, pages 88–104, 2003.
2. C. Baier, F. Ciesinski, and M. Groesser. Quantitative analysis of distributed randomized protocols. In *Proc. of the tenth International Workshop on Formal Methods for Industrial Critical Systems (FMICS 05)*, 2005.
3. C. Baier, P. D’Argenio, and M. Größer. Partial order reduction for probabilistic branching time. In *Proc. QAPL*, 2005.
4. C. Baier, M. Größer, and F. Ciesinski. Partial order reduction for probabilistic systems. In *QEST 2004* [30], pages 230–239.
5. Christel Baier and Marielle Stoelinga. Norm functions for probabilistic bisimulations with delays. In *Proc. FOSSACS 2000*, volume 1784 of *LNCS*, pages 1–16, 2000.
6. Dimitri P. Bertsekas. *Dynamic Programming and Optimal Control*. Athena Scientific, Belmont, Massachusetts, 1995.
7. A. Bianco and L. de Alfaro. Model checking of probabilistic and nondeterministic systems. In *Proc. Foundations of Software Technology and Theoretical Computer Science (FST & TCS)*, volume 1026 of *Lecture Notes in Computer Science*, pages 499–513, 1995.
8. E. Clarke, E. Emerson, and A. Sistla. Automatic verification of finite-state concurrent systems using temporal logic specifications. *ACM Transactions on Programming Languages and Systems*, 8(2):244–263, April 1986.
9. P.R. D’Argenio and P. Niebert. Partial order reduction on concurrent probabilistic programs. In *QEST 2004* [30], pages 240–249.
10. L. de Alfaro. *Formal Verification of Probabilistic Systems*. PhD thesis, Stanford University, Department of Computer Science, 1997.
11. L. de Alfaro. Temporal logics for the specification of performance and reliability. In *Proc. STACS*, volume 1200 of *Lecture Notes in Computer Science*, pages 165–179, 1997.

12. L. de Alfaro. How to specify and verify the long-run average behavior of probabilistic systems. In *Proc. 13th Annual IEEE Symposium on Logic in Computer Science (LICS)*, IEEE Press, pages 454–465, 1998.
13. R. Gerth, R. Kuiper, D. Peled, and W. Penczek. A partial order approach to branching time logic model checking. In *Proc. 3rd Israel Symposium on the Theory of Computing Systems (ISTCS'95)*, pages 130–139. IEEE Press, 1995.
14. P. Godefroid. *Partial Order Methods for the Verification of Concurrent Systems: An Approach to the State Explosion Problem*, volume 1032 of *Lecture Notes in Computer Science*. Springer-Verlag, 1996.
15. P. Godefroid, D. Peled, and M. Staskauskas. Using partial-order methods in the formal validation of industrial concurrent programs. In *Proc. International Symposium on Software Testing and Analysis*, pages 261–269. ACM Press, 1996.
16. D. Griffioen and F. Vaandrager. Normed simulations. In *Proc. 10th International Computer Aided Verification Conference*, volume 1427 of *LNCSS*, pages 332–344, 1998.
17. H. Hansson. *Time and Probability in Formal Design of Distributed Systems*. Series in Real-Time Safety Critical Systems. Elsevier, 1994.
18. H. Hansson and B. Jonsson. A logic for reasoning about time and reliability. *Formal Aspects of Computing*, 6(5):512–535, 1994.
19. A. Hinton, M. Kwiatkowska, G. Norman, and D. Parker. PRISM: A tool for automatic verification of probabilistic systems. In *Proc. 12th International Conference on Tools and Algorithms for the Construction and Analysis of Systems (TACAS'06)*, 2006. To appear.
20. G. Holzmann. *The SPIN Model Checker, Primer and Reference Manual*. Addison Wesley, 2003.
21. C. Jones. *Probabilistic Non-Determinism*. PhD thesis, University of Edinburgh, 1990.
22. B. Jonsson and K. Larsen. Specification and refinement of probabilistic processes. In *Proc. LICS*, pages 266–277. IEEE CS Press, 1991.
23. R. Kurshan, V. Levin, M. Minea, D. Peled, and H. Yenig. Static partial order reduction. In *Proc. Tools and Algorithms for Construction and Analysis of Systems (TACAS)*, volume 1384 of *Lecture Notes in Computer Science*, pages 345–357, 1998.
24. K. Namjoshi. A simple characterization of stuttering bisimulation. In *Proc. FSTTCS*, volume 1346 of *Lecture Notes in Computer Science*, pages 284–296, 1997.
25. N. Pekergin and Sana Younes. Stochastic model checking with stochastic comparison. In *Proc. EPEWWS-FM*, volume 3670 of *Lecture Notes in Computer Science*, pages 109–123, 2005.
26. D. Peled. All from one, one for all: On model checking using representatives. In *Proc. 5th International Computer Aided Verification Conference (CAV)*, volume 697 of *Lecture Notes in Computer Science*, pages 409–423, 1993.
27. D. Peled. Partial order reduction: Linear and branching time logics and process algebras. In [28], pages 79–88, 1996.
28. D. Peled, V. Pratt, and G. Holzmann, editors. *Partial Order Methods in Verification*, volume 29(10) of *DIMACS Series in Discrete Mathematics and Theoretical Computer Science*. American Mathematical Society, 1997.
29. M. L. Puterman. *Markov Decision Processes—Discrete Stochastic Dynamic Programming*. John Wiley & Sons, Inc., New York, 1994.
30. *Proceedings of the 1st International Conference on Quantitative Evaluation of Systems (QEST 2004)*. Enschede, the Netherlands. IEEE Computer Society Press, 2004.
31. R. Segala. *Modeling and Verification of Randomized Distributed Real-Time Systems*. PhD thesis, Massachusetts Institute of Technology, 1995.
32. R. Segala and N. Lynch. Probabilistic simulations for probabilistic processes. *Nordic Journal of Computing*, 2(2):250–273, 1995.
33. A. Valmari. Stubborn set methods for process algebras. In [28], pages 79–88, 1996.

## A Appendix: Proof sketch of Lemma 2

Our goal is to show that normed reward bisimulation equivalent MDPs fulfil the same  $\text{PCTL}_r$  formulae. For this, we have to show that if  $\mathcal{M} \approx_{nr} \mathcal{N}$  then there is a transformation

“scheduler  $A$  for  $\mathcal{M} \mapsto$  scheduler  $B$  for  $\mathcal{N}$ ”

such that  $A$  and  $B$  are equivalent for  $\text{PCTL}_r$ -properties, i.e., they have the same expected reward for all  $\text{PCTL}_r$ -formulae  $\Phi$  and the same probabilities for the path formulae  $\Phi_1 U_I \Phi_2$ . We assume here starting states  $s$  and  $s'$  for  $A$  and  $B$ , respectively, such that  $(s, s') \in R$  for some normed reward bisimulation for  $(\mathcal{M}, \mathcal{N})$ . In particular, this applies to the initial states  $s = s_{\text{init}}^{\mathcal{M}}$  and  $s' = s_{\text{init}}^{\mathcal{N}}$ . By symmetry, we obtain the equivalence of  $\mathcal{M}$  and  $\mathcal{N}$  for  $\text{PCTL}_r$ .

Let  $(R, \eta_1, \eta_2, \eta_1^-, \eta_2^-)$  be a normed reward bisimulation for  $(\mathcal{M}, \mathcal{N})$ . For simplicity, let us assume that  $A$  is deterministic. If we are given a pair  $(s, s') \in R$  such that  $A$  chooses action  $\alpha$  for  $s$  then the corresponding scheduler  $B$  chooses action  $\alpha$  for  $s'$  in case (N2.1). In case (N2.2)  $B$  “waits” until  $A$  has generated a finite path  $s = s_0 \xrightarrow{\alpha} s_1 \xrightarrow{\alpha_1} \dots \xrightarrow{\alpha_{n-1}} s_n$  consisting of non-probabilistic stutter actions with reward 0 where case (N2.2) applies to the pairs  $(s_i, s')$  and the chosen action  $\alpha_i$ ,  $i = 1, \dots, n-1$ , while for  $(s_n, s')$  and the chosen action  $\alpha_n$  we are in case (N2.1) or (N2.3). (Note that in (N2.2),  $\eta_1$  is strictly decreasing. Hence, it is impossible to generate an infinite  $A$ -path by (N2.2) only.) In the former case,  $B$  selects action  $\alpha_n$  for  $s'$ . If case (N2.3) applies to  $(s_n, s')$  and  $\alpha_n$  then  $B$  generates a finite path  $s' = s'_0 \xrightarrow{\beta_0} s'_1 \xrightarrow{\beta_1} \dots \xrightarrow{\beta_{m-1}} s'_m$  consisting of non-probabilistic stutter actions with reward 0 where (N2.3) holds for the pairs  $(s_n, s'_j)$  and the chosen action  $\alpha_n$  for  $s_n$ ,  $j = 1, \dots, m-1$ , while for  $(s_n, s'_m)$  and the chosen action  $\alpha_n$  we are in case (N2.1) or (N2.3). We may repeat this argument by alternating between the cases (N2.2) and (N2.3). There are two possibilities:

*Case 1:* Eventually case (N2.1) applies. Thus, the above technique yields a finite  $A$ -path  $\sigma$  with last state  $t$  and a  $B$ -path ending in state  $t'$ , each of them consisting of non-probabilistic stutter actions with reward 0 such that  $(t, t') \in R$  and case (N2.1) applies for the chosen action  $\gamma$  under  $A$ . Hence,  $B$  can make the same choice and performs action  $\gamma$  in  $t'$ . The weight function condition in (N2.1) guarantees that  $\mathbf{P}_{\mathcal{M}}(t, \gamma, \cdot) \sqsubseteq_R \mathbf{P}_{\mathcal{N}}(t', \gamma, \cdot)$ .

*Case 2:* (N2.1) never applies. Thus alternating between the cases (N2.2) and (N2.3) forever, each scheduler  $A$  and  $B$  creates a single path consisting of non-probabilistic stutter actions with zero reward.

We then may continue this technique of simulating  $\mathcal{M}$ 's behaviour under  $A$  in  $\mathcal{N}$  with a scheduler  $B$  using the concept of weight function on the level of paths. The formal arguments are rather technical and very similar to the techniques worked out by Segala [31] for other (weak) bisimulation-relations. We skip the technical details of this scheduler-transformation, but observe that the given scheduler  $A$  for  $\mathcal{M}$  and the derived scheduler  $B$  for  $\mathcal{N}$  generate with the same probability infinite paths of the form

$$u_0 \xrightarrow{*}_{C_1} \tilde{u}_1 \rightarrow u_2 \xrightarrow{*}_{C_2} \tilde{u}_2 \rightarrow u_3 \xrightarrow{*}_{C_3} \tilde{u}_3 \rightarrow \dots$$

where  $C_1, C_2, \dots$  are equivalence classes under  $\approx_{nr b}$  (note that  $C_i = C_{i+1}$  is possible). For  $C$  to be an  $\approx_{nr b}$ -equivalence class  $u \rightarrow_C^* \tilde{u}$  denotes a finite path built out of non-probabilistic stutter actions with reward 0 such that  $u, \tilde{u}$  and all intermediate states of that path belong to  $C$ . Independent of their length, the cumulative reward of these path fragments is 0.  $u_0$  denotes the starting state under consideration ( $u_0 = s$  for  $\mathcal{M}$  and  $u_0 = s'$  for  $\mathcal{N}$ ).

By induction on length of PCTL<sub>r</sub>-formulae, it can now be shown that normed reward bisimulation equivalent MDPs fulfil the same PCTL<sub>r</sub>-formulae.