



AVACS – Automatic Verification and Analysis of Complex
Systems

REPORTS

of SFB/TR 14 AVACS

Editors: Board of SFB/TR 14 AVACS

Optimal Schedulers for Time-Bounded Reachability in CTMDPs

by
Markus Rabe Sven Schewe

Publisher: Sonderforschungsbereich/Transregio 14 AVACS
(Automatic Verification and Analysis of Complex Systems)
Editors: Bernd Becker, Werner Damm, Martin Fränzle, Ernst-Rüdiger Olderog,
Andreas Podelski, Reinhard Wilhelm
ATRs (AVACS Technical Reports) are freely downloadable from www.avacs.org

Copyright © October 2009 by the author(s)
Author(s) contact: Markus Rabe (rabe@cs.uni-saarland.de).

Optimal Schedulers for Time-Bounded Reachability in CTMDPs

Markus Rabe¹ and Sven Schewe²

¹ Saarland University*
rabe@cs.uni-saarland.de

² University of Liverpool
sven.schewe@liverpool.ac.uk

Abstract. We study time-bounded reachability in continuous-time Markov decision processes for various scheduler classes. Such reachability problems play a paramount rôle in dependability analysis and the modelling of manufacturing and queueing systems. Consequently, their analysis has been studied intensively, and techniques for the approximation of optimal control are well understood. From a mathematical point of view, however, the question of approximation is secondary compared to the fundamental question whether or not optimal control exists.

We demonstrate the existence of optimal schedulers for all commonly considered scheduler classes. For scheduler classes without full access to time we provide an effective technique to determine simple optimal schedulers that converge to an easy-to-compute memoryless scheduling policy after a finite number of steps.

1 Introduction

Markov decision processes (MDPs) are a framework that incorporates both nondeterministic and probabilistic choices. They are used in a variety of applications such as the control of manufacturing processes [11, 5] or queueing systems [13]. We study a real time version of MDPs, continuous-time Markov decision processes (CTMDPs), which are a natural formalism for modelling in scheduling [4, 11] and stochastic control theory [5]. CTMDPs can also be seen as a unified framework for different stochastic model types used in dependability analysis [12, 11, 8, 6, 9].

The analysis of CTMDPs usually concerns the different possibilities to resolve the nondeterminism by means of a scheduler (also called strategy). Typical questions cover qualitative as well as quantitative properties, such as: “Can the nondeterminism be resolved by a scheduler such that a predefined property holds?” or respectively “Which scheduler optimises a given objective function?”.

In this paper, we study the *maximal time-bounded reachability problem* [11, 3, 15, 9, 10] for CTMDPs. Time-bounded reachability is the standard control problem to construct a scheduler that controls the Markov decision process such that the likelihood of reaching a goal region within a given time bound is maximised, and to determine the probability. For CTMDPs, the answer to both questions naturally depends on the power

* This project was supported by the Saarbrücken Graduate School of Computer Science funded by the Initiative for Excellence of the German federal and state governments

a scheduler has to observe the run of the system—in particular if it can observe time—and on its ability to store and process this information. For the common classes of schedulers, research has focused on efficient approximation techniques [3, 9, 10], while the existence of optimal schedulers has remained open.

Overview. Given its practical importance, the bounded reachability problem for Markov decision processes has been intensively studied [2, 3, 15, 9, 10]. While previous research focused on *approximating* of optimal scheduling policies [3, 10], we prove the *existence* of optimal schedulers in this paper.

Unlike for discrete time MDPs, the power of observing the elapsing of time is an important aspect in the design of schedulers for CTMDPs, and various classes of schedulers that differ in their observational power have been discussed in the literature [9, 3]. Intuitively, the differences in these classes concern the ability to store information, and to measure time.

Figure 1 shows a comparison between the commonly considered scheduler classes, where schedulers that can store the history, its length, or nothing at all are marked H (for history-dependent), C (for hop-counting), and P (for positional), respectively. Schedulers that can observe time are marked with a T (timed), and with TT (total time) if they have the power to revoke their decision.

Revoking decisions is a concept first discussed in [9] that extends schedulers on a different level than their observational power: while traditional scheduler classes require the schedulers to fix their decisions as soon as they enter a location, TT schedulers may change their decision over time, even while residing in a location.

The arrows denote inclusions between classes, which are direct implications of their definitions. The classes depicted in Figure 1 are ordered top down by their maximal reachability probabilities as known from the literature [3, 9].

In principle, approximating optimal schedulers is simple for all scheduler classes. For schedulers that can observe time, it suffices to discretise time and to increase the sample rate [10], and for time-abstract schedulers, it suffices to optimise the reachability within a bounded number of steps and to let this bound grow to infinity [3].

Efficient techniques to determine these rates have, for example, been discussed for uniform CTMDPs—CTMDPs with a constant transition rate—by Baier, Hermanns, Katon, and Haverkort [3].

Contribution. This paper has contributions on two levels: The clean result on the technical level is a proof that optimal schedulers exist for all commonly considered scheduler classes, but we deem the simple insights on the conceptual level that led to these results to be of similar importance.

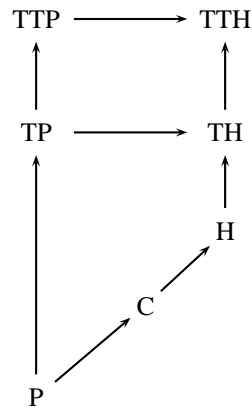


Fig. 1. Scheduler hierarchy

Markov + Time = Markov. Markov processes are mathematical models for the random evolution of *memoryless* systems, that is, systems for which the likelihood of future events, at any given moment, depends only on their present state, and not on the past. We observe that *continuous-time Markov chains and decision processes remain Markovian if we add the time that has passed to the state space.* We use this observation in Section 4 to introduce *time-extended* CTMDPs, which contain the time that has passed as part of their state space. This approach has an immediate implication for all time-dependent scheduler classes: It implies without further ado that the scheduler classes TP and TH as well as the classes TTP and TTH coincide, because optimal scheduler decisions in a Markovian system (with simple objectives like time-bounded reachability) cannot depend on the history. As a result, the description of optimal time-dependent schedulers in Section 4 is simple.

Reasoning about time-abstract scheduler classes is slightly more involved, because time-abstract schedulers do not have access to the precise time that remains for reaching the goal region. Phrased in terms of time-extended CTMDPs, these schedulers do not know precisely in which state of the time-extended CTMDP they are, but they can infer a distribution over the states in which they could potentially be. While this argument is not used explicitly in Section 3, it was the driving factor in our research that led to the construction of optimal time-abstract schedulers. It also provides quick and intuitive alternative proofs for the traditional result [3] that counting and history-dependent schedulers provide the same time-bounded reachability probability for uniform CTMDPs, but different ones for non-uniform CTMDPs: while the distribution over the states of the time-extended CTMDP coincides in the first case, it differs in the latter.

Optimal Schedulers. Our technical contributions are simple algorithms for the construction of optimal time-abstract (Section 3) and time-dependent (Section 4) schedulers.

The algorithmic solution for time-abstract schedulers builds on the observation that, if time has almost run out, we can use a *greedy strategy* that optimises our chances to reach our goal in a single step. Reaching it in more steps is then used as a tie-break criterion with decreasing power for increasing distance. We show that such a memoryless greedy scheduler exists and is indeed optimal after a certain step bound.

As a small side-result, we also extended the result that allowing for randomisation does not increase the time-bounded reachability probability to all scheduler classes.

2 Continuous-Time Markov Decision Processes

A *continuous-time Markov decision process* \mathcal{M} is a tuple $(L, Act, \mathbf{R}, \nu, B)$ with a finite set of locations L , a finite set of actions Act , a rate matrix $\mathbf{R} : (L \times Act \times L) \rightarrow \mathbb{Q}_{\geq 0}$, an initial distribution $\nu \in Dist(L)$, and a goal region $B \subseteq L$. We define the total exit rate for a location l and an action a as $\mathbf{R}(l, a, L) = \sum_{l' \in L} \mathbf{R}(l, a, l')$. For a CTMDP we require that, for all locations $l \in L$, there must be an action $a \in Act$ such that $\mathbf{R}(l, a, L) > 0$, and we call such actions *enabled*. We define $Act(l)$ to be the set of enabled actions in location l . If there is only one enabled action per location, a CTMDP \mathcal{M} is a continuous-time Markov chain [7]. If multiple actions are available, we need to resolve the nondeterminism by means of a scheduler (also called strategy or scheduling policy). As usual, we

assume the goal region to be absorbing, and we use $\mathbf{P}(l, a, l') = \frac{\mathbf{R}(l, a, l')}{\mathbf{R}(l, a, L)}$ to denote the time-abstract transition probability.

Note, that we explicitly distinguish between *locations* and *states*. We consider a state to be a location at a certain point of time. This notion will prove to be helpful when considering time-dependent schedulers in Section 4.

Uniform CTMDPs. We call a CTMDP uniform with rate λ if, for every location l and action $a \in \text{Act}(l)$, the total exit rate $\mathbf{R}(l, a, L)$ is λ . In this case the probability $p_{\lambda}(n)$ that there are exactly n discrete events (transitions) in time t is Poisson distributed: $p_{\lambda}(n) = e^{-\lambda t} \cdot \frac{(\lambda t)^n}{n!}$.

We define the *uniformisation* \mathcal{U} of a CTMDP \mathcal{M} as the uniform CTMDP obtained by creating a copy $l_{\mathcal{U}}$ for every location l . We call the new copies unobservable, and all locations $l \in L \subset L_{\mathcal{U}}$ observable. Let λ be the maximal total exit rate in \mathcal{M} . The new rate matrix $\mathbf{R}_{\mathcal{U}}$ extends \mathbf{R} by first adding the rate $\mathbf{R}_{\mathcal{U}}(l, a, l_{\mathcal{U}}) = \lambda - \mathbf{R}(l, a, L)$ for every location $l \in L$ and action $a \in \text{Act}$ of \mathcal{M} , and by then copying the outgoing transitions from every observable location l to its unobservable counterpart $l_{\mathcal{U}}$, while the other components remain untouched. The intuition behind this uniformisation technique is that it enables us to distinguish whether a step would have occurred in the original automaton or not.

Paths. A *timed path* π in CTMDP \mathcal{M} is a finite sequence in $(L \times \text{Act} \times \mathbb{R}_{\geq 0})^* \times L = \text{Paths}(\mathcal{M})$. We write

$$l_0 \xrightarrow{a_0, t_0} l_1 \xrightarrow{a_1, t_1} \dots \xrightarrow{a_{n-1}, t_{n-1}} l_n$$

for a sequence π , and we require $t_{i-1} < t_i$ for all $i < n$. The t_i denote the system's time when the events happen. The corresponding *time-abstract path* is defined as $l_0 \xrightarrow{a_0} l_1 \xrightarrow{a_1} \dots \xrightarrow{a_{n-1}} l_n$. We use $\text{Paths}_{\text{abs}}(\mathcal{M})$ to denote the set of all such projections and $|\cdot|$ to count the number of actions in a path. Concatenation of paths π, π' will be written as $\pi \circ \pi'$ if the last state of π is the first state of π' .

Schedulers. The system's behaviour is not completely defined by the CTMDP, but also by a scheduler that resolves the nondeterminism. When analysing properties of a CTMDP, such as the reachability probability, we usually quantify over a class of schedulers. We restrict all scheduler classes to those schedulers creating a measurable probability space (cf. [14]), and we consider the following common classes, which differ in their power to observe events and to revoke their decisions:

- *Total time history-dependent* (TTH) schedulers $\text{Paths}(\mathcal{M}) \times \mathbb{R}_{\geq 0} \rightarrow D$
that map timed paths and the elapsed time to decisions.
- *Total time positional* (TTP) schedulers $L \times \mathbb{R}_{\geq 0} \rightarrow D$
that map locations and the elapsed time to decisions.
- *Timed history* (TH) schedulers $\text{Paths}(\mathcal{M}) \rightarrow D$
that map timed paths to decisions.
- *Timed positional* (TP) schedulers $L \times \mathbb{R}_{\geq 0} \rightarrow D$
that map locations and the time until the last state change to decisions.

- *Time-abstract history-dependent* (H) schedulers $Paths_{abs}(\mathcal{M}) \rightarrow D$
that map time-abstract paths to decisions.
- *Time-abstract hop-counting* (C) schedulers $L \times \mathbb{N} \rightarrow D$
that map locations and the number of hops (length of the path) to decisions.
- *Positional* (P) or memoryless schedulers $L \rightarrow D$
that map locations to decisions.

Decisions D are either randomised (R), in which case $D = Dist(Act)$ is the set of distributions over enabled actions, or are restricted to deterministic (D) choices, that is $D = Act$. Where it is necessary to distinguish randomised and deterministic versions we will add a postfix to the scheduler class, for example HD and HR.

Induced Probability Space. We build our probability space in the natural way: we first define the probability measure for cylindric sets of paths that start with

$$l_0 \xrightarrow{a_0, t_0} l_1 \xrightarrow{a_1, t_1} \dots \xrightarrow{a_{n-1}, t_{n-1}} l_n,$$

with $t_j \in I_j$ for all $j < n$, and for non-overlapping open intervals I_0, I_1, \dots, I_{n-1} , to be the usual probability that a path starts with these actions for a randomised scheduler \mathcal{S} that may not revoke its decisions, and such that $\mathcal{S}(l_0 \xrightarrow{a_0, t_0} \dots \xrightarrow{a_{i-1}, t_{i-1}} l_i)$ is equivalent for all $(t_0, \dots, t_{i-1}) \in I_0 \times \dots \times I_{i-1}$:

$$\int_{t_0 \in I_0, t_1 \in I_1, \dots, t_{n-1} \in I_{n-1}} \prod_{i=0}^{n-1} \mathcal{S}(l_0 \xrightarrow{a_0, t_0} \dots \xrightarrow{a_{i-1}, t_{i-1}} l_i)(a_i) \cdot \mathbf{R}(l_i, a_i, l_{i+1}) \cdot e^{-\mathbf{R}(l_i, a_i, L)(t_i - t_{i-1})},$$

assuming $t_{-1} = 0$.

From this basic building block, we build our probability measure for measurable sets of paths and measurable schedulers in the usual way (cf. [14]). The similar space for TT schedulers, which may revoke their decisions, is described in Appendix B.

Time-Bounded Reachability Probability. For a given CTMDP $\mathcal{M} = (L, Act, \mathbf{R}, \nu, B)$ and a given measurable scheduler \mathcal{S} that resolves the nondeterminism, we use the following notations for the probabilities:

- $Pr_{\mathcal{S}}^{\mathcal{M}}(l, t)$ is the probability of reaching the goal region B in time t when starting in location l ,
- $Pr_{\mathcal{S}}^{\mathcal{M}}(t) = \sum_{l \in L} \nu(l) Pr_{\mathcal{S}}^{\mathcal{M}}(l, t)$ denotes the probability of reaching the goal region B in time t ,
- $Pr_{\mathcal{S}}^{\mathcal{M}}(t; k)$ denotes the probability of reaching the goal region B in time t and in at most k discrete steps, and
- $PR_{\mathcal{S}}^{\mathcal{M}}(\pi, t)$ is the probability to traverse the time-abstract path π within time t .

As usual, the supremum of the time-bounded reachability probability over a particular scheduler class is called the time-bounded reachability of \mathcal{M} for this scheduler class, and we use ‘max’ instead of ‘sup’ to indicate that this value is taken for some *optimal scheduler* \mathcal{S} of this class.

Step Probability Vector. Given a scheduler \mathcal{S} and a location l for a CTMDP \mathcal{M} , we define the *step probability vector* $d_{l,\mathcal{S}}$ of infinite dimension. An entry $d_{l,\mathcal{S}}[i]$ for $i \geq 0$ denotes the probability to reach goal region B in up to i steps from location l (not considering any time constraints).

3 Time-Abstract Schedulers

In this section, we show that *optimal* schedulers exist for all natural time-abstract classes, that is, for CD, CR, HD, and HR. Moreover, we show that there are optimal schedulers that become positional after a small number of steps, which we compute with a simple algorithm. We also show that randomisation does not yield any advantage: deterministic schedulers are as good as randomised ones. Our proofs are constructive, and thus allow for the construction of optimal schedulers. This also provides the first procedure to precisely determine the time-bounded reachability probability, because we can now reduce this problem to solving the time-bounded reachability problem of continuous-time Markov chains [2].

Our proof consists of two parts. We first consider the class of uniform CTMDPs, which are much simpler to treat in the time-abstract case, because we can use Poisson distributions to describe the number of steps taken within a given time bound. For uniform CTMDPs it is already known that the supremum over the bounded reachability collapses for all time-abstract scheduler classes from CD to HR [3]. It therefore suffices to show that there is a CD scheduler which takes this value.

We then show that a similar claim holds for CD and HD scheduler in the general class of not necessarily uniform CTMDPs. In this case, it also holds that there are simple optimal schedulers that converge against a positional scheduler after a finite number of steps, and that randomisation does not improve the time-bounded reachability probability. However, in the non-uniform case the time-abstract path contains more information about the remaining time than its length only, and bounded reachability of history-dependent and counting schedulers usually deviate (see [3] for a simple example).

We start this section with the introduction of *greedy schedulers*, HD schedulers that favour reachability in a small number of steps over reachability with a larger number of steps; the positional schedulers against which the CD and HD schedulers converge are such greedy schedulers.

3.1 Greedy Schedulers

The natural objective when seeking optimal schedulers is to maximise time-bounded reachability $Pr_S^{\mathcal{M}}(l, t)$ for every location l with respect to a particular scheduler class such as HD. Unfortunately, this optimisation problem is comparably complex.

However, when the remaining time t is close to 0, then increasing the likelihood of reaching the goal region in few steps dominates the impact of reaching it later. While we have no direct access to the remaining time in the time-abstract case, we can infer the distribution over the remaining time from the time-abstract history (or its length).

Since the expected remaining time converges to 0 when the number of transitions goes to infinity, we can argue in a way similar to the time-dependent case.

This motivates the introduction of greedy schedulers. Schedulers are called greedy, if they (greedily) look for short-term gain, and favour it over any long-term effect. Greedy schedulers that optimise the reachability within the first k steps have been exploited in the efficient analysis of CTMDPs [3]. We call such schedulers *k-optimal* (because they are the optimal schedulers for $k = \infty$).

To understand the principles of optimal control, a simpler form greediness proves to be more appropriate: We call an HD scheduler *greedy* if it maximises the step probability vector of every location l with respect to the lexicographic order (for example $(0, 0.2, 0.3, \dots) >_{lex} (0, 0.1, 0.4, \dots)$). To prove the existence of greedy schedulers, we draw from the fact that the supremum $d_l = \sup_{S \in HD} d_{l,S}$ obviously exists, where the supremum is to be read as a supremum with respect to the lexicographic order. An action $a \in Act(l)$ is called *greedy* for a location $l \notin B$ if it satisfies $shift(d_l) = \sum_{l' \in L} \mathbf{P}(l, a, l') d_{l'}$, where $shift(d_l)$ shifts the vector by one position (that is, $shift(d_l)[i] = d_l[i + 1] \forall i \in \mathbb{N}$). For locations l in the goal region B , all enabled actions $a \in Act(l)$ are greedy.

Lemma 1. *Greedy schedulers exist, and they can be described as the class of schedulers that choose a greedy action upon every reachable time-abstract path.*

Proof. It is plain that, for every non-goal location $l \notin B$, $shift(d_l) \geq \sum_{l' \in L} \mathbf{P}(l, a, l') d_{l'}$ holds for every action a , and that equality must hold for some.

For a scheduler S that always chooses greedy actions, a simple inductive argument shows that $d_l[i] = d_{l,S}[i]$ holds for all $i \in \mathbb{N}$, while it is easy to show that $d_l > d_{l,S}$ holds if S deviates from greedy decisions upon a path that is possible under its own scheduling policy and does not contain a goal location. \square

This allows in particular to fix a positional *standard greedy scheduler* by fixing an arbitrary greedy action for every location.

To determine the set of greedy actions, let us consider a deterministic scheduler S that starts in a location l with a non-greedy action a . Then $shift(d_{l,S}) \leq \sum_{l' \in L} \mathbf{P}(l, a, l') d_{l'}$ holds true, where the sum $\sum_{l' \in L} \mathbf{P}(l, a, l') d_{l'}$ corresponds to the scheduler choosing the non-greedy action a at location l and acting greedy in all further steps. Let $d_{l,a} + \sum_{l' \in L} \mathbf{P}(l, a, l') d_{l'}$ denote the step probability vector of such schedulers.

We know that $d_{l,S} \leq d_{l,a} < d_l$. Hence, there is not only a difference between $d_{l,S}$ and d_l , this difference will not occur at a higher index as the first difference between the newly defined $d_{l,a}$ and d_l . The finite number of locations and actions thus implies the existence of a bound k on the occurrence of this first difference between $d_{l,a}$ and d_l as well as $d_{l,S}$ and d_l . While the existence of such a k suffices to demonstrate the existence of optimal schedulers, we show in Appendix A that this constant $k < |L|$ is smaller than the CTMDP itself.

Having established such a bound k , it suffices to compare schedulers up to this bound. This provides us with the greedy actions, and also with the initial sequence $d_{l,a}[0], d_{l,a}[1], \dots, d_{l,a}[k]$ for all locations l and actions a . Consequently, we can determine a positive lower bound $\mu > 0$ for the first non-zero entry of the vectors $d_l - d_{l,S}$. We call this lower bound μ the *discriminator* of the CTMDP. Intuitively, the discriminator μ represents the minimal advantage of greedy strategy over non-greedy strategies.

3.2 Uniform CTMDPs

In this subsection, we show that every CD or HD scheduler for a uniform CTMDP can be transformed into a scheduler that converges to this standard greedy scheduler.

In the quest for an optimal CD scheduler, it is useful to consider the fact that the maximal reachability probability can be computed using the step probability vector, because the likelihood that a particular number of steps happen in time t is independent of the scheduler:

$$Pr_S^{\mathcal{M}}(t) = \sum_{l \in L} v(l) \sum_{i=0}^{\infty} d_{l,S}[i] \cdot p_{\lambda_{\mathcal{M}}}(i). \quad (1)$$

Moreover, the Poisson distribution $p_{\lambda_{\mathcal{M}}}$ has the useful property that the probability of taking k steps is falling very fast. We define the *greed bound* $n_{\mathcal{M}}$ to be a natural number, for which

$$\mu p_{\lambda_{\mathcal{M}}}(n) \geq \sum_{i=1}^{\infty} p_{\lambda_{\mathcal{M}}}(n+i) \quad \forall n \geq n_{\mathcal{M}} \quad (2)$$

holds true. It suffices to choose $n_{\mathcal{M}} \geq \frac{2\lambda_{\mathcal{M}}}{\mu}$ since it implies $\mu p_{\lambda_{\mathcal{M}}}(n) \geq 2p_{\lambda_{\mathcal{M}}}(n+1)$, $\forall n > n_{\mathcal{M}}$ (which yields (2) by simple induction). Such a greed bound implies that the decrease in likelihood of reaching the goal region in few steps caused by making a non-greedy decision after the greed bound dwarfs any potential later gain. We use this observation to improve any given CD or HD scheduler \mathcal{S} that makes a non-greedy decision after $\geq n_{\mathcal{M}}$ steps by replacing the behaviour after this history by a greedy scheduler. Finally, we use the interchangeability of greedy schedulers to introduce a scheduler $\bar{\mathcal{S}}$ that makes the same decisions as \mathcal{S} on short histories and follows the standard greedy scheduling policy once the length of the history reaches the greed bound. For this scheduler, we show that $Pr_{\bar{\mathcal{S}}}^{\mathcal{M}}(t) \geq Pr_{\mathcal{S}}^{\mathcal{M}}(t)$ holds true.

Theorem 1. *For uniform CTMDPs, there is an optimal scheduler for the classes CD and HD that converges to the standard greedy scheduler after $n_{\mathcal{M}}$ steps.*

Proof. Let us consider any HD scheduler \mathcal{S} that makes a non-greedy decision after a time-abstract path π of length $|\pi| \geq n_{\mathcal{M}}$ with last location l . If the path ends in, or has previously passed, the goal region, or if the probability of the history π is 0, that is, if it cannot occur with the scheduling policy of \mathcal{S} , then we can change the decision of \mathcal{S} on every path starting with π arbitrarily—and in particular to the standard greedy scheduler—without altering the reachability probability.

If $PR_{\mathcal{S}}^{\mathcal{M}}(\pi, t) > 0$, then we change the decisions of the scheduler \mathcal{S} for paths with prefix π such that they comply with the standard greedy scheduler. We call the resulting HD scheduler \mathcal{S}' and analyse the change in reachability probability using Equation (1):

$$Pr_{\mathcal{S}'}^{\mathcal{M}}(t) - Pr_{\mathcal{S}}^{\mathcal{M}}(t) = PR_{\mathcal{S}}^{\mathcal{M}}(\pi, t) \cdot \sum_{i=0}^{\infty} (d_l[i] - d_{l, \mathcal{S}_{\pi}}[i]) \cdot p_{\lambda_{\mathcal{M}}}(|\pi| + i),$$

where $\mathcal{S}_{\pi} : \pi' \mapsto \mathcal{S}(\pi \circ \pi')$ is the HD scheduler which prefixes its input with the path π and then calls the scheduler \mathcal{S} . The greedy criterion implies $d_l > d_{l, \mathcal{S}_{\pi}}$ with respect to the lexicographic order, and we can apply Equation 2 to deduce that the difference $Pr_{\mathcal{S}'}^{\mathcal{M}}(t) - Pr_{\mathcal{S}}^{\mathcal{M}}(t)$ is non-negative.

Likewise, we can concurrently change the scheduling policy to the standard greedy scheduler for all paths of length $\geq n_{\mathcal{M}}$ for which the scheduler \mathcal{S} makes non-greedy decisions. In this way, we obtain a scheduler \mathcal{S}'' that makes non-greedy decisions only in the first $n_{\mathcal{M}}$ steps, and yields a (not necessarily strictly) better time-bounded reachability probability than \mathcal{S} .

Since all greedy schedulers are interchangeable without changing the time-bounded reachability probability (and even without altering the step probability vector), we can modify \mathcal{S}'' such that it follows the standard greedy scheduling policy after $\geq n_{\mathcal{M}}$ steps, resulting in a scheduler $\bar{\mathcal{S}}$ that comes with the same time-bounded reachability probability as \mathcal{S}'' . Note that $\bar{\mathcal{S}}$ is counting if \mathcal{S} is counting.

Hence, the supremum over the time-bounded reachability of all CD/HD schedulers is equivalent to the supremum over the bounded reachability of CD/HD schedulers that deviate from the standard greedy scheduler only in the first $n_{\mathcal{M}}$ steps. This class is finite, and the supremum over the bounded reachability is therefore the maximal bounded reachability obtained by one of its representatives. \square

Hence, we have shown the existence of a—simple—optimal time-bounded CD scheduler. Using the fact that the suprema over the time-bounded reachability probability coincide for CD, CR, HD, and HR scheduler [3], we can infer that such a scheduler is optimal for all of these classes.

Corollary 1. $\max_{\mathcal{S} \in CD} Pr_{\mathcal{S}}^{\mathcal{M}}(t) = \max_{\mathcal{S} \in HR} Pr_{\mathcal{S}}^{\mathcal{M}}(t)$ holds for all uniform CTMDPs \mathcal{M} . \square

3.3 Non-uniform CTMDPs

Reasoning over non-uniform CTMDPs is harder than reasoning over uniform CTMDPs, because the likelihood of seeing exactly k steps does not adhere to the simple Poisson distribution, but depends on the precise history. Even if two paths have the same length, they may imply different probability distributions over the time passed so far. Knowing the time-abstract history therefore provides a scheduler with more information about the system's state than merely its length. As a result, it is simple to construct example CTMDPs, for which history-dependent and counting schedulers can obtain different time-bounded reachability probabilities [3].

In this subsection, we extend the results from the previous subsection to general CTMDPs. We show that simple optimal CD/HD scheduler exist, and that randomisation does not yield an advantage:

$$\max_{\mathcal{S} \in CD} Pr_{\mathcal{S}}^{\mathcal{M}}(t) = \max_{\mathcal{S} \in CR} Pr_{\mathcal{S}}^{\mathcal{M}}(t) \quad \text{and} \quad \max_{\mathcal{S} \in HD} Pr_{\mathcal{S}}^{\mathcal{M}}(t) = \max_{\mathcal{S} \in HR} Pr_{\mathcal{S}}^{\mathcal{M}}(t).$$

To obtain this result, we work on the uniformisation \mathcal{U} of \mathcal{M} instead of working on \mathcal{M} itself. We argue that the behaviour of a general CTMDP \mathcal{M} can be viewed as the observable behaviour of its uniformisation \mathcal{U} , using a scheduler that does not *see* the new transitions and locations. Schedulers from this class can then be replaced by (or viewed as) schedulers that do not *use* the additional information. And finally, we can approximate schedulers that do not use the additional information by schedulers that do not use it initially, where initially means until the number of visible steps—and hence

in particular the number of steps—exceeds the greed bound $n_{\mathcal{U}}$ of the uniformisation \mathcal{U} of \mathcal{M} . Comparable to the argument from the proof of Theorem 1, we show that we can restrict our attention to the standard greedy scheduler after this initial phase, which leads again to a situation where considering a finite class of schedulers suffices to obtain the optimum.

Lemma 2. *The greedy decisions and the step probability vector coincide for the observable and unobservable copy of each location in the uniformisation \mathcal{U} of any CTMDP \mathcal{M} .*

Proof. The observable and unobservable copy of each location reach the same successors under the same actions with the same transition rate. \square

We can therefore choose a positional *standard greedy scheduler* whose decisions coincide for the observable and unobservable copy of each location.

For the *uniformisation* \mathcal{U} of a CTMDP \mathcal{M} , we define the function $vis : Paths_{abs}(\mathcal{U}) \rightarrow Paths_{abs}(\mathcal{M})$ that maps a path π of \mathcal{U} to the corresponding path in \mathcal{M} , the *visible path*, by deleting all unobservable locations and their preceding transitions from π . (Note that all paths in \mathcal{U} start in an observable location.) We call a scheduler *n-visible* if its decisions only depend on the visible path and coincide for the observable and unobservable copy of every location for all paths containing up to n visible steps. We call a scheduler *visible* if it is n -visible for all $n \in \mathbb{N}$.

We call a HD/HR scheduler an (n -)visible HD/HR scheduler if it is (n -)visible, and we call an (n -)visible HD/HR scheduler a visible CD/CR scheduler if its decisions depend only on the length of the visible path, and an n -visible CD/CR scheduler if its decisions depend only on the length of the visible path for all paths containing up to n visible steps. The respective classes are denoted with according prefixes, for example, n -vCD. Note that (n -)visible counting schedulers are not counting.

It is a simple observation that we can study visible CD, CR, HD, and HR schedulers on the uniformisation \mathcal{U} of a CTMDP \mathcal{M} instead of studying CD, CR, HD, and HR schedulers on \mathcal{M} .

Lemma 3. $S \mapsto S \circ vis$ is a bijection from visible CD, CR, HD, or HR schedulers for the uniformisation \mathcal{U} of a CTMDP \mathcal{M} onto CD, CR, HD, or HR schedulers, respectively, of \mathcal{M} that preserves the time-bounded reachability probability: $Pr_S^{\mathcal{U}}(t) = Pr_{S \circ vis}^{\mathcal{M}}(t)$. \square

At the same time, copying the argument from the proof of Theorem 1, an $n_{\mathcal{U}}$ -visible CD or HD scheduler S can be adjusted to the $n_{\mathcal{U}}$ -visible CD or HD scheduler \bar{S} that deviates from S only in that it complies with the standard greedy scheduler for \mathcal{U} after $n_{\mathcal{U}}$ visible steps, without decreasing the time-bounded reachability probability. These schedulers are visible schedulers from a finite sub-class, and hence some representative of this class takes the optimal value.

Lemma 4. *The following equations hold for the uniformisation \mathcal{U} of a CTMDP \mathcal{M} :*

$$\max_{S \in n_{\mathcal{U}}\text{-vCD}} Pr_S^{\mathcal{U}}(t) = \max_{S \in \text{vCD}} Pr_S^{\mathcal{U}}(t) \quad \text{and} \quad \max_{S \in n_{\mathcal{U}}\text{-vHD}} Pr_S^{\mathcal{U}}(t) = \max_{S \in \text{vHD}} Pr_S^{\mathcal{U}}(t).$$

Proof. We have shown in Theorem 1 that turning to the standard greedy scheduling policy after $n_{\mathcal{U}}$ or more steps can only increase the time-bounded reachability probability. This implies that we can turn to the standard greedy scheduler after $n_{\mathcal{U}}$ visible steps.

The scheduler resulting from this adjustment does not only remain $n_{\mathcal{U}}$ -visible, it becomes a visible CD and HD scheduler, respectively. Moreover, it is a scheduler from the finite subset of CD or HD schedulers, respectively, whose behaviour may only deviate from the standard scheduler within the first $n_{\mathcal{U}}$ visible steps. \square

We can therefore construct optimal CD and HD scheduler for every CTMDP \mathcal{M} . To prove that optimal CD and HD scheduler are also optimal CR and HR scheduler, respectively, we first prove the simpler lemma that this holds for k -bounded reachability.

Lemma 5. *k -optimal CD or HD schedulers are also k -optimal CR or HR schedulers, respectively.*

Proof. For a CTMDP \mathcal{M} we can turn an arbitrary CR or HR scheduler \mathcal{S} into a CD or HD scheduler \mathcal{S}' with a time and k -bounded reachability probability that is at least as good as the one of \mathcal{S} by first determinising the scheduler decisions from the $(k+1)$ st step onwards—this has obviously no impact on k -bounded reachability—and then determinising the remaining randomised choices.

Replacing a single randomised decision on a path π (for history-dependent schedulers) or on a set of paths Π (for counting schedulers) that end(s) in a location l is safe, because the time and k -bounded reachability probability of a scheduler is an affine combination—the affine combination defined by $\mathcal{S}(\pi)$ and $\mathcal{S}(|\pi|, l)$, respectively—of the $|\text{Act}(l)|$ schedulers resulting from determinising this single decision. Hence, we can pick one of them whose time and k -bounded reachability probability is at least as high as the one of \mathcal{S} .

As the number of these randomised decisions is finite ($\leq k|L|$ for CR, and $\leq k^{|L|}$ for HR schedulers), this results in a deterministic scheduler after a finite number of improvement steps. \square

Theorem 2. *Optimal CD schedulers are also optimal CR schedulers.*

Proof. First, the probability that the goal region B is reached in more than k steps converges to 0, independent of the scheduler. Together with Lemma 5, this implies

$$\sup_{\mathcal{S} \in CR} Pr_{\mathcal{S}}^{\mathcal{M}}(t) = \lim_{n \rightarrow \infty} \sup_{\mathcal{S} \in CR} Pr_{\mathcal{S}}^{\mathcal{M}}(t; k) = \lim_{n \rightarrow \infty} \sup_{\mathcal{S} \in CD} Pr_{\mathcal{S}}^{\mathcal{M}}(t; k) \leq \max_{\mathcal{S} \in CD} Pr_{\mathcal{S}}^{\mathcal{M}}(t),$$

while \geq is implied by $CD \subseteq CR$. \square

Analogously, we can prove the similar theorem for history-dependent schedulers:

Theorem 3. *Optimal HD schedulers are also optimal HR schedulers.* \square

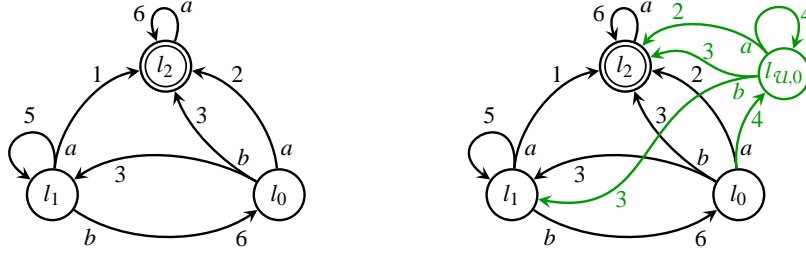


Fig. 2. The example CTMDP \mathcal{M} (left) and the reachable part of its uniformisation \mathcal{U} (right).

3.4 Example

To exemplify our proposed construction, let us consider the example CTMDP \mathcal{M} depicted in Figure 2. As \mathcal{M} is not uniform, we start with constructing the uniformisation \mathcal{U} of \mathcal{M} (cf. Figure 2).

\mathcal{U} has the uniform transition rate $\lambda = 6$. Independent of the initial distribution of \mathcal{M} , the unobservable copies of l_1 and l_2 are not reachable in \mathcal{U} , because the initial distribution of a uniformisation assigns all probability weight to observable locations, and the transition rate of all enabled actions in l_1 and l_2 in \mathcal{M} is already λ . (Unobservable copies of a location l are only reachable from the observable and unobservable copy of l upon enabled actions a with non-maximal exit rate $\mathbf{R}(l, a, L) \neq \lambda$.)

Disregarding the unreachable part of \mathcal{U} , there are only 8 positional schedulers for \mathcal{U} , and only 4 of them are visible (that is, coincide on l_0 and $l_{\mathcal{U},0}$). They can be characterised by $\mathcal{S}_1 = \{l_0 \mapsto a, l_1 \mapsto a\}$, $\mathcal{S}_2 = \{l_0 \mapsto a, l_1 \mapsto b\}$, $\mathcal{S}_3 = \{l_0 \mapsto b, l_1 \mapsto a\}$, and $\mathcal{S}_4 = \{l_0 \mapsto b, l_1 \mapsto b\}$. In order to determine a greedy scheduler, we first determine step probability vectors:

For l_0 : $d_{l_0, \mathcal{S}_1} = d_{l_0, \mathcal{S}_2} = (\frac{1}{3}, \frac{5}{9}, \frac{19}{27}, \dots)$, $d_{l_0, \mathcal{S}_3} = (\frac{1}{2}, \frac{7}{12}, \frac{43}{72}, \dots)$, $d_{l_0, \mathcal{S}_4} = (\frac{1}{2}, \frac{1}{2}, \frac{3}{4}, \dots)$.

For l_1 : $d_{l_1, \mathcal{S}_1} = d_{l_1, \mathcal{S}_3} = (\frac{1}{6}, \frac{7}{36}, \frac{71}{216}, \dots)$, $d_{l_1, \mathcal{S}_2} = (0, \frac{1}{3}, \frac{5}{9}, \dots)$, $d_{l_1, \mathcal{S}_4} = (0, \frac{1}{2}, \frac{1}{2}, \dots)$.

Note that, in the given example, it suffices to compute the step probability vector for a single step to determine that \mathcal{S}_3 is optimal (w.r.t. the greedy optimality criterion); in general, it suffices to consider as many steps as the CTMDP has locations. Since deviating from \mathcal{S}_3 decreases the chance to reach the goal location l_2 in a single step by $\frac{1}{6}$ both from l_0 and l_1 , the discriminator $\mu = \frac{1}{6}$ is easy to compute.

Our coarse estimation provides a greed bound of $n_{\mathcal{U}} = \lceil 72 \cdot t \rceil$, where t is the time bound, but $n_{\mathcal{U}} = \lceil 42 \cdot t \rceil$ suffices to satisfy Equation (2).

When seeking optimal schedulers from any of the discussed classes, we can focus on the finite set of those schedulers that comply with \mathcal{S}_3 after $n_{\mathcal{U}}$ (visible) steps. In the following subsection, we describe how the precise model checking technique of Aziz et al. [2] can be exploited to turn the existence proof into an effective technique for the construction of optimal schedulers.

3.5 Constructing Optimal Schedulers

The techniques discussed in the previous subsections prove the existence of optimal schedulers by showing that it suffices to consider deterministic schedulers that deviate

from a standard greedy scheduler only in the first $n_{\mathcal{M}}$ (or $n_{\mathcal{U}}$) steps. The combination of each of these schedulers with the respective CTMDP can be viewed as a *finite* continuous-time Markov *chain* (CTMC).

Aziz et al. [2] have shown that the time-bounded reachability probability of CTMCs are computable (and comparable) finite sums $\sum_{i \in I} \eta_i e^{\delta_i}$, where the individual η_i and δ_i are algebraic numbers. This provides us with a constructive extension of the results developed in this section:

Corollary 2. *We can effectively construct optimal CD, CR, HD, and HR schedulers. \square*

Corollary 3. *We can compute the time-bounded reachability probability of optimal schedulers as finite sums $\sum_{i \in I} \eta_i e^{\delta_i}$, where the η_i and δ_i are algebraic numbers. \square*

These corollaries, however, lean on the precise CTMC model checking approach of Aziz et al. [2], which only demonstrates the effective decidability of this problem. We deem it unlikely that a complexity for finding optimal strategies can be provided prior to determining the respective CTMC model checking complexity.

4 Time-Dependent Schedulers

In this section we make use of a simple but illuminative shift in our view on the control problem for a CTMDP \mathcal{M} : We consider the time that has passed as part of the state-space³ of a time-extended CTMDP (tCTMDP), turning the time-bounded reachability problem to reach B in time t_0 into an ordinary reachability problem to reach $B \times [0, t_0]$ in a tCTMDP \mathcal{M}_{t_0} .

This extension has obviously no effect on the Markovian character of the tCTMDP. In particular for a TT scheduler \mathcal{S} , which can revoke its decisions, the probability $Pr_{\mathcal{S}}^{\mathcal{M}}((l, t))$ to reach the goal region $B \times [0, t_0]$ from a state (l, t) in the time extended CTMDP \mathcal{M}_{t_0} is independent of the history.

For a traditional time-dependent schedulers \mathcal{S} , the probability $Pr_{\mathcal{S}}^{\mathcal{M}}((l, t))$ to reach a location in the goal region from a state (l, t) is memoryful in general, as the decisions made by the scheduler \mathcal{S} depend on the time that the location l was entered. However, the behaviour becomes memoryless if we focus on the points of time at which a discrete transition took place.

In both cases it is simple to translate positional schedulers for the resulting time-extended CTMDP \mathcal{M}_{t_0} to equivalent TTP/TP schedulers for the original CTMDP \mathcal{M} .

For every TP scheduler \mathcal{S} , we have $Pr_{\mathcal{S}}^{\mathcal{M}_{t_0}}((l, t)) = 1$ for all goal states $(l, t) \in B \times [0, t_0]$, as we have reached the goal region in time in this case, and $Pr_{\mathcal{S}}^{\mathcal{M}_{t_0}}((l, t)) = 0$ for all locations $l \in L$ and all $t > t_0$, because the goal region cannot be reached in time any longer if it has not been visited before.

For a measurable deterministic positional scheduler \mathcal{S} and a non-goal locations $l' \notin B$ and time $t \in [0, t_0]$, we will reach the goal region in time (provided we have not

³ Adding time to the state-space leads to a construction that resembles the semantics of timed automata [1], but with a simpler treatment of time (only one clock and no resets).

reached it before), if we reach it in time with or after the following transition. Hence,

$$Pr_S^{\mathcal{M}_0}((l,t)) = \sum_{l' \in L} \mathbf{R}(l, \mathcal{S}((l,t)), l') \int_t^\infty Pr_S^{\mathcal{M}_0}(l', \tau) e^{-\mathbf{R}(l, \mathcal{S}((l,t)), L)\tau} d\tau$$

holds true, where $Pr_S^{\mathcal{M}_0}((l,t))$ denotes the probability of reaching the goal region when the location l is *entered* at time t . Different to tCTMDPs for TT schedulers (cf. Appendix B), tCTMDPs for traditional T schedulers therefore have a discrete flavour.

Naturally, this shift in our way of looking at the problem has no influence on the probability of reaching our objective, and the following equations must hold:

$$\sup_{S \in TP} Pr_S^{\mathcal{M}}(l, t_0 - t) = \sup_{S \in P} Pr_S^{\mathcal{M}_0}((l,t)) \quad \text{and} \quad \sup_{S \in TP} Pr_S^{\mathcal{M}}(t_0) = \sum_{l \in L} v(l) \sup_{S \in P} Pr_S^{\mathcal{M}_0}((l,0)).$$

For both T and TT schedulers, the hard part is to show that an optimal measurable scheduler exists. We prove that there are optimal schedulers in the class of randomised history and time-dependent schedulers that are deterministic and positional.

Theorem 4. $\max_{S \in TP} Pr_S^{\mathcal{M}}(t) = \sup_{S \in TH} Pr_S^{\mathcal{M}}(t)$, and randomisation does not improve the result.

Proof. The formulas given above for positional schedulers, as well as similar formulas for history-dependent schedulers, are clearly dominated by the functions defined by

$$Pr_P^{\mathcal{M}_0}((l,t)) = \max_{a \in Act(l)} \sum_{l' \in L} \mathbf{R}(l, a, l') \int_t^\infty Pr_P^{\mathcal{M}_0}(l', \tau) e^{-\mathbf{R}(l, a, L)\tau} d\tau.$$

For an extension to randomised schedulers, the maximum over the actions needs to be replaced by a supremum over the distributions in an intermediate step, but as suprema over affine combinations over a finite set of values are taken in one of these values, the same function is dominating the functions for measurable history-dependent randomised schedulers as well.

The hard part is to show that a measurable scheduler exists that takes these maxima, that is, that no non-measurable change between different actions is required. To prove this, we show how to construct a measurable deterministic positional scheduler \mathcal{S} that always chooses actions a that take the maximum value.

To determine suitable scheduler decisions for a location l for such a scheduler \mathcal{S} , we disintegrate $[0, t_0]$ into measurable sets $\{T_a \mid a \in Act(l)\}$, such that \mathcal{S} only makes decisions that maximise $\sum_{l' \in L} \mathbf{R}(l, a, l') \int_t^\infty Pr_P^{\mathcal{M}_0}(l', \tau) e^{-\mathbf{R}(l, a, L)\tau} d\tau$. (For positions outside of $[0, t_0]$, that is, for times behind our time bound t_0 , the behaviour of the scheduler does not matter, and $\mathcal{S}(l, t)$ can be fixed to any constant decision $a \in Act(l)$ for all $t > t_0$.)

We start with fixing an arbitrary order \succ on the actions in $Act(l)$, and introduce, for each point $t \in [0, t_0]$, an order \succ_t on the actions determined by the value of $\sum_{l' \in L} \mathbf{R}(l, a, l') \int_t^\infty Pr_P^{\mathcal{M}_0}(l', \tau) e^{-\mathbf{R}(l, a, L)\tau} d\tau$, using \succ as a tie-breaker.

1. For the action a in $Act(l)$ that is minimal with respect to \succ , we start by fixing the open set $O_a = [0, t_0]$ of points in time where the scheduler does not make a decision $a' \succ a$ (where open set in this proof refers to sets open in $[0, t_0]$).

2. We then define the set T_a as the points $t \in O_a$ in time, for which the action a is maximal with respect to \succ_t .
 T_a is an open measurable set with a countable fringe, and for all points $t \in \overline{T_a} \setminus T_a$ it holds that a maximises $\sum_{l' \in L} \mathbf{R}(l, b, l') \int_t^\infty Pr_P^{\mathcal{M}_{l_0}}((l', \tau)) e^{-\mathbf{R}(l, b, L)\tau} d\tau$ among all actions $b \in Act(l)$, though not strictly. (A detailed description why the continuity of $\sum_{l' \in L} \mathbf{R}(l, b, l') \int_t^\infty Pr_P^{\mathcal{M}_{l_0}}((l', \tau)) e^{-\mathbf{R}(l, b, L)\tau} d\tau$ for all actions $b \in Act(l)$ implies that T_a is open, measurable, and has a countable fringe is supplied in Appendix B.)
3. We fix $\mathcal{S}(l, t) = a$ for all $t \in \overline{T_a} \cap O_a$.
4. If there is a next smaller (with respect to \succ) action $a' = \max\{a'' \prec a\}$, then we fix the new open set $O_{a'} = O_a \setminus \overline{T_a}$ for a' , and proceed with Step 2.

Repeating this for all non-goal locations $l \notin B$, and fixing arbitrary decisions for the goal locations (independent of the time passed) provides the sought measurable deterministic time-dependent positional scheduler that dominates all history-dependent randomised time-dependent schedulers. \square

The proof for a similar theorem for TT schedulers is moved to Appendix B due to space restrictions.

Theorem 5. $\max_{\mathcal{S} \in TTP} Pr_S^{\mathcal{M}}(t) = \sup_{\mathcal{S} \in TTH} Pr_S^{\mathcal{M}}(t)$, and randomisation does not improve the result.

References

1. Rajeev Alur and David L. Dill. A Theory of Timed Automata. *Theoretical Computer Science*, 126(2):183–235, 1994.
2. Adnan Aziz, Kumud Sanwal, Vigyan Singhal, and Robert Brayton. Model-checking continuous-time Markov chains. *Transactions on Computational Logic*, 1(1):162–170, 2000.
3. Christel Baier, Holger Hermanns, Joost-Pieter Katoen, and Boudewijn R. Haverkort. Efficient computation of time-bounded reachability probabilities in uniform continuous-time Markov decision processes. *Theoretical Computer Science*, 345(1):2–26, 2005.
4. J. Bruno, P. Downey, and G. N. Frederickson. Sequencing Tasks with Exponential Service Times to Minimize the Expected Flow Time or Makespan. *Journal of the ACM*, 28(1):100–113, 1981.
5. Eugene A. Feinberg. Continuous Time Discounted Jump Markov Decision Processes: A Discrete-Event Approach. *Mathematics of Operations Research*, 29(3):492–524, 2004.
6. H. Hermanns. *Interactive Markov Chains and the Quest for Quantified Quality*. LNCS 2428. Springer-Verlag, 2002.
7. Vidyadhar G. Kulkarni. *Modeling and Analysis of Stochastic Systems*. Chapman & Hall, Ltd., London, UK, 1995.
8. M. A. Marsan, G. Balbo, G. Conte, S. Donatelli, and G. Franceschinis. Modelling with Generalized Stochastic Petri Nets. *SIGMETRICS Performance Evaluation Review*, 26(2):2, 1998.
9. Martin R. Neuhäüßer, Mariëlle Stoelinga, and Joost-Pieter Katoen. Delayed Nondeterminism in Continuous-Time Markov Decision Processes. In *Proceedings of FOSSACS '09*, pages 364–379, 2009.
10. Martin R. Neuhäüßer and Lijun Zhang. Time-Bounded Reachability in Continuous-Time Markov Decision Processes. Technical report, 2009.

11. Martin L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. Wiley-Interscience, April 1994.
12. William H. Sanders and John F. Meyer. Reduced Base Model Construction Methods for Stochastic Activity Networks. In *Proceedings of PNPM'89*, pages 74–84, 1989.
13. Linn I. Sennott. *Stochastic Dynamic Programming and the Control of Queueing Systems*. Wiley-Interscience, 1999.
14. Nicolás Wolovick and Sven Johr. A Characterization of Meaningful Schedulers for Continuous-Time Markov Decision Processes. In *Proceedings of FORMATS'06*, pages 352–367, 2006.
15. L. Zhang, H. Hermanns, E. M. Hahn, and B. Wachter. Time-bounded model checking of infinite-state continuous-time Markov chains. In *Proceedings of ACSD'08*, pages 98–107, 2008.

Appendix

A Greedy Schedulers

The proof of the existence of a positional greedy scheduler in Subsection 3.1 is not constructive, because we do not provide means to compute a constant k that satisfies $d_l \neq d_{l,a} \Rightarrow \exists k' \leq k. d_l[k] > d_{l,a}[k]$. Moreover, suitable constants μ and $n_{\mathcal{M}}$ can only be computed once k is known. Without being able to provide a suitable constant, we could therefore provide an algorithm that converges to optimal time-abstract schedulers, but we would not be able to determine whether an optimal solution has already been reached.

The techniques we exploit in this appendix to show that $k < |L|$ is in fact smaller than the uniform CTMDP \mathcal{M} itself draw from linear algebra, and are, while simple, a bit unusual in this context. We first turn to the simpler notion of Markov chains, resolving the nondeterminism in accordance with the positional standard greedy scheduler \mathcal{S} whose existence was shown in Subsection 3.1.

We first lift the step probability vector from locations to distributions, where $d_v = \sum_{l \in L} v(l) d_l$ is, for a distribution $v : L \rightarrow [0, 1]$, the normal affine combination of the step probability vectors of the individual locations. We call two distributions $v, v' : L \rightarrow [0, 1]$ equivalent, if their step probability vectors $d_v = d_{v'}$ are equal, and i -step equivalent if they are equal up to position i ($\forall j \leq i. d_v[j] = d_{v'}[j]$). We immediately extend these definitions to arbitrary vectors $v : L \rightarrow \mathbb{R}$.

We then define the vector spaces D_i of multitudes of differences $v - v'$ of i -step equivalent distributions $v, v' : L \rightarrow [0, 1]$. That is, the differences between i -step equivalent distributions form a spanning set of D_i .

As a basis for our construction, D_0 is simple to construct: it is the vector space that contains the multitudes of differences $v - v'$ of distributions $v, v' : L \rightarrow [0, 1]$ that are equally likely in the goal region. The first implies the restriction $\sum_{l \in L} \delta(l) = 0$, and the latter the restriction $\sum_{l \in B} \delta(l) = 0$ for all $\delta \in D_0$. On the other hand, all vectors that satisfy these side conditions are clearly multitudes of differences $v - v'$ of 0-step equivalent distributions $v, v' : L \rightarrow [0, 1]$, such that D_0 is completely described by this description.

Once we have computed the vector space D_i , we can compute the vector space O_i that contains a vector δ if it is a multitude of differences $v - v'$ of distributions $v, v' : L \rightarrow [0, 1]$, such that $\text{shift}(d_v)$ and $\text{shift}(d_{v'})$ are i -step equivalent. The transition from step probability vectors to the *shift* of them is a simple linear operation, which essentially transforms the distributions according to the transition matrix, but drops all weight that is already accumulated in the goal region. Hence, we can obtain O_i from D_i by a simple linear transformation of the vector space.

Now, D_{i+1} obviously consists of those vectors in D_i that are also in O_i , and $D_{i+1} = D_i \cap O_i$ can be obtained by an intersection of two vector spaces.

Now $D_0 \supseteq D_1 \supseteq \dots \supseteq D_{|L|-2} \supseteq \dots$ obviously holds, and $D_i = D_{i+1}$ implies $O_i = O_{i+1}$, and hence $D_i = D_\infty$. As D_0 is a $|L| - 2$ dimensional vector space, and inequality implies the loss of at least one dimension, a fixed point is reached after at most $|L| - 2$ steps. That is, two distributions are equivalent, if, and only if, they are $(|L| - 2)$ -step equivalent.

Having established this, we apply it on the distribution $v_{l,a}$ obtained in one step from a position $l \notin B$ when choosing the action a , as compared to the distribution v_l obtained when choosing the action $S(l)$ defined by the positional greedy scheduler.

Now, $d_l > d_{l,a}$ holds if, and only if $\text{shift}(d_l) = d_{v_l} > d_{v_{l,a}} = \text{shift}(d_{l,a})$, which implies $d_{v_l}[k'] > d_{v_{l,a}}[k']$ for some $k' \leq |L| - 2$, and hence $d_l[k] > d_{l,a}[k]$ for some $k < |L|$.

B TT Schedulers

In this appendix we describe the small differences that occur when we allow for schedulers that have the capability to revoke their decisions.

Probability Space. As a first adjustment, we have to build a probability space that covers this generalisation. Such spaces are not hard to build (cf. [9, 10] for locally uniform CTMDPs): We can simply define measures for simple types of these schedulers, and complete the measure space in the usual way.

That is, we start with defining the probability measure for sets of paths

$$I_0 \xrightarrow{a_0, t_0} I_1 \xrightarrow{a_1, t_1} \dots \xrightarrow{a_{n-1}, t_{n-1}} I_n$$

with $0 < t_0 < t_1 < t_2 < \dots < t_{n-1}$, such that $t_0 \in I_0, t_1 \in I_1, \dots, t_{n-1} \in I_{n-1}$, for disjoint open intervals I_0, I_1, \dots, I_{n-1} , and schedulers that revoke their decisions in finitely many points r_1, r_2, \dots, r_m , but whose decisions do not depend on the times t_0, t_1, \dots, t_{n-1} .

For such simple sets of paths and schedulers, we can compute the probability to obtain a path in this cylindrical set as

$$\int_{t_0 \in I_0, t_1 \in I_1, \dots, t_{n-1} \in I_{n-1}} \prod_{i=0}^{n-1} \mathbf{R}(l_i, a_i, l_{i+1}) \prod_{i=1}^{m+n} e^{-\mathbf{R}(l'_i, a'_i, L)(t'_i - t'_{i-1})},$$

where

- $t'_0 = 0$,
- $t'_1 < t'_2 < \dots < t'_{m+n}$ is the chain of points in time that contains $t_0 < t_1 < t_2 < \dots < t_{n-1}$ and r_1, r_2, \dots, r_m ,
- l'_i is the location the CTMDP is in for the interval (t'_{i-1}, t'_i) , and
- a'_i is the decision the scheduler would make in the time interval (t'_{i-1}, t'_i) , which is also the decision it makes at the times t_j of the discrete transitions.

The extension to randomised schedulers is trivial.

These probabilities for cylindrical sets then become the basic building blocks of our probability space: As usual, we can build a σ -algebra over these sets, and complete the resulting simple measure space. Note that this definition does not raise the requirement of locally uniform schedulers that was considered necessary previously (cf. [9, 10]), although using locally uniform schedulers admittedly simplifies representing the measure of these cylindrical sets of traces to the same integral used in Section 2.

Optimal TT Schedulers. Based on the resulting probabilistic space, we argue as in Section 4 that we can consider tCTMDPs instead of the standard ones, and that the resulting tCTMDPs remain Markovian. This suggests a proof for the existence of optimal TT schedulers comparable to the prove for time-dependent schedulers that cannot revoke their decisions.

The main difference to the proof in Section 4 is that TT schedulers can revoke their decision in any point of time, and the resulting tCTMDP $\mathcal{M}_{t_0}^{\mathcal{M}_0}$ is Markovian in any state, rather than only in any discrete entry point. This takes away the discrete flavour from the T scheduler case. Comparable to the case of T schedulers, we know that

- $Pr_S^{\mathcal{M}_{t_0}^{\mathcal{M}_0}}((l, t)) = 1$ holds for all goal states $l \in B$ and all $t \leq t_0$,
- $Pr_S^{\mathcal{M}_{t_0}^{\mathcal{M}_0}}((l, t)) = 0$ holds for all locations $l \in L$ and all $t > t_0$, and
- $Pr_S^{\mathcal{M}_{t_0}^{\mathcal{M}_0}}((l, t_0)) = 0$ holds for all non-goal locations $l \notin B$.

holds for every scheduler. For a measurable positional scheduler \mathcal{S} , we now have that

$$Pr_S^{\mathcal{M}_{t_0}^{\mathcal{M}_0}}((l, t)) = \sum_{l' \in L} \mathbf{R}(l, \mathcal{S}((l, t)), l') \cdot \left(Pr_S^{\mathcal{M}_{t_0}^{\mathcal{M}_0}}((l, t)) - Pr_S^{\mathcal{M}_{t_0}^{\mathcal{M}_0}}((l', t)) \right)$$

for all non-goal locations $l \notin B$, and all $t \in [0, t_0]$ holds true, where $\dot{Pr}_S^{\mathcal{M}_{t_0}^{\mathcal{M}_0}}((l, t))$ is the derivation of $Pr_S^{\mathcal{M}_{t_0}^{\mathcal{M}_0}}((l, t))$ to the second argument, that is, to the time.

Naturally, our shift in the way we look at the problem has again no influence on the probability of reaching our objective, and the following equations must hold:

$$\sup_{\mathcal{S} \in TTP} Pr_S^{\mathcal{M}}(l, t_0 - t) = \sup_{\mathcal{S} \in P} Pr_S^{\mathcal{M}_{t_0}^{\mathcal{M}_0}}((l, t)) \quad \text{and} \quad \sup_{\mathcal{S} \in TTP} Pr_S^{\mathcal{M}}(t_0) = \sum_{l \in L} v(l) \sup_{\mathcal{S} \in P} Pr_S^{\mathcal{M}_{t_0}^{\mathcal{M}_0}}((l, 0)).$$

The hard part is again to show that an optimal measurable scheduler exists.

Theorem 5. $\max_{\mathcal{S} \in TTP} Pr_S^{\mathcal{M}}(t_0) = \sup_{\mathcal{S} \in THH} Pr_S^{\mathcal{M}}(t_0)$, and randomisation does not improve the result.

Proof. The formulas discussed above provide us with simple differential equations, and the functions that we yield for positional schedulers, as well as those that we would get for history-dependent schedulers, are clearly dominated by the functions defined by the differential equation

$$\dot{Pr}_S^{\mathcal{M}_{t_0}^{\mathcal{M}_0}}((l, t)) = \min_{a \in Act(l)} \sum_{l' \in L} \mathbf{R}(l, a, l') \cdot \left(Pr_S^{\mathcal{M}_{t_0}^{\mathcal{M}_0}}((l, t)) - Pr_S^{\mathcal{M}_{t_0}^{\mathcal{M}_0}}((l', t)) \right) \quad \text{for all } t \in [0, t_0].$$

For an extension to randomised schedulers, the minimum over the actions needs to be replaced by an infimum over the distributions in an intermediate step, but as the infima over affine combinations of a finite set of values takes its minimum in one of these values, the same differential equations defines a dominating function.

Just like in the proof of Theorem 4, the hard part of the proof is to show that there is a measurable scheduler \mathcal{S} that always chooses an action a that minimises this value. This

guarantees $\sum_{l \in L} \nu(l) Pr_S^{\mathcal{M}}(l, 0) = \sup_{S \in TH} Pr_S^{\mathcal{M}}(t_0)$. We can construct such a scheduler similarly to the construction of an optimal scheduler in the proof of Theorem 4.

To construct the scheduler decisions for a location l for a measurable scheduler \mathcal{S} , we disintegrate $[0, t_0]$ into measurable sets $\{T_a \mid a \in Act(l)\}$, such that \mathcal{S} only makes decisions that minimise $\sum_{l' \in L} \mathbf{R}(l, a, l') \cdot (Pr_S^{\mathcal{M}_0}((l, t)) - Pr_S^{\mathcal{M}_0}((l', t)))$. (For positions outside of $[0, t_0]$, that is, for times behind the time bound t_0 , the behaviour of the scheduler does not matter and $\mathcal{S}(l, t)$ can be fixed to any constant decision $a \in Act(l)$ for all $t > t_0$.)

We start with fixing an arbitrary order \succ on the actions in $Act(l)$, and introduce, for each point $t \in [0, t_0]$, an order \succ_t on the actions determined by the value of $\sum_{l' \in L} \mathbf{R}(l, a, l') \cdot (Pr_S^{\mathcal{M}_0}((l, t)) - Pr_S^{\mathcal{M}_0}((l', t)))$, using \succ as a tie-breaker.

1. For the action a in $Act(l)$ that is maximal with respect to \succ , we start by fixing the open set $O_a = [0, t_0]$ of points in time where the scheduler does not make a decision $a' \succ a$ (where open set in this proof refers to sets open in $[0, t_0]$).
2. We then define the set T_a as the points $t \in O_a$ in time, for which the action a is minimal with respect to \succ_t .

Being minimal with respect to \succ_t requires the value $\sum_{l' \in L} \mathbf{R}(l, a, l') \cdot (Pr_S^{\mathcal{M}_0}((l, t)) - Pr_S^{\mathcal{M}_0}((l', t)))$ to be strictly smaller for a compared to the respective value of all other actions $a' \prec a$, which implies that this sum is also strictly smaller for all t' in some ε -environment of t in $[0, t_0]$.

As O_a is open, such an ε -environment is a subset of O_a , which implies that T_a is open.

We now first fix one such ε -environment $\mathcal{E}_t \subseteq O_a$ for all $t \in T_a$, and then fix an arbitrary sequence $\mathcal{E}'_0, \mathcal{E}'_1, \dots$ from these open sets such that, for all $i \in \mathbb{N}$, there is no $t \in T_a$ such that $\int_{\mathcal{E}_t \setminus \bigcup_{j < i} \mathcal{E}'_j} d\tau > 2 \int_{\mathcal{E}'_i \setminus \bigcup_{j < i} \mathcal{E}'_j} d\tau$.

Now, the open set $\mathcal{E}_a = \bigcup_{i \in \mathbb{N}} \mathcal{E}'_i$ is measurable, and since obviously $\lim_{i \rightarrow \infty} \int_{\mathcal{E}'_i \setminus \bigcup_{j < i} \mathcal{E}'_j} = 0$ holds true, $\int_{\mathcal{E}_t \setminus \bigcup_{i \in \mathbb{N}} \mathcal{E}'_i} d\tau = 0$ follows for all $t \in T_a$, which implies $\mathcal{E}_t \subseteq \overline{\mathcal{E}_a}$ (because the measurable open set $\mathcal{E}_a \setminus \overline{\mathcal{E}_a}$ is a 0 set, and hence empty).

Since \mathcal{E}_a is a countable union of open intervals, $|\overline{\mathcal{E}_a} \setminus \mathcal{E}_a| \leq |\mathbb{N}|$ is countable, and hence a 0 set. Furthermore, $\mathcal{E}_a \subseteq T_a \subseteq \overline{\mathcal{E}_a}$ holds true, and consequently $\overline{T_a} = \overline{\mathcal{E}_a}$ is measurable.

3. Moreover, for the points of time in the fringe of T_a , the continuity of $\sum_{l' \in L} \mathbf{R}(l, a', l') \cdot (Pr_S^{\mathcal{M}_0}((l, t)) - Pr_S^{\mathcal{M}_0}((l', t)))$ for every action $a' \in Act(l)$ guarantees that a still minimises this value, although not necessarily strictly. Hence, we can fix $\mathcal{S}(l, t) = a$ for all $t \in \overline{T_a} \cap O_a$.
4. If there is a next smaller (with respect to \succ) action $a' = \max\{a'' \prec a\}$, we fix the new open set $O_{a'} = O_a \setminus \overline{T_a}$ for a' , and proceed with step 2.

This way, we can construct a *measurable* scheduler that provides the optimal solution. \square