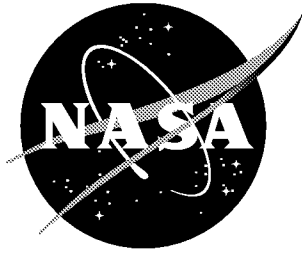


NASA / CR-2000-210296



Analysis of Phase-Type Stochastic Petri Nets With Discrete and Continuous Timing

Robert L. Jones
ASRC Aerospace Corporation, Greenbelt, Maryland

November 2000

The NASA STI Program Office ... in Profile

Since its founding, NASA has been dedicated to the advancement of aeronautics and space science. The NASA Scientific and Technical Information (STI) Program Office plays a key part in helping NASA maintain this important role.

The NASA STI Program Office is operated by Langley Research Center, the lead center for NASA's scientific and technical information. The NASA STI Program Office provides access to the NASA STI Database, the largest collection of aeronautical and space science STI in the world. The Program Office is also NASA's institutional mechanism for disseminating the results of its research and development activities. These results are published by NASA in the NASA STI Report Series, which includes the following report types:

- **TECHNICAL PUBLICATION.** Reports of completed research or a major significant phase of research that present the results of NASA programs and include extensive data or theoretical analysis. Includes compilations of significant scientific and technical data and information deemed to be of continuing reference value. NASA counterpart of peer-reviewed formal professional papers, but having less stringent limitations on manuscript length and extent of graphic presentations.
- **TECHNICAL MEMORANDUM.** Scientific and technical findings that are preliminary or of specialized interest, e.g., quick release reports, working papers, and bibliographies that contain minimal annotation. Does not contain extensive analysis.
- **CONTRACTOR REPORT.** Scientific and technical findings by NASA-sponsored contractors and grantees.

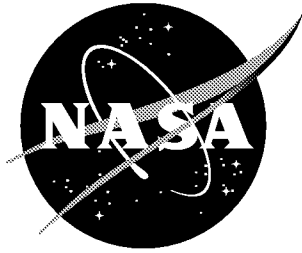
- **CONFERENCE PUBLICATION.** Collected papers from scientific and technical conferences, symposia, seminars, or other meetings sponsored or co-sponsored by NASA.
- **SPECIAL PUBLICATION.** Scientific, technical, or historical information from NASA programs, projects, and missions, often concerned with subjects having substantial public interest.
- **TECHNICAL TRANSLATION.** English-language translations of foreign scientific and technical material pertinent to NASA's mission.

Specialized services that complement the STI Program Office's diverse offerings include creating custom thesauri, building customized databases, organizing and publishing research results ... even providing videos.

For more information about the NASA STI Program Office, see the following:

- Access the NASA STI Program Home Page at <http://www.sti.nasa.gov>
- E-mail your question via the Internet to help@sti.nasa.gov
- Fax your question to the NASA STI Help Desk at (301) 621-0134
- Phone the NASA STI Help Desk at (301) 621-0390
- Write to:
NASA STI Help Desk
NASA Center for Aerospace Information
7121 Standard Drive
Hanover, MD 21076-1320

NASA / CR-2000-210296



Analysis of Phase-Type Stochastic Petri Nets With Discrete and Continuous Timing

Robert L. Jones
ASRC Aerospace Corporation, Greenbelt, Maryland

National Aeronautics and
Space Administration

Langley Research Center
Hampton, Virginia 23681-2199

Prepared for Langley Research Center
under Contract NAS1-99124

November 2000

Available from:

NASA Center for AeroSpace Information (CASI)
7121 Standard Drive
Hanover, MD 21076-1320
(301) 621-0390

National Technical Information Service (NTIS)
5285 Port Royal Road
Springfield, VA 22161-2171
(703) 605-6000

Contents

1	Introduction	4
1.1	Motivation	4
1.2	Objective	5
1.3	Organization and Assumptions	5
1.4	Notation	6
2	Background	7
2.1	Petri Nets	7
2.2	Markov Models	12
2.2.1	Discrete-Time Markov Chains	13
2.2.2	Continuous-Time Markov Chains	15
2.3	Phase-Type Models	21
2.4	Semi-Markov Models	25
2.5	Semi-Regenerative Models	28
2.6	Generalized Semi-Markov Models	35
3	Proposed Research	39
4	Preliminary Research	41
4.1	Analyzing the Underlying Stochastic Process	41
4.1.1	Theory Applied to PDPNs in General	44
4.1.2	Isochronous PDPNs	46
4.1.3	Synchronous PDPNs	48
4.1.4	Mixed PDPNs	49
4.1.5	Asynchronous PDPNs	51
4.2	PDPN Stationary Solution Algorithm	53
4.2.1	Embedding with Elimination	56
4.2.2	Conversion Matrix	64
4.2.3	Presentation of the Algorithm	65
4.3	Complexity Analysis	73
4.4	Summary of Current Accomplishments	81
5	Future Work	82
	References	83

List of Figures

2.1	Example Petri net model.	8
2.2	Reachability graph of example PN model.	10
2.3	Example continuous-time phase-type (PH) random variables.	21
2.4	Example discrete-time phase-type (DPH) random variables with step τ	21
2.5	Markovianized process of the example SPN of Figure 2.1.	24
2.6	Transition probabilities for a portion of Figure 2.5.	25
2.7	Examples of possible Petri net and stochastic confusion.	26
2.8	Markov renewal process.	27
2.9	Semi-regenerative process sample path.	28
2.10	Semi-regenerative process for generally distributed t_3 and all others Expo.	33
2.11	Two-Level Semi-regenerative process.	34
2.12	GSMP where t_1, t_2 are generally distributed and all others Expo.	36
4.1	Markovianized process of the example PDPN.	42
4.2	Example non-Markovian process when $t_2 \sim \text{Geom}(q, 3)$ and $t_3 \sim \text{Const}(2)$	43
4.3	Studying a PDPN regeneration period.	44
4.4	(a) Characterization and (b) sample path of the Isochronous PDPN.	47
4.5	(a) Characterization and (b) sample path of the Synchronous PDPN.	48
4.6	(a) Characterization and (b) sample path of the Mixed PDPN.	50
4.7	(a) Characterization and (b) sample path of the Asynchronous PDPN.	52
4.8	Example model.	59
4.9	Single-step EMC model.	59
4.10	Reducing the EMC.	60
4.11	Multi-step observation example.	62
4.12	Multi-step EMC model.	63
4.13	Multi-step example with two Geom transitions.	63
4.14	Performance comparison.	73
4.15	Reduced effort for $\rho_1 = \rho_2 = 0.9$	77
4.16	Reduced effort for $\rho_1 = \rho_2 = 0.2$	78
4.17	Reduced effort for $\rho_1 = 0.9, \rho_2 = 0.8$	79
4.18	Reduced effort for $\rho_1 = 0.8, \rho_2 = 0.9$	80

List of Algorithms

2.2.1 Extended uniformization algorithm	20
4.2.1 PDPN stationary solution algorithm	66
4.2.2 Regeneration period solution for DTMC only case	67
4.2.3 Regeneration period solution for CTMC only case	67
4.2.4 General solution procedure when both DTMC and CTMC are active	69
4.2.5 Enhanced solution procedure when both DTMC and CTMC are active and age DPH transitions are not enabled.	72

Chapter 1

Introduction

1.1 Motivation

The Petri net formalism has proven its usefulness in modeling discrete-state systems that move between states in discrete or continuous time and that may be characterized as sequential or concurrent, synchronous or asynchronous, deterministic or stochastic, or any combination for that matter. Through higher levels of abstraction, a Petri net (PN) permits a compact specification of the underlying mathematical model (usually a stochastic process) that is amenable to computer analysis. But applying PN modeling to arbitrarily complex systems requires the solution of difficult problems from a computational point of view.

A common problem is the computational complexity often required to solve stochastic PN (SPN) models with realistic assumptions about the logical and timed behavior, as opposed to simple behavior that restricts the applicability to toy models. SPNs are most easily solved when firing delays of transitions are either exponentially or geometrically distributed. Then, the underlying stochastic process is a continuous-time Markov chain in the former and a discrete-time Markov chain in the latter. Efficient solution techniques for such Markov chains are well known. However, such modeling assumptions may be unrealistic for many systems, leading to inaccurate results when adopted.

A more general approach is to allow the transition firing delays to have general distributions. In their full generality, such SPNs are classified as non-Markovian and specify generalized semi-Markov processes. Unfortunately, in the absence of any restrictions the study of such nets by analytical or numerical means is so computationally expensive that simulation becomes the only practical means to a solution. With some conveniently-chosen restrictions, however, the solution to the model can be made efficient in some cases while still lending itself to useful applications. The investigation of efficient exact and approximate solutions to such classes of non-Markovian PNs is the focus of our research.

Much of our research finds its foundation in the early work of Molloy [1], Bobbio and Cumani [2], and Marsan et al. [3, 4]: Molloy introducing execution policies, the combination of deterministic and discrete-time random behavior, and the expansion of phase-type firing delays at the state-space level, and Bobbio, Marsan, et al. doing the same for continuous time. Complete specifications of the semantics of a stochastic Petri net requires consideration about how (and in what order) enabled transitions are selected to fire and what happens to the remaining firing time of other enabled transitions when another transition fires. An

execution policy formally defines such semantics by specifying the policy used to select the enabled transition that fires and the way memory is kept of the past history of the net [4].

In [1], Molloy provided a comprehensive overview of the SPN semantics and behavior as a function of the execution policies. He also showed that when the probability distributions of SPN transitions are discrete phase-type, an otherwise non-Markovian underlying process can be transformed into a homogeneous discrete-time Markov chain defined over an expanded state space. Particularly relevant to this work is that Molloy showed how transitions with deterministic firing delays can be combined with geometric firing-delay transitions with the restriction that the deterministic delays are equal to the *basic step* of the geometric distributions. These ideas were brought up-to-date by Ciardo [5] while introducing the discrete deterministic and stochastic PN (DDSPN) formalism that allows firing delays with discrete phase-type distributions, all sharing a basic step, and having an underlying discrete-time Markov chain. The notion of a “basic step” is important to our proposed research as well. The basic step period, denoted by τ , is defined as the sojourn time in each state of the underlying discrete-time Markov chain.

In [3], Marsan and Chiola introduced the deterministic and stochastic PN (DSPN), and marks the first time deterministic behavior was integrated with continuous-time random behavior. Bobbio and Cumani in [2] showed how continuous-time, phase-type firing delays can be expanded at the state-space level to form a continuous-time Markov chain. This was discussed again by Marsan et al. in [4] while providing an extensive discussion of execution policies for generally distributed firing delays.

1.2 Objective

Our research involves the formal development of a new class of non-Markovian SPNs based on phase-type firing delays in both discrete and continuous time, present simultaneously in the *same* model. We build upon the extended SPN formalism, one that includes constructs that increase its modeling power, from a logical point of view, to that of a Turing machine as well as features that provide modeling conveniences. Such modern extensions include inhibitor arcs, transition priorities, transition-enabling guards, marking-dependent arc multiplicities, and marking-dependent execution policies.

When possible, efficient and exact solution algorithms will be developed with certain restrictions; otherwise, approximate solution algorithms will be investigated. In this way, we anticipate that this new SPN formalism may also prove to be useful in many modeling problems while still affording an efficient solution.

1.3 Organization and Assumptions

Relevant background material is provided in Chapter 2, which provides the foundation for our chosen approach. Topics include discrete- and continuous-time Markov chains, semi-Markov chains, semi-regenerative processes, generalized semi-Markov processes, characteristics specific to each, and known solution methods. These are the classes of underlying stochastic processes for popular SPN formalisms used today, and the classification of these formalisms is closely related to the stochastic process which they can specify. The generality and so-

lution complexity associated with these stochastic processes determine the modeling power and efficiency, and therefore, the SPN’s applicability and practical usefulness. Therefore, by also discussing the complexity issues germane to the solution methods, the background chapter also serves to motivate the proposed approach towards our objective: developing a SPN formalism that lends itself to useful modeling applications *and* efficient numerical analysis.

Chapter 3 provides an outline of the proposed research. Preliminary research results are provided in Chapter 4, which includes the formalization of the new SPN class and analysis theory, culminating into an exact, stationary solution algorithm. Time-dependent analysis is shown to be difficult except for a special case, making a strong argument for approximate solutions, which is planned for later. The chapter ends with a comparative analysis of the new SPN formalism with other noteworthy extensions in terms of modeling power and solution complexity. The preliminary results are followed by the plan towards completing the research in Chapter 5.

1.4 Notation

For the definitions and methods that follow, we restrict ourselves to homogeneous (time invariant) models, and we assume that a race execution policy is employed; that is, the transition with the earliest firing time is selected to fire next. Nevertheless, a preselection execution policy can be modeled by our formalism since a mechanism to resolve “contemporary” events is still required. We assume that transition firing events are atomic (no time elapses) and always sequential even if the firing times are contemporary. From a modeling perspective, contemporary firings still occur at the same time. It is only that we choose to impose a sequential ordering policy so that new markings can be unambiguously determined. As for timing, we may sometimes allow the probability distribution functions associated with transition firings to be marking dependent.

Sets are denoted with calligraphic letters. Vectors and matrices (usually lower and upper case letters, respectively) are denoted with bold text and the corresponding elements are denoted with (usually subscripted) plain text. The notion of “state” for the models and stochastic processes presented herein is actually a vector, dimensioned on the set of natural numbers \mathbb{N} and sometimes paired with supplementary information, also vectors dimensioned on \mathbb{N} or the set of real numbers \mathbb{R} . But our model solutions take the form of probability distribution vectors that either provide the state-occupancy probabilities at certain times or at steady state, or provide cumulative probabilities of occupying states over intervals of time. In either case, each (vector) state $\mathbf{i} \in \mathbb{N}^n$ must be mapped to some index $i \in \mathbb{N}$ that is associated with the state’s lexicographic position in the solution vector $\mathbf{p} = [p_i] \in \mathbb{R}^{|\mathcal{S}|}$ on the complete set of states \mathcal{S} , also known as the state space. When we refer to state “ \mathbf{i} ” in bold text, we mean its vector form, and when we refer to state “ i ” in plain text, we are referring to its lexicographic index, unless otherwise stated.

Chapter 2

Background

2.1 Petri Nets

A *Petri net*, such as the one pictured in Figure 2.1, is a directed bipartite graph described by the tuple $\mathbf{PN} = (\mathcal{P}, \mathcal{T}, \mathcal{A}, A^-, A^+, A^\circ, g, \prec, \mathbf{m}_0)$ with finite vertex sets \mathcal{P} (*places*) and \mathcal{T} (*transitions*) and a finite set of arcs (drawn as directed line segments), $\mathcal{A} \subseteq \mathcal{P} \times \mathcal{T} \cup \mathcal{T} \times \mathcal{P}$. Places (drawn as circles) can contain an integer number of *tokens* (drawn as dots or denoted by a number inside the place). We denote the marking of the net by a row vector $\mathbf{m} \in \mathbb{N}^{|\mathcal{P}|}$ that contains as entries the number of tokens in each place. Hence, m_p denotes the number of tokens in place p of marking $\mathbf{m} = [m_1, m_2, \dots, m_{|\mathcal{P}|}]$. The vector \mathbf{m}_0 denotes the initial marking.

Markings can be altered when enabling rules are satisfied at transitions (drawn as rectangles), permitting one or more transitions to *fire*, thereby removing tokens from input places and depositing them to output places according to the connecting arc multiplicities. *Arc multiplicities* are defined on arcs as either nonnegative integer constants or marking dependent functions that return a nonnegative integer; the semantics depend on the type of arc. For *input arcs*, the multiplicity, denoted by $A_{tp}^-(\mathbf{m})$, specifies the minimum number of tokens needed in place p before transition t can become enabled in marking \mathbf{m} ; this number of tokens is then removed if the transition is indeed chosen to fire. Special input arcs, called *inhibitor arcs* (drawn as directed lines with a circle at the end instead of an arrow), have a complementary effect on the enabling of transitions. For inhibitor arcs, the multiplicity, denoted by $A_{tp}^\circ(\mathbf{m})$, is the minimum number of tokens needed in place p to disable transition t in marking \mathbf{m} . For *output arcs*, the multiplicity, denoted by $A_{tp}^+(\mathbf{m})$, specifies the number of tokens that will be deposited in place p when transition t fires in marking \mathbf{m} .

Guards, denoted by g and defined for transitions, are functions $\mathbb{N}^{|\mathcal{P}|} \rightarrow \{true, false\}$ that conveniently specify additional firing rules on transitions. Given some marking \mathbf{m} , $g_t(\mathbf{m})$ must return *true* to enable transition t .

The \mathbf{PN} component $\prec \subset \mathcal{T} \times \mathcal{T}$ specifies an acyclic, preselection, priority relation, which can resolve conflicts between competing transitions attempting to fire.

Either inhibitor arcs, guards, priorities, or marking-dependent multiplicities alone increases the modeling power of a PN to that of a Turing machine (so we can represent any computational model), and hence are referred to as Turing extensions [6, 7]. Including all four Turing extensions merely provides additional modeling conveniences since the modeling

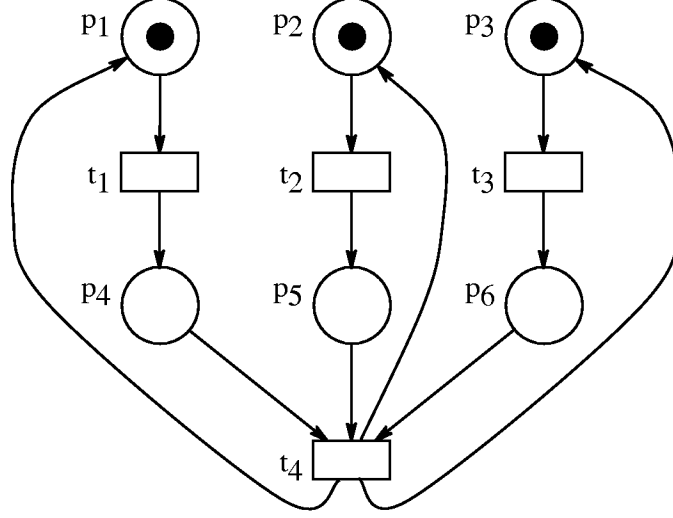


Figure 2.1: Example Petri net model.

power can no longer increase.

A transition $t \in \mathcal{F}(\mathbf{m})$, the *set of enabled transitions* in marking \mathbf{m} , if all of the following hold:

1. $g_t(\mathbf{m}) = \text{true}$
2. all of its input places p contain at least as many tokens as the corresponding input arc multiplicity $A_{tp}^-(\mathbf{m})$:

$$\forall p \in \mathcal{P}, m_p \geq A_{tp}^-(\mathbf{m})$$

3. all of its inhibitor arc places p contain fewer tokens than the arc multiplicity $A_{tp}^\circ(\mathbf{m})$:

$$\forall p \in \mathcal{P}, m_p < A_{tp}^\circ(\mathbf{m})$$

4. no other transition with higher priority \prec is enabled:

$$\forall u \in \mathcal{T}, u \not\prec t \text{ or } u \notin \mathcal{F}(\mathbf{m})$$

The firing of a transition is assumed to be atomic, consuming zero time. And, timing constraints aside, the PN defined above can evolve through markings originating from the initial marking by firing enabled transitions in any order. A transition $t \in \mathcal{F}(\mathbf{m}_0)$ can fire thereby changing the marking, $\mathbf{m}_0 \xrightarrow{t} \mathbf{m}_1$, where \mathbf{m}_1 is obtained by consuming tokens from input places and depositing tokens to output places according to the input and output arc multiplicities $A_{t\bullet}^-$ and $A_{t\bullet}^+$, respectively. By treating the $A_{t\bullet}^-$ and $A_{t\bullet}^+$ as vectors, we can write an equation for the next marking as

$$\begin{aligned} \mathbf{m}_1 &= \mathbf{m}_0 + A_{t\bullet}^+(\mathbf{m}_0) - A_{t\bullet}^-(\mathbf{m}_0) \\ &= \mathbf{m}_0 + \mathbf{1}_t \mathbf{A}(\mathbf{m}_0) \end{aligned}$$

where $\mathbf{A}(\mathbf{m}) = \mathbf{A}^+(\mathbf{m}) - \mathbf{A}^-(\mathbf{m})$ is called the *incidence matrix* and $\mathbf{1}_t$ is a unit vector with a 1 at the t^{th} position and 0 everywhere else. It is convenient to extend this next-marking computation to one that takes a sequence $s \in \mathcal{T}^*$ of transition firings as input, where \mathcal{T}^* denotes the set of transition sequences obtained by concatenating zero or more transitions from \mathcal{T} . We do this by defining the *next-marking* function $M : \mathcal{T}^* \times \mathbb{N}^{|\mathcal{P}|} \rightarrow \mathbb{N}^{|\mathcal{P}|}$ to operate on a transition firing sequence $s = (t_1, t_2, \dots, t_n) \in \mathcal{T}^n$ and a marking \mathbf{m} and return a new marking. With $\epsilon = ()$ denoting the *null* sequence, the next marking function M is defined recursively as

$$\begin{aligned} M(\epsilon, \mathbf{m}) &= \mathbf{m}, \\ M((t_1, t_2, \dots, t_n), \mathbf{m}) &= M((t_2, \dots, t_n), \mathbf{m} + \mathbf{1}_{t_1} \mathbf{A}(\mathbf{m})) \quad \text{if } t_1 \in \mathcal{F}(\mathbf{m}), \end{aligned}$$

and is undefined otherwise.

The PN behavior characterized by the set of markings reachable from the initial marking and the transition firings that cause the net to enter one marking from another can be represented as a directed graph with vertices corresponding to markings and arcs corresponding to the firing of transitions, completely constructed using M . Such a graph is called the *reachability graph*. The *reachability set*, \mathcal{R} , the set of reachability graph vertices, is the set of all markings reachable by a sequence of transition firings starting from the initial marking \mathbf{m}_0 :

$$\mathcal{R} = \{ \mathbf{m} : \exists s \in \mathcal{T}^*, \mathbf{m} = M(s, \mathbf{m}_0) \}$$

The reachability graph of the example PN model is portrayed in Figure 2.2 assuming for the moment that the net is untimed. The possible state space is subject to the number of tokens that can reside in each place p_1, p_2, \dots, p_6 and the possible sequence of transition firings that move the net between markings starting from the initial marking $(m_1 m_2 \dots m_6) = (111000)$. Without timing constraints, transitions t_1, t_2 , and t_3 are concurrent with each other and each can fire asynchronously. However, synchronization is imposed after these three transitions fire before transition t_4 can become enabled and fire, returning the net to the initial marking. Because all transitions have a fair chance of firing, the reachability graph contains all possible markings and all possible transition firing sequences.

Petri nets as defined are useful in the study of many types of systems, with or without concurrency, with or without synchronization. But without the inclusion of time, we are limited to the qualitative analysis of properties like liveness, deadlock, boundedness, and invariants [8]. To broaden the applicability of PNs, the notion of time has been incorporated into the Petri net by various researchers with various generalities by requiring that an enabled transition *delay* some amount of time before firing. Ultimately, the specification captured by the Petri net must be transformed into an (underlying) mathematical model that can be solved to obtain quantitative measures. When the *firing delays* are specified as random variables (or even if deterministic but contemporary transition firings are allowed), the underlying model is a stochastic process, the solution of which governs the overall complexity of the model solution. As one would expect, the tractability of the solution decreases as the generality of the model increases.

Extended Petri nets with the most convenient stochastic models include those with geometrically distributed (Geom) firing delays [1] having an underlying discrete-time Markov

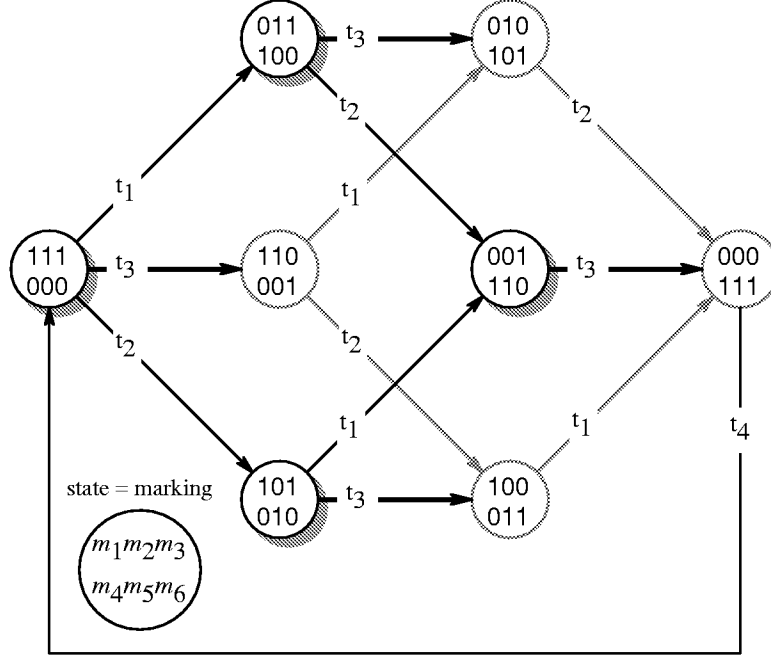


Figure 2.2: Reachability graph of example PN model.

chain and those with exponentially distributed (Expo) firing delays [9] having an underlying continuous-time Markov chain. These Markovian extensions have proven useful in the years for studying discrete-event systems with random behavior. But with the usefulness of these models to more complex and realistic systems in question, more recent extensions have tried to incorporate non-Markovian behavior. Noteworthy extensions and associated underlying processes are the phase-type SPN with an underlying, expanded continuous-time Markov chain [2], the extended SPN (ESPN) [10] with an underlying semi-Markov chain, the deterministic and stochastic PN (DSPN) [3] and the Markov regenerative SPN (MRSPN) [11, 12, 13] with an underlying semi-regenerative process, and the discrete deterministic and stochastic PN (DDSPN) [5] with an underlying, expanded discrete-time Markov chain. However, the complexity of solving more general underlying stochastic processes oftentimes limits in practice their usefulness to problems with small dimensions.

Consider now the SPN where the states reachable from the initial state are subject to the possible state transitions under the constraints imposed by the *timed* execution. Sometimes the reachability graph of the PN is isomorphic to the underlying stochastic process that models its timed execution. It is important to know when this property holds because it determines when and how the reachability graph of markings and the underlying stochastic (marking) process can be constructed. When they are isomorphic, there is a one-to-one correspondence between states in the reachability set \mathcal{R} and the state space of the stochastic process, denoted by \mathcal{S} . This is the case for SPNs with Geom or Expo firing delays, which have underlying Markov chains. In such cases, it may be convenient to build the reachability graph first and then construct the (matrix) equations for the stochastic process from it second [14]. Also when they are isomorphic, one may wish to perform the qualitative analysis by operating

on the incidence matrix or reachability graph separately from studying the stochastic marking process [8]. If we know that isomorphism is not guaranteed or does not exist, then \mathcal{S} may be a strict subset of \mathcal{R} . Consequently, the results from reasoning about the reachability set \mathcal{R} independently from the timing specification is less meaningful, and possibly misleading.

The semantics of a SPN model also depends on the chosen execution policy. The execution policy specifies two things: 1) how the next transition to fire is selected among those enabled and 2) how memory is kept regarding the *remaining firing time* (RFT), or similarly the “age”, of transitions with non-memoryless distributions.

In regard to selecting transitions to fire, the policy most frequently assumed is the *race* policy where the transition selected to fire is the one with the minimum remaining firing time over all enabled transitions. But it is also possible to perform the selection on the basis of additional specifications that do not depend on the duration of the activities associated with the transitions that are enabled. One such policy is called *preselection* where transitions are selected to fire among the enabled set according to *a priori* information, independent of the firing delay distributions. For example, a probability mass function (pmf) can be defined over the set of enabled transitions in a given marking, and used to choose the transition that fires next. This preselection can be done globally, defined for all markings of the net, which implies serialization of all activities. Alternatively, the preselection can be done locally, defined over transition groupings, not necessarily disjoint, within which a preselection policy is applied. Local preselection can be done in concert with the race policy as follows. In a marking that enables transitions belonging to these groups, the next transition to fire is identified by selecting first, with time independent criteria, an enabled transition (if one exists) from each of the groups, and then by choosing among the preselected transitions the one whose firing delay is minimal [4].

We allow a combination of race policy with pre- and post-selection in terms of priorities. Under the race policy, we select among the enabled transition whose sampled firing delay is (statistically) the shortest. This policy provides very useful models of systems that exhibit concurrency where multiple activities compete such that the first to finish determines the change in system state. We assume that *immediate transitions*, those that fire in zero time, have a higher priority of firing over timed transitions. Thus, we implicitly employ a preselection policy between timed and immediate transitions. Other execution policies like pre- or post-selection among timed transitions are discussed at length in [4] and may be required as well to resolve confusion, which is discussed later.

In regard to how memory is kept about the age, or RFT, of transitions, the memory policy is only meaningful for transitions with non-memoryless firing-delay distributions since these are the only transition that can “age”. Transitions that have memoryless distributions (the Expo in continuous time and the Geom in discrete time, the Const(0) is a special case of the Geom) are not affected since their firing delays can equivalently be sampled after every transition firing. Memory policies are as important to the semantics of non-Markovian SPNs as the net topology and the policy used to select the next transition to fire. There are three policies that are most-frequently used in modeling applications, namely, *resampling*, *enabling memory*, and *age memory* [1, 4]. The chosen memory policy need not be global; a different policy can be associated with different transitions and be marking dependent [15]. A resampling policy requires that transitions obtain a new firing delay, sampled from the respective distribution functions, after some transition fires, including itself. Since each

transition firing causes a state change, a resampling policy for all transitions results in a semi-Markov process, which enjoys an absence of memory after each state change. Such a policy (with race execution) is useful for modeling competing activities (modeled of course by transitions where the amount of work performed is represented by the firing delay) in which the next state of the system is decided by the activity that finishes first. Consequently, the work performed by the losing transitions is lost. Alternatively, an enabling-memory policy causes the firing delays to be resampled only when a transition becomes enabled again after being disabled. The enabling memory policy is useful in modeling activities where work is performed until either completion or termination by another event causing the work performed to be lost. Finally, the age-memory policy causes the firing delay to be resampled only after the transition itself fires, even though it may have been disabled and re-enabled many times. Thus, the work performed by age-memory activities is never lost, which makes such transitions useful in modeling tasks that can be preempted and then resumed at the same point.

2.2 Markov Models

When the SPN firing delays are defined by random variables, the SPN provides a compact specification of an underlying stochastic process. Hereafter, we will denote this stochastic process as $\{X(\theta) : \theta \geq 0\}$: a collection of random variables defined over the same probability space, indexed by a time parameter θ , and taking on values in a state space \mathcal{S} , which may be finite or infinite as well as continuous or discrete. A stochastic process that has a discrete state space is called a *chain*. The index (time) parameter θ can also be continuous or discrete. Hereafter, we will denote discrete-time processes by $\{X_\theta : \theta \in \mathbb{N}\}$.

A stochastic process (chain) $\{X(\theta) : \theta \geq 0\}$ with the property

$$\begin{aligned} \Pr\{X(v + \theta) = j \mid X(v) = i, X(\zeta) = x(\zeta), 0 \leq \zeta < v\} = \\ \Pr\{X(v + \theta) = j \mid X(v) = i\} \end{aligned}$$

$\forall i, j, x(\zeta) \in \mathcal{S}, \theta \geq 0, v \geq 0$, is called a *Markov process (chain)* and the property is referred to as the “memoryless” or *Markov property*. Thus, the determination of the state the process will transition to next depends solely on the current state and not the past history of the process. When the value of that conditional probability is independent of v , the process is said to be *homogeneous* or time invariant; this is one of our assumptions. Also, because of the discrete nature of the PN markings, we will mostly concern ourselves with discrete-state processes or chains. The acronyms DTMC and CTMC are used for discrete-time and continuous-time Markov chains, respectively. There will be occasions, presented later, when the “state” is supplemented with additional, continuous-valued information for modeling purposes thereby giving rise to a continuous-state process.

The amount of time that a process spends in a state is referred to as its *sojourn time*. Clearly, if the evolution of a Markov process depends only on the current state, it should not matter how long the process remains in the current state before making a transition. Thus, the sojourn time is geometrically distributed (Geom) for DTMCs and exponentially distributed (Expo) for CTMCs. This is expected since the Geom and Expo random variables are the only ones that exhibit the “memoryless property” for discrete and continuous random

variables, respectively. Without limiting ourselves to only Expo and Geom firing delays, the efficient analysis of SPNs using DTMCs and CTMCs will be the main focus of our research.

It follows from the memoryless property of Markov chains, letting

$$P_{ij}(\theta) = \Pr\{X(\theta) = j \mid X(0) = i\},$$

that

$$P_{ij}(\theta + v) = \sum_{k \in \mathcal{S}} P_{ik}(\theta) P_{kj}(v).$$

This is known as the Chapman-Kolmogorov equation for Markov chains and is key in formulating the analytical solutions of Markovian models.

2.2.1 Discrete-Time Markov Chains

Consider a DTMC defined by its transition matrix $\mathbf{\Pi} = [\Pi_{ij}]$, $i, j \in \mathcal{S}$, where

$$\Pi_{ij} = \Pr\{X_1 = j \mid X_0 = i\}$$

gives the conditional transition probabilities between states in one step or *jump*. It is often the case with DTMC models that the time spent in each state is of no concern, only the states that can be occupied after a given number of “jumps” is of interest. But we can also imagine that the DTMC remains in each state a fixed amount of time τ , referred to as its *basic step* time.

A fundamental property of DTMCs is that the Chapman-Kolmogorov equation takes the form

$$P_{ij}(\theta) = [\mathbf{\Pi}^\theta]_{ij}$$

where $\theta \in \mathbb{N}$. However, we do not have to, nor would we want to, compute the matrix $\mathbf{\Pi}^\theta$. Instead, we can compute the unconditional probability vector

$$\mathbf{x}^{(\theta)} = [\Pr\{X_\theta = i\}] = \mathbf{x}^{(0)} \mathbf{\Pi}^\theta$$

iteratively using the recursive relation

$$\mathbf{x}^{(\theta)} = \mathbf{x}^{(\theta-1)} \mathbf{\Pi}$$

known as the *power method* where $\mathbf{x}^{(0)} = [\Pr\{X_0 = i\}] \in \mathbb{R}^{|\mathcal{S}|}$ is given by the initial probability distribution.

We can also compute the cumulative probability vector $\mathbf{y}^{(\theta)} = [y_j^{(\theta)}] = \int_0^\theta \mathbf{x}^{(\lfloor u \rfloor)} du$, defined as

$$y_j^{(\theta)} = \mathbb{E}[\text{number of visits to } j \in \mathcal{S} \text{ until time } \theta \mid X_0 = i]$$

with an *extended power method*:

$$\begin{aligned} \mathbf{y}^{(n)} &= \mathbf{y}^{(n-1)} + \mathbf{x}^{(n-1)} \\ \mathbf{x}^{(n)} &= \mathbf{x}^{(n-1)} \mathbf{\Pi} \end{aligned}$$

for $n \leftarrow 1$ to θ with initial condition $\mathbf{y}^{(0)} = \mathbf{0}$ (the vector of all zeros).

If the DTMC contains transient states that lead to strongly-connected recurrent states (including absorbing states), then we can partition the state space into \mathcal{S}_T and \mathcal{S}_R , the set of transient and recurrent states, respectively, such that $\mathcal{S} = \mathcal{S}_T \cup \mathcal{S}_R$. Then, by defining a new matrix $\tilde{\mathbf{\Pi}}$ by restricting $\mathbf{\Pi}$ to states in \mathcal{S}_T only, we can compute the cumulative probability vector $\tilde{\mathbf{y}} = [\tilde{y}_j] \in \mathbb{R}^{|\mathcal{S}_T|}$ defined as

$$\tilde{y}_j = \lim_{\theta \rightarrow \infty} y_j^{(\theta)} = \mathbb{E}[\text{number of visits to } j \in \mathcal{S}_T \text{ until absorption} \mid X_0 = i]$$

with the same extended power method except that matrix $\tilde{\mathbf{\Pi}}$ is used instead of $\mathbf{\Pi}$, $\tilde{\mathbf{x}} \in \mathbb{R}^{|\mathcal{S}_T|}$ is used, and we stop when probability mass remaining in \mathcal{S}_T , determined from the vector norm $\|\tilde{\mathbf{x}}\|_1$, becomes small enough.

The power method can also be used to compute the stationary or steady-state solution $\mathbf{x} = [x_i]$, $i \in \mathcal{S}$, of DTMCs where $x_i = \lim_{\theta \rightarrow \infty} \Pr\{X_\theta = i\}$, by iterating long enough for the sequence $\{x^{(n)}\}_{n=0}^{\infty}$ to converge to \mathbf{x} .

For the limiting measures where $\theta \rightarrow \infty$, convergence utilizing the power method may take a long time, making for a poor method in practice. Alternatively, we can observe that a stationary solution \mathbf{x} , if it exists, would satisfy the equation $\mathbf{x}\mathbf{\Pi} = \mathbf{x}$. This is just the case for ergodic DTMCs (those that are irreducible, aperiodic, and positive recurrent). Fortunately, ergodic DTMCs have a *unique* stationary solution that satisfies the system of equations

$$\mathbf{x}\mathbf{\Pi} = \mathbf{x} \quad \text{subject to} \quad \sum_{i \in \mathcal{S}} x_i = 1. \quad (2.1)$$

Of course, we could utilize direct methods like standard Gaussian elimination or LU decomposition to solve for \mathbf{x} . But because the coefficient matrix based on $\mathbf{\Pi}$ would be modified and susceptible to fill-in and because $\mathbf{\Pi}$ is typically very large for realistic models, direct methods are rarely used in practice. Even though the DTMC specified by $\mathbf{\Pi}$ (and most any stochastic model for that matter) is *large*, it is at least *sparse* in general. So iterative methods like Gauss-Seidel and successive overrelaxation (SOR), which have much faster convergence than the iterative power method, are usually employed to compute \mathbf{x} . These methods do not modify the iteration matrix based on $\mathbf{\Pi}$ and if sparse matrix storage is used, the time complexity is $O(N\eta)$ where N is the number of iterations needed for convergence and η is the number of nonzero entries in matrix $\mathbf{\Pi}$. Unfortunately, N may be unbounded since these iterative methods do not guarantee convergence for all initial guesses for \mathbf{x} . This is because $\mathbf{\Pi}$ is a stochastic matrix (each row sum is one), which makes the spectral radius (the largest eigenvalue) equal to one. Iterative methods have guaranteed convergence for any initial guess only when the spectral radius is less than one [16].

With these iterative methods, we can also solve the following system of linear equations for the cumulative probability vector $\tilde{\mathbf{y}}$ when \mathcal{S} contains transient states:

$$\tilde{\mathbf{y}} = \tilde{\mathbf{x}}^{(0)} + \tilde{\mathbf{y}}\tilde{\mathbf{\Pi}}$$

or equivalently

$$\tilde{\mathbf{y}}(\mathbf{I} - \tilde{\mathbf{\Pi}}) = \tilde{\mathbf{x}}^{(0)}. \quad (2.2)$$

While the use of iterative methods to solve Equation 2.1 may not converge for all initial guesses, convergence is instead guaranteed for Equation 2.2, since $(\mathbf{I} - \tilde{\mathbf{\Pi}})$ is an M-matrix, (nonsingular, elements are less than or equal to zero, and having a nonnegative inverse [16]), which always has a spectral radius less than one. Therefore, solving Equation 2.2 with iterative methods has guaranteed convergence, regardless of the initial guess [17].

2.2.2 Continuous-Time Markov Chains

Consider now a CTMC, where for all continuous points in time, state transitions can occur and the process is memoryless. Let the rate at which the process transitions from state i to state j be denoted by $\lambda_{ij} \in \mathbb{R}^+$, $i, j \in \mathcal{S}$. The sojourn time in each state i is exponentially distributed, so by letting $\lambda_i = \sum_{j \in \mathcal{S}} \lambda_{ij}$, we can obtain the expected sojourn time from

$$\mathbb{E}[\text{sojourn time in state } i] = \lambda_i^{-1}.$$

The interpretation of the rates is such that

$$\lim_{\theta \rightarrow 0} \frac{P_{ij}(\theta)}{\theta} = \lambda_{ij}, \quad i \neq j \quad (2.3)$$

$$\lim_{\theta \rightarrow 0} \frac{1 - P_{ii}(\theta)}{\theta} = \lambda_i. \quad (2.4)$$

By observing the CTMC just after each state transition, we can construct a DTMC consisting of the possible sequence of states the process can move between over time θ . If only limiting measures where $\theta \rightarrow \infty$ are sought then we can also compute the probability of transitioning between state i and state j in the DTMC from the ratio λ_{ij}/λ_i since the sojourn times are exponentially distributed. Let $\text{diag}(\lambda_{ii})$ be matrix with λ_i along the diagonal $\forall i \in \mathcal{S}$ and zero elsewhere. Then the embedded DTMC matrix $\mathbf{\Pi} \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}|}$ constructed from the CTMC is defined as

$$\mathbf{\Pi} = \text{diag}(\lambda_i)^{-1} [\lambda_{ij}]_{i \neq j} \quad (2.5)$$

Constructing a DTMC by observing a stochastic process (a CTMC in this case) at times when the Markov property holds is called *embedding*. Then steady-state, state-occupancy measures such as the stationary distribution or time-to-absorption (TTA) can be computed from the embedded DTMC (EMC) using Equations 2.1 and 2.2, respectively. But, since the quantities \mathbf{x} and \mathbf{y} from Equations 2.1 and 2.2 for the EMC are interpreted as or are based on the “number of visits” to states, we must convert these measures back to the original CTMC by appropriately scaling them. The needed “scaling factors” come from the knowledge that, with each visit, the expected sojourn time in each state i is just λ_i^{-1} . For example, the CTMC stationary distribution $\mathbf{p} = [p_i]$, $p_i = \lim_{\theta \rightarrow \infty} \Pr\{X(\theta) = i\}$, can be computed by using the embedded DTMC matrix (2.5) and then computing

$$\mathbf{x}\mathbf{\Pi} = \mathbf{x} \quad \text{subject to} \quad \sum_{i \in \mathcal{S}} x_i = 1, \quad \tilde{x}_i = x_i \lambda_i^{-1}, \quad p_j = \frac{\tilde{x}_j}{\sum_{k \in \mathcal{S}} \tilde{x}_k}$$

where the last step is needed to re-normalize the distribution so that it sums to one once again.

Embedding a CTMC is applicable for steady-state solutions, not time-dependent solutions. With simple embedding, the CTMC is observed only at times when state transitions occur. Consequently, information concerning how long the process sojourns in states is lost making the EMC time-dependent solutions useless. For time-dependent analysis, we must once again make use of the Chapman-Kolmogorov equation. For steady-state analysis, embedding will once again be useful for more complicated stochastic process than CTMCs, as discussed in more detail in later sections.

By manipulating the Chapman-Kolmogorov equation:

$$\begin{aligned} P_{ij}(\theta + v) &= \sum_{k \in \mathcal{S}} P_{ik}(\theta) P_{kj}(v) \\ P_{ij}(\theta + v) - P_{ij}(\theta) &= \sum_{k \in \mathcal{S}} P_{ik}(\theta) P_{kj}(v) - P_{ij}(\theta) \\ P_{ij}(\theta + v) - P_{ij}(\theta) &= \sum_{k \neq i} P_{ik}(\theta) P_{kj}(v) - (1 - P_{jj}(v)) P_{ij}(\theta) \end{aligned}$$

then dividing by v and taking the limit as $v \rightarrow 0$:

$$\lim_{v \rightarrow 0} \frac{P_{ij}(\theta + v) - P_{ij}(\theta)}{v} = \lim_{v \rightarrow 0} \left\{ \frac{\sum_{k \neq i} P_{ik}(\theta) P_{kj}(v) - (1 - P_{jj}(v)) P_{ij}(\theta)}{v} \right\}$$

and finally substituting 2.3 and 2.4, we get what is known as *Kolmogorov's forward equation*:

$$\frac{d}{d\theta} P_{ij}(\theta) = \sum_{k \neq i} P_{ik}(\theta) \lambda_{kj} - P_{ij}(\theta) \lambda_j. \quad (2.6)$$

The interchange of the summation and the limit, needed to obtain Equation 2.6, is not always justified, but it does hold for most models including those with finite state spaces, as is the case here [18]. Similarly, we can derive the *Kolmogorov's backward equation*,

$$\frac{d}{d\theta} P_{ij}(\theta) = \sum_{k \neq i} \lambda_{ik} P_{kj}(\theta) - \lambda_i P_{ij}(\theta), \quad (2.7)$$

by looking backwards in time from a given state.

By defining what is called the *infinitesimal generator matrix*, $\mathbf{Q} = [Q_{ij}] \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}|}$, as

$$\mathbf{Q} = [\lambda_{ij}]_{i \neq j} - \text{diag}(\lambda_i)$$

and letting $\mathbf{P}(\theta) = [P_{ij}(\theta)] \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}|}$, Kolmogorov's forward and backward equations can be rewritten in matrix form as

$$\frac{d}{d\theta} \mathbf{P}(\theta) = \mathbf{P}(\theta) \mathbf{Q} \quad (2.8)$$

and

$$\frac{d}{d\theta} \mathbf{P}(\theta) = \mathbf{Q} \mathbf{P}(\theta) \quad (2.9)$$

respectively. These differential equations can sometimes be solved with conventional, direct or indirect means (such as Runge-Kutta or Laplace transforms) but this restricts the usefulness of CTMC models to small problems. Instead, we could make use of the known solution

$$\mathbf{P}(\theta) = e^{\mathbf{Q}\theta},$$

which is a matrix exponential computed from

$$e^{\mathbf{Q}\theta} = \sum_{n=0}^{\infty} \frac{(\mathbf{Q}\theta)^n}{n!}. \quad (2.10)$$

However, the matrix exponential method is susceptible to subtractive cancellation errors due to the positive and negative entries in \mathbf{Q} , which makes the method unstable.

In practice, the time-dependent solution of the CTMC is usually computed using Jensen's method, also known as *uniformization*. The basic idea behind uniformization is to perform time-dependent analysis on a DTMC constructed from the CTMC in a way similar to "embedding" except that all states are forced to have the same expected sojourn time by imposing self-loops on states where necessary. By uniformizing the CTMC, the CTMC is observed at times of state transitions, when they occur naturally, and at more frequent times, when self-transitions occur. Hence, information is retained about how long the CTMC occupies each state, unlike the embedding method.

The basic uniformization algorithm defines a DTMC matrix, $\mathbf{A} \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}|}$, as

$$\mathbf{A} = q^{-1}\mathbf{Q} + \mathbf{I} \quad (2.11)$$

where $q \geq \max_i Q_{ii}$ is chosen as the "sampling" rate, at least as large as the maximum outgoing rate of all CTMC states. The sampling rate q is normally chosen to be slightly *larger* than the maximum outgoing rate to ensure that the DTMC will not be periodic. Also, all \mathbf{A} entries are non-negative. So, substituting

$$\mathbf{Q} = (\mathbf{A} - \mathbf{I})q,$$

derived from 2.11, into the matrix exponential 2.10 provides an efficient and much more stable computation of $\mathbf{P}(\theta)$, for a given initial probability distribution $\mathbf{x}^{(0)} = [\Pr\{X_0 = i\}] \in \mathbb{R}^{|\mathcal{S}|}$, that makes use of the Poisson random variable:

$$\begin{aligned} \mathbf{x}^{(0)}\mathbf{P}(\theta) &= \mathbf{x}^{(0)}e^{\mathbf{Q}\theta} \\ &= \mathbf{x}^{(0)}e^{(\mathbf{A}-\mathbf{I})q\theta} \\ &= \mathbf{x}^{(0)}e^{\mathbf{A}q\theta}e^{-\mathbf{I}q\theta} \\ &= \mathbf{x}^{(0)}\sum_{n=0}^{\infty} \frac{\mathbf{A}^n(q\theta)^n}{n!} \cdot \sum_{n=0}^{\infty} \frac{\mathbf{I}^n(-q\theta)^n}{n!} \\ &= \mathbf{x}^{(0)}\sum_{n=0}^{\infty} \frac{\mathbf{A}^n(q\theta)^n}{n!} \cdot \sum_{n=0}^{\infty} \frac{(-q\theta)^n}{n!} \\ &= \mathbf{x}^{(0)}\sum_{n=0}^{\infty} \frac{(q\theta)^n e^{-q\theta}}{n!} \mathbf{A}^n. \end{aligned}$$

The summation can be easily computed with

$$\sum_{n=0}^{\infty} \mathbf{x}^{(n)} \text{Poiss}(n, q\theta) \quad (2.12)$$

where

$$\mathbf{x}^{(n)} = \mathbf{x}^{(n-1)} \mathbf{A}$$

and

$$\text{Poiss}(n, q\theta) = \text{Poiss}(n-1, q\theta) \frac{q\theta}{n}, \quad \text{Poiss}(0, q\theta) = e^{-q\theta}.$$

So we see that uniformization is another version of the power method, only extended for suitability in studying CTMCs. Essentially, uniformization subordinates the uniformized process with a Poisson birth process. This has a nice interpretation. The Poisson process determines the probability that the uniformized process (the DTMC) can make n jumps within fixed time θ . Given n jumps, \mathbf{A}^n determines the set of states that can be occupied conditioned on the initial state. By doing this for all possible n (and summing the probabilities), we can uncondition on n and obtain our desired solution for the original CTMC process.

In practice, computing $\mathbf{x}^{(0)} e^{\mathbf{Q}\theta}$ requires that we truncate the infinite series 2.12 and sum the remaining terms in a numerically stable way. To do this, we employ the Fox-Glynn algorithm [19], which defines left and right truncation points, L_1 and R_1 , respectively, so that

$$\mathbf{x}^{(0)} e^{\mathbf{Q}\theta} \approx \sum_{n=L_1}^{R_1} \mathbf{x}^{(n)} \text{Poiss}(n, q\theta)$$

and the error is bounded by 10^{-d} (d digits of precision) if

$$L_1 = \max_{k \in \mathbb{N}} \left\{ \sum_{n=0}^k \text{Poiss}(n, q\theta) \leq \frac{10^{-d}}{2} \right\}$$

$$R_1 = \min_{k \in \mathbb{N}} \left\{ 1 - \sum_{n=L_1}^k \text{Poiss}(n, q\theta) \leq 10^{-d} \right\}.$$

Note that although the Poisson probabilities are only computed in the range $L_1 \leq n \leq R_1$, the vector-matrix multiplications must be done for $0 \leq n \leq R_1$.

Just as in the power method for DTMCs, we can compute the same $\mathbf{y}^{(\theta)}$ and $\tilde{\mathbf{y}}$, defined in the previous section as

$$\mathbf{y}^{(\theta)} = \int_0^\theta \mathbf{x}^{(u)} du = \mathbf{x}^{(0)} \int_0^\theta e^{\mathbf{Q}u} du$$

$$\tilde{\mathbf{y}} = \lim_{\theta \rightarrow \infty} \mathbf{y}^{(\theta)}$$

at the same time $\mathbf{x}^{(0)}e^{\mathbf{Q}\theta}$ is computed [20]. Substitution of the uniformization computation for $e^{\mathbf{Q}u}$ yields

$$\mathbf{x}^{(0)} \int_0^\theta e^{\mathbf{Q}u} du = \int_0^\theta \left(\sum_{n=0}^{\infty} \mathbf{x}^{(n)} \cdot \text{Poiss}(n, qu) \right) du$$

and after factoring the summation series over n , we have

$$\sum_{n=0}^{\infty} \mathbf{x}^{(n)} \int_0^\theta \text{Poiss}(n, qu) du,$$

which can be written equivalently as

$$\frac{1}{q} \sum_{n=0}^{\infty} \mathbf{x}^{(n)} \left(1 - \sum_{\ell=0}^n \text{Poiss}(\ell, q\theta) \right).$$

using integration by parts. Although we must begin at $n = L_2 = 0$ here, a right truncation point R_2 can be found with bounded error. Since the total sojourn in all states over an interval $[0, \theta]$ must be θ , so that

$$\left\| \mathbf{x}^{(0)} \int_0^\theta e^{\mathbf{Q}u} du \right\|_1 = \theta,$$

we can stop when the difference between θ and $(1 - \sum_{\ell=0}^n \text{Poiss}(\ell, q\theta))$ becomes small. Thus, we have the right truncation point

$$R_2 = \min_{k \in \mathbb{N}} \left\{ \theta - \frac{1}{q} \sum_{n=0}^k \left(1 - \sum_{\ell=0}^n \text{Poiss}(\ell, q\theta) \right) \right\} \leq 10^{-d}$$

When computing both $\mathbf{x}^{(0)} \int_0^\theta e^{\mathbf{Q}u} du$ and $\mathbf{x}^{(0)}e^{\mathbf{Q}\theta}$, the smallest left truncation point $L = 0$ and largest right truncation point $R = \max(R_1, R_2)$ are chosen.

In case the CTMC with generator \mathbf{Q} is ergodic, we can also test for stationary conditions to check whether θ is large enough for the DTMC to have reached steady state. Detecting stationary conditions that occur before the right truncation is reached, and halting, can result in significant performance gains.

But if stationary solutions are sought then there is a better way just as in the DTMC case. When stationary or steady-state equilibrium is reached, the change in probability mass between states becomes zero. So we can set the derivative in Equation 2.8 to zero and obtain global balance equations for the CTMC:

$$\mathbf{x}\mathbf{Q} = \mathbf{0} \tag{2.13}$$

where \mathbf{x} satisfies the stationary solution $\lim_{\theta \rightarrow \infty} \mathbf{x}^{(\theta)}$ independent of the initial probability distribution when the CTMC is ergodic. The same argument used in the previous section against employing direct methods apply here as well; iterative methods are better even though convergence is not guaranteed for all initial guesses.

With iterative methods, we can also solve the following system of linear equations for the cumulative probability vector $\tilde{\mathbf{y}} = \lim_{\theta \rightarrow \infty} \mathbf{y}^{(\theta)}$:

$$\tilde{\mathbf{y}}\tilde{\mathbf{Q}} = -\tilde{\mathbf{x}}^{(0)}. \quad (2.14)$$

Like in Equation 2.2 for the DTMC case, the use of iterative methods for Equation 2.2 enjoys guaranteed convergence since $-\tilde{\mathbf{Q}}$ is also an M-matrix. Because the use of Equations 2.13 and 2.14 has the same complexity as Equations 2.1 and 2.2, respectively, steady-state solutions are computed directly from the CTMC in practice.

Because the uniformization algorithm is integral to the solution algorithms developed later, we provide it here as Algorithm 2.2.1, without steady-state detection. The algorithm computes the transient probability vector $\boldsymbol{\pi}^{(\theta)} \in \mathbb{R}^{|\mathcal{S}_i|}$ on the CTMC state space \mathcal{S}_i originating from the initial state i where $\pi_j = \Pr\{X(\theta) = j \mid X(0) = i\}$ and the cumulative probability vector $\boldsymbol{\sigma} = \int_0^\theta \boldsymbol{\pi}^{(u)} du$. The algorithm assumes that $\boldsymbol{\pi}$ contains the initial probability distribution when the algorithm is invoked.

Algorithm 2.2.1 Extended uniformization algorithm

- 1: Given the solution time θ , generator $\mathbf{Q} \in \mathbb{R}^{|\mathcal{S}_i| \times |\mathcal{S}_i|}$, and initial probability vector $\boldsymbol{\pi} \in \mathbb{R}^{|\mathcal{S}_i|}$,
 - 2: Let $q \leftarrow 1.02 \cdot \max_i |Q_{ii}|$ and $\mathbf{A} \leftarrow q^{-1}\mathbf{Q} + \mathbf{I}$
 - 3: $\hat{\mathbf{x}} \leftarrow \boldsymbol{\pi}$
 - 4: $\boldsymbol{\pi} \leftarrow \mathbf{0}$
 - 5: $\boldsymbol{\sigma} \leftarrow \mathbf{0}$
 - 6: $s \leftarrow 1$
 - 7: Choose L and R for desired precision 10^{-d}
 - 8: Compute $\text{Poiss}(n), \forall n, L \leq n \leq R$ using Fox-Glynn algorithm
 - 9: **for** $n \leftarrow 0$ to $L - 1$ **do**
 - 10: $\boldsymbol{\sigma} \leftarrow \boldsymbol{\sigma} + \hat{\mathbf{x}}$
 - 11: $\hat{\mathbf{x}} \leftarrow \hat{\mathbf{x}}\mathbf{A}$
 - 12: **end for**
 - 13: **for** $n \leftarrow L$ to R **do**
 - 14: $s \leftarrow s - \text{Poiss}(n)$
 - 15: $\boldsymbol{\pi} \leftarrow \boldsymbol{\pi} + \text{Poiss}(n) \cdot \hat{\mathbf{x}}$
 - 16: $\boldsymbol{\sigma} \leftarrow \boldsymbol{\sigma} + s \cdot \hat{\mathbf{x}}$
 - 17: $\hat{\mathbf{x}} \leftarrow \hat{\mathbf{x}}\mathbf{A}$
 - 18: **end for**
 - 19: $\boldsymbol{\sigma} \leftarrow \boldsymbol{\sigma}/q$
 - 20: **return** solutions $\boldsymbol{\pi}, \boldsymbol{\sigma}$
-

2.3 Phase-Type Models

Phase-type random variables are defined as the time-to-absorption of Markov chains with at least one absorbing state. An absorbing CTMC (via a rate matrix) is used for *continuous phase-type distributions* (let PH denote this family hereafter) and an absorbing DTMC (via a stochastic matrix) is used for *discrete phase-type distributions* (denoted hereafter by DPH). Each must also include a specification of the initial state occupancy probabilities that the absorbing Markov chain assumes when a new random variable is “sampled”. Special cases of PH, as shown in Figure 2.3, include: exponential (Expo), Erlang, hyper-exponential (Hyper), and hypo-exponential (Hypo). Special cases of DPH, as shown in Figure 2.4, include: geometric (Geom), constant integer multiples of τ (Const), and discrete uniform (Equiprob).

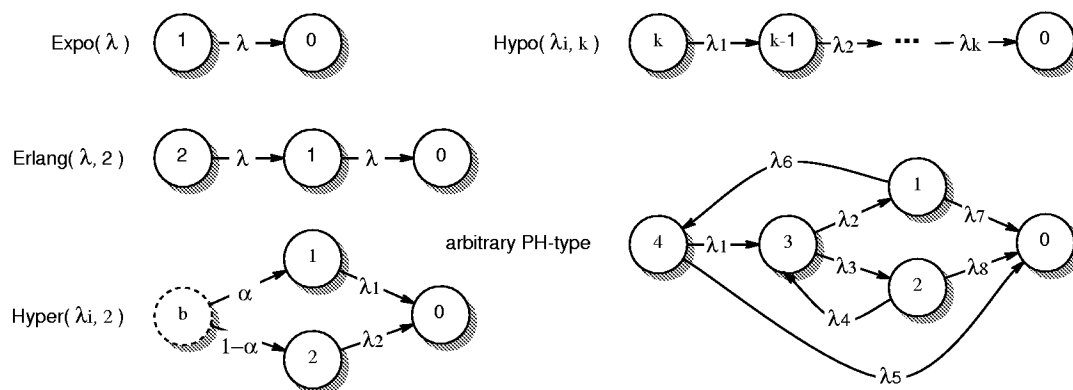


Figure 2.3: Example continuous-time phase-type (PH) random variables.

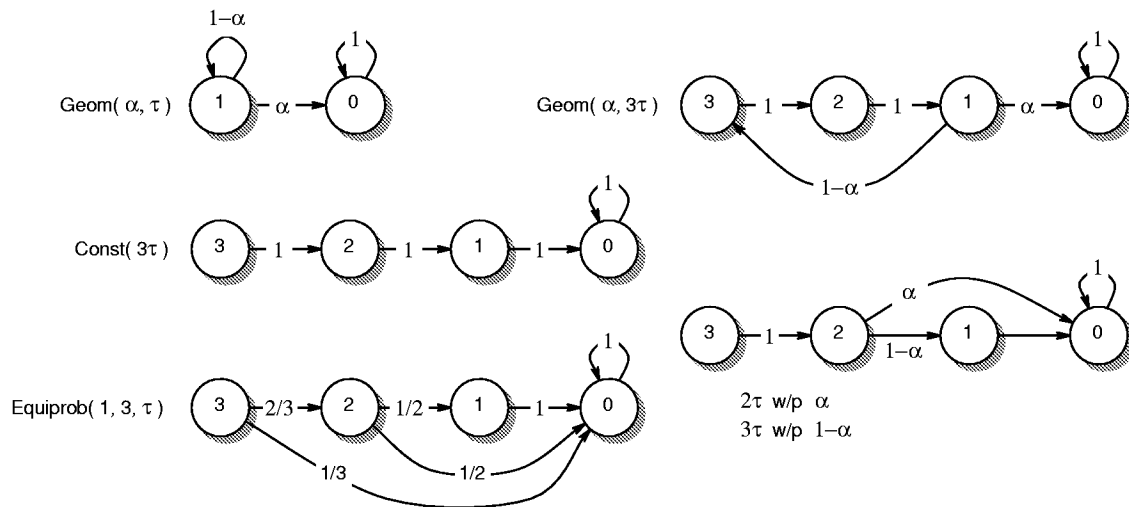


Figure 2.4: Example discrete-time phase-type (DPH) random variables with step τ .

The reachability graph and corresponding state space of a SPN with PH or DPH firing delays is constructed by *expanding* the state space and state transitions so as to remember the RFT of each enabled transition until one or more transitions can fire. The RFT for phase-type firing delays is naturally discretized, and so the pairing of each possible phase vector ϕ together with each possible, discrete PN marking vector \mathbf{m} produces a state (\mathbf{m}, ϕ) within an expanded Markov chain. By including enough information in the current state about the past evolution of the process, the past can be forgotten, essentially *Markovianizing* the process. Because the idea of supplemented states is important to our research, we will soon revisit this topic in greater detail. After firing a transition t , the execution policies that define what happens to the RFT of transitions still enabled are then applied to create a new phase vector ϕ' paired with the new marking $M(t, \mathbf{m})$.

If the transition t that fires is once again enabled in the new marking, then a new firing delay is sampled from its PH or DPH distribution function (an absorbing Markov chain) and included in ϕ' of the new state. If the firing transition t is not enabled in the marking, then its phase component, ϕ'_t , within vector ϕ' is unspecified, a “don’t care”. This procedure continues until no new state is discovered or until no transitions can fire (this event results in an absorbing state in the reachability graph).

For all phase-type transitions t and all reachable markings \mathbf{m} , the possible combination of phases can be obtained by performing the Cartesian product of the *phase space*, denoted by \mathcal{D}^t , of each absorbing Markov chain that specifies the phase-type firing delay: a stochastic matrix $\mathbf{D}^t(\mathbf{m})$ for DPH phases or a rate matrix $\mathbf{E}^t(\mathbf{m})$ for PH phases. We allow the firing delays to be marking dependent by letting these matrices be functions of the marking. For DPH phases, the result is a *Cartesian product* of the constituent phase spaces and an *arithmetic product* of all corresponding one-step probabilities. For PH phases, the result is a similar Cartesian product and an *arithmetic sum* of all corresponding rates. Hereafter, we let \mathcal{D} denote the potential phase space in the expanded model. When discrete-time and continuous-time models are considered separately, an expanded DTMC results from DPH phases and an expanded CTMC results from PH phases.

For DPH models, the expanded DTMC states and transition probabilities are specified formally using the *Kronecker multiplication*. Letting $\mathbf{A} \in \mathbb{R}^{r \times r}$ and $\mathbf{B} \in \mathbb{R}^{s \times s}$, the *Kronecker product* \otimes is defined as

$$\mathbf{A} \otimes \mathbf{B} = \begin{bmatrix} a_{11}\mathbf{B} & a_{12}\mathbf{B} & \cdots & a_{1r}\mathbf{B} \\ a_{21}\mathbf{B} & a_{22}\mathbf{B} & \cdots & a_{2r}\mathbf{B} \\ \vdots & \vdots & \ddots & \vdots \\ a_{r1}\mathbf{B} & a_{r2}\mathbf{B} & \cdots & a_{rr}\mathbf{B} \end{bmatrix}$$

the result being of order rs [21]. Then, for each DPH transition t , the expanded DTMC on the total, *potential* phase space \mathcal{D} is given by

$$\mathbf{D}(\mathbf{m}) = \bigotimes_t \mathbf{D}^t(\mathbf{m}) \tag{2.15}$$

with order $\prod_t |\mathcal{D}^t|$. The stochastic matrix $\mathbf{D}(\mathbf{m})$ completely specifies all conditional next phases and associated probabilities that can occur in one basic step. By referring to a

“potential” phase space, we emphasize that not all combinations of discrete phases given by the matrix equation 2.15 is reachable.

For PH models, the expanded CTMC states and transition rates are specified formally using the *Kronecker addition*. Using \otimes and the identity matrices \mathbf{I}_r and \mathbf{I}_s of order r and s , respectively, to obtain the correct dimension for summing \mathbf{A} and \mathbf{B} , the *Kronecker sum* \oplus is defined as

$$\mathbf{A} \oplus \mathbf{B} = \mathbf{A} \otimes \mathbf{I}_s + \mathbf{I}_r \otimes \mathbf{B}$$

Then, for each PH transition t , the expanded CTMC on the total, potential phase space \mathcal{D} is given by

$$\mathbf{E}(\mathbf{m}) = \bigoplus_t \mathbf{E}^t(\mathbf{m}) \quad (2.16)$$

with order $\prod_t |\mathcal{D}^t|$. The rate matrix $\mathbf{E}(\mathbf{m})$ completely specifies all conditional next phases and associated transition rates that can occur in continuous time.

An example reachability graph isomorphic to an expanded DTMC is shown in Figure 2.5 assuming DPH transition timing: $t_1 \sim \text{Geom}(p, 1)$, $t_2 \sim \text{Geom}(q, 3)$, $t_3 \sim \text{Const}(2)$, and $t_4 \sim \text{Const}(1)$. The states are composed of both marking and RFT information corresponding to the discrete, firing delay phases of each enabled transition. Phases of transitions not enabled are of no importance and consequently are indicated by “•” symbols.

Note that the state transitions associated with “ ϵ ” are those where no transition actually fires, but merely update the phase information. The one-step transition probabilities originating from state $(110001, 11 \bullet \bullet)$ are shown in Figure 2.6. Here we see the presence of the simultaneous firing of t_1 and t_2 between states $(110001, 11 \bullet \bullet)$ and $(000111, \bullet \bullet \bullet 1)$. Unlike the continuous-time PH models where the probability of any two transition having the same firing time is zero, such contemporary firings are possible, indeed likely, with DPH models.

Contemporary firings not only have the potential of making the reachability graph more dense (more state transitions) than one with PH transitions, but can create *confusion*. Two transition are said to be *concurrent* when each can fire and the firing of one does not affect the firing of the other. Two transitions are said to be in *conflict* when the firing of one prevents the firing of the other. When we have both concurrency and conflict, we *may* have confusion [8].

Consider the examples of possible confusion in Figure 2.7 where contemporary firings are possible and where there exists a mix of concurrency and conflict. In both examples, transitions a and c are concurrent, and transitions b and c are in conflict. In the bottom example, there is also conflict between transitions a and b . Even though contemporary firings are possible, a firing sequence must be chosen and applied to the function M to determine the next marking. PN confusion can occur when the next marking depends on the order in which transition firings are chosen. When considering the marking

$$(m_1 m_2 m_3 m_4 m_5) = (10100)$$

shown in the top example, we only see the concurrency between transitions a and c , but if we choose to fire transition a first, and move to the next marking

$$(m_1 m_2 m_3 m_4 m_5) = (01100),$$

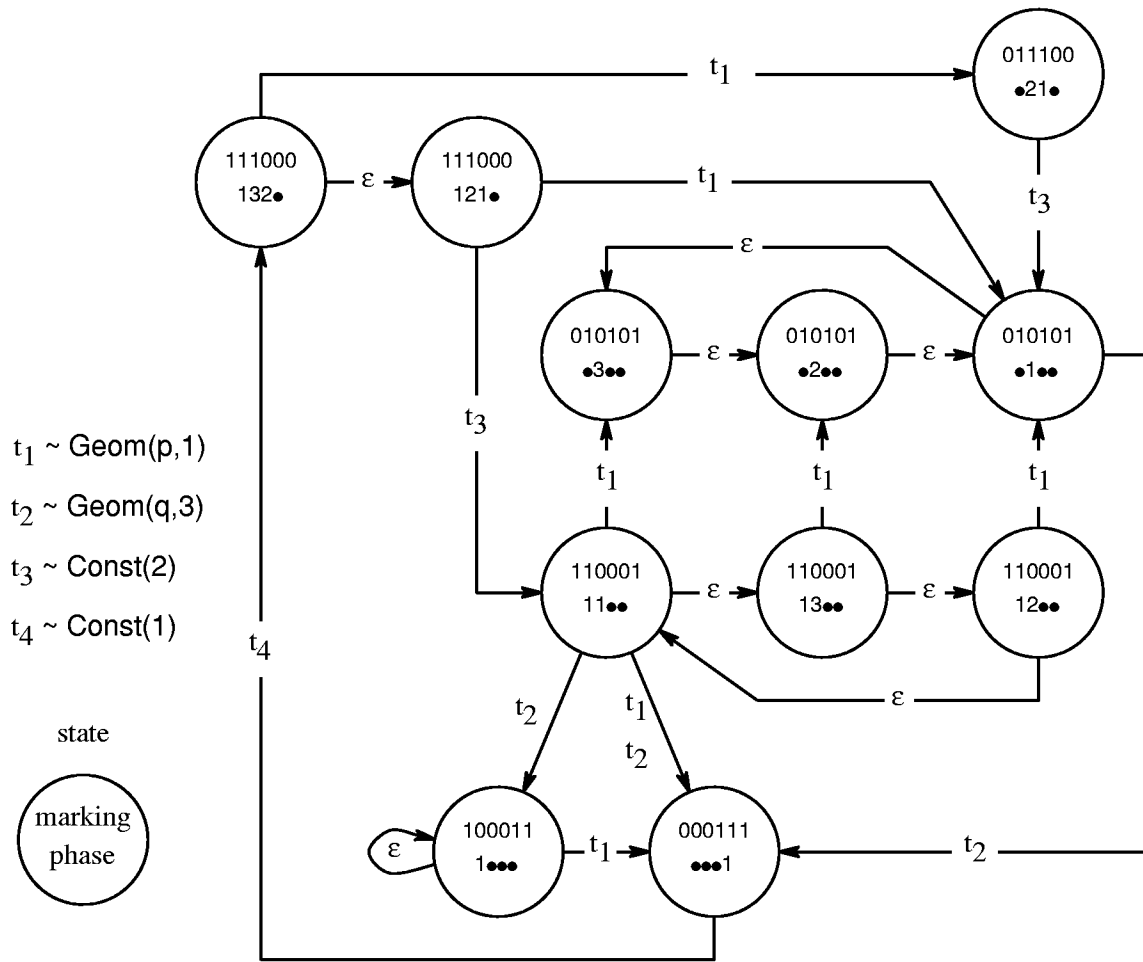


Figure 2.5: Markovianized process of the example SPN of Figure 2.1.

we encounter a conflict between transitions b and c . If we instead choose to fire transition c first, and move to the next marking

$$(m_1 m_2 m_3 m_4 m_5) = (10001),$$

then transition b is prevented from firing whether or not transition a fires next in the same marking. Thus, different outcomes arise depending on which transition, a or c , is chosen first as shown by the reachability graph in the top-right corner of Figure 2.7. The confusion about which transition to fire first must be resolved by either avoiding the confusion in the first place, by preventing the bothersome enabling of transitions, or by forcing a particular firing sequence. Either fix can be accomplished by employing a different net specification, guards, or preselection priority. Preselection priorities can resolve conflicts in untimed PNs or within timed PNs for immediate transitions that can fire in zero time. Using preselection priorities for example, different probabilities can be assigned to each sequence leading to the three new states, such that all probabilities sum to one. This results in a stochastic model even though the PN may be untimed.

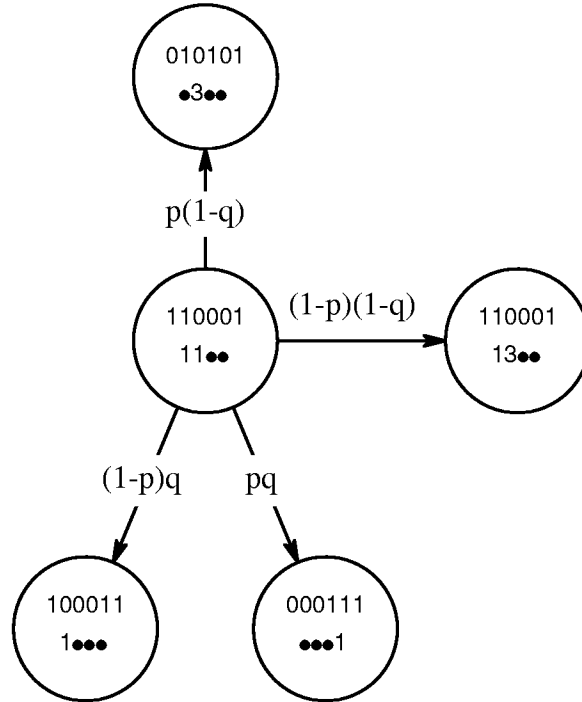


Figure 2.6: Transition probabilities for a portion of Figure 2.5.

When analyzing timed PNs, we may also have stochastic confusion. In the bottom example of Figure 2.7, PN confusion does not exist since the net reaches the same marking no matter which sequence is chosen. However, if impulse rewards* $\gamma : \mathcal{T} \rightarrow \mathbb{R}$ are used such that $\gamma(a) + \gamma(c) \neq \gamma(b)$ then stochastic confusion results since the reward measure depends on the chosen sequence. If the stochastic outcome differs depending on the order in which transitions in a contemporary firing sequence are selected to fire, then the model is not “well defined”. Such confusion may be resolved the same way as with PN net confusion or with the addition of *postselection* priorities. Rather than preventing the simultaneous enabling of transitions that may lead to confusion, we may instead leave them be, let them delay, and if their simultaneous firing leads to confusion, decide then which transition gets to fire first using postselection priorities. The understanding of this problem and possible solutions, including the use of postselection priorities, have already been discussed in [22, 23, 24] and can be brought to bear on the research proposed here.

2.4 Semi-Markov Models

If the sojourn times in states are Expo or Geom random variables, we have a Markov chain: a CTMC in the former, which is memoryless for all time, and a DTMC in the latter, which is memoryless at times multiple of the basic step. Alternatively, if the sojourn times in states are *all* equal to the *same* constant $\tau \in \mathbb{N}$, we still have a DTMC. From the previous

*Each time a state transition occurs due to the firing of t , a reward of $\gamma(t)$ is accumulated.

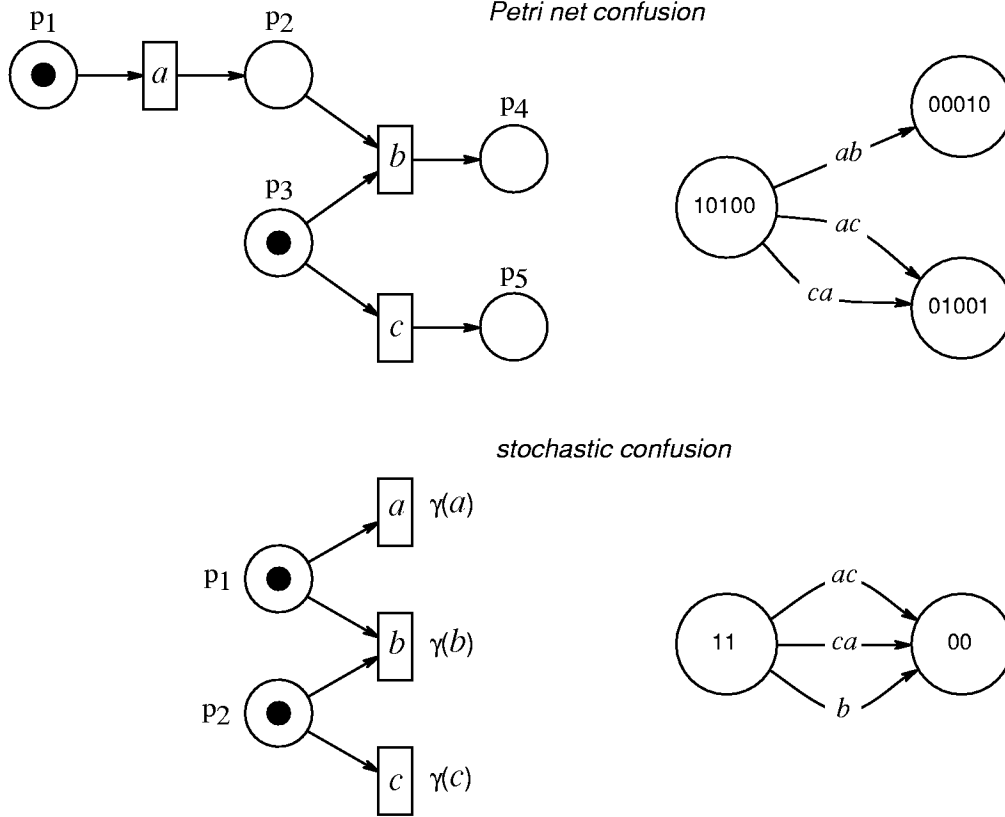


Figure 2.7: Examples of possible Petri net and stochastic confusion.

section, we know that although such models are restricted, they lend themselves to efficient time-dependent and stationary analysis.

In more general cases, semi-Markov processes satisfy the Markov property at times of jumps when state transitions occur, but not necessarily between jumps. Consequently, the sojourn times for semi-Markov processes can be arbitrary, nonmemoryless random variables. Satisfying the Markov property in at least a “semi” way affords efficient, stationary analysis just as strict Markov chains, but time-dependent analysis becomes difficult. The following theory used to study semi-Markov chains is provided below in preparation for the more general theory needed to study semi-regenerative processes, of which Markov and semi-Markov chains are special cases.

Consider a random variable X_n , defined for each $n \in \mathbb{N}$ and taking values from the state space \mathcal{E} , and a random variable T_n , likewise defined for each n but taking values in \mathbb{R}^+ such that $T_0 = 0$ and $T_n \leq T_{n+1}$, $\forall n \in \mathbb{N}$. The process $\{(X_n, T_n) : n \in \mathbb{N}\}$ is called a *Markov renewal process* (MRP) with state space \mathcal{E} if the following holds $\forall n \in \mathbb{N}$, $\forall j \in \mathcal{E}$:

$$\Pr\{X_{n+1} = j, T_{n+1} - T_n \leq \theta \mid X_0, X_1, \dots, X_n; T_0, T_1, \dots, T_n\} = \Pr\{X_{n+1} = j, T_{n+1} - T_n \leq \theta \mid X_n\} = \Pr\{X_1 = j, T_1 \leq \theta \mid X_0\}$$

The sequence $\{X_n : n \in \mathbb{N}\}$ is a DTMC. An example MRP sample path is portrayed in Figure 2.8.

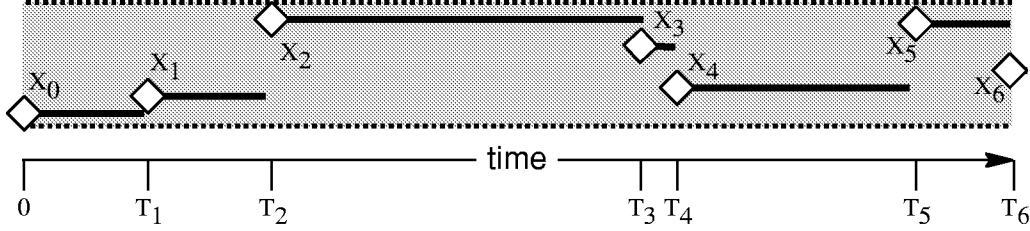


Figure 2.8: Markov renewal process.

Consider a stochastic process $\{X(\theta) : \theta \geq 0\}$ with state space \mathcal{S} that has an *embedded* MRP with state space $\mathcal{E} \subseteq \mathcal{S}$. That is, observing $X(\theta)$ at certain, random times T_n and recording the state X_n occupied at those times produces a MRP. Analogous to the CTMC embedding discussed in the previous section, if we can determine the transition probability matrix \mathbf{P} of the embedded DTMC (EMC) and its stationary solution $x_i = \lim_{n \rightarrow \infty} \Pr\{X_n = i\}$, $i \in \mathcal{E}$, then we can easily compute the stationary distribution $p_j = \lim_{\theta \rightarrow \infty} \Pr\{X(\theta) = j\}$, $j \in \mathcal{S}$, of the complete process. For a semi-Markov process, one that enjoys a *renewal* after every state transition, as portrayed in Figure 2.8, we have $\mathcal{S} = \mathcal{E}$ and $E[T_1 | X_0 = i] = E[\text{sojourn in } i]$, $i \in \mathcal{E}$, and

$$\tilde{x}_i = x_i \cdot E[\text{sojourn in } i], \quad p_j = \frac{\tilde{x}_j}{\sum_{k \in \mathcal{S}} \tilde{x}_k}. \quad (2.17)$$

This well known technique of “embedding” is based on the following reasoning. The stationary probability distribution can be interpreted as the fraction of time the process resides in each state. For there to be a unique stationary solution, the EMC must be ergodic; i.e., it is aperiodic, positive recurrent, and irreducible. The aperiodic property ensures that we can compute an unique stationary solution, given that the other two properties also hold. The positive recurrent property means that the process after leaving some state will eventually return to the same state in some finite time. The irreducible property means that every state can reach every other state. So to determine the expected cycle time of the stationary process, we need only pick a single reference state. Let state k be this reference state. Then after determining the stationary solution of the EMC, $x_i, \forall i \in \mathcal{E}$, we can interpret the ratio x_i/x_k as the expected number of visits to state i between two visits to state k . The expected sojourn time in state i within a stationary cycle is just $E[\text{sojourn in } i] \cdot x_i/x_k$ and the expected cycle time is just $\sum_{j \in \mathcal{E}} E[\text{sojourn in } j] \cdot x_j/x_k$. Equations 2.17 then follows from the interpretation that the stationary probability distribution is the fraction of time the process resides in each state within the expected cycle time.

The embedding method makes intuitive sense as well when we think of it as just scaling the stationary probability distribution of the EMC according to the expected sojourn times and then normalizing so that the new distribution sums to one. This idea of “scaling” is useful in understanding the theory that follows.

2.5 Semi-Regenerative Models

For a *semi-regenerative process*, one that becomes a probabilistic replica of itself at certain random times T_n , given the same state X_n , the stationary solution is computed in a way similar to semi-Markov processes except that $\mathbf{\Pi}$ and $E[\text{sojourn in } k \text{ during } [0, T_1) \mid X_0 = i]$, $i \in \mathcal{E}$, $k \in \mathcal{S}_i$, must be computed by studying the *subordinate process*, the process with state space \mathcal{S}_i that evolves between renewals [25]:

$$\tilde{x}_k = \sum_{i \in \mathcal{E}} x_i \cdot E[\text{sojourn in } k \text{ during } [0, T_1) \mid X_0 = i], \quad p_j = \frac{\tilde{x}_j}{\sum_{\ell \in \mathcal{S}} \tilde{x}_\ell}. \quad (2.18)$$

An example sample path of a semi-regenerative process is portrayed in Figure 2.9.

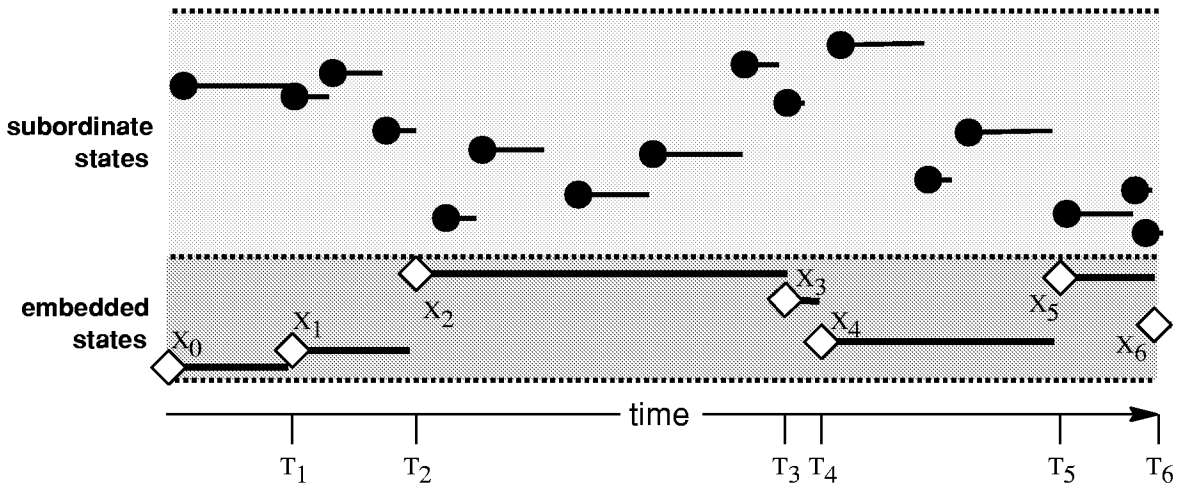


Figure 2.9: Semi-regenerative process sample path.

The problem of solving a semi-regenerative process using Markov renewal theory is reduced to studying the stochastic process between the T_n points in time when the process enjoys an absence of memory concerning its past. Unlike the semi-Markov process, the semi-regenerative process can occupy many states (within the subordinate process) between renewals, referred to hereafter as *regeneration* times. So the embedding technique used for semi-Markov processes by “scaling” is applicable to semi-regenerative processes except that we proportionally redistribute the EMC stationary distribution over all the states visited in the subordinate process between regenerations, followed by normalization. For semi-regenerative processes, the distribution, scaling, and normalization required to compute the stationary solution from the EMC is referred to hereafter as *conversion*. We will return to the subject of embedding a semi-regenerative process after we present the key aspects of Markov renewal theory, which is needed to compute the EMC and the conversion factors.

In general, Markov renewal theory depends on the specification of the following two quantities:

$$\begin{aligned} G_{ij}(\theta) &= \Pr\{X_{n+1} = j, T_{n+1} - T_n \leq \theta \mid X_n = i\} \\ &= \Pr\{X_1 = j, T_1 \leq \theta \mid X_0 = i\}, \\ H_{ik}(\theta) &= \Pr\{X(\theta) = k, T_1 > \theta \mid X_0 = i\} \quad i, j \in \mathcal{E}, k \in \mathcal{S}_i \end{aligned}$$

where the elimination of n is justified by assuming homogeneity. $H_{ik}(\theta)$ gives the transient probability of occupying state $k \in \mathcal{S}_i \subseteq \mathcal{S}$ at time θ before the next regeneration given the initial state $i \in \mathcal{E}$ entered at the last regeneration. The subordinate process evolves during the interval $[T_n, T_{n+1})$ or equivalently $[0, T_1)$ when $X_0 = X_n$. $G_{ij}(\theta)$ gives the state transition probability of the EMC between two consecutive regenerations jointly with the distribution of the regeneration period T_1 .

The quantities G and H can be recursively combined similar to the Chapman-Kolmogorov equation to obtain

$$P_{ij}(\theta) = H_{ij}(\theta) + \sum_{k \in \mathcal{E}} \int_0^\theta dG_{ik}(v) P_{kj}(\theta - v) \quad i, j \in \mathcal{S} \quad (2.19)$$

known as the *Markov renewal equation*. Markov renewal theory is essentially the application of this equation to aid the study of semi-regenerative processes.

If \mathcal{E} is finite, the Markov renewal equation is satisfied by the unique solution

$$P_{ij}(\theta) = \sum_{k \in \mathcal{E}} \int_0^\theta dR_{ik}(v) H_{kj}(\theta - v), \quad (2.20)$$

where R is the *Markov renewal function* [25]. The Markov renewal function $R_{ij}(\theta)$ is defined as the expected number of renewals observed at a fixed state $j \in \mathcal{E}$ starting from state $i \in \mathcal{E}$ within a fixed interval $[0, \theta]$:

$$R_{ij}(\theta) = \sum_{n=0}^{\infty} \Pr\{X_n = j, T_n \leq \theta \mid X_0 = i\} = \sum_{n=0}^{\infty} G_{ij}^n(\theta).$$

where $G_{ij}^n(\theta)$ is the n -fold convolution of $G_{ij}(\theta)$ with itself.

Finding the transient probability distribution $P_{ij}(\theta)$ that satisfies Equations 2.19 or 2.20 is not a trivial task in general. For models with large dimensions, direct solution in the time domain is expensive and is susceptible to numerical difficulties or instabilities. Alternatively, Equation 2.19 can be solved as a linear system of equations in the s -domain by utilizing the Laplace-Stieltjes transform [26]. However, the numerical inversion necessary to obtain the time-domain solution afterwards is also complex if numerical instabilities are to be avoided. So it would seem that the difficulties in obtaining an exact, time-dependent solution of semi-regenerative processes makes a good case for either reasonable, simplifying restrictions that make R easier to compute, or approximations [27].

Fortunately for stationary analysis, the method of embedding indirectly produces a solution that satisfies

$$\lim_{\theta \rightarrow \infty} P_{ij}(\theta) = \lim_{\theta \rightarrow \infty} \sum_{k \in \mathcal{E}} \int_0^\theta dR_{ik}(v) H_{kj}(\theta - v)$$

for any initial state $i \in \mathcal{S}$ (assumed to be an embedded state at time 0) with much less difficulty than solving Equation 2.19 or 2.20 directly [25]. To do this, we determine the EMC transition matrix from

$$H_{ij} = \lim_{\theta \rightarrow \infty} G_{ij}(\theta) \quad i, j \in \mathcal{E} \quad (2.21)$$

where G itself is determined from H by observing the subordinate process at regeneration times. At the same time, the conversion factors, denoted hereafter by $\mathbf{h} = [h_{ik}] \in \mathbb{R}^{|\mathcal{E}| \times |\mathcal{S}|}$, are computed $\forall i \in \mathcal{E}, k \in \mathcal{S}_i$, from

$$h_{ik} = \mathbb{E}[\text{sojourn in } k \text{ during } [0, T_1) \mid X_0 = i] = \int_0^\infty H_{ik}(\theta) d\theta \quad (2.22)$$

from which we can also compute the expected value of a typical regeneration period

$$\mathbb{E}[T_1 \mid X_0 = i] = \sum_{k \in \mathcal{S}_i} h_{ik}.$$

Consequently, the complexity of the method is dictated by the complexity of studying the subordinate process and the EMC. The solution complexity of these two subproblems depends to a large degree on the regeneration points that are sampled—which ones and how many. The PDPN solution algorithms we propose in Chapter 4 are developed with these considerations in mind.

We now focus on the application of Markov renewal theory to the deterministic and stochastic PN (DSPN) [3] and its generalization, the Markov regenerative SPN (MRSPN) [12]. We assume for clarity that the probability distribution function or PDF (also known as the cumulative distribution function or CDF) for the firing delay of a generally distributed transition t ,

$$F^t(\theta) = \Pr\{\text{transition } t \text{ firing delay} \leq \theta\},$$

is not marking dependent and that t cannot be preempted by the firing of another transition. While marking dependent PDFs and preemption policies have been addressed, in [28, 15, 29, 30] for example, such situations introduce an unnecessary complication to the following discussion. Although we assume marking independent PDFs for our work, we do allow certain kinds of preemption and so this topic is discussed later in Chapter 4.

Starting with $G_{ij}(\theta)$ and conditioning on the events $\{T_1 = v\}$ and $\{X(v) = k\}$, we can equivalently write, for $i, j \in \mathcal{E}$,

$$G_{ij}(\theta) = \sum_{k \in \mathcal{S}_0} \int_0^\theta \Pr\{X_1 = j \mid X(v) = k\} \Pr\{X(v) = k \mid X_0 = i\} d\Pr\{T_1 \leq v \mid X_0 = i\}. \quad (2.23)$$

For MRSPNs, which are restricted so that at most one generally distributed transition is enabled in any marking, and with our assumptions, regeneration points are observed at times when generally distributed transitions either become enabled or fire. When only Expo transitions are enabled, regeneration points are observed just after each state transition as a

result of Expo transition firings. This means that the subordinate process is a general CTMC at worst, when a generally distributed transition is enabled, and a single state CTMC at best, when no generally distributed transitions are enabled.

Consider a regeneration period starting with a known embedded state $i \in \mathcal{E}$. For the simple case when only Expo transitions are enabled, let λ_i^{-1} denote the expected sojourn time in state i and let λ_{ij} be the rate that the exponentially distributed process transitions between state i and some next state j . Then Π_{ij} is given in closed form as $\lambda_{ij}\lambda_i^{-1}$.

For the more complex case, when a generally distributed transition t is enabled, we can substitute the PDF for transition t , $F^t(v)$, in place of $\Pr\{T_1 \leq v \mid X_0 = i\}$ and recognize that $\Pr\{X(v) = k \mid X_0 = i\}$ is just the solution of the subordinate CTMC at time v when t may fire. The quantity $\Pr\{X_1 = j \mid X(v) = k\}$ is simply the probability of entering marking j when t fires in marking k , possibly followed by a firing sequence of immediate transitions, all occurring in zero time. We denote this *switching probability* with Δ_{kj}^t . If no immediate transition firing can occur, $\Delta_{kj}^t = 1$ if $k = j$ and 0 otherwise. These substitutions into Equation 2.23 yields

$$G_{ij}(\theta) = \sum_{k \in \mathcal{S}_0} \int_0^\theta [e^{\mathbf{Q}_i v}]_{ik} dF^t(v) \Delta_{kj}^t \quad i, j \in \mathcal{E} \quad (2.24)$$

for the MRSPN where \mathbf{Q}_i denotes the CTMC generator with state space \mathcal{S}_i and initial state $i \in \mathcal{E}$.

When only exponential transitions are enabled in state $i \in \mathcal{E}$, we know that the expected sojourn time in i is λ_i^{-1} from Equation 2.4, and so $H_{ii}(\theta) = e^{-\lambda_i \theta}$. For the more complex case, we can rewrite $H_{ik}(\theta)$ equivalently as

$$\begin{aligned} H_{ik}(\theta) &= \Pr\{X(\theta) = k, T_1 > \theta \mid X_0 = i\} \\ &= \Pr\{X(\theta) = k \mid T_1 > \theta, X_0 = i\} \Pr\{T_1 > \theta \mid X_0 = i\} \quad i \in \mathcal{E}, k \in \mathcal{S}_i \end{aligned}$$

and again for the MRSPN when a generally distributed transition t is enabled in state $i \in \mathcal{E}$, we can substitute $\Pr\{X(\theta) = k \mid T_1 > \theta, X_0 = i\}$ with the transient solution of the subordinate CTMC at time θ and obtain $\Pr\{T_1 > \theta \mid X_0 = i\}$ from the PDF of t to get

$$H_{ik}(\theta) = [e^{\mathbf{Q}_i \theta}]_{ik} (1 - F^t(\theta)). \quad (2.25)$$

The DSPN is a special case of the MRSPN where the non-exponential transitions t have deterministic firing delays, specified by $\text{Const}(\tau_t)$. With the substitution $F(v)^t = \text{Const}(\tau_t)$ in Equations 2.24 and 2.25 we can obtain the Equations 2.26 and 2.27 for the EMC $\mathbf{\Pi} = [\Pi_{ij}]$ and conversion factors $\mathbf{h} = [h_{ik}]$, respectively:

$$\begin{aligned} \Pi_{ij} &= \lim_{\theta \rightarrow \infty} G_{ij}(\theta) \\ &= \sum_{k \in \mathcal{S}_0} \int_0^\infty [e^{\mathbf{Q}_i v}]_{ik} dF^t(v) \Delta_{kj}^t \\ &= \sum_{k \in \mathcal{S}_0} \int_0^\infty [e^{\mathbf{Q}_i v}]_{ik} \delta(\tau_t) d\tau_t \Delta_{kj}^t \\ &= \sum_{k \in \mathcal{S}} [e^{\mathbf{Q}_i \tau_t}]_{ik} \Delta_{kj}^t \end{aligned} \quad (2.26)$$

and

$$\begin{aligned}
h_{ik} &= \int_0^\infty H_{ik}(\theta) d\theta \\
&= \int_0^\infty [e^{\mathbf{Q}_i\theta}]_{ik} (1 - F^t(\theta)) d\theta \\
&= \int_0^\infty [e^{\mathbf{Q}_i\theta}]_{ik} (1 - I(\theta - \tau_t)) d\theta \\
&= \int_0^{\tau_t} [e^{\mathbf{Q}_i\theta}]_{ik} d\theta
\end{aligned} \tag{2.27}$$

where $\delta(\cdot)$ is the unit impulse function and $I(\cdot)$ is the unit step function. We will find Equations 2.26 and 2.27 useful in the analysis of the non-Markovian SPN proposed in the remaining chapters.

We end this section by providing in Figure 2.10 the underlying semi-regenerative process for our running example when all transitions but t_3 are exponentially distributed. The semi-regenerative process *graph* shown here is the same in structure, but not timing, of course, whether we consider the model to be a MRSPN or a DSPN by assuming transition t_3 to be generally distributed or deterministic, respectively. States in the EMC are shadowed to distinguish them from the states in the subordinate Markov chain that evolve due to the firing of transitions t_1 and t_2 while t_3 is enabled. Once t_3 has fired, all states visited are considered embedded states until t_3 becomes enabled once again. The initial state (111000) is considered to be both an embedded state and subordinate state by definition. Note that the semi-regenerative process graph is isomorphic to the PN reachability graph given in Figure 2.2. This is not happenstance. This is always the case for DSPNs or MRSPNs when at most one deterministic or general transition, respectively, is enabled in any marking.

A *sufficient* condition for isomorphism between the PN reachability set \mathcal{R} and the underlying stochastic state space \mathcal{S} is if all firing delays have continuous distributions with infinite support $[0, \infty)$. But this is not a necessary condition as shown by the fact that DSPNs, with finite-support, deterministic firing delays, have reachability sets that *are* isomorphic to the underlying stochastic state space. The necessary conditions for isomorphism ensure that the timing specification allow every enabled transition to have a chance of firing. That is, $\forall \mathbf{m} \in \mathcal{R}$:

1. The quantity $\Pr\{t \text{ fires, firing delay} \leq \theta \mid \text{history}\}$, at the present, is uniquely determinable by observing the past evolution constituting the *history*,
2. $\sum_{t \in \mathcal{F}(\mathbf{m})} \lim_{\theta \rightarrow \infty} \Pr\{t \text{ fires, firing delay} \leq \theta \mid \text{history}\} = 1$, and
3. $\forall t \in \mathcal{F}(\mathbf{m}), \Pr\{t \text{ fires} \mid \text{history}\} = \lim_{\theta \rightarrow \infty} \Pr\{t \text{ fires, firing delay} \leq \theta \mid \text{history}\} > 0$

are true [4]. Simply put, (1) and (2) ensure unambiguous determination of future states given the past even when timing constraints are imposed, and (3) ensures that every transition has an *opportunity* to fire in every marking that enables it. When these criteria are met, the determination of \mathcal{R} is independent of the firing-delay distributions. These conditions are satisfied by DSPNs and MRSPNs since at most one deterministic or general, respectively, transition t is enabled in any given marking. Because all other enabled (Expo) transitions

have continuous distribution functions with infinite support starting at 0, transition t is able to fire before any other enabled transition, and the other transitions are able to fire before transition t . Hence, every enabled transition has an opportunity to fire.

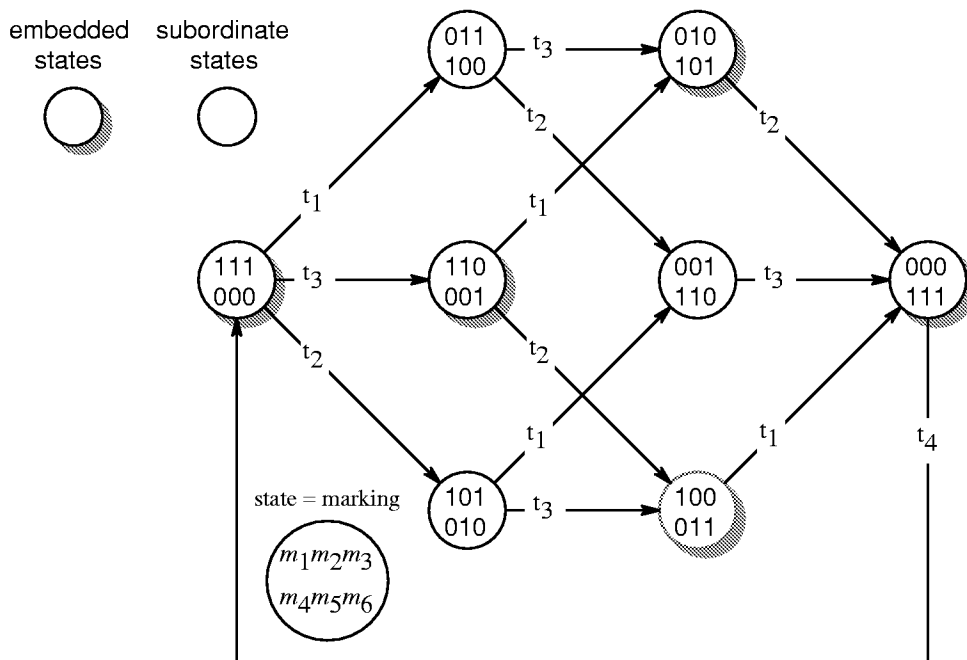


Figure 2.10: Semi-regenerative process for generally distributed t_3 and all others Expo.

The complexity of studying semi-regenerative processes increases with the complexity of the subordinate process. Models with at most one generally distributed transition enabled are convenient because they restrict the subordinate process to a CTMC. But when we allow the simultaneous enabling of multiple, generally distributed transitions, the subordinate process is more complicated. For example, the MRSPNs was extended in [29] by allowing multiple general transitions to be simultaneously enabled, provided that only one of the general transitions defines the next regeneration point. But the subordinate process becomes a semi-Markov chain, which is more difficult to solve in the transient than a CTMC. With unrestricted models in general, the subordinate process may be a semi-Markov chain or worse, perhaps even a semi-regenerative process by itself.

For example, consider our running example model with two generally distributed and two exponentially distributed transitions: $t_1 \sim F^{t_1}(\theta)$, $t_2 \sim F^{t_2}(\theta)$, $t_3 \sim \text{Expo}(\lambda)$, and $t_4 \sim \text{Expo}(\mu)$. Since t_1 and t_2 are simultaneously enabled in markings (111000) and (110001) , we cannot be sure, in general, that the underlying process is semi-regenerative. But if, for instance, the PDFs were such that t_1 *always delays longer than* t_2 then, because they both become enabled simultaneously, we know that once transition t_1 fires, everything about the past since t_1 and t_2 became enabled can be forgotten. Because the stochastic marking process regenerates itself at this point in time, and depends on the state reached when t_1 fires, the process is semi-regenerative.

Starting in the initial state (111000) , assumed to be a regeneration point, the next regeneration point would coincide with the firing of transition t_1 and the subordinate process

that evolves in between (due to the firing of t_2 and t_3) is itself a semi-regenerative process. So we have a two-level semi-regenerative process hierarchy—regeneration when t_2 fires at level 2 and regeneration when t_1 fires at level 1. This underlying process is portrayed in Figure 2.11 where the embedded states at level 1 (the ones we wish to place in \mathcal{E}) are shadowed.

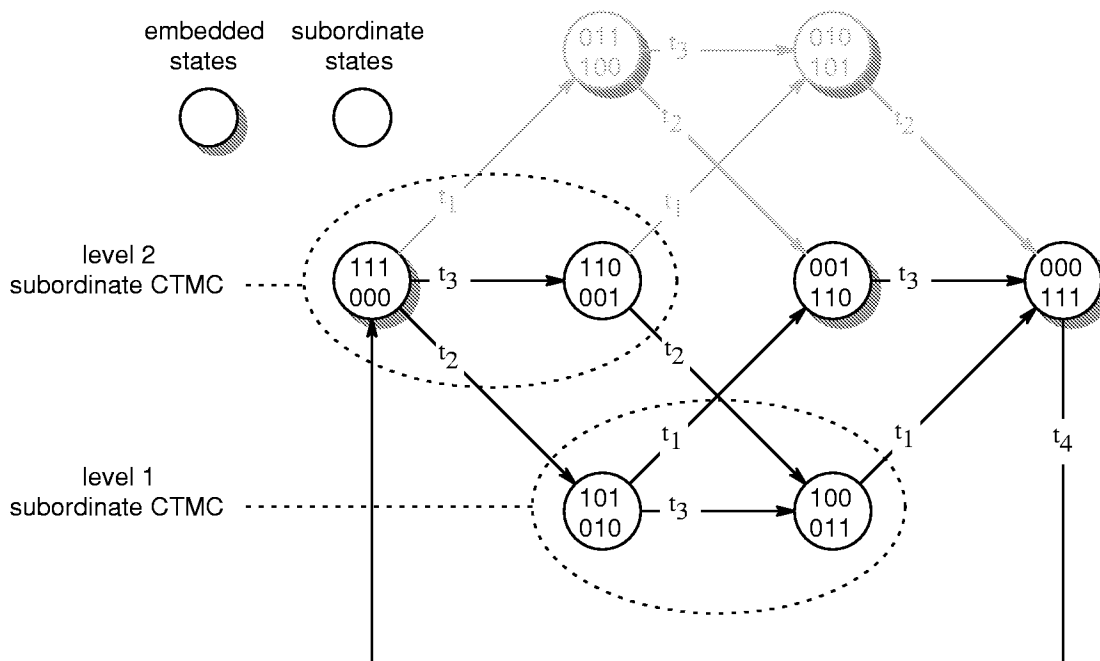


Figure 2.11: Two-Level Semi-regenerative process.

The subordinate processes for levels 1 and 2 are partitioned into two groups. The level 2 subordinate process, a 2-state CTMC, has initial state (111000) with probability one; it is also an embedded state by definition. The level 1 subordinate process, also a 2-state CTMC, has an initial probability distribution subject to the transient analysis at level 2, which gives the probability of entering either of the two states (101010) or (100011) when t_2 fires. The embedded state entered at level 1 depends on the subordinate state occupied when t_1 finally fires. The grayed portion of the graph indicates the states that are no longer reachable from (111000) because of the imposed timing, i.e., transition t_1 must fire after t_2 . Unlike the previous examples, the reachability graph is obviously not isomorphic to the stochastic process and $\mathcal{S} \subset \mathcal{R}$.

In cases like this, studying the subordinate process between regeneration points becomes as difficult as studying the actual process. Moreover, under different assumptions, regeneration points may be rare or the process may not even be semi-regenerative. Without simplifying assumptions, the underlying process of an SPN may be a generalized semi-Markov process.

2.6 Generalized Semi-Markov Models

By definition, the future markings of an untimed PN depends only on the current marking, not the past markings, and so the Markov property holds. It is the nature of the timed marking process, a stochastic process, that complicates matters when any transition t has a PDF, $F^t(\cdot)$, that is not memoryless, and hence the Markov property does not hold. Without restrictions, the underlying process of an SPN is known as a *generalized semi-Markov process* (GSMP).

By including in the GSMP state the RFT information for each transition along with the current marking, everything about the past that is needed to determine the future evolution of the stochastic process is contained in the current state. Hence, the past can be forgotten. Alternatively, we could instead augment markings with *age* information, which records the times since transitions became enabled without firing. Either way, by effectively *Markovianizing* the process, customary Markov techniques can be applied, albeit to a larger, possibly continuous, state space and, perhaps, a more complicated reachability graph. This “Markovianizing” of the underlying process, which can then be solved using a generalization of Kolmogorov’s forward (or backward) equations, is called the *method of supplementary variables*, which was first proposed by Cox in [31]. However, the new, Markovianized process is a *continuous-state* process in general because of the continuous nature of the age or RFT information that are augmented with the discrete marking information.

The method of supplementary variables has been applied to DSPNs and MRSPNs in [32] as an alternative to Markov renewal theory. Unlike Markov renewal theory, this method is still applicable when the process is a GSMP. A fourth-order, stationary solution algorithm has been proposed in [33] for DSPN and MRSPN models. Unfortunately, the solution of such system of equations is usually too numerically challenging for anything other than models with very small dimensions. We will back up this claim while describing the method of supplementary variables using our running example.

When employing the method of supplementary variables, there is some freedom of choice in how the generalized Kolmogorov’s equations are constructed. We can choose to look forward or backward, use age or RFT variables, and specify differential or integral formulas. The published literature regarding the solution of extended DSPNs and MRSPNs with the method of supplementary variables tend to use forward equations and *age* variables, and we choose to do so here as well.

Generalized Kolmogorov’s equations, forward or backward, require the use of a stationary probability density function (pdf) on the state and age variable jointly:

$$q_k(v) = \lim_{\theta \rightarrow \infty} \frac{d}{dv} \Pr\{ X(\theta) = k, a \leq v \}$$

where $v \in \mathbb{R}^+$, $k \in \mathcal{S}$, and $a \in \mathbb{R}^+$ denotes the age variable. A multidimensional pdf, $q_k(a_1, a_2, \dots, a_n)$, is used for states that enable more than one generally distributed transition. Instead of the constant rates given by the infinitesimal generator matrix, the generalized Kolmogorov’s equations also require the use of *instantaneous rate functions*, defined as

$$\lambda^t(v) = \frac{f^t(v)}{1 - F^t(v)}$$

where $f^t(t) = \frac{d}{dv} F^t(v)$ is the pdf of the firing delay for transition t . Since $\lambda^t(\cdot)$ is a conditional probability function (not a PDF however), this allows the age information a to be taken into consideration when computing the state transition rates. That is, if a denotes the age of transition t then $\lambda^t(a)da$ is interpreted as the conditional probability that t fires in the next da interval given that it has not fired in time a since becoming enabled. Of course, $\lambda^t(\cdot) = \lambda$, a constant, if t is exponentially distributed with rate λ . So when only exponentially distributed transitions are enabled, the forward equations will degenerate to the familiar system of equations (Equation 2.8) presented earlier for CTMCs.

As a consequence of the supplementary variable conditioning, we must be careful of the initial value and boundary conditions when constructing the state equations. Moreover, the Kolmogorov's forward equations must be *partial differential equations* when more than one transition with general firing delays are enabled simultaneously, coinciding with states having multiple age variables.

Consider the last example model used at the end of the previous section but without any special assumptions about $F^{t_1}(\cdot)$ and $F^{t_2}(\cdot)$. In this case, the underlying process is a GSMP. Let $a_1 \in \mathbb{R}$ and $a_2 \in \mathbb{R}$ be the *ages* of transitions t_1 and t_2 , respectively. Of course, age information for transitions t_3 and t_4 can be omitted since the Expo PDF is the same when conditioned with age (or RFT) information; hence, nothing concerning elapsed time needs to be remembered. Only states that enable t_1 , t_2 , or both are supplemented with a_1 , a_2 , or both, respectively. The GSMP for our example is shown in Figure 2.12 while also introducing a different depiction of “supplemented states”, enumerated by $k \in \{1, 2, \dots, 8\} = \mathcal{S}$. That is, the supplemented states are portrayed as two circles: the smaller, raised one contains the marking and the larger one contains the age information. This depiction will be used in Chapter 4 when presenting the solution algorithms for our new SPN class, except that instead of continuous age, discrete information concerning the RFT of phase-type transitions will be recorded. Referring to Figure 2.12, the state equations using the method of supplementary “age” variables are constructed as follows.

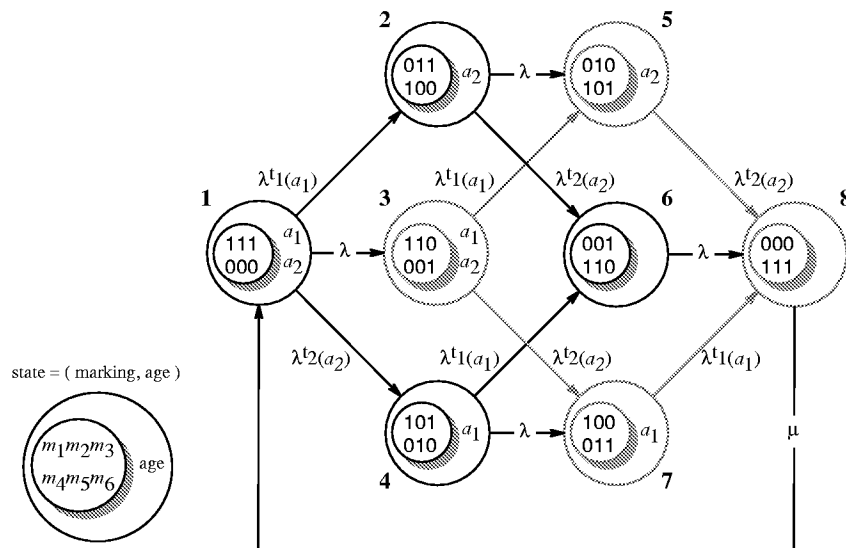


Figure 2.12: GSMP where t_1, t_2 are generally distributed and all others Expo.

First, we must consider the initial values for the age variables, which are both reset to zero only upon entering state 1 from state 8. We assume that $F^{t_1}(\cdot)$ and $F^{t_2}(\cdot)$ do not have mass at the origin and so t_1 and t_2 must delay some positive value before firing, and so no other events causing outflow of probability mass needs capturing. Consequently, the initial value condition is simply

$$q_1(0, 0) = q_8\mu.$$

Second, we have the state equations for positive-valued age variables, constructed in the spirit of Kolmogorov's forward equation:

$$\begin{aligned} \frac{d}{da_1}q_4(a_1) &= \int_0^\infty q_1(a_1, a_2)\lambda^{t_2}(a_2) da_2 - q_4(a_1)\lambda, & a_1 > 0 \\ \frac{d}{da_1}q_7(a_1) &= \int_0^\infty q_3(a_1, a_2)\lambda^{t_2}(a_2) da_2 + q_4(a_1)\lambda, & a_1 > 0 \\ \frac{d}{da_2}q_2(a_2) &= \int_0^\infty q_1(a_1, a_2)\lambda^{t_1}(a_1) da_1 - q_2(a_2)\lambda, & a_2 > 0 \\ \frac{d}{da_2}q_5(a_2) &= \int_0^\infty q_3(a_1, a_2)\lambda^{t_1}(a_1) da_1 + q_2(a_2)\lambda, & a_2 > 0 \\ \left(\frac{\partial}{\partial a_1} + \frac{\partial}{\partial a_2}\right)q_3(a_1, a_2) &= q_1(a_1, a_2)\lambda, & a_1 > 0, a_2 > 0 \end{aligned}$$

The integrals are needed to uncondition on the age variable associated with the transition that fires, thereby causing the state transition of interest.

Third, we have equations for states without age variables, which resemble the familiar flow-balance equations for CTMCs:

$$\begin{aligned} \int_0^\infty q_4(a_1)\lambda^{t_1}(a_1) da_1 + \int_0^\infty q_2(a_2)\lambda^{t_2}(a_2) da_2 - q_6\lambda &= 0 \\ \int_0^\infty q_7(a_1)\lambda^{t_1}(a_1) da_1 + \int_0^\infty q_5(a_2)\lambda^{t_2}(a_2) da_2 + q_6\lambda - q_8\mu &= 0 \end{aligned}$$

Fourth, since we are only interested in the stationary probability distribution of markings, we also need equations that eliminate the age variables:

$$p_k = \begin{cases} \iint_0^\infty q_k(a_1, a_2) da_1 da_2 & : k \in \{1, 3\} \\ \int_0^\infty q_k(a_2) da_2 & : k \in \{2, 5\} \\ \int_0^\infty q_k(a_1) da_1 & : k \in \{4, 7\} \\ q_k & : k \in \{6, 8\} \end{cases}$$

where $p_k = \lim_{\theta \rightarrow \infty} \Pr\{X(\theta) = k\}$, $k \in \mathcal{S}$.

And finally, we need one last equation to normalize the solution:

$$\sum_{k \in \mathcal{S}} p_k = 1.$$

Of course, all of the above equations need to be satisfied simultaneously. Numerical solution of such system of equations ultimately requires the discretization of the continuous

variable differential-integro equations so that finite difference equations can be solved instead. We must also assume that $F^{t_1}(\cdot)$ and $F^{t_2}(\cdot)$ have finite support so that a finite dimension mesh can be constructed. But even then, computing the solution is numerically challenging. We do not intend to solve models this way. The real purpose of this exercise was to show that even with this small, simple model, solving general SPN models with underlying GSMPs is computationally challenging and costly.

Alternative solution techniques for GSMP models have been investigated. One of these, presented in [34], observes the process at fixed intervals and records the marking and RFT at these times. This procedure gives rise to an embedded general state-space Markov chain from which state equations can be written and then transformed into a system of Volterra equations. These Volterra equations permit the specification of state transitions subject to the clock readings at the equidistant time intervals and can be simpler to solve numerically than the method of supplementary variables.

When faced with practical considerations, the modeler is usually limited to simulation. SPN-specific simulation techniques have also been investigated. In [35], structural properties and conditions imposed on the SPN model are exploited that ensures the underlying process is regenerative so that faster regenerative simulation can be employed. A regenerative process is in fact a special case of semi-regenerative processes, one that regenerates itself probabilistically but also independently to the state entered at each regeneration time, thus it is more restricted. For example, a semi-regenerative process with only a single state is a regenerative process. A different simulation technique can be found in [36] where time-averaged statistics can be obtained using methods based on standardized time series, in particular, the method of batch means with the number of batches fixed. This method can analyze simulated output where the regeneration methods are not applicable.

The disadvantages to using simulation for general SPN models is that the results can be inaccurate, limited to confidence intervals, and require long simulation times. Also, the state space is never explored exactly, a shortfall that may preclude logical analysis and model-based formal verification. That is, just because some state specific property (such as deadlock) or sequence-specific property (such as an undesirable chain of events) is not observed in the simulation run does not mean that such properties do not exist in the model. It could simply mean that the simulation was not run long enough.

Chapter 3

Proposed Research

The concept essential to studying a non-Markovian process using Markovian techniques is the inclusion of additional information in the state for the purpose of remembering the time each enabled transition has to delay before its scheduled firing. This time is initially sampled from its probability distribution function when a transition first becomes enabled. Compared with general firing delays, phase-type firing delays have the advantage that the extra RFT information included in the state is conveniently *discretized* into states of an absorbing Markov chain, greatly simplifying the solution when either PH or DPH are used alone.

PH or DPH random variables can approximate any general random variable arbitrarily well. The use of phase-type firing delays has had success in the past as a way of broadening the applicability of SPN models. Indeed, continuous PH random variables like Erlang can even approximate discrete, constant random variables by including enough stages, since the variance diminishes as the number of stages increases. Similarly, the discrete Geom random variable can approximate the continuous Expo random variable arbitrarily well by reducing the basic step size. But approximating general timing behavior by allowing either PH or DPH alone in a given model offers less modeling convenience and, possibly, requires more computational work.

Limiting a model to either PH or DPH timing may also require more computational work. Capturing a deterministic activity into a model with only PH-type timing requires an approximation using an Erlang random variable and may require many stages, and hence adding to the state space. possibly, requires more computational work.

Fortunately, such restrictions are unnecessary.

We propose to extend the SPN definition of the previous chapter to include non-Markovian timing that may prove useful to many problems while still affording an efficient, numerical solution. To this end, we elect to extend the phase-type SPN formalism for use in both discrete and continuous time, present simultaneously in the *same* model. This research effort will develop a new class of SPN that permits transition firing delays with both PH and DPH distributions. We call this new formalism a *Phased Delay Petri Net* (PDPN). The contribution herein includes the first time that PH behavior has been combined with DPH behavior in the same model, thereby extending the aforementioned research that considered each separately. Alone, PH and DPH models enjoy the “memoryless” property of underlying CTMCs and DTMCs, respectively, making efficient solutions possible. Together, the reasoning about the combined behavior becomes complicated. The proposed research will

show that the underlying process is semi-regenerative at best (perhaps degenerating to a DTMC, CTMC, or semi-Markov chain at times) and a generalized semi-Markov process at worst [25, 37]. So, our contribution also includes a formalized understanding of the stochastic process underlying this new class of SPN and the development of efficient algorithms for its solution.

After presenting the PDPN in a general setting, we intend to show how certain simplifying assumptions restrict the underlying stochastic process to one that is manageable with efficient solutions. Investigations will be conducted into techniques that offer efficiencies in exact solutions when practical. Otherwise, approximate solution algorithms will be sought that give heuristically good results with high fidelity and efficiency. Although these assumptions may also restrict the set of problems that can be modeled, we anticipate nevertheless that this will not preclude the usefulness of the PDPN to many real-world applications.

While the PDPN approach has the benefit in fidelity that comes from mixing PH and DPH behavior, the expansion of the state space required by this approach will undoubtedly compound the well known “state-space explosion problem”. However, this extra burden on memory can be alleviated by utilizing advanced data structures and manipulation algorithms. For instance, decision diagrams proposed in [38, 39] offer very compact storage with a fraction of the memory requirements as conventional sparse storage structures where the memory usage grows linearly with the number of states. These advanced data structures should be just as applicable to storing the PDPN reachability graph (for performance analysis) as well as the reachability set (for logical analysis).

We can also exploit the convenient formulation of the PH and DPH state space expansion using Kronecker addition and multiplication, respectively. By employing data structures and manipulation algorithms similar to the aforementioned ones, we can take advantage of the Kronecker representations to efficiently store the expanded state space implicitly. That is, the sparse PH and DPH Markov chains can be stored in isolation and the expanded state space and matrix entries can be constructed as needed using the Kronecker operators. Such techniques have already been investigated in [40, 41, 42], and we feel that these recent advancements in compact state-space and matrix storage techniques make phase-expansion approaches worth revisiting.

Chapter 4

Preliminary Research

Before efficient analysis techniques can be investigated, the PDPN model must first be formalized, its inherent properties understood, and interesting subclasses defined to aid the investigation of solution algorithms. In this section, we discuss the underlying stochastic process of the PDPN, how to study it, and at the same time define three subclasses to the general PDPN base class. The subclasses restrict the underlying process to one that can be more easily studied. We discuss the findings of our preliminary research and the implications to solution complexity and applicability of PDPNs. We end this chapter with a proposed stationary solution algorithm and its analysis.

4.1 Analyzing the Underlying Stochastic Process

For convenience, we partition the set of transitions \mathcal{T} into the set \mathcal{T}_C having PH distributions, the set \mathcal{T}_D having DPH distributions, and the set \mathcal{T}_Z of immediate transitions having $\text{Const}(0)$ distributions. Although the “immediate” transitions in \mathcal{T}_Z are, in fact, special cases of DPH transitions, we consider them separately since they are given higher priority in firing over the timed transitions in $\mathcal{T}_C \cup \mathcal{T}_D$, and hence are usually handled separately during the analysis.

As already stated, when PH (DPH) is used alone (possibly in conjunction with immediate transitions), an otherwise non-Markovian process becomes a CTMC (DTMC). When PH and DPH are allowed to coexist, we have complicated matters, but not as badly as allowing completely general distributions. Towards constructing a PDPN reachability graph and the corresponding stochastic process specification, we must determine the possible combination of phases that can occur, ultimately leading to a phase that allows some transition t to fire. While at least one $t \in \mathcal{T}_D$ is enabled, phase changes occur at discrete instants of time $n\tau$, where $n \in \mathbb{N}$ and τ denotes the basic step interval common to all DPH transitions, referred to hereafter as the *clock*. Because PH distributions are continuous, the probability of observing a PH phase change at any particular point in time, in particular at times $n\tau$, is zero. Therefore, we can consider the DPH phase changes separately from PH phase changes. The possible combination of DPH and PH phases for each marking \mathbf{m} can be obtained from Equations 2.15 and 2.16, respectively. Hereafter, we let \mathcal{D} denote the potential phase space. The actual state space of the stochastic process will be denoted by $\mathcal{S} \subseteq \mathcal{R} \times \mathcal{D}$.

An example reachability graph assuming that transition t_3 is either $\text{Erlang}(\cdot, 2) \in \text{PH}$ or $\text{Const}(2) \in \text{DPH}$ and all others are $\text{Expo} \in \text{PH}$ is given in Figure 4.1. The states are

shadowed in a way that distinguishes the RFT information of transition t_3 corresponding to its firing delay *phase*.

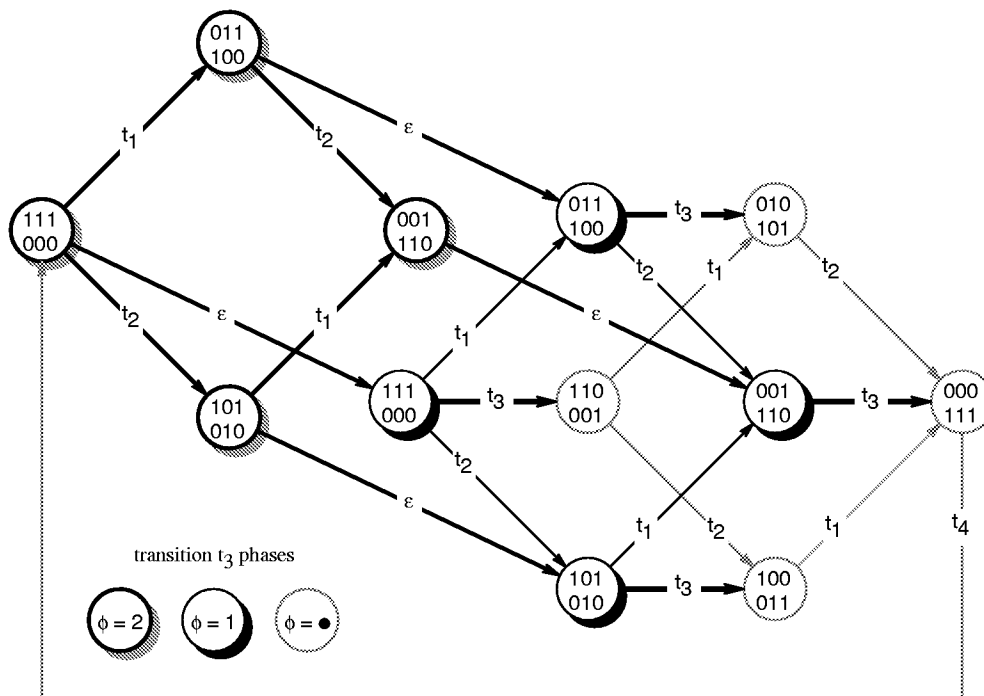


Figure 4.1: Markovianized process of the example PDPN.

PDPNs do not satisfy the conditions for isomorphism between the reachability graph and state space because multiple \mathcal{T}_D transitions can be enabled and fired simultaneously, a sequence $s \in (\mathcal{T}_D \cup \mathcal{T}_Z)^*$, with the possibility of disabling any or all co-enabled transitions in a way that prevents their firing in certain states. Therefore, the construction of the PDPN reachability graph and the underlying stochastic process must be done in concert. As an example of how the timing constraints can restrict the reachability graph, consider the case where $t_2 \sim \text{Geom}(q, 3)$, $t_3 \sim \text{Const}(2)$, and the other transitions are Expo. The resulting stochastic marking process is portrayed in Figure 4.2 where the grayed portion indicates the states that are no longer accessible due to the time constraints. Whereas the sojourn times in markings (100011) and (000111) are exponentially distributed and therefore memoryless, this is not the case for the other markings.

In general, the underlying stochastic process of the PDPN is a generalized semi-Markov process. Although state equations can be constructed, as discussed in Chapter 2, and the method of supplementary variables applied, the method requires the solution of partial differential equations, which is computationally intensive in general. Recall that for MRSPNs, the restriction that at most one generally distributed transition is enabled in any marking simplifies the model to one consisting of ordinary differential equations, which are easier to solve. But, since this restricted marking process is a semi-regenerative process, the supplementary variable can be eliminated altogether by constructing the solution algorithm around Markov renewal theory.

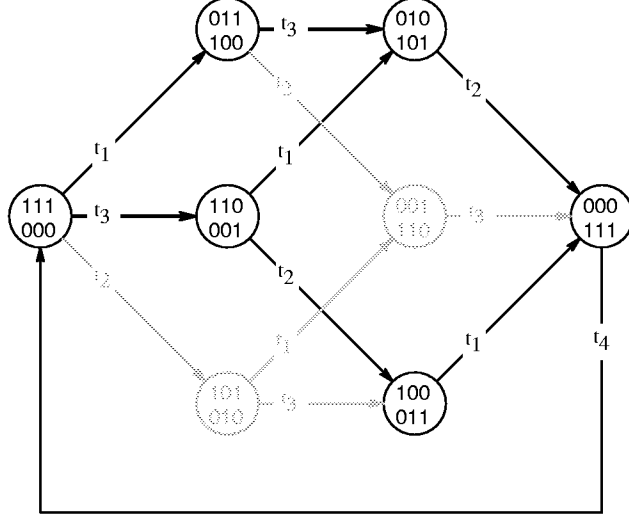


Figure 4.2: Example non-Markovian process when $t_2 \sim \text{Geom}(q, 3)$ and $t_3 \sim \text{Const}(2)$.

For PDPNs, we can eliminate the clock variables from the state, thereby reducing the state-space memory along with the computational costs if we restrict the PDPN so that *synchronization* is maintained among all enabled \mathcal{T}_D transitions. By synchronizing the \mathcal{T}_D transitions, the analysis is simplified, subject to only a single clock, thus requiring the storage of just one clock variable in each state. Even better, the clock variable can be eliminated altogether for this restricted PDPN by recognizing the underlying stochastic process as a semi-regenerative process. The conditions necessary for such synchronization in any marking $\mathbf{m} \in \mathcal{R}$ are:

1. transitions $t \in \mathcal{T}_D$ are never enabled by a firing sequence $s \in \mathcal{T}_C \mathcal{T}_Z^*$ except when $\mathcal{F}(\mathbf{m}) \cap \mathcal{T}_D = \emptyset$, and
2. if a firing sequence $s \in \mathcal{T}_C \mathcal{T}_Z^*$ resets the current phase of a transition $t \in \mathcal{F}(\mathbf{m}) \cap \mathcal{T}_D$ then it must do so for all transitions in $\mathcal{F}(\mathbf{m}) \cap \mathcal{T}_D$.

Theorem 4.1.1 *The stochastic process underlying a synchronized PDPN is semi-regenerative.*

PROOF. Let $X = \{X(\theta) : \theta \geq 0\}$ denote the underlying stochastic process. Clearly, if only \mathcal{T}_D transitions are enabled, X is a DTMC, and if only \mathcal{T}_C transitions are enabled, X is a CTMC, both of which are special cases of the semi-regenerative class. Consider periods when both \mathcal{T}_D and \mathcal{T}_C transitions are enabled. Because the firing sequences $s \in (\mathcal{T}_C \cup \mathcal{T}_Z)^*$ are expanded into Expo and Const(0) state transitions, which are memoryless for all time, we need only observe the successive times when the \mathcal{T}_D transitions become enabled, fire, or undergo phase advancements. Restrictions (1) and (2) ensure that all such events for different \mathcal{T}_D are synchronized and therefore occur at successive jump times $T_n = T_{n-1} + \tau$ at which time state X_n is entered. The sequence of states $\{X_n : n \geq 0\}$ form a DTMC, and together the sequence $\{(X_n, T_n) : n \geq 0\}$ forms a Markov renewal process. Therefore it follows, by definition, that X is a semi-regenerative process. \square

The process shown in Figure 4.2 is semi-regenerative since t_2 and t_3 maintain synchronization with respect to the basic clock advancements. As such, Markov renewal theory can be applied to solve the model. While it has been shown that solution algorithms based on Markov renewal theory has the same asymptotic costs as the method of supplementary variables [11], we believe Markov renewal theory can be more intuitive in some ways, interpretation of the results to the original process can be preserved, and opportunities to eliminate phase information from the state can be exploited. But to make Markov renewal theory applicable, we must impose the above conditions on the PDPN model to ensure that the underlying stochastic process is semi-regenerative.

4.1.1 Theory Applied to PDPNs in General

We already know that, separately, PH and DPH based SPNs enjoy the efficient solution of underlying CTMCs and DTMCs, respectively. We also know that mixing PH and DPH behavior requires the solution of a semi-regenerative process, thereby complicating matters. However, by assuming synchronization between \mathcal{T}_D transition when enabled and with the aid of Markov renewal theory, we have reduced the analysis problem to one of studying multiple CTMCs (the subordinate processes), one for each embedded state in \mathcal{E} , and one DTMC (the EMC). To this end, we must simultaneously study the evolution of both the DTMC and each CTMC in turn as they interact with one another.

Figure 4.3 shows a sample path observation of a typical PDPN regeneration period aided by Markov renewal theory. Because isomorphism between the timed and untimed PDPN reachability graph is not guaranteed and because of the potential interaction between the DTMC and CTMC models, we take the approach of constructing the stochastic matrix, $\mathbf{\Pi}$, of the EMC, one row at a time.

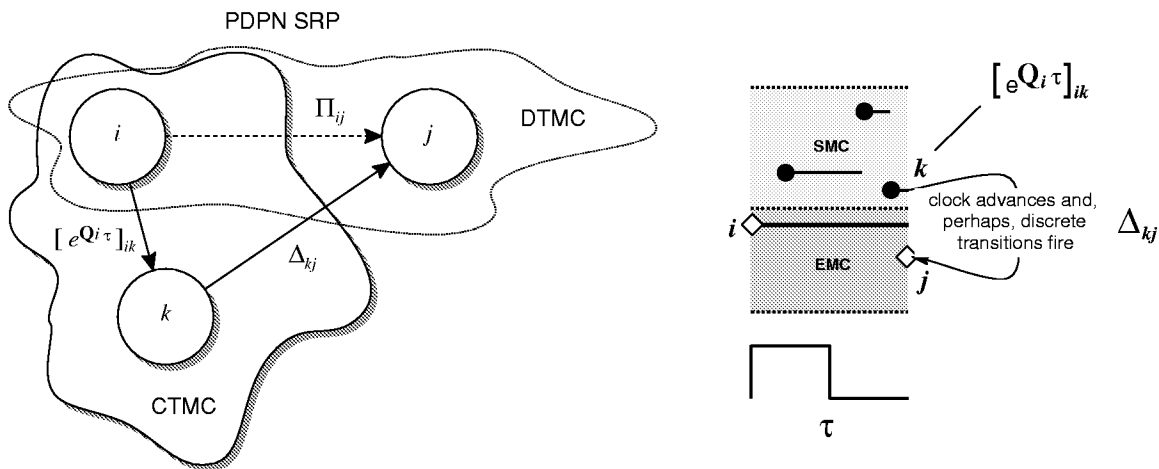


Figure 4.3: Studying a PDPN regeneration period.

The basic algorithm to construct $\mathbf{\Pi}$ can be as follows. Starting with a known embedded state $i \in \mathcal{E}$, we observe the subordinate CTMC (SMC) up to the next clock advance, a period of *at most* length τ . The regeneration period T_1 defined here is actually a random variable over the range $(0, \tau]$. T_1 is exactly τ if at least one $t \in \mathcal{T}_D$ remains enabled during the entire period. However, if all transitions in \mathcal{T}_D become disabled or they all are forced to simultaneously reset to a new firing delay (resample) due to the firing of a transition from $\mathcal{T}_C \cup \mathcal{T}_Z$, then T_1 will be less than τ . Assuming that $T_1 = \tau$, we simply solve the SMC (with generator \mathbf{Q}_i and state space \mathcal{S}_i originating from the embedded state i) at time τ . The state occupied at time τ , when the clock advance occurs, is applied to the next state *switching matrix* $\mathbf{\Delta}$ that computes the set of states reachable after some $s \in (\mathcal{T}_D \cup \mathcal{T}_Z)^*$ following the clock advance. This set of next states are new embedded states and are added to \mathcal{E} . This procedure repeats until no new embedded states are found. Formally, the analysis is expressed by the application of Markov renewal theory with the PDPN properties in mind:

$$H_{ik}(\theta) = \left[e^{\mathbf{Q}_i \theta} \right]_{ik} \quad (4.1)$$

$$\Pi_{ij} = \sum_{k \in \mathcal{S}_i} \left[e^{\mathbf{Q}_i \tau} \right]_{ik} \Delta_{kj} \quad (4.2)$$

where, of course, $e^{\mathbf{Q}_i \theta}$ is the transient solution of the SMC at time θ , one solution for each $i \in \mathcal{E}$. Notice that these equations are the same as for DSPNs except that for PDPNs we have a fixed deterministic delay τ . The following equations that are needed for the stationary solution of PDPNs are the same as well and are repeated here for convenience.

We can distinguish the two cases of $\{T_1 = \tau\}$ and $\{T_1 < \tau\}$ by appropriately constructing the SMC. That is, the CTMC states reached that also coincide with \mathcal{T}_D simultaneous disabling or resampling are made absorbing and are regarded as embedded states in \mathcal{E} . The set of absorbing states, which are formed in this way, will be denoted hereafter as the set \mathcal{E}_i . In this way, we *trap* such events and associated probability mass in these absorbing states when solving the CTMC at time τ . If the total probability mass absorbed is α , then because $\Pr\{T_1 > \tau\} = 0$, we know that $\Pr\{T_1 < \tau\} = \alpha$ and $\Pr\{T_1 = \tau\} = 1 - \alpha$. For stationary analysis where we are interested in constructing the EMC matrix $\Pi_{ij} = \lim_{\theta \rightarrow \infty} G_{ij}(\theta)$, the exact value of T_1 is of no consequence, only its expected value is important. With the appropriately constructed absorbing CTMC, the expected value of T_1 is determined from the cumulative probabilities

$$h_{ik} = \mathbb{E}[\text{sojourn in } k \text{ during } [0, T_1) \mid X_0 = i] = \left[\int_0^\tau e^{\mathbf{Q}_i v} dv \right]_{ik} \quad i \in \mathcal{E}, k \in \mathcal{S}_i \setminus \mathcal{E}_i \quad (4.3)$$

in each CTMC state, which are computed anyway to obtain the necessary conversion factors. Then

$$\mathbb{E}[T_1 \mid X_0 = i] = \sum_k h_{ik}$$

for each $i \in \mathcal{E}$ and for all $k \in \mathcal{S}_i \setminus \mathcal{E}_i$.

The stationary solution $\mathbf{x} = [x_i] \in \mathbb{R}^{|\mathcal{E}|}$ of the EMC satisfies the set of balance equations

$$\mathbf{x}\mathbf{\Pi} = \mathbf{x} \quad \text{subject to} \quad \sum_{i \in \mathcal{E}} x_i = 1. \quad (4.4)$$