

AperTO - Archivio Istituzionale Open Access dell'Università di Torino

Phase Type and Matrix Exponential Distributions in Stochastic Modeling

This is the author's manuscript

Original Citation:

Availability:

This version is available <http://hdl.handle.net/2318/1645687> since 2017-07-25T16:13:54Z

Publisher:

Springer International Publishing

Published version:

DOI:10.1007/978-3-319-30599-8_1

Terms of use:

Open Access

Anyone can freely access the full text of works made available as "Open Access". Works made available under a Creative Commons license can be used according to the terms and conditions of said license. Use of all other works requires consent of the right holder (author or publisher) if not exempted from copyright protection by the applicable law.

(Article begins on next page)

Metadata of the chapter that will be visualized in SpringerLink

Book Title	Principles of Performance and Reliability Modeling and Evaluation	
Series Title		
Chapter Title	Phase Type and Matrix Exponential Distributions in Stochastic Modeling	
Copyright Year	2016	
Copyright HolderName	Springer International Publishing Switzerland	
Corresponding Author	Family Name	Horvath
	Particle	
	Given Name	Andras
	Prefix	
	Suffix	
	Division	Dipartimento di Informatica
	Organization	Università degli Studi di Torino
	Address	C.so Svizzera 185, 10149, Turin, Italy
	Email	horvath@di.unito.it
Author	Family Name	Scarpa
	Particle	
	Given Name	Marco
	Prefix	
	Suffix	
	Division	Dipartimento di Ingegneria
	Organization	Università degli Studi di Messina
	Address	C.da di Dio, 98166, Messina, Italy
	Email	mscarpa@unime.it
Author	Family Name	Telek
	Particle	
	Given Name	Miklos
	Prefix	
	Suffix	
	Division	Department of Networked Systems and Services, MTA-BME Information Systems Research Group
	Organization	Budapest University of Technology and Economics
	Address	Magyar Tudosok krt. 2, Budapest, 1117, Hungary
	Email	telek@hit.bme.hu

Abstract Since their introduction, properties of Phase Type (PH) distributions have been analyzed and many interesting theoretical results found. Thanks to these results, PH distributions have been profitably used in many modeling contexts where non-exponentially distributed behavior is present. Matrix Exponential (ME) distributions are distributions whose matrix representation is structurally similar to that of PH distributions but represent a larger class. For this reason, ME distributions can be usefully employed in modeling contexts in place of PH distributions using the same computational techniques and similar algorithms, giving rise to new opportunities the fact, they are able to represent different dynamics, e.g., faster dynamics, or the same dynamics but at lower computational cost. In this work, we deal with the

characteristics of PH and ME distributions, and their use in stochastic analysis of complex systems. Moreover, the techniques used in the analysis to take advantage of them are revised.

Phase Type and Matrix Exponential Distributions in Stochastic Modeling

Andras Horvath, Marco Scarpa and Miklos Telek

1 **Abstract** Since their introduction, properties of Phase Type (PH) distributions have
2 been analyzed and many interesting theoretical results found. Thanks to these results,
3 PH distributions have been profitably used in many modeling contexts where non-
4 exponentially distributed behavior is present. Matrix Exponential (ME) distributions
5 are distributions whose matrix representation is structurally similar to that of PH
6 distributions but represent a larger class. For this reason, ME distributions can be
7 usefully employed in modeling contexts in place of PH distributions using the same
8 computational techniques and similar algorithms, giving rise to new opportunities
9 the fact, they are able to represent different dynamics, e.g., faster dynamics, or the
10 same dynamics but at lower computational cost. In this work, we deal with the
11 characteristics of PH and ME distributions, and their use in stochastic analysis of
12 complex systems. Moreover, the techniques used in the analysis to take advantage
13 of them are revised.

14 1 Introduction

15 Stochastic modeling has been used for performance analysis and optimization of
16 computer systems for more than five decades [19]. The main analysis method behind
17 this effort was the continuous time Markov chains (CTMC) description of the sys-

A. Horvath (✉)

Dipartimento di Informatica, Università degli Studi di Torino, C.so Svizzera 185,
10149 Turin, Italy
e-mail: horvath@di.unito.it

M. Scarpa

Dipartimento di Ingegneria, Università degli Studi di Messina, C.da di Dio,
98166 Messina, Italy
e-mail: mscarpa@unime.it

M. Telek

Department of Networked Systems and Services, MTA-BME Information
Systems Research Group, Budapest University of Technology and Economics,
Magyar Tudosok krt. 2, Budapest 1117, Hungary
e-mail: telek@hit.bme.hu

© Springer International Publishing Switzerland 2016

L. Fiondella and A. Puliafito (eds.), *Principles of Performance and Reliability
Modeling and Evaluation*, Springer Series in Reliability Engineering,
DOI 10.1007/978-3-319-30599-8_1

tem behavior and the CTMC-based analysis of the performance measures of interest. With the evolution of computing devices, model description languages (e.g., queueing systems, Petri nets, process algebras), and model analysis techniques (a wide range of software tools with efficient analysis algorithm using adequate data representation and memory management) the analysis of more and more complex systems has become possible. One of main modeling limitations of the CTMC-based approach is the limitation on the distribution of the random time durations, which is restricted to be exponentially distributed. Unfortunately, in a wide range of practical applications, the empirical distribution of field data differs significantly from the exponential distribution. The effort to relax this restriction of the CTMC-based modeling on exponentially distributed durations resulted in the development of many alternative stochastic modeling methodologies (semi-Markov and Markov regenerative processes [11], analysis with the use of continuous system parameters [8]), yet all of the alternative modeling methodologies suffer from infeasible computational complexity very quickly when the complexity of the systems considered increases beyond basic examples.

It remains a significant research challenge to relax the modeling restriction of the exponentially distributed duration time and still evaluate complex model behaviors. To this end, one of the most promising approaches is the extension of CTMC-based analysis to non-exponentially distributed durations. Initial steps in this direction date back to the activity of A.K. Erlang in the first decades of the twentieth century as reported in [10]. These initial trials were referred to as the method of phases, which influenced later terminology. M.F. Neuts characterized a set of distributions which can be incorporated into CTMC-based analysis by introducing the set of phase type (PH) distributions [16].

The extension of CTMC-based analysis (where the durations are exponentially distributed) with PH distributed durations requires the generation of a large CTMC, referred to as extended Markov chain (EMC), which combines the system behavior with the description of the PH distributions. In this chapter, we summarize the basics of EMC-based stochastic analysis and provide some application examples. Finally, we note that in this work we restrict our attention to continuous time stochastic models, but that the same approach applies for discrete time stochastic models as well.

1.1 Structure of the Chapter

The next two sections, Sects. 2 and 3, summarize the basic information on PH and ME distributions, respectively. The following two sections, Sects. 4 and 5, discuss the analysis procedure for complex stochastic systems with PH and ME distributed durations, respectively. The tools available to support EMC-based analysis of stochastic systems is presented in Sect. 6. Numerical examples demonstrate the modeling and analysis capabilities of the approach are discussed in Sect. 7 and the main findings and conclusions are given in Sect. 8.

2 PH Distributions and Their Basic Properties

2.1 Assumed Knowledge

Transient behavior of a finite state Markov chain with generator \mathbf{Q} and initial distribution π , specifically, the transient probability vector $p(t)$, satisfies the ordinary differential equation

$$\frac{d}{dt} p(t) = p(t)\mathbf{Q}, \text{ with initial condition } p(0) = \pi,$$

whose solution is a matrix exponential function

$$p(t) = \pi e^{\mathbf{Q}t}, \quad (1)$$

where the matrix exponential term is defined as

$$e^{\mathbf{Q}t} = \sum_{i=0}^{\infty} \frac{t^i}{i!} \mathbf{Q}^i.$$

The properties of generator \mathbf{Q} and initial distribution π are as follows. The elements of π are probabilities, i.e., nonnegative numbers not greater than one. The off-diagonal elements of \mathbf{Q} are transition intensities, i.e., nonnegative numbers. The diagonal elements of \mathbf{Q} are such that each row sum is zero, i.e., the diagonal elements are non-positive. The elements of π sum to one, that is $\sum_i \pi_i = \pi \mathbf{1} = 1$. Each row of a generator matrix sums to zero, that is $\sum_j Q_{ij} = 0$, or equivalently, in vector form, we can write $\mathbf{Q}\mathbf{1} = \mathbf{0}$, where $\mathbf{1}$ is a column vector of ones and $\mathbf{0}$ is a column vector of zeros. Hereafter, the sizes of vector $\mathbf{1}$ and $\mathbf{0}$ are defined by the context such that the dimensions in the vector expressions are compatible.

The stationary distribution of an irreducible finite state Markov chain with generator \mathbf{Q} , $p \triangleq \lim_{t \rightarrow \infty} p(t)$, can be computed as the unique solution of the linear system of equations

$$p\mathbf{Q} = \mathbf{0}, \quad p\mathbf{1} = 1. \quad (2)$$

In this chapter, we focus on the computation of the initial distribution and the generator matrix of the EMC and do not discuss the efficient solution methods for solving (1) and (2).

2.2 Phase Type Distributions

PH distributions are defined by the behavior of a Markov chain, which is often referred to as the background Markov chain behind a PH.

Let $X(t)$ be a Markov chain with n transient and one absorbing states, meaning that the absorbing state is reachable (by a series of state transitions) from all transient states, but when the Markov chain moves to the absorbing state it remains there forever. Let π be the initial distribution of the Markov chain, that is $\pi_i = P(X(0) = i)$. Without loss of generality, we number the states of the Markov chain such that state $1, \dots, n$ are transient states and state $n + 1$ is the absorbing state. The generator matrix of such a Markov chain has the following structure

$$\mathbf{Q} = \begin{bmatrix} \mathbf{A} & \mathbf{a} \\ \mathbf{0} & 0 \end{bmatrix},$$

where \mathbf{A} is a square matrix of size n and \mathbf{a} is a column vector of size n . Since the rows of the generator matrix sum to zero, the elements of \mathbf{a} can be computed from \mathbf{A} , that is $\mathbf{a} = -\mathbf{A}\mathbf{1}$. Similarly, the first n elements of the initial vector π , denoted by $\boldsymbol{\alpha}$, completely defines the initial vector, since the $(n + 1)$ st element of π is $1 - \boldsymbol{\alpha}\mathbf{1}$. We note that $\boldsymbol{\alpha}$ defines the initial probabilities of the transient states. With the help of this Markov chain, we are ready to define PH distributions.

Definition 1 The time to reach the absorbing state of a Markov chain with a finite number of transient and an absorbing state

$$T = \min\{t : X(t) = n + 1, t \geq 0\},$$

is phase type distributed.

Throughout this document, we assume that the Markov chain starts from one of the transient states and consequently $\boldsymbol{\alpha}\mathbf{1} = 1$, i.e., there is no probability mass at zero and T has a continuous distribution on \mathbb{R}^+ . Since the time to reach the absorbing state is a transient measure of the Markov chain, we can evaluate the distribution of random variable T , based on the transient analysis of the Markov chain with initial distribution π and and generator matrix \mathbf{Q}

$$F_T(t) = P(T < t) = P(X(t) = n + 1) = \pi e^{\mathbf{Q}t} e_{n+1},$$

where e_{n+1} is the $(n + 1)$ st unit vector (the column vector with zero elements except in position $n + 1$ which is one).

This straight forward description of the distribution of T is not widely used due to the redundancy of matrix \mathbf{Q} and vector π . Indeed, matrix \mathbf{A} and the initial vector associated with the transient states, $\boldsymbol{\alpha}$, define all information about the distribution of T and the analytical description based on $\boldsymbol{\alpha}$ and \mathbf{A} is much simpler to use in more complex stochastic models. To obtain the distribution based on $\boldsymbol{\alpha}$ and \mathbf{A} , we carry on the block structure of matrix \mathbf{Q} in the computation.

$$\begin{aligned}
 121 \quad F_T(t) &= P(T < t) = P(X(t) = n + 1) = 1 - \sum_{i=1}^n P(X(t) = n + 1) \\
 122 \quad &= 1 - [\alpha, 0] e^{Q^t} \begin{bmatrix} \mathbf{1} \\ 0 \end{bmatrix} = 1 - [\alpha, 0] \sum_{i=0}^{\infty} \frac{t^i}{i!} \begin{bmatrix} \mathbf{A} & \mathbf{a} \\ \mathbf{0} & 0 \end{bmatrix}^i \begin{bmatrix} \mathbf{1} \\ 0 \end{bmatrix} \\
 123 \quad &= 1 - [\alpha, 0] \sum_{i=0}^{\infty} \frac{t^i}{i!} \begin{bmatrix} \mathbf{A}^i & \bullet \\ \mathbf{0} & 0 \end{bmatrix} \begin{bmatrix} \mathbf{1} \\ 0 \end{bmatrix} = 1 - \alpha \sum_{i=0}^{\infty} \frac{t^i}{i!} \mathbf{A}^i \mathbf{1} = 1 - \alpha e^{A^t} \mathbf{1}, \\
 124 \quad &
 \end{aligned}$$

125 where \bullet indicates irrelevant matrix block whose elements are multiplied by zero.
 126 The PDF of T can be obtained from the derivative of its CDF.

$$\begin{aligned}
 127 \quad f_T(t) &= \frac{d}{dt} F_T(t) = \frac{d}{dt} \left(1 - \alpha \sum_{i=0}^{\infty} \frac{t^i}{i!} \mathbf{A}^i \mathbf{1} \right) = -\alpha \sum_{i=0}^{\infty} \frac{d}{dt} \frac{t^i}{i!} \mathbf{A}^i \mathbf{1} \\
 128 \quad &= -\alpha \sum_{i=1}^{\infty} \frac{t^{i-1}}{(i-1)!} \mathbf{A}^{i-1} \mathbf{A} \mathbf{1} = -\alpha e^{A^t} \mathbf{A} \mathbf{1} = \alpha e^{A^t} \mathbf{a}, \\
 129 \quad &
 \end{aligned}$$

130 where we used $\mathbf{a} = -\mathbf{A} \mathbf{1}$ in the last step.

131 Before computing the remaining properties of PH distributions we need to classify
 132 the eigenvalues of \mathbf{A} . The i, j element of matrix e^{A^t} contains the probability that
 133 starting from transient state i the Markov chain is in transient state j at time t . If
 134 states $1, \dots, n$ are transient states then as t tends to infinity e^{A^t} tends to zero, which
 135 means that the eigenvalues of \mathbf{A} have negative real part and, as a consequence, \mathbf{A} is
 136 non-singular.

137 The Laplace transform of T , $E(e^{-sT})$, can be computed as

$$\begin{aligned}
 138 \quad f_T^*(s) &= E(e^{-sT}) = \int_{t=0}^{\infty} e^{-st} f_T(t) dt = \int_{t=0}^{\infty} e^{-st} \alpha e^{A^t} \mathbf{a} dt \\
 139 \quad &= \alpha \int_{t=0}^{\infty} e^{(-s\mathbf{I} + \mathbf{A})t} dt \mathbf{a} = \alpha (s\mathbf{I} - \mathbf{A})^{-1} \mathbf{a}, \\
 140 \quad &
 \end{aligned}$$

141 where we note that the integral surely converges for $\mathcal{R}(s) \geq 0$ because in this case
 142 the eigenvalues of $-s\mathbf{I} + \mathbf{A}$ also possess a negative real part.

143 To compute the k th moment of T , $E(T^k)$, we need the following integral relation

$$144 \quad [t^k e^{A^t}]_0^{\infty} = \int_{t=0}^{\infty} k t^{k-1} e^{A^t} dt + \int_{t=0}^{\infty} t^k e^{A^t} A dt,$$

146 whose left-hand side is zero because the eigenvalues of \mathbf{A} possess a negative real
 147 part. Multiplying both side with $(-\mathbf{A})^{-1}$ we get

$$\int_{t=0}^{\infty} t^k e^{At} dt = k \int_{t=0}^{\infty} t^{k-1} e^{At} dt (-A)^{-1}.$$

Using this relation, the k th moment of T is

$$\begin{aligned} E(T^k) &= \int_{t=0}^{\infty} t^k f_T(t) dt = \alpha \int_{t=0}^{\infty} t^k e^{At} dt (-A)^{-1} \mathbf{1} = k\alpha \int_{t=0}^{\infty} t^{k-1} e^{At} dt \mathbf{1} \\ &= k(k-1)\alpha \int_{t=0}^{\infty} t^{k-2} e^{At} dt (-A)^{-1} \mathbf{1} = \dots = k! \alpha (-A)^{-k} \mathbf{1}. \end{aligned}$$

These four properties of PH distributions (CDF, PDF, Laplace transform, and moments) have several interesting consequences and some of which we summarize below.

- Matrix $(-A)^{-1}$ has an important stochastic meaning. Let T_{ij} be the time spent in transient state j before moving to the absorbing state when the Markov chain starts from state i . For $E(T_{ij})$, we have

$$E(T_{ij}) = \frac{\delta_{ij}}{-A_{ii}} + \sum_{k, k \neq i} \frac{A_{ik}}{-A_{ii}} E(T_{kj}),$$

where δ_{ij} is the Kronecker delta symbol. The first term of the left-hand side is the time spent in state j while the Markov chain is in the initial state, and the second term is the time spent in state j during later visits to j . Multiplying both sides by $-A_{ii}$ and adding $E(T_{ij}) A_{ii}$ gives

$$0 = \delta_{ij} + \sum_k A_{ik} E(T_{kj}),$$

whose matrix form is

$$\mathbf{0} = \mathbf{I} + A\bar{\mathbf{T}} \quad \longrightarrow \quad \bar{\mathbf{T}} = (-A)^{-1},$$

where $\bar{\mathbf{T}}$ is the matrix composed of the elements $E(T_{ij})$. Consequently, the (ij) element of $(-A)^{-1}$ is $E(T_{ij})$, which is a nonnegative number.

- $f_T^*(s)$ is a rational function of s whose numerator is at most order $n-1$ and denominator is at most order n . This is because

$$\begin{aligned} f_T^*(s) &= \alpha (s\mathbf{I} - A)^{-1} \mathbf{a} = \sum_i \sum_j \alpha_i (s\mathbf{I} - A)^{-1}_{ij} \mathbf{a}_j \\ &= \sum_i \sum_j \alpha_i \left[\frac{\det_{ji}(s\mathbf{I} - A)}{\det(s\mathbf{I} - A)} \right] \mathbf{a}_j = \frac{\sum_i \sum_j \alpha_i \mathbf{a}_j \det_{ji}(s\mathbf{I} - A)}{\det(s\mathbf{I} - A)}. \end{aligned}$$

175 $\det_{ji}(\mathbf{M})$ denotes the determinant of the matrix obtained by removing row j and
 176 column i of matrix \mathbf{M} . The denominator of the last expression is an order n
 177 polynomial of s , while the numerator is the sum of order $n - 1$ polynomials,
 178 which is at most an order $n - 1$ polynomial of s .

179 • This rational Laplace transform representation indicates that a PH distribution
 180 with n transient state can be represented by $2n - 1$ independent parameters. A
 181 polynomial of order n is defined by $n + 1$ coefficients, and a rational function of
 182 order $n - 1$ numerator, and order n denominator is defined by $2n + 1$ parameters.
 183 Normalizing the denominator such that the coefficient of s^n is 1 and considering
 184 that $\int_0^\infty f_T(t)dt = \lim_{s \rightarrow 0} f_T^*(s) = 1$ adds two constraints for the coefficients, from
 185 which the number of independent parameters is $2n - 1$.

186 • The PDF of a PH distribution is the sum of exponential functions. Let $\mathbf{A} = \mathbf{B}^{-1} \Delta \mathbf{B}$
 187 be the Jordan decomposition¹ of \mathbf{A} and let $\mathbf{u} = \alpha \mathbf{B}^{-1}$ and $\mathbf{v} = \mathbf{B} \mathbf{a}$. Then,

$$188 \quad f_T(t) = \alpha e^{\mathbf{A}t} \mathbf{a} = \alpha \mathbf{B}^{-1} e^{\Delta t} \mathbf{B} \mathbf{a} = \mathbf{u} e^{\Delta t} \mathbf{v}.$$

190 At this point, we distinguish two cases.

191 – The eigenvalues of \mathbf{A} are different and Δ is a diagonal matrix. In this case, $f_T(t)$
 192 is a sum of exponential functions because

$$193 \quad f_T(t) = \mathbf{u} e^{\Delta t} \mathbf{v} = \sum_i u_i v_i e^{\lambda_i t} = \sum_i c_i e^{\lambda_i t},$$

195 where $c_i = u_i v_i$ is a constant coefficient of the exponential function.

196 Here the eigenvalues (λ_i) as well as the associated coefficients (c_i) can be real
 197 or complex conjugate pairs. For a complex conjugate pair of eigenvalues, we
 198 have

$$199 \quad c_i e^{\lambda_i t} + \bar{c}_i e^{\bar{\lambda}_i t} = 2|c_i| e^{\mathcal{R}(\lambda_i)t} \cos(\mathcal{I}(\lambda_i)t - \varphi_i),$$

200 where $c_i = |c_i| e^{i\varphi_i}$, $\mathcal{R}(\lambda_i)$ and $\mathcal{I}(\lambda_i)$ are the real and the imaginary part of λ_i
 201 and i is the imaginary unit.

202 – There are eigenvalues of \mathbf{A} with higher multiplicity and Δ contains real Jordan
 203 blocks. The matrix exponent of a Jordan block is

$$204 \quad \exp \left[\left(\begin{array}{cccc} \lambda & 1 & & \\ & \lambda & 1 & \\ & & \ddots & \ddots \\ & & & \lambda \end{array} \right) t \right] = \begin{pmatrix} e^{\lambda t} & t e^{\lambda t} & \frac{1}{2!} t^2 e^{\lambda t} & \frac{1}{3!} t^3 e^{\lambda t} \\ & e^{\lambda t} & t e^{\lambda t} & \frac{1}{2!} t^2 e^{\lambda t} \\ & & \ddots & \ddots \\ & & & e^{\lambda t} \end{pmatrix}.$$

¹The case of different Jordan blocks with identical eigenvalue is not considered here, because it cannot occur in non-redundant PH representations.

Consequently, the density function takes the form

$$f_T(t) = \sum_{i=1}^{\#\lambda} \sum_{j=1}^{\#\lambda_i} c_{ij} t^{j-1} e^{\lambda_i t},$$

where $\#\lambda$ is the number of different eigenvalues and $\#\lambda_i$ is the multiplicity of λ_i . Similar to the previous case, the eigenvalues (λ_i) as well as the associated coefficients ($c_{i,j}$) can be real or complex conjugate pairs. For a complex conjugate pair of eigenvalues, we have

$$c_{i,j} t^{j-1} e^{\lambda_i t} + \bar{c}_{i,j} t^{j-1} e^{\bar{\lambda}_i t} = 2|c_{i,j}| t^{j-1} e^{\mathcal{R}(\lambda_i)t} \cos(\mathcal{I}(\lambda_i)t - \varphi_{i,j}),$$

where $c_{i,j} = |c_{i,j}| e^{i\varphi_{i,j}}$.

As a result of all of these cases, the density function of a PH distribution possesses the form

$$f_T(t) = \sum_{i=1}^{\#\lambda_R} \sum_{j=1}^{\#\lambda_i^R} c_{ij} t^{j-1} e^{\lambda_i^R t} + \sum_{i=1}^{\#\lambda_C} \sum_{j=1}^{\#\lambda_i^C} 2|c_{i,j}| t^{j-1} e^{\mathcal{R}(\lambda_i^C)t} \cos(\mathcal{I}(\lambda_i^C)t - \varphi_{i,j}) \quad (3)$$

where $\#\lambda_R$ is the number of different real eigenvalues and $\#\lambda_C$ is the number of different complex conjugate eigenvalue pairs.

• In general, infinitely many Markov chains can represent the same PH distribution.

– The following *similarity transformation* generates representations with identical size.

Let \mathbf{T} be a non-singular matrix with unit row sums ($\mathbf{T}\mathbf{1} = \mathbf{1}$). The vector–matrix pairs (α, \mathbf{A}) and $(\alpha\mathbf{T}, \mathbf{T}^{-1}\mathbf{A}\mathbf{T})$ are two different vector–matrix representations of the same PH distribution, since

$$F_T(t) = 1 - \alpha\mathbf{T}e^{\mathbf{T}^{-1}\mathbf{A}\mathbf{T}t}\mathbf{1} = 1 - \alpha\mathbf{T}\mathbf{T}^{-1}e^{\mathbf{A}t}\mathbf{T}\mathbf{1} = 1 - \alpha e^{\mathbf{A}t}\mathbf{1}.$$

– Representations with different sizes can be obtained as follows.

Let matrix \mathbf{V} of size $m \times n$ be such that $\mathbf{V}\mathbf{1} = \mathbf{1}$.

The vector–matrix pairs (α, \mathbf{A}) of size n and (γ, \mathbf{G}) of size m are two different vector–matrix representations of the same PH distribution if $\mathbf{A}\mathbf{V} = \mathbf{V}\mathbf{G}$ and $\alpha\mathbf{V} = \gamma$ because

$$F_T(t) = 1 - \gamma e^{\mathbf{G}t}\mathbf{1} = 1 - \alpha\mathbf{V}e^{\mathbf{G}t}\mathbf{1} = 1 - \alpha e^{\mathbf{A}t}\mathbf{V}\mathbf{1} = 1 - \alpha e^{\mathbf{A}t}\mathbf{1}$$

in this case.

3 Matrix Exponential Distributions and Their Basic Properties

In the definition of PH distributions, vector α is a probability vector with nonnegative elements and matrix A is a generator matrix with negative diagonal and nonnegative off-diagonal elements. Relaxing these sign constraints for the vector and matrix elements and maintaining the matrix exponential distribution (and density) function results in the set of matrix exponential (ME) distributions.

Definition 2 Random variable T with distribution function

$$F_T(t) = 1 - \alpha e^{At} \mathbf{1},$$

where α is a finite real vector and A is a finite real matrix, is matrix exponentially distributed.

The size of α and A plays the same role as the number of transient states in case of PH distributions. By definition, the set of PH distributions with a given size is a subset of the set of PH distributions with the same size.

ME distributions share the following basic properties with PH distributions: matrix exponential distribution function, matrix exponential density function, moments, rational Laplace transform, the same set of functions as in (3), and non-unique representation. The main difference between the matrix exponential and the PH classes comes from the fact that the sign constraints on the elements of generator matrixes restrict the eigenvalue structure of such matrixes, while such restrictions do not apply in case of ME distributions. For example, the eigenvalues of an order three PH distribution with dominant eigenvalue θ satisfy $\mathcal{R}(\lambda_i) \leq \theta$ and $|\mathcal{I}(\lambda_i)| \leq (\theta - \mathcal{R}(\lambda_i))/\sqrt{3}$, while the eigenvalues of an order three ME distribution with dominant eigenvalue θ satisfy $\mathcal{R}(\lambda_i) \leq \theta$ only. This flexibility of the eigenvalues has significant consequence on the flexibility of the set of order three PH and ME distributions. For example, the minimal squared coefficient of variation among the order three PH and ME distributions are 1/3 and 0.200902, respectively.

The main difficulty encountered when working with ME distributions is that a general vector–matrix pair does not always define a nonnegative density function, while a vector–matrix pair with the sign constraints of PH distributions does. Efficient numerical methods have been proposed recently to check the nonnegativity of a matrix exponential function defined by a general vector–matrix pair, but general symbolic conditions are still missing.

4 Analysis of Models with PH Distributed Durations

If all durations (service times, interarrival times, repair times, etc.) in a system are distributed according to PH distributions, then its overall behavior can be captured by a continuous time Markov chain, referred to as extended Markov chain (EMC).

274 In this section, we show how to derive the infinitesimal generator of this EMC using
 275 Kronecker operations. The methodology here described has been originally presented
 276 in the case of Discrete PHs in [17] and more recently in the case of Continuous PHs
 277 in [13]

278 To this end we first introduce the notation used to describe the model. By \mathcal{S} , we
 279 denote the set of states and by $N = |\mathcal{S}|$ the number of states. The states themselves
 280 are denoted by s_1, s_2, \dots, s_N . The set of activities is denoted by \mathcal{A} and the set of those
 281 that are active in state s_i is denoted by \mathcal{A}_i . The activities are denoted by a_1, a_2, \dots, a_M
 282 with $M = |\mathcal{A}|$. When activity a_i is completed in state s_j then the system moves from
 283 state s_j to state $n(j, i)$, i.e., n is the function that provides the next state. We assume
 284 that the next state is a deterministic function of the current state and the activity that
 285 completes. We further assume that there does not exist a triple, k, i, j , for which
 286 $s_k \in \mathcal{S}$, $a_i \in \mathcal{A}$, $a_j \in \mathcal{A}$ and $n(k, i) = n(k, j)$. These two assumptions, which make
 287 the formulas simpler, are easy to relax in practice. There can be activities that end
 288 when the system moves from state s_i to state s_j even if they do not complete and are
 289 active both in s_i and in s_j . These activities are collected in the set $e(i, j)$. The PH
 290 distribution that is associated with activity a_i is characterized by the initial vector
 291 α_i and matrix A_i . As before, we use the notation $\mathbf{a}_i = -A_i \mathbf{1}$ to refer to the vector
 292 containing the intensities that lead to completion of activity a_i . The number of phases
 293 of the PH distribution associated with activity a_i is denoted by n_i .

294 *Example 1 PH/PH/1/K queue with server break-downs.* As an example, we consider,
 295 using the above-described notation, a queue in which the server is subject to failure
 296 only if the queue is not empty. The set of states is $\mathcal{S} = \{s_1, s_2, \dots, s_{2K+1}\}$ where s_1
 297 represents the empty queue, s_{2i} with $1 \leq i \leq K$ represents the state with i clients
 298 in the queue and the server up, and s_{2i+1} with $1 \leq i \leq K$ represents the state with
 299 i clients and the server down. There are four activities in the system: a_1 represents
 300 the arrival activity, a_2 the service activity, a_3 the failure activity,² and a_4 the repair
 301 activity. The vectors and matrices that describe the associated PH distributions are
 302 $\alpha_1, \alpha_2, \alpha_3, \alpha_4$ and A_1, A_2, A_3, A_4 . In this example, we assume that the arrival
 303 activity is active if the system is not full and it is inactive if the system is full.
 304 The service and the failure activities are active if the queue is not empty and the
 305 server is up. The repair activity is active if the queue is not empty and the server is
 306 down. Accordingly, we have $\mathcal{A}_1 = \{a_1\}$, $\mathcal{A}_{2i} = \{a_1, a_2, a_3\}$ for $1 \leq i \leq K - 1$,
 307 $\mathcal{A}_{2i+1} = \{a_1, a_4\}$ for $1 \leq i \leq K - 1$, $\mathcal{A}_{2K} = \{a_2, a_3\}$, and $\mathcal{A}_{2K+1} = \{a_4\}$. The next
 308 state function is as follows: for arrivals we have $n(1, 1) = s_2$ and $n(i, 1) = s_{i+2}$ with
 309 $2 \leq i \leq 2K - 1$; for services $n(2, 2) = s_1$ and $n(2i, 2) = s_{2i-2}$ with $2 \leq i \leq K$;
 310 for failures $n(2i, 3) = s_{2i+1}$ with $1 \leq i \leq K$; for repairs $n(2i + 1, 4) = s_{2i}$ with
 311 $1 \leq i \leq K$. We assume that the failure activity ends every time when a service
 312 activity completes, i.e., failure is connected to single jobs and not to the aging of the
 313 server. Other activities end only when they complete or when such a state is reached
 314 in which they are not active. Accordingly, $e(2i, 2i - 2) = \{a_3\}$ for $2 \leq i \leq K$.

²Failure is more like an event than an activity but, in order to keep the discussion clearer, we refer to it as failure activity.

Based on the description of the ingredients of the model, it is possible to derive blocks of the initial probability vector and the blocks of the infinitesimal generator of the corresponding CTMC. Let us start with the infinitesimal generator, which we denote by \mathbf{Q} , composed of $N \times N$ blocks. The block of \mathbf{Q} that is situated in the i th row of blocks and in the j th column of blocks is denoted by \mathbf{Q}_{ij} . A block in the diagonal, \mathbf{Q}_{ii} describes the parallel execution of the activities that are active in s_i . The parallel execution of CTMCs can be captured by the Kronecker-sum operator (\oplus), and thus we have

$$\mathbf{Q}_{ii} = \bigoplus_{j: s_j \in \mathcal{A}_i} \mathbf{A}_j .$$

An off-diagonal block, \mathbf{Q}_{ij} , is not a zero matrix only if there exists an activity whose completion moves the system from state s_i to state s_j . Let us assume that the completion of activity a_k moves the system from state s_i to state s_j , i.e., $n(i, k) = s_j$. The corresponding block, \mathbf{Q}_{ij} , must

- reflect the fact that activity a_k completes and restarts if a_k is active in s_j ,
- reflect the fact that activity a_k completes and does not restart if a_k is not active in s_j ,
- end activities that are active in s_i but not in s_j ,
- start those activities that are not active in s_i but are active in s_j ,
- end and restart those activities that are active both in s_i and in s_j but are in $e(i, j)$,
- and maintain the phase of those that are active both in s_i and in s_j and are not in $e(i, j)$.

The joint treatment of the above cases can be carried out by the Kronecker-product operator and thus we have

$$\mathbf{Q}_{ij} = \bigotimes_{l: 1 \leq l \leq M} \mathbf{R}_l$$

with

$$\mathbf{R}_l = \begin{cases} \mathbf{a}_k & \text{if } l = k \text{ and } k \notin \mathcal{A}_j \\ \mathbf{a}_k \boldsymbol{\alpha}_k & \text{if } l = k \text{ and } k \in \mathcal{A}_j \\ \mathbf{1}_{n_i} & \text{if } l \neq k \text{ and } k \in \mathcal{A}_i \text{ and } k \notin \mathcal{A}_j \\ \boldsymbol{\alpha}_l & \text{if } l \neq k \text{ and } k \notin \mathcal{A}_i \text{ and } k \in \mathcal{A}_j \\ \mathbf{1}_{n_i} \boldsymbol{\alpha}_l & \text{if } l \neq k \text{ and } k \in \mathcal{A}_i \text{ and } k \in \mathcal{A}_j \text{ and } k \in e(i, j) \\ \mathbf{I}_{n_i} & \text{if } l \neq k \text{ and } k \in \mathcal{A}_i \text{ and } k \in \mathcal{A}_j \text{ and } k \notin e(i, j) \\ \mathbf{1} & \text{otherwise} \end{cases}$$

where the subscripts to $\mathbf{1}$ and \mathbf{I} indicate their size.

The initial probability vector of the CTMC, $\boldsymbol{\pi}$, is a row vector composed of N blocks which must reflect the initial probabilities of the states of the system and the initial probabilities of the PH distributions of the active activities. Denoting by π_i the initial probability of state s_i , the i th block of the initial probability vector, $\boldsymbol{\pi}_i$, is given as

$$\pi_i = \bigotimes_{j:s_j \in \mathcal{A}_i} \alpha_j .$$

Example 2 For the previous example, the diagonal blocks, which must reflect the ongoing activities, are the following:

$$\begin{aligned} \mathcal{Q}_{1,1} &= A_1, \quad \mathcal{Q}_{2i,2i} = A_1 \oplus A_2 \oplus A_3, \quad \mathcal{Q}_{2i+1,2i+1} = A_1 \oplus A_4, \\ \mathcal{Q}_{2K,2K} &= A_2 \oplus A_3, \quad \mathcal{Q}_{2K+1,2K+1} = A_4 \quad \text{with } 1 \leq i \leq K-1 \end{aligned}$$

Arrival in state s_1 takes the system to state s_2 . The corresponding block must complete and restart the arrival activity and must restart both the service and the failure activity:

$$\mathcal{Q}_{12} = a_1 \alpha_1 \otimes \alpha_2 \otimes \alpha_3 \quad (4)$$

Arrival in state s_{2i} (server up) takes the system to state s_{2i+2} . If the system does not become full then the corresponding block must complete and restart the arrival activity and must maintain the phase of both the service and the failure activity. If the system becomes full, the arrival activity is not restarted. Accordingly, we have

$$\begin{aligned} \mathcal{Q}_{2i,2i+2} &= a_1 \alpha_1 \otimes \mathbf{I}_{n_2} \otimes \mathbf{I}_{n_3} \quad \text{with } 1 \leq i \leq K-2 \\ \mathcal{Q}_{2K-2,2K} &= a_1 \otimes \mathbf{I}_{n_2} \otimes \mathbf{I}_{n_3} \end{aligned}$$

An arrival in state s_{2i+1} (server down) takes the system to state s_{2i+3} . If the system does not become full then the corresponding block must complete and restart the arrival activity and must maintain the phase of the repair activity. If the system becomes full, the arrival activity is not restarted. Accordingly, we have

$$\begin{aligned} \mathcal{Q}_{2i+1,2i+3} &= a_1 \alpha_1 \otimes \mathbf{I}_{n_4} \quad \text{with } 1 \leq i \leq K-2 \\ \mathcal{Q}_{2K-1,2K+1} &= a_1 \otimes \mathbf{I}_{n_4} \end{aligned}$$

Service completion can take place in three different situations. If the system becomes empty then the phase of the arrival activity is maintained, the service activity is completed and the failure activity is put to an end. If the system neither becomes empty nor was full then the phase of the arrival activity is maintained, the service activity is completed and restarted, and the failure activity ends and restarts. Finally, if the queue was full then the arrival activity is restarted, the service activity is completed and restarted, and the failure activity is put to an end and restarted. Accordingly, we have

$$\mathcal{Q}_{2,1} = \mathbf{I}_{n_1} \otimes a_2 \otimes \mathbf{1}_{n_3}$$

$$\begin{aligned}
 380 \quad Q_{2i,2i-1} &= \mathbf{I}_{n_1} \otimes a_2 \alpha_2 \otimes \mathbf{I}_{n_3} \alpha_3 \quad \text{with } 1 < i < K \\
 381 \quad Q_{2K,2K-2} &= \alpha_1 \otimes a_2 \alpha_2 \otimes \mathbf{I}_{n_3} \alpha_3 \quad (5) \\
 382
 \end{aligned}$$

383 The failure activity can be completed in two different situations. If the system is not
 384 full, then the phase of the arrival activity is maintained. If the system is full then
 385 the arrival activity is not active. In both cases, the service activity ends, the failure
 386 activity is completed and the repair activity is initialized.

$$\begin{aligned}
 387 \quad Q_{2i,2i+1} &= \mathbf{I}_{n_1} \otimes \mathbf{I}_{n_2} \otimes a_3 \otimes \alpha_4 \quad \text{with } 1 \leq i < K \\
 388 \quad Q_{2K,2K+1} &= \mathbf{I}_{n_2} \otimes a_3 \otimes \alpha_4 \\
 389
 \end{aligned}$$

390 Similarly to the failure activity, also the repair activity can be completed in two
 391 different situations because the arrival activity can be active or inactive. In both
 392 cases, the service activity and the failure activity must be initialized and the repair
 393 activity completes.

$$\begin{aligned}
 394 \quad Q_{2i+1,2i} &= \mathbf{I}_{n_1} \otimes \alpha_2 \otimes \alpha_3 \otimes a_4 \quad \text{with } 1 \leq i < K \\
 395 \quad Q_{2K+1,2K} &= \alpha_2 \otimes \alpha_3 \otimes a_4 \\
 396
 \end{aligned}$$

397 5 Analysis of Stochastic Systems with ME Distributed 398 Durations

399 The most important observation to take from this section is that all steps of the method
 400 of EMCs (as explained in the previous section) remain directly applicable in case
 401 of ME distributed durations (where the (α_i, A_i) vector–matrix pairs describe ME
 402 distributions). In that case, the only difference is that the signs of the vector and matrix
 403 elements are not restricted to be nonnegative in case of the vector elements and off-
 404 diagonal matrix elements and to be negative in case of the diagonal matrix elements.
 405 Consequently, the model description does not allow a probabilistic interpretation via
 406 Markov chains.

407 This general conclusion was obtained through serious research efforts. Following
 408 the results in [12], it was suspected that in a stochastic model ME distributions could
 409 be used in place of PH distributions and several results would carry over, but it was
 410 not easy to prove these conjectured results in the general setting because probabilistic
 411 arguments associated with PH distributions no longer hold. In [1], it was shown that
 412 matrix geometric methods can be applied for quasi-birth–death processes (QBDs)
 413 with rational arrival processes (RAPs) [3], which can be viewed as an extension of
 414 ME distributions to arrival processes. To prove that the matrix geometric relations
 415 hold, the authors of [1] use an interpretation of RAPs proposed in [3]. However,
 416 the models considered are limited to QBDs. For the model class of SPNs with ME

417 distributed firing times, the applicability of the EMC-like analysis was proved in [2]
 418 and refined for the special case when the ME distribution has no PH representation
 419 in [4].

420 6 Analysis tools

421 Based on the common representation of the EMC through the Kronecker algebra,
 422 smart algorithms have been developed recently to optimize memory usage. These
 423 algorithms build the EMC in a completely symbolic way, both at the process state
 424 space level and at the expanded state space level, as deeply explained in [13] that we
 425 use as reference.

426 The algorithm presented in [13] is based on two high level steps:

- 427 1. to generate the reachability graph of the model (which collects the system states
 428 in a graph according to their reachability from an initial set of states) using a
 429 symbolic technique;
- 430 2. to enrich the symbolically stored reachability graph with all the necessary infor-
 431 mation to evaluate Kronecker expressions representing the expanded state space.

432 Step 1 is performed using symbolic technique based on complex data structures like
 433 Multi-Valued Decision Diagram (MDD) [18] to encode the model state space; step
 434 2 adds information related to each event memory policy to the encoded state space.
 435 In manner it is possible to use on the fly expressions introduced in Sects. 4 and 5 to
 436 compute various probability measures of the model.

437 6.1 Symbolic Generation of Reachability Graph

438 Both traditional performance or dependability evaluation techniques and more recent
 439 model checking-based approaches are grounded in the knowledge of the set of states
 440 that the system considered can reach starting from a particular initial state (or in
 441 general from a set of initial states). Symbolic techniques [5] focus on generating
 442 a compact representation of huge state spaces by exploiting a model's structure
 443 and regularity. A model has a structure when it is composed of K sub-models, for
 444 some $K \in \mathbb{N}$. In this case, a global system state can be represented as a K -tuple
 445 (q^1, \dots, q^K) , where q^k is the local state of sub-model k (having some finite size n^k).

446 The use of (MDDs) for the encoding of model state spaces was introduced by
 447 Miner and Ciardo in [14]. MDDs are rooted, directed, acyclic graphs associated with
 448 a finite ordered set of integer variables. When used to encode a state space, an MDD
 449 has the following structure:

- 450 • nodes are organized into $K + 1$ levels, where K is the number of sub-models;
- 451 • level K contains only a single non-terminal node, the root, whereas levels $K - 1$
- 452 through 1 contain one or more non-terminal nodes;
- 453 • a non-terminal node at level k has n^k arcs pointing to nodes at level $k - 1$;

454 A state $s = (q^1, \dots, q^K)$ belongs to S if and only if a path exists from the root node
 455 to the terminal node 1 such that, at each node, the arc corresponding to the local
 456 state q^k is followed. In [6], and then in [7], Ciardo et al. proposed the *Saturation*
 457 algorithm for the generation of reachability graphs using MDDs. Such an iteration
 458 strategy improves both memory and execution time efficiency.

459 An efficient encoding of the reachability graph is built in the form of a set of
 460 Kronecker matrices $\mathbf{W}_{e,k}$ with $e \in \mathcal{A}$ and $k = 1, \dots, K$, where \mathcal{A} is the set collecting
 461 all the system events or activities. $\mathbf{W}_{e,k}[i_k, j_k] = 1$ if state j_k of sub-model k is
 462 reachable from state i_k due to event e . According to such a definition, the next
 463 state function of the model can be encoded as the incidence matrix given by the
 464 boolean sum of Kronecker products $\sum_{e \in \mathcal{A}} \otimes_{K \geq k \geq 1} \mathbf{W}_{e,k}$. As a consequence, the
 465 matrix representation \mathbf{R} of the reachability graph of the model can be obtained by
 466 filtering the rows and columns of such a matrix corresponding to the reachable global
 467 states encoded in the MDD and replacing each non-null element with the labels of
 468 the events that cause the corresponding state transition.

469 Saturation Unbound is a very effective way to represent the model state space
 470 and the related reachability graph of a model. In any case, the methodology we are
 471 dealing with is not strictly dependent on any particular algorithm to efficiently store
 472 the reachability graph. We refer to the *Saturation Unbound* algorithm simply because
 473 its efficiency is well known [7].

474 6.2 Annotating the Reachability Graph

475 The use of Saturation together with the Kronecker representation presented in previ-
 476 ous sections enable solution of the derived stochastic process. However, knowledge
 477 of the reachability graph of the untimed system as produced by Saturation is not
 478 sufficient to manage the infinitesimal generator matrix \mathbf{Q} on the fly according to the
 479 symbolic representation. Considering that the information about the enabled events
 480 for all the system states is contained in the high level description of the model and it
 481 can be evaluated on the fly when needed with a negligible overhead, the only addi-
 482 tional information needed is knowledge about the sets of active but not enabled events
 483 in each state s ($T_a^{(s)}$). Using Saturation for the evaluation of the reachability graph
 484 requires an additional analysis step for the computation of such an information and
 485 use of a different data structure for storage. Multi Terminal Multi-Valued Decision
 486 Diagram (MTMDD) [15] is used for this purpose.

487 The main differences with respect to MDDs are that: (1) more than two terminal
 488 nodes are present in an MTMDD and (2) such nodes can be labeled with arbitrary
 489 integer values, rather than just 0 and 1. An MTMDD can efficiently store both the

490 system state space \mathcal{S} and the sets $T_a^{(s)}$ of active but not enabled events for all $s \in \mathcal{S}$;
 491 this is necessary, in our approach, to correctly evaluate non-null blocks of \mathbf{Q} matrix.
 492 In fact, while an MDD is only able to encode a state space, an MTMDD is also
 493 able to associate an integer to each state. Thus, the encoding of sets $T_a^{(s)}$ can be
 494 done associating to each possible set of events an integer code that unambiguously
 495 represents it. Let us associate to each event an unique index n such that $1 \leq n \leq \|\mathcal{A}\|$.
 496 Then the integer value associated to one of the possible sets $T_a^{(s)}$ is computed starting
 497 from the indices associated with the system events that belong to it in the following
 498 way:

$$b_M \cdot 2^A + \dots + b_n \cdot 2^n + \dots + b_1 \cdot 2^1 + 1 = \sum_{i=1}^M b_i 2^i + 1$$

500 where

$$b_i = \begin{cases} 1, & \text{if event } e_i \in T_a^{(s)} \\ 0, & \text{otherwise} \end{cases}$$

502 In this manner all the necessary information to apply the Kronecker-based expres-
 503 sions on the fly are provided; the only remaining need is a method to evaluate the set
 504 $T_a^{(s)}$ given a referring state s .

505 In [13], the following theorem has been proved.

506 **Theorem 1** *Given a model \mathcal{M} , a state $s_0 \in \mathcal{S}$ and an event $\bar{e} \in \mathcal{A}$ with an age*
 507 *memory policy associated, then $\bar{e} \in T_a^{(s_0)}$ iff $\bar{e} \notin T_e^{(s_0)}$ and one of the following*
 508 *statements holds:*

- 509 1. $\exists s_1 \in \mathcal{S}, \exists e_1 \in \mathcal{A}, s_1 \neq s_0, e_1 \neq \bar{e} \mid s_0 \in \mathcal{N}_{e_1}(s_1) \wedge \bar{e} \in T_e^{(s_1)}$
- 510 2. $\exists s_1 \in \mathcal{S}, s_1 \neq s_0 \mid s_0 \in \mathcal{N}(s_1) \wedge \bar{e} \in T_a^{(s_1)}$

511 where \mathcal{N}_{e_1} is the next state function associated to event e_1 .

512 Note that function \mathcal{N} is the equivalent to the $n(\cdot, \cdot)$ defined in Sect. 4; function \mathcal{N}_e
 513 instead differs for the restriction to the firing of a specific event e . We use this notation
 514 because it is less cumbersome in this specific context.

515 Theorem 1 gives a way to evaluate if an event e belongs to the set $T_a^{(s_0)}$ or not.
 516 In fact, according to the statements proved, it is possible to characterize a state
 517 s_0 with respect to the system event memory policies by exploring its reachability
 518 graph. Exploration can be performed using classical bread-first search and depth-
 519 first search algorithms, easily applicable to an explicitly stored reachability graph; it
 520 is more complicated to apply classical search algorithms when the graph is stored in
 521 implicit manner as is the case when MTMDD data structures are used.

522 In this case, a different approach can be used by resorting to Computational Tree
 523 Logic (CTL) formulas that have been shown to be very efficient for data structures
 524 like MDD and MTMDD. The use of CTL formulas to evaluate sets $T_a^{(s)}$ is justified
 525 by a theorem introduced in [13]. Before recalling this theorem, we need to introduce
 526 a CTL operator.

527 **Definition 3** Let $s_0 \in \mathcal{S}$ be a state of a discrete state process with state space \mathcal{S} , and
 528 let p and q be two logical conditions on the states. Let also $\mathcal{F}(s) \subseteq \mathcal{N}(s) \cup \mathcal{N}^{-1}(s)$ be
 529 a reachability relationship between two states in \mathcal{S} that defines a desired condition

530 over the paths. Then s_0 satisfies the formula $E_{\mathcal{F}}[pUq]$, and we will write $s_0 \models$
 531 $E_{\mathcal{F}}[pUq]$, iff $\exists n \geq 0, \exists s_1 \in \mathcal{F}(s_0), \dots, \exists s_n \in \mathcal{F}(s_{n-1}) \mid (s_n \models q) \wedge (\forall m <$
 532 $n, s_m \models p)$.

533 In definition above, we used the path quantifier E with the meaning *there exists a*
 534 *path* and the tense operator U with the meaning *until*, as usually adopted in CTL
 535 formulas.

536 Given Definition 3, the following theorem holds:

537 **Theorem 2** *An event $\bar{e} \in \mathcal{E}$, with an age memory policy associated, belongs to*
 538 $T_a^{(s_0)}$, *with $s_0 \in \mathcal{S}$, iff $s_0 \models E_{\mathcal{F}}[pUq]$ over a path at least one long, where*
 539 *p and q are the statements “ \bar{e} is not enabled” and “ \bar{e} is enabled,” respectively,*
 540 *and $\mathcal{F}(s) = \mathcal{N}^{-1}(s) \setminus \mathcal{N}_{\bar{e}}^{-1}(s)$.*

541 Thanks to Theorem 2, evaluation of the CTL formula $E_{\mathcal{F}}[pUq]$ makes possible
 542 to evaluate whether an event \bar{e} is active but not enabled in state s_0 or not by setting
 543 condition p as \bar{e} is not enabled and q as \bar{e} is enabled. This is the last brick to build
 544 an algorithm able to compute state probabilities of a model, where the event are PH
 545 or ME distributed; in fact, it is possible to characterize all the active and/or enabled
 546 events in all the different states and to apply the Kronecker expressions with this
 547 information to solve the derived EMC.

548 7 Examples

549 In this section, we present two examples where non-exponentially distributed dura-
 550 tions are present. In the first example, these durations are approximated by PH dis-
 551 tributions, while in the second example they are described by ME distributions.

552 7.1 Reliability Model of Computer System

553 We introduce a reliability model where we use PH distributions as failure times.
 554 The model is specified according to the Petri net depicted in Fig. 1, where the usual
 555 graphical notation for the places, transitions, and arcs has been adopted.

556 The system under study is a distributed computing system composed of a cluster
 557 of two computers. Each of them has three main weak points: the motherboard, CPU,
 558 and disk. Interconnections inside the cluster are provided by a manager in such a
 559 way that the overall system is seen as a single unit. In the distributed system, the two
 560 computers work independently, driven by the manager that acts as a load balancer to
 561 split the work between them. Since the manager represents a single point of failure, a
 562 second instance is deployed for redundancy in the system; this latter instance operates
 563 in cold standby when the main computer manager works and it is powered on when
 564 it fails.

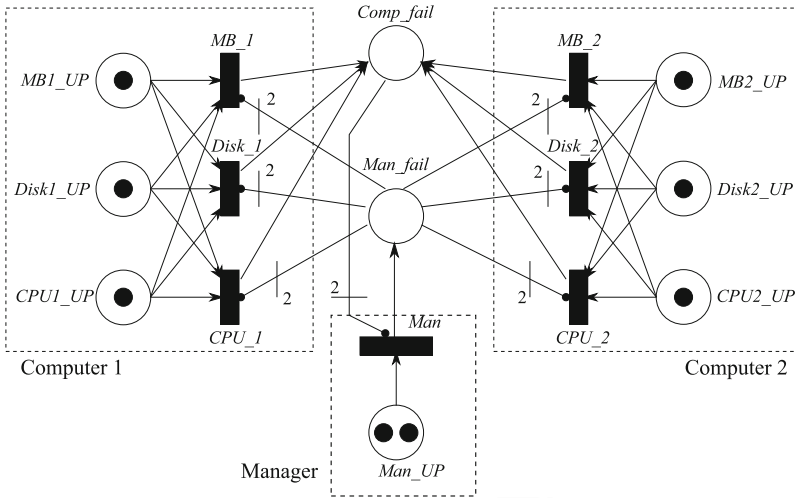


Fig. 1 Computer system reliability model

565 Due to this configuration, the distributed system works when at least one of the two
 566 computers works and the computer manager properly operates. The main components
 567 of each computational unit (CPU, motherboard, and disk) may fail rendering the
 568 unit inoperable. In the Petri net model, faults in the CPU, motherboard, and disk
 569 are modeled by the timed transitions MB_i , $Disk_i$, and CPU_i whose firing
 570 represents the respective faulty event in the i -th Computer; the operating conditions
 571 of components are represented by a token in the places $CPUi_UP$, MBi_UP , and
 572 $Disk_UP$. When one of the transitions above fires a token is flushed out of the
 573 place and a token is put in the place $Comp_fail$. At the same time, all the other
 574 transitions related to the faulty events in the same unit become disabled because
 575 the unit is considered down and thus no more faults can occur. Two tokens in the
 576 place $Comp_fail$ means that the two computational units are both broken and the
 577 overall distributed system is not operational. Similarly, transition Man models the
 578 fault of a manager unit. Its firing flushes a token out of the place Man_UP and
 579 puts a token in the place Man_fail . Thanks to the redundancy, the first manager unit
 580 fault is tolerated whereas the system goes down when a second fault occurs. This
 581 state is represented in the Petri net by two tokens in the place Man_fail . In both
 582 faulty states, all the transitions are disabled and an absorbing state is reached. In
 583 terms of Petri net objects, the not operational condition is expressed by the following
 584 statement:

$$585 \quad (\#Comp_fail = 2) \vee (\#Man_fail = 2), \quad (6)$$

586 where the symbol $\#P$ states the number of token in place P .

Table 1 Failure time distribution parameters

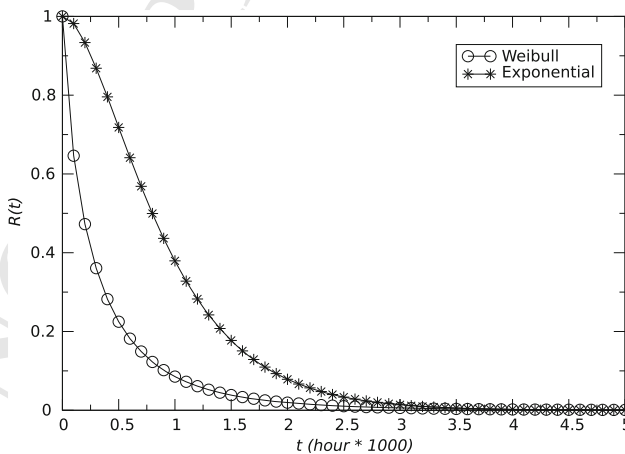
Transition	Weibull			
	β_f	η_f	E	λ
MB_1, MB_2	0.5965	1.20	1.82	0.55
$Disk_1, Disk_2$	0.5415	1.00	1.71	0.59
CPU_1, CPU_2	0.399	1.46	3.42	0.29
Man	0.5965	1.20	1.82	0.55

As usual in reliability modeling, the time to failure of the components has been modeled using Weibull distributions whose cumulative distribution function is

$$F(t) = 1 - e^{-(t/\eta_f)^{\beta_f}}.$$

This choice has been also supported by measures done on real systems such as those analyzed in [9]. The parameters of the Weibull distributions used for the Petri net transitions of Fig. 1 are reported in Table 1.

Weibull distributions have been introduced in the model through the use of 10-phase PH distributions, approximating them by evaluating the formula (6). The results obtained are depicted in Fig. 2. To better highlight the usefulness of the modeling approach presented here, the Petri net model was solved by imposing exponential distributions as transition firing times. In fact, the use of exponential distributions is quite common to obtain a more tractable model. The value of the parameters λ used in this second run was computed as the reciprocal of the expected value, E , of the corresponding Weibull distributions (listed in Table 1). The result obtained are also depicted in Fig. 2. As can be easily noted, the use of exponential distributions produces optimistic results compared to the use of Weibull distributions, making the system appear more reliable than it is in reality.

**Fig. 2** Computer system reliability $R(t)$

623

$$\begin{array}{cccccc}
 0 & 0 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0 & 0 \\
 a_2 \otimes I_3 & 0 & 0 & 0 & 0 \\
 0 & a_2 \otimes A_5 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0 & 0 \\
 0 & 0 & a_2 \otimes I_4 \otimes A_5 & 0 & 0 \\
 I_3 \otimes a_4 & 0 & a_3 \otimes I_4 \otimes A_5 & 0 & 0 \\
 A_3 & a_3 \otimes A_5 & 0 & 0 & 0 \\
 0 & A_5 & 0 & 0 & a_5 \\
 0 & a_4 \otimes I_5 & A_4 \oplus A_5 & I_4 \otimes a_5 & 0 \\
 0 & 0 & 0 & A_4 & a_4 \\
 0 & 0 & 0 & 0 & 0
 \end{array}$$

624

The vector that provides the initial configuration is $|A_1 \otimes A_2, 0, \dots, 0|$.

625

In order to illustrate a feature of ME distributions that cannot be exhibited by PH distributions, we applied an ME distribution with “oscillating” PDF to describe the duration of activities 1, 2, 4, and 5. The vector–matrix pair of this ME distribution is

626

628

$$A_1 = A_2 = A_4 = A_5 = |1.04865, -0.0340166, -0.0146293|,$$

629

630

$$A_1 = A_2 = A_4 = A_5 = \begin{vmatrix} -1 & 0 & 0 \\ 0 & -1 & -20 \\ 0 & 20 & -1 \end{vmatrix},$$

631

and its PDF is depicted in Fig. 5. The duration of the remaining activity, namely activity 3, is distributed according to an Erlang distribution with four phases and average execution time equal to 1, i.e.,

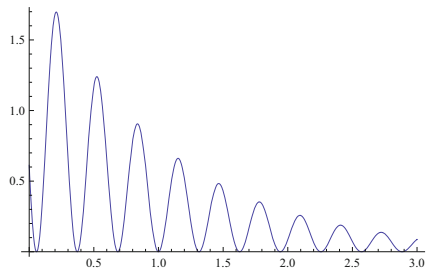
632

633

634

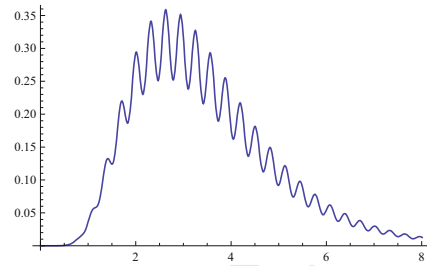
$$A_3 = |1, 0, 0, 0|, \quad A_3 = \frac{1}{4} \begin{vmatrix} -1 & 1 & 0 & 0 \\ 0 & -1 & 1 & 0 \\ 0 & 0 & -1 & 1 \\ 0 & 0 & 0 & -1 \end{vmatrix}.$$

Fig. 5 Oscillating activity duration pdf



REV IN PRINT

Fig. 6 Overall accomplishment time pdf



635 The model was then used to characterize the PDF of the time that is needed to
 636 accomplish the whole mission. The resulting PDF is shown in Fig.6 and one can
 637 observe that the oscillating nature of the distribution of the activity durations carries
 638 over into the overall completion time distribution.

639 8 Conclusions

640 While the evolution of computing devices and analysis methods resulted in a sharp
 641 increase in the complexity of computable CTMC models, CTMC-based analysis
 642 had been restricted to the analysis of stochastic models with exponentially distrib-
 643 uted duration times. A potential extension of CTMC-based analysis is the inclusion
 644 of PH distributed duration times, which enlarges the state space, but still has a fea-
 645 sible computational complexity. We surveyed the basics of PH distributions and the
 646 analysis approach to generate the EMC.

647 A more recent development in this field is the extension of the EMC-based analy-
 648 sis with ME distributed duration times. With respect to the steps of the analysis
 649 method, the EMC-based analysis and its extension with ME distributions are identi-
 650 cal. However, because ME distributions are more flexible than the PH distributions
 651 (more precisely, the set of PH distributions of a given size is a subset of the set of
 652 ME distributions of the same size) this extension increases the modeling flexibility
 653 of the set of models which can be analyzed with a given computational complexity.

654 Apart of the steps of the EMC-based analysis method, we discussed the tool sup-
 655 port available for the automatic execution of the analysis method. Finally, application
 656 examples demonstrate the abilities of the modeling and analysis methods.

657 References

- 658 1. Asmussen S, Bladt M (1999) Point processes with finite-dimensional conditional probabilities.
 659 Stoch Process Appl 82:127–142
- 660 2. Bean NG, Nielsen BF (2010) Quasi-birth-and-death processes with rational arrival process
 661 components. Stoch Models 26(3):309–334

- 662 3. Buchholz P, Horvath A, Telek M (2011) Stochastic Petri nets with low variation matrix exponentially distributed firing time. *Int J Perform Eng* 7:441–454, 2011 (Special issue on performance and dependability modeling of dynamic systems)
- 663
- 664 4. Buchholz P, Telek M (2010) Stochastic petri nets with matrix exponentially distributed firing times. *Perform Eval* 67:1373–1385
- 665
- 666 5. Burch JR, Clarke EM, McMillan KL, Dill DL, Hwang LJ (1990) Symbolic model checking: 1020 states and beyond. In: Fifth annual IEEE symposium on logic in computer science, 1990. LICS '90, Proceedings, pp 428–439
- 667
- 668 6. Ciardo G, Luttgen G, Siminiceanu R (2001) Saturation: an efficient iteration strategy for symbolic state space generation. In: Proceedings of tools and algorithms for the construction and analysis of systems (TACAS), LNCS 2031. Springer, pp 328–342
- 669
- 670 7. Ciardo G, Marmorstein R, Siminiceanu R (2003) Saturation unbound. In: Proceedings of TACAS. Springer, pp 379–393
- 671
- 672 8. Cox DR (1955) The analysis of non-markovian stochastic processes by the inclusion of supplementary variables. *Proc Cambridge Philos Soc* 51(3):433–441
- 673
- 674 9. Distefano S, Longo F, Scarpa M, Trivedi KS (2014) Non-markovian modeling of a blade-center chassis midplane. In: Computer performance engineering, vol 8721 of Lecture Notes in Computer Science. Springer International Publishing, pp 255–269
- 675
- 676 10. Kleinrock L (1975) Queuing systems, vol 1: theory. Wiley Interscience, New York
- 677
- 678 11. Kulkarni VG (1995) Modeling and analysis of stochastic systems. Chapman & Hall
- 679
- 680 12. Lipsky L (2008) Queueing theory: a linear algebraic approach. Springer
- 681
- 682 13. Longo F, Scarpa M (2015) Two-layer symbolic representation for stochastic models with phase-type distributed events. *Int J Syst Sci* 46(9):1540–1571
- 683
- 684 14. Miner AS, Ciardo G (1999) Efficient reachability set generation and storage using decision diagrams. In: Application and Theory of Petri Nets 1999 (Proceedings 20th international conference on applications and theory of Petri Nets. Springer, pp 6–25)
- 685
- 686 15. Miner A, Parker D (2004) Symbolic representations and analysis of large state spaces. In: Validation of stochastic systems, LNCS 2925, Dagstuhl (Germany). Springer, pp 296–338
- 687
- 688 16. Neuts M (1975) Probability distributions of phase type. In: Amicorum L, Florin EH (eds) University of Louvain, pp 173–206
- 689
- 690 17. Scarpa M, Bobbio A (1998) Kronecker representation of stochastic petri nets with discrete ph distributions. In: Proceedings of IEEE international computer performance and dependability symposium, 1998. IPDS'98. pp 52–62
- 691
- 692 18. Srinivasan A, Ham T, Malik S, Brayton RK (1990) Algorithms for discrete function manipulation. In: IEEE international conference on computer-aided design, 1990. ICCAD-90. Digest of technical papers, pp 92–95
- 693
- 694 19. Trivedi K (1982) Probability and statistics with reliability, queueing and computer science applications. Prentice-Hall, Englewood Cliffs
- 695
- 696
- 697
- 698
- 699

Author Queries

Chapter 1

Query Refs.	Details Required	Author's response
	No queries.	

UNCORRECTED PROOF

MARKED PROOF

Please correct and return this set

Please use the proof correction marks shown below for all alterations and corrections. If you wish to return your proof by fax you should ensure that all amendments are written clearly in dark ink and are made well within the page margins.

<i>Instruction to printer</i>	<i>Textual mark</i>	<i>Marginal mark</i>
Leave unchanged	... under matter to remain	Ⓟ
Insert in text the matter indicated in the margin	∧	New matter followed by ∧ or ∧ [Ⓢ]
Delete	/ through single character, rule or underline or ┌───┐ through all characters to be deleted	Ⓞ or Ⓞ [Ⓢ]
Substitute character or substitute part of one or more word(s)	/ through letter or ┌───┐ through characters	new character / or new characters /
Change to italics	— under matter to be changed	↙
Change to capitals	≡ under matter to be changed	≡
Change to small capitals	≡ under matter to be changed	≡
Change to bold type	~ under matter to be changed	~
Change to bold italic	≈ under matter to be changed	≈
Change to lower case	Encircle matter to be changed	≡
Change italic to upright type	(As above)	⊕
Change bold to non-bold type	(As above)	⊖
Insert 'superior' character	/ through character or ∧ where required	Υ or Υ under character e.g. Υ or Υ
Insert 'inferior' character	(As above)	∧ over character e.g. ∧
Insert full stop	(As above)	⊙
Insert comma	(As above)	,
Insert single quotation marks	(As above)	Ƴ or ƴ and/or Ƶ or ƶ
Insert double quotation marks	(As above)	ƴ or ƶ and/or Ƶ or Ʒ
Insert hyphen	(As above)	⊞
Start new paragraph	┌	┌
No new paragraph	┐	┐
Transpose	└┐	└┐
Close up	linking ○ characters	Ⓞ
Insert or substitute space between characters or words	/ through character or ∧ where required	Υ
Reduce space between characters or words		↑