# Product-Form Approximation of Tandem Queues via Matrix Geometric Methods

Giuliano Casale, Peter G. Harrison, Maria Grazia Vigliotti

Department of Computing

Imperial College London

g.casale@imperial.ac.uk, {pgh, mgv98}@doc.ic.ac.uk

*Abstract*—We introduce a product-form approximation for tandem networks with Poisson arrivals and non-exponential service times. The proposed technique perturbs the model state space to match the sufficient conditions for product-form solution provided by the Reversed Compound Agent Theorem (RCAT). After characterizing the relationship between RCAT product-forms and matrix geometric solutions, we develop an algorithm based on nonlinear programming that automatically searches for an approximating product-form model.

## I. INTRODUCTION

We consider a tandem pair of queues processing requests that arrive at the first queue according to a Poisson process with rate $\lambda$. Both queues have infinite capacity and schedule requests according to a first-come first-served policy. Service times are phase-type distributed, so they can fit a wide range of distributions including exponential, hyper-exponential, hypo-exponential, and approximate certain heavy-tail distributions. Phase-type distributions are defined using the method of phases, so that time delays are modeled in terms of passage times in continuous-time Markov chains; upon activation of a tagged transition, a service time sample is generated as the cumulative time elapsed from the last activation of a tagged transition.

In this paper, we represent phase-type distributed service times using the Markovian arrival process (MAP) notation. Denote by $\mu_{k,k'}$ the rate of a tagged transition between states $k$ and $k'$, a MAP is described by a pair of matrices $(\boldsymbol{D}_0, \boldsymbol{D}_1)$ where $\mu_{k,k'}$ is the entry in row $k$ and column $k'$ of $\boldsymbol{D}_1$, while $\boldsymbol{D}_0 = \boldsymbol{Q} - \boldsymbol{D}_1$. $\boldsymbol{Q}$ is the infinitesimal generator of the underlying continuous-time Markov chain that does not distinguish between tagged and untagged transitions. Let us call queue $a$ the first queue to receive jobs, and let queue $b$ be the second queue. We denote by $(\boldsymbol{D}_0^a, \boldsymbol{D}_1^a)$ and $(\boldsymbol{D}_0^b, \boldsymbol{D}_1^b)$ the service processes of the first and second queue respectively and assume that their numbers of phases are $K$ and $H$.

## II. RCAT CONDITIONS FOR TANDEM QUEUES

Product-form approximations of state probabilities are obtained in the present work by applying the RCAT result [4]. RCAT provides sufficient conditions for communicating Markov processes to enjoy a product-form solution in their

joint state space and has been shown to unify existing product-form results for queueing networks, stochastic Petri nets, and to drive the construction of new product-forms [1]. In RCAT, a feed-forward tandem network is represented by two communicating Markov processes with infinitesimal generators $\boldsymbol{Q}^a$ : $(n, k) \rightarrow (n', k')$ for queue $a$ and $\boldsymbol{Q}^b$ : $(m, h) \rightarrow (m', h')$ for queue $b$, where $n$ is the number of jobs in queue $a$ and $k$ is the active phase in its service process, the population $m$ and phase $h$ are similarly defined for queue $b$. The equilibrium probabilities of $\boldsymbol{Q}^a$ and $\boldsymbol{Q}^b$ are denoted respectively by the vectors $\boldsymbol{\alpha}_n$, $n \geq 0$, and $\boldsymbol{\beta}_m$, $m \geq 0$, where $\alpha_{n,k}$ is the $k$th element of $\boldsymbol{\alpha}_n$, the equilibrium probability of state $(n, k)$; similarly, $\beta_{m,h}$ in $\boldsymbol{\beta}_m$ is the probability of state $(m, h)$ in queue $b$. The tandem network state space is readily obtained by composing $\boldsymbol{Q}^a$ and $\boldsymbol{Q}^b$, see [4] for related material.

RCAT's conditions for a product form solution in the joint process given by the composition of $\boldsymbol{Q}^a$ and $\boldsymbol{Q}^b$ are as follows [4]:

C1) queue $b$ can accept an incoming job in any of its states $(m, h)$ – such local state-transitions are called *passive*;

C2) each state $(n, k)$ of queue $a$ has an incoming transition for each tagged event corresponding to a job arrival at queue $b$ – such local state-transitions are called *active*;

C3) the sum of the reversed rates of all active transitions incoming to any state $(n, k)$ is constant (within each class of transition in a multi-class model).

We here consider the formulation of C3 proposed in [7]. If all the above conditions are simultaneously satisfied, then the equilibrium distribution of the model becomes

$$\pi(n, k, m, h) = \alpha_{n,k}\beta_{m,h} \qquad (1)$$

where $\beta_{m,h}$ is determined after setting to an appropriate value the arrival rate at the second queue $b$ according to RCAT's rate equations [4]. We refer to (1) as RCAT product-form since other product-form conditions also exist, e.g., ERCAT product-forms [1].

In the tandem models we consider, condition C1 is always satisfied since queue $b$ has infinite buffer size. We propose to modify the state space of $\boldsymbol{Q}^a$ to match RCAT conditions C2 and C3 and define a new approximating model that has product-form solution. Due to the complexity of the problem, we focus here on models with a single-class of transitions

synchronizing between $\boldsymbol{Q}^a$ and $\boldsymbol{Q}^b$, i.e., all active transitions (tagged events) are marked with an identical label [4]. While the results provided here are general under this assumption, the single-class case does restrict, in practice, the range of systems where we can find an approximating model that admits an RCAT product-form. We plan to characterize, in future work, the class of models that admit such single-class RCAT product-forms as well as to extend our approach to multiple classes of transitions. In particular, it would be interesting to assess the applicability to nonrenewal MAP service times which are hard to approximate using a single-class of transitions. Under the single-class assumption, condition C3 is equivalent to verifying the existence of a constant $\bar{q}$ such that

$$\sum_{(n',k')} \frac{\alpha_{n',k'} q((n',k') \to (n,k))}{\alpha_{n,k}} = \bar{q}, \qquad \forall (n,k) \quad (2)$$

where the summation considers only active transition rates $q((n',k') \to (n,k))$ between states $(n',k')$ and $(n,k)$ in $\boldsymbol{Q}^a$ that result in an arrival to queue $b$.

The basic idea of the product-form approximation we propose is to introduce perturbations on the structure and on the rates of the $\boldsymbol{Q}^a$ process to make the tandem network a product-form model. Throughout the next sections, we take the same qualitative approach proposed in [5] for $M/E_2/1 \to -/M/1$ networks, where $E_2$ is an Erlang-2 service process, but greatly extend its scope to queues with phase-type distributed service times. This leads to substantial differences in the approach since, without the special structure of an $E_2$ process, it is no longer tractable to symbolically solve for $\bar{q}$ in (2). Determining the $\bar{q}$ value is fundamental in verifying RCAT's conditions and in determining $\boldsymbol{\beta}_m$ using RCAT's rate equations [4], [5].

### A. Matrix Geometric Solutions and RCAT Product-Forms

We begin with observing that since the behavior of the first queue is independent of the second queue, we can solve for the equilibrium state probabilities of $\boldsymbol{Q}^a$ using the quasi-birth-death (QBD) process

$$\boldsymbol{Q}^a = \begin{bmatrix} \boldsymbol{L}_0 & \boldsymbol{F}_0 & 0 & 0 & 0 & \cdots \\ \boldsymbol{B}_0 & \boldsymbol{L} & \boldsymbol{F} & 0 & 0 & \cdots \\ 0 & \boldsymbol{B} & \boldsymbol{L} & \boldsymbol{F} & 0 & \cdots \\ 0 & 0 & \boldsymbol{B} & \boldsymbol{L} & \boldsymbol{F} & \cdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \end{bmatrix}$$

where the arrival and service rates define the forward transition matrix $\boldsymbol{F} = \lambda \boldsymbol{I}$, the local transition matrix $\boldsymbol{L} = -\lambda \boldsymbol{I} + \boldsymbol{D}_0^a$, and the backward transition matrix $\boldsymbol{B} = \boldsymbol{D}_1^a$. The transition matrices $\boldsymbol{L}_0$, $\boldsymbol{B}_0$, and $\boldsymbol{F}_0$ describe boundary conditions for states where the first queue is empty. The aim of this section is to develop a characterization of the relationship between RCAT condition C3 and matrix geometric solutions of the above QBD. We find in particular that an RCAT product-form imposes special properties on the equilibrium of the QBD.

We begin by writing the global balance equations for states where queue $a$ has population of two or more jobs, i.e.,

$$\boldsymbol{\alpha}_{n-1} \boldsymbol{F} + \boldsymbol{\alpha}_n \boldsymbol{L} + \boldsymbol{\alpha}_{n+1} \boldsymbol{B} = \boldsymbol{0}, \qquad n \geq 2. \quad (3)$$

We now make the fundamental observation that, for a product-form model that satisfies C3, inserting (2) in the global balance equations of $\boldsymbol{Q}^a$ for state $(n,k)$ always cancels out the probability flux of the incoming transitions from the states $(n',k')$ contributing to the summation in (2). This can be equivalently expressed in matrix notation by first noting that C3 implies

$$\boldsymbol{\alpha}_n \hat{\boldsymbol{L}} = \boldsymbol{\alpha}_{n+1} \boldsymbol{B}, \qquad n \geq 1 \quad (4)$$

in which we define

$$\hat{\boldsymbol{L}} = \begin{bmatrix} \bar{q}\delta_1 & 0 & 0 & \cdots & 0 \\ 0 & \bar{q}\delta_2 & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \cdots & \vdots \\ 0 & 0 & 0 & \cdots & \bar{q}\delta_K \end{bmatrix}$$

where $\delta_k$ is 1 if the $k$th column of $\boldsymbol{D}_1^a$ has at least one nonzero element, and 0 otherwise. Then inserting (4) into (3), C3 implies that

$$\boldsymbol{\alpha}_{n-1} \boldsymbol{F} + \boldsymbol{\alpha}_n (\boldsymbol{L} + \hat{\boldsymbol{L}}) = \boldsymbol{0}, \qquad n \geq 2 \quad (5)$$

Since $\boldsymbol{L}$ has an inverse, also $\boldsymbol{L} + \hat{\boldsymbol{L}}$ has an inverse and so

$$\boldsymbol{\alpha}_n = \boldsymbol{\alpha}_{n-1} \boldsymbol{H}, \qquad n \geq 2 \quad (6)$$

where $\boldsymbol{H} = -\boldsymbol{F}(\boldsymbol{L} + \hat{\boldsymbol{L}})^{-1}$. Note that this is also the solution of the matrix equation $\boldsymbol{F} + \boldsymbol{H}\boldsymbol{L} + \boldsymbol{H}\hat{\boldsymbol{L}} = \boldsymbol{0}$. However, the matrix geometric solution requires that there exist a rate matrix $\boldsymbol{R}$ such that

$$\boldsymbol{\alpha}_n = \boldsymbol{\alpha}_{n-1} \boldsymbol{R}, \qquad n \geq 2 \quad (7)$$

where $\boldsymbol{R}$ is the non-negative minimum norm solution of $\boldsymbol{F} + \boldsymbol{R}\boldsymbol{L} + \boldsymbol{R}^2 \boldsymbol{B} = \boldsymbol{0}$. Thus, we find that if the QBD process satisfies C3 there exist a rate matrix

$$\boldsymbol{H} = -\boldsymbol{F}(\boldsymbol{L} + \hat{\boldsymbol{L}})^{-1} = \boldsymbol{R}(\boldsymbol{L} + \boldsymbol{R}\boldsymbol{B})(\boldsymbol{L} + \hat{\boldsymbol{L}})^{-1},$$

that provides an additional geometric relation (6) for the computation of the equilibrium solution. Note that $\boldsymbol{H}$ is different from $\boldsymbol{R}$ whenever $\boldsymbol{R}\boldsymbol{B} \neq \hat{\boldsymbol{L}}$, yet (6) and (7) clearly imply dependencies between the spectral properties of $\boldsymbol{H}$ and $\boldsymbol{R}$. In fact, we note that in the limit $n \to \infty$, $\boldsymbol{\alpha}_n = k\eta^j \boldsymbol{\alpha}_{n-j}$ for some constant $k$, where $\eta \in [0,1]$ is the largest eigenvalue of $\boldsymbol{R}$ that describes the asymptotic decay rate of the queue-length probabilities. However, by (6), $\eta$ must also be the largest eigenvalue of $\boldsymbol{H}$. Thus, $\boldsymbol{R}$ and $\boldsymbol{H}$ must have the same spectral radius

$$\rho(\boldsymbol{R}) = \rho(\boldsymbol{H}) = \eta \quad (8)$$

which imposes identical stability conditions on the QBD, i.e., $\rho(\boldsymbol{R}) < 1$. The importance of (8) is that one can compute $\bar{q}$ by searching for the $\hat{\boldsymbol{L}}$ matrix that satisfies $\rho(\boldsymbol{R}) = \rho(\boldsymbol{H})$ and then extract $\bar{q}$ directly from its main diagonal.

To complete the characterization, we are left with the implications of C3 on the global balance conditions for the vectors $\boldsymbol{\alpha}_0$ and $\boldsymbol{\alpha}_1$; such global balance equations are not considered in (3). Following the same argument used before,

we see that C3 implies the following conditions

$$\boldsymbol{\alpha}_0 \boldsymbol{F}_0 + \boldsymbol{\alpha}_1 \boldsymbol{L} + \boldsymbol{\alpha}_1 \boldsymbol{R} \boldsymbol{B} = \mathbf{0} \tag{9}$$

$$\boldsymbol{\alpha}_0 \boldsymbol{L}_0 + \boldsymbol{\alpha}_1 \boldsymbol{B}_0 = \mathbf{0} \tag{10}$$

$$\boldsymbol{\alpha}_0 \mathbf{1} + \boldsymbol{\alpha}_1 (\boldsymbol{I} - \boldsymbol{R})^{-1} \mathbf{1} = 1 \tag{11}$$

$$\boldsymbol{\alpha}_0 \hat{\boldsymbol{L}}_0 - \boldsymbol{\alpha}_1 \boldsymbol{B}_0 = \mathbf{0} \tag{12}$$

$$\boldsymbol{\alpha}_1 \hat{\boldsymbol{L}} - \boldsymbol{\alpha}_1 \boldsymbol{R} \boldsymbol{B} = \mathbf{0} \tag{13}$$

where we have used the fact that $\boldsymbol{\alpha}_2 = \boldsymbol{\alpha}_1 \boldsymbol{R}$, $\mathbf{1}$ is a vector of ones of proper size, (11) normalizes the probabilities, and $\hat{\boldsymbol{L}}_0$ is defined similarly to $\hat{\boldsymbol{L}}$ but relative to the states in which the QBD is empty. Note that in (9)-(13) only the last two equations are specific to C3, since (9)-(11) hold for all QBDs.

## III. PRODUCT-FORM APPROXIMATION

In order to obtain a product-form model that behaves as closely as possible to the original tandem network, we propose to introduce two perturbations in the state space. The first assumes that the rates in $\boldsymbol{F}_0$ can be perturbed in order to satisfy condition C3 on the entire state space. The second perturbation involves adding a self-looping active transition to each state $(n, k)$ in $\boldsymbol{Q}^a$ that does not satisfy the RCAT condition C2; such transitions are referred to as *invisible* transitions and to match condition C3 must have rate $\bar{q}$. Note that the activation of an invisible transition does not change the steady state distribution of queue $a$, and its only effect is to have a new job arriving to queue $b$ but that is *not* due to a departure from queue $a$. Following an argument similar to [5], the invisible transitions are found to describe the action of an exogenous Interrupted Poisson Process (IPP) that is superposed with the departure flow of the QBD. Thus, the modified $\boldsymbol{Q}^a$ process also communicates, without changing state, with $\boldsymbol{Q}^b$ by enabling job arrivals from the IPP process[1].

Since C2 can easily be imposed, the main challenge of the proposed approximation is to search for appropriate rates in $\boldsymbol{F}_0$ that ensure an equilibrium probability distribution satisfying (4) and (9)-(13) which are equivalent to C3. In order to meet the infinite conditions (4) we first observe that any set of $K$ powers of $\boldsymbol{R}$ must be linearly dependent, $\boldsymbol{R}$ being a square matrix of order $K$. Thus, C3 is satisfied if (9)-(13) and (4), for $n = 1, \ldots, K$, hold true. To impose this, we search over the space of the rates in $\boldsymbol{F}_0$ using a nonlinear optimization program. After initial computation of $\boldsymbol{R}$, at each iteration we determine the values of $\boldsymbol{\alpha}_0$ and $\boldsymbol{\alpha}_1$ for the current assignment of the rates in $\boldsymbol{F}_0$; this may be done using tools such as [2] which automatically determine $\boldsymbol{\alpha}_0$ and $\boldsymbol{\alpha}_1$ satisfying conditions (9)-(11). The values of $\boldsymbol{\alpha}_0$, $\boldsymbol{\alpha}_1$, and $\boldsymbol{F}_0$ for the current iteration can then be used to drive the minimization of the objective function

$$f_{obj} = ||\boldsymbol{\alpha}_0 \hat{\boldsymbol{L}}_0 - \boldsymbol{\alpha}_1 \boldsymbol{B}_0||_2 + \sum_{i=1}^{K} ||\boldsymbol{\alpha}_1 \boldsymbol{R}^{i-1} \hat{\boldsymbol{L}} - \boldsymbol{\alpha}_1 \boldsymbol{R}^i \boldsymbol{B}||_2$$

which is zero if C3 holds on all states of the QBD. We have observed that the above formulation often avoids stagnation of

---

[1]Since the above approximation only injects new load into the system, the resulting product-form model may be used as a pessimistic estimate on the performance of the tandem network, see [3], [6] and references therein for related techniques.

the iteration that arise when searching directly for a feasible solution to the nonlinear system of equations (4) and (9)-(13) in the unknowns $\boldsymbol{F}_0$, $\boldsymbol{\alpha}_0$, and $\boldsymbol{\alpha}_1$. A summary of the proposed technique to impose C2 and C3 is given in the following pseudocode:

1) Determine $\boldsymbol{R}$ for given $\boldsymbol{F}$, $\boldsymbol{L}$, $\boldsymbol{B}$
2) Compute the largest eigenvalue $\eta \in [0, 1]$ of $\boldsymbol{R}$
3) Determine $\bar{q}$ in $\hat{\boldsymbol{L}}$ such that

$$\rho(-\boldsymbol{F}(\boldsymbol{L} + \hat{\boldsymbol{L}})^{-1}) = \eta$$

4) Using the rate $\bar{q}$ found in the previous step, search for values of the entries in $\boldsymbol{F}_0$ satisfying (4) and (9)-(13). At each iteration, $\boldsymbol{\alpha}_0$, $\boldsymbol{\alpha}_1$ are immediately determined by the linear conditions (9)-(13) for assigned $\boldsymbol{F}_0$.
5) Add self-looping transitions with rate $\bar{q}$ to states that do not match C2
6) The rate $\bar{q}$ provides the unknown rate of arrivals to queue $b$ to be used to determine the probability distribution $\boldsymbol{\beta}_n$ in $\boldsymbol{Q}_b$. See [5], [4] for additional details.

## IV. CASE STUDY

We illustrate our product-form approximation on a case study. We consider a tandem network $M/Hypo_3/1 \rightarrow -/E_2/1$, where $E_2$ denotes an Erlang-2 distribution and $Hypo_3$ indicates a hypo-exponential process composed of three sequential stages with rates $\mu_{1,1} = 1, \mu_{2,2} = 2, \mu_{3,3} = 3$, where $\mu_{1,1}, \mu_{2,2} \in \boldsymbol{D}_0$ and $\mu_{3,3} \in \boldsymbol{D}_1$. The arrival process to the first queue is Poisson with rate $\lambda = 0.25$. The QBD for $\boldsymbol{Q}^a$ is then defined by

$$\boldsymbol{B} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 3 & 0 & 0 \end{bmatrix}, \quad \boldsymbol{L} = \begin{bmatrix} -1.25 & 1 & 0 \\ 0 & -2.25 & 2 \\ 0 & 0 & -3.25 \end{bmatrix},$$

$$\boldsymbol{F} = \begin{bmatrix} 0.25 & 0 & 0 \\ 0 & 0.25 & 0 \\ 0 & 0 & 0.25 \end{bmatrix}, \quad \boldsymbol{B}_0 = \begin{bmatrix} 0 \\ 0 \\ 3 \end{bmatrix}$$

The rate matrix $\boldsymbol{R}$ is computed using the tool presented in [2] and its eigenvalues are found to be $eig(\boldsymbol{R}) = 0.3522, 0.0856 + 0.008i, 0.0856 - 0.008i$, thus $\eta = \rho(\boldsymbol{R}) = 0.3522$. Using this value, we determine numerically the minimum of $|\rho(-\boldsymbol{F}(\boldsymbol{L} + \hat{\boldsymbol{L}})^{-1}) - \eta|$; this function is plotted in Figure 1 for different values of the $\bar{q}$ element in $\hat{\boldsymbol{L}}$. The minimization finds that a product form requires a reversed rate $\bar{q} = 0.54$, so that

$$\hat{\boldsymbol{L}} = \begin{bmatrix} 0.54 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

where the zero diagonal values are due to zero columns in $\boldsymbol{B}$. We now define $\boldsymbol{F}_0 = [\lambda + x_0, x_1, \ldots, x_K]$, where $x_j$ are the perturbations we introduce in $\boldsymbol{F}_0$ to obtain a product-form. We set up a non-linear optimization program searching for the values of $x_0, x_1, \ldots, x_K$ to generate a distribution that satisfies (4) and (9)-(13), and use MATLAB's fmincon function to solve this and obtain the following result: $x_0 = 0$, $x_1 = 0.1623$, $x_2 = 0.1278$. Thus

$$\boldsymbol{F}_0 = \begin{bmatrix} 0.2500 & 0.1623 & 0.1278 \end{bmatrix}, \quad \boldsymbol{L}_0 = \begin{bmatrix} -0.5401 \end{bmatrix},$$

Fig. 1. Plot of the absolute spectral radius difference between $\boldsymbol{H}$ and $\boldsymbol{R}$. The minimum is for $\bar{q} = 0.54$.
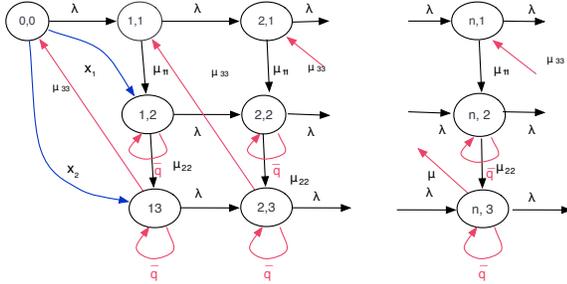


Fig. 2. The single-class RCAT product-form approximation is based on the (invisible) self-looping transitions in phases 2 and 3 and on the new transitions $x_1$ and $x_2$ (highlighted in blue) introduced in $\boldsymbol{F}_0$. All active transitions (highlighted in red) generate arrivals to queue $b$.

| state | $\bar{q}$ value | state | $\bar{q}$ value |
|---|---|---|---|
| $(0, 0)$ | 0.54012951 | $(10, 1)$ | 0.54013125 |
| $(1, 1)$ | 0.54013191 | $(11, 1)$ | 0.54013125 |
| $(2, 1)$ | 0.54013167 | $(12, 1)$ | 0.54013125 |
| $(3, 1)$ | 0.54013140 | $(13, 1)$ | 0.54013125 |
| $(4, 1)$ | 0.54013131 | $(14, 1)$ | 0.54013125 |
| $(5, 1)$ | 0.54013125 | $(15, 1)$ | 0.54013125 |
| $(6, 1)$ | 0.54013125 | $(16, 1)$ | 0.54013125 |
| $(7, 1)$ | 0.54013125 | $(17, 1)$ | 0.54013125 |
| $(8, 1)$ | 0.54013125 | $(18, 1)$ | 0.54013125 |
| $(9, 1)$ | 0.54013125 | $\ldots$ | $\ldots$ |

TABLE II
MARGINAL PROB. OF QUEUE $b$ IN A $M/Hypo_3/1 \rightarrow -/E_2/1$ MODEL

| population | original model | product-form approx. |
|---|---|---|
| 0 | 0.7525 | 0.3997 |
| 1 | 0.2075 | 0.2790 |
| 2 | 0.0347 | 0.1529 |
| 3 | 0.0046 | 0.0782 |
| 4 | 0.0007 | 0.0425 |
| 5 | 0.0001 | 0.0233 |
| 6 | 0.0000 | 0.0121 |
| 7 | 0.0000 | 0.0058 |
| 8 | 0.0000 | 0.0040 |
| 9 | 0.0000 | 0.0018 |

where $\boldsymbol{L}_0$ is defined by the element $-\lambda - x_0 - x_1 - x_2 = -0.5401$ that appears on the main diagonal of the QBD; note that this is equal to $-\bar{q}$ similarly to what found in [5] for the Erlang-2 case.

For validation purpose, we compute numerically the entire probability distribution $\boldsymbol{\alpha}_n$, $n \geq 0$, which gives the reversed rate $\bar{q}$ in (2) that needs to be the same at all instances in order to satisfy conditions C3 of RCAT. Table I shows that this is the case with good precision, hence we can reasonably condition C3 to hold. Based on this result, the state space of the first queue is corrected as shown in Figure 2, where the red transitions represent marked events that produce an arrival to the second queue. Note in particular that red self-loops are used to denote arrivals to the second queue due to the extra IPP flow with rate $\bar{q}$.

We conclude the case study by considering the marginal probabilities of queue $b$ in the above case study. Probability values are obtained by simulation initializing both queues in the empty state. Table II reports simulation results; similar differences are observed if the Erlang-2 distribution is replaced by an exponential or hyper-exponential distribution. In this experiment, the utilization of queue $a$ is $\rho_a = 0.458$ in the original model and $\rho'_a = 0.462$ after the modification of $\boldsymbol{F}_0$. For queue $b$ it is $\rho_b = 0.25$ in the original model and becomes $\rho'_b = 0.60$ in the product-form approximation. Thus, we see that the additional jobs injected in the system by the IPP flow and by the $x_j$ perturbations result in an increased load at queue $b$. In fact, all product-form values for states where the queue is busy are greater than the corresponding values for the original model. Indeed, depending on the service characteristics at the second queue, there exist perturbations that result in saturation at the second queue, in this case the approximation does not return a valid product-form model. We plan to characterize conditions under which these cases arise in future work.

## V. CONCLUSION

We have proposed a new approximation method for tandem networks of queues. We have found that product-form conditions applied to QBDs reveal a novel matrix geometric relationship between state probabilities. This allows us to integrate the RCAT conditions into an optimization program that constructs an approximate product-form.

Possible generalization of the above results include considering generalized arrival processes and nonrenewal service. Similarly, network structures that are more general than tandem have for long been accommodated by RCAT, and these provide another obvious research direction.

## REFERENCES

[1] S. Balsamo, P.G. Harrison, A. Marin. A unifying approach to product-forms in networks with finite capacity constraints, to appear in *Proc. of ACM SIGMETRICS 2010*.

[2] D.A.Bini, B.Meini, S. Steffé, B. Van Houdt. Structured Markov chains solver: software tools. in *Proc. of SMCTOOLS*, 2006.

[3] P. Buchholz. Bounding stationary results of Tandem networks with MAP input and MAP service time distributions. *Proc. of ACM SIGMETRICS*, 191–202,2006.

[4] P.G. Harrison. Turning back time in Markovian process algebra. *Theor. Comp. Sci.*. 290(3):1947–1986, 2003.

[5] P.G. Harrison and M.G. Vigliotti. Perturbation of a non-product-form network into a new product-form: equilibrium state probabilities and response time density. *Proc. of VALUETOOLS*, 2009.

[6] A. Heindl. Decomposition of general queueing networks with MMPP inputs and customer losses. *Perf. Eval.*, 51(2–4), 117–136, Feb 2003.

[7] A. Marin and M.G. Vigliotti. A general result for deriving product-form solutions of Markovian models. *Proc. of WOSP/SIPEW*, 2010.