

Markovian approximations for a grid computing network with a ring structure

J. F. Pérez and B. Van Houdt

Performance Analysis of Telecommunication Systems Research Group,
Department of Mathematics and Computer Science,
University of Antwerp - IBBT,
email: {juanfernando.perez, benny.vanhoudt}@ua.ac.be

Abstract

Optical Grid networks allow many computing sites to share their resources by connecting them through high-speed links, providing a more efficient use of the resources and a timely response for incoming jobs. These jobs originate from users connected to each of the sites and, in contrast to traditional queueing networks, a particular job does not have to be processed in a predefined site. Furthermore, a job is always processed locally if there is an available local server. In this paper we propose two different methods to approximate the performance of an optical Grid network with a ring topology. The first method is based on approximating the inter-overflow time process, while the second separately characterizes the periods where jobs are overflowed and the periods where they are served locally. Both approaches rely on a marked Markovian representation of the overflow process at each station and on reducing this representation by moment-matching methods. The results show that the methods accurately approximate the rate of locally processed jobs, one of the main performance measures.

1 Introduction

Optical grid networks offer the possibility of sharing resources between several computing sites/stations by connecting them through high-speed links. Each of these sites is typically connected to a number of users who request the processing of tasks/jobs in any of its servers. The site always tries to process the jobs locally but, if all its servers are busy, the jobs can be sent to any of the other sites. Furthermore, there is no preference about where the jobs should be processed [10, 11]. Thus, the site where a particular job is served depends not only on a predefined routing matrix, but also on the server availability, the scheduling algorithm and the topology of the network. In this paper we focus on a grid network with a ring topology where each station is physically connected to only one other site through an unidirectional fiber link. If an incoming job finds all the local servers busy, the job is sent to the next station in the ring where it is served by an available server. If no servers are available in that station, the job is sent again to the next site. In this manner the job goes from one station to the next trying to find an idle server. After trying all the stations once, the job is dropped from the network. This description reveals two main characteristics of the grid ring network: first, the jobs transmitted from one site to another form an overflow process; second, every station in the ring receives both newly-generated and overflowing jobs, which may come from any of the other sites.

A grid network typically consists of many stations and, thus, the scalability of the methods to analyze these networks is of major importance. In this work we decompose the network

in single nodes to perform the analysis of each station separately. This technique has been widely used in the analysis of queueing networks that do not have a product-form solution. Particularly, moment-matching methods have played a key role in approximating the arrival, departure and overflow processes of the single nodes in the network, see [25, 26, 37] and references therein. This approach has been used recently to approximate the departure process of queues with Markovian input [19, 20, 21, 22]. Since the traffic between stations in the ring forms an overflow process, the analysis of teletraffic networks with alternate routing is particularly relevant [24, 25, 30, 31, 38]. In those networks the incoming traffic is first offered to a trunk group with a finite buffer. If the buffer has no available space, the traffic is sent to an alternative trunk group, forming an overflow process. To compute the performance measures for each of the flows in the network separately, the methods in [25, 31] rely on representing each of these flows as an independent process. This approach is well-suited for general networks with different routing patterns, but dimensionality problems arise when analyzing even small networks. In [30] the authors propose to group all but one overflow process into a single description in order to model the arrival process at each station with one Poisson and two overflow processes. With this method the dimensionality problems are avoided, but an additional approximation is introduced.

In this work we exploit the ring structure to approximate the overflow process at each station by means of a marked Markovian point process. The markings are used to differentiate each of the flows in the network, leading to a more compact representation of the overflow process. This representation has to be further reduced to make the analysis of the network feasible. We propose two different methods to build an approximate reduced representation of the overflow process. The first approach aims at approximating the stationary inter-overflow time process by matching a set of higher (joint) moments that uniquely characterize the reduced point process. The second approach separately characterizes the periods where jobs are overflowed and the periods where they are served locally, combining their reduced representation into a single one. The first approach is stepwise; we start by approximating the stationary inter-overflow time distribution with a renewal process that matches three first moments relying on Phase-Type distributions. We then analyze the effect of capturing more moments by relying on the class of Matrix-Exponential distributions. Additionally, we include information of the joint moments of successive inter-overflow times by using recent matching methods based on the marked and unmarked versions of the Rational Arrival Process. These approximations are tested for different network parameters with particular emphasis on computing the rate of locally processed jobs (local rate) and the total traffic matrix. This is of vital importance when dimensioning the link capacities interconnecting the different sites [10, 12]. The results show that the methods are able to adequately approximate these measures, being especially accurate for the local rate.

The paper is organized as follows. Section 2 provides a description of the grid network as well as a discussion of some characteristics common to both approximation methods. The approximation method based on the inter-overflow time process is presented in Section 3, while the second method, called ON-OFF, is introduced in Section 4. Section 5 compares the performance of both methods for various realistic network configurations.

2 The grid network

We consider a grid network that consists of N nodes arranged in a ring topology. The job arrivals at node i are represented by a Poisson process with rate λ_i (this assumption will be relaxed to allow for more general arrival processes as described at the end of this section).

The Poisson assumption is based on Grid level measurements [9] and has been employed in previous works on Grid dimensioning [10, 12]. When a job arrives at node i , it is served by any of the C_i servers in that node, if at least one is available. In case all of them are busy, the job is sent to the next node in the ring. The job jumps from one node to the next until it finds an idle server. If a job originating in node i arrives at station $i - 1$ and there are no available servers, it must be dropped. Thus, a job that has tried all the stations once is dropped in the station before the one where it was generated. The service time of a job in station i is assumed to be exponentially distributed with rate μ_i . This means that the service time depends on the station serving the job, but not on the station where it entered the grid (in other words, all sites generate similar jobs).

Even in the case where each type of job comes from a Poisson process, the actual flow that arrives to a particular station is a complex mixture of all the job types present in the network. Hence we will approximate the input process¹ at station i as a marked Markovian Arrival Process [16] (MMAP[N]) characterized by a set of $m_i \times m_i$ matrices $\{D_0^i, D_1^i, \dots, D_N^i\}$. These matrices will be built through an iterative process in order to incorporate information about the overflowing jobs along the network into the analysis of each station. The MMAP[N] is a point process well suited for modeling purposes since it admits general inter-arrival times and a correlation structure. The process is driven by an underlying Continuous-Time Markov Chain (CTMC) with generator matrix $D^i = \sum_{j=0}^N D_j^i$. Each of its transitions might generate zero or one arrival. The matrix D_0^i describes phase transitions without arrivals, while the matrices $\{D_j^i, j = 1, \dots, N\}$ describe the transition rates related to an arrival of type (originated in station) j . These matrices as well as the off-diagonal elements of D_0^i are non-negative, while the diagonal elements of D_0^i are negative and such that $D^i \mathbf{1} = \mathbf{0}$, where $\mathbf{1}$ and $\mathbf{0}$ are column vectors with all its entries equal to 1 and 0, respectively.

Let $\{N_i(t), t \geq 0\}$ (resp. $\{J_i(t), t \geq 0\}$) be the number of busy servers (resp. the phase of the arrival process) at station i at time t . Then $\{(N_i(t), J_i(t)), t \geq 0\}$ is a CTMC on the state space $\{(k, l), 0 \leq k \leq C_i, 1 \leq l \leq m_i\}$, that describes the state of the station i , based on the approximated arrival process. Its generator matrix is given by

$$Q^i = \begin{bmatrix} D_0^i & D_+^i & 0 & \dots & 0 & 0 \\ \mu_i I & D_0^i - \mu_i I & D_+^i & \dots & 0 & 0 \\ 0 & 2\mu_i I & D_0^i - 2\mu_i I & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & D_0^i - (C_i - 1)\mu_i I & D_+^i \\ 0 & 0 & 0 & \dots & C_i \mu_i I & D^i - C_i \mu_i I \end{bmatrix}, \quad (1)$$

where $D_+^i = \sum_{j=1}^N D_j^i$ and $D^i = D_0^i + D_+^i$.

If the matrices $\{D_0^i, D_1^i, \dots, D_N^i\}$ are specified, the steady state probability vector can be easily computed by exploiting the structure of matrix Q^i as a finite Quasi-Birth-and-Death (QBD) process [14, 27, 32]. The time and memory complexity of the algorithm in [14] to compute the steady-state distribution are $O(C_i m_i^3)$ and $O(C_i m_i^2)$, respectively. Nonetheless, the arrival process to a particular station is determined by the superposition of the new arrivals arriving at that station and the overflow process coming from the previous station in the network. Even though these could be exactly represented by including in the arrival process the state of all the stations in the network, the size of such a representation is huge even for a small number of sites and servers per site. Previous work by Meier-Hellstern [31] considers the problem of approximating the arrival process at each station by representing each

¹In Section 3.1.2 we will relax this approximation to a marked Rational Arrival Process (MRAP[N])

of the flows in the network as a Markov Modulated Poisson Process (MMPP) and combining them for each station. An MMPP can be seen as an MMAP[1] where only the matrix D_1 is specified and it has zero off-diagonal elements. In [31] the algorithm of Heffes [18] is used to reduce each of the MMPPs to be of size 2. Even though this approximation is well suited to represent different routing strategies, it does not scale well with the number of stations because the input process to a particular station is of size 2^{N-1} . Our approximation methods make use of the ring structure to avoid this large size while keeping meaningful information about the overflow process at each site.

The overflow process at station i can be represented as an MMAP[N] with parameters $\{E_0^i, E_1^i, \dots, E_N^i\}$ given by

$$E_0^i = \begin{bmatrix} D_0^i & D_+^i & 0 & \dots & 0 & 0 \\ \mu_i I & D_0^i - \mu_i I & D_+^i & \dots & 0 & 0 \\ 0 & 2\mu_i I & D_0^i - 2\mu_i I & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & D_0^i - (C_i - 1)\mu_i I & D_+^i \\ 0 & 0 & 0 & \dots & C_i \mu_i I & D_0^i + D_s^i - C_i \mu_i I \end{bmatrix}, \quad (2)$$

$$E_j^i = \begin{bmatrix} 0 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & D_j^i \end{bmatrix}, j = 1, \dots, N, j \neq s,$$

and $E_s^i = 0$, where s denotes the next station after node i and 0 is a zero matrix of appropriate dimension. This representation reflects the fact that the jobs originating from station s must be dropped if they can not be served in node i . Even though this representation is exact, it is not useful for practical implementation, because the size of the matrices that describe the overflow process becomes extremely large after a few stations. Thus, it is useful to reduce the size of these matrices such that the reduced representation keeps some characteristics of the original and can be used as part of the arrival process at the next station. We propose two different methods to find an approximate representation of the overflow process, which result in an MMAP[N] with matrices $\{\bar{D}_0^i, \bar{D}_1^i, \dots, \bar{D}_N^i\}$. The first one is based on the approximation of the inter-overflow time process. The second method divides the representation of the overflow process into two sets: one characterizes the time periods where any arriving job is sent to the next station (ON period), while the other captures the behavior of the site when there are servers available to process an incoming job (OFF period). Even though each of the approximation methods relies on a different type of information, some main features are alike. For example, both methods represent the overflow process at station i by means of a reduced MMAP[N] (or a generalization of it) with parameters $\{\bar{D}_0^i, \bar{D}_1^i, \dots, \bar{D}_N^i\}$. Other common features are presented in the next section.

2.1 Common characteristics of the approximation methods

Iterative approach: The ring topology implies that all the nodes receive both newly generated and overflow jobs, meaning that the arrival process at each station depends on the analysis of the previous one. To contemplate this, both methods consider an iterative strategy in which an arbitrary site is first analyzed considering only its own traffic. The next station in the ring is then considered, including both newly generated and overflow jobs (from the previous station only). The analysis continues at each station and, after considering all the

stations once, the overflow process contains jobs of (possibly) all types. Then, it is possible to reanalyze the “first” station, but now the arrivals include both newly generated jobs as well as jobs coming from (possibly) all the other sites in the ring. This sequential analysis is performed several times for each site until the traffic matrix (that contains the amount of traffic between each pair of sites) changes less than a predetermined value ϵ , e.g., $\epsilon = 10^{-8}$.

Moment matching: Another relevant issue for both methods is the reduction of the representation of an inter-event process by means of the moment-matching methods introduced in [7, 33, 34, 35, 37]. The reduction of the process always implies the reduction of the inter-event time distribution by matching some of its moments with a distribution with smaller representation (if the moments are attainable by the matching distribution). Some of these methods provide closed-form formulas for the parameters of the matching distribution [33, 37], making them well suited for iterative procedures as those presented in the next sections. The other methods require more computational effort, but this is still negligible compared to the computation of the moments themselves. This is due to the large number of phases in the exact representation of the inter-event time distribution. The direct formulas to compute the moments [27] have a time and memory complexity of $O(C_i^3 m_i^3)$ and $O(C_i^2 m_i^2)$, respectively. However, the representation of the inter-event time distribution has a block-tridiagonal structure as the one shown in Equation (1). To exploit this structure we make use of the algorithm introduced in [14] to compute the first two moments of the first-passage time distribution to higher levels in a finite QBD, where level k corresponds to the set of states with k busy servers: $\{(k, l), 1 \leq l \leq m_i\}$. Using the generating function of the first-passage times described in [14] it is possible to determine higher moments of this distribution to use them as input for the moment-matching methods. Appendix A discusses why the inter-event times used by the approximation methods in Sections 3 and 4 can be interpreted as first-passage times to higher levels in a finite QBD. It also describes an algorithm based on [14] to compute any number of moments of this distribution. This is particularly useful for the method in Section 4, where the size of the arrival process at each station grows linearly with the number of stations. The time and memory complexity of the algorithm to compute these moments are thereby reduced to $O(C_i m_i^3)$ and $O(C_i m_i^2)$, respectively.

Arrival process: To characterize the arrival process at each station, the overflow coming from the previous station and the new arrivals must be combined. Under the assumption of Poisson arrivals, the new incoming jobs at station s and the overflow from station i can be combined as an MMAP[N] with parameters

$$D_0^s = \bar{D}_0^i - \lambda_s I, \quad D_j^s = \bar{D}_j^i, \quad j = 1, \dots, N, j \neq s, \quad D_s^s = \lambda_s I. \quad (3)$$

It will become clear that the approximation introduced in Section 3 can deal not only with Poisson arrivals, but with more general point processes as well. Let the arrivals at station s be described by a Markovian Arrival Process (MAP) characterized by $\{B_0^s, B_1^s\}$ [27, 29]. A MAP process is a point process that can represent general correlated inter-arrival times. It can be seen as a special case of the MMAP[N] process with a single arrival type. Assuming this arrival process, the combined stream at station s can be represented as an MMAP[N] with parameters

$$D_0^s = \bar{D}_0^i \oplus B_0^s, \quad D_j^s = \bar{D}_j^i \otimes I, \quad j = 1, \dots, N, j \neq s, \quad D_s^s = I \otimes B_1^s,$$

where \otimes and \oplus stand for Kronecker product and sum [5], respectively. Even though the second approximation (Section 4) can in principle include this more general process, the size

of the representation of the arrival process increases exponentially with the dimension of the MAP at each station. For the technique in Section 3, we do not encounter such an exponential increase.

3 An approximation based on inter-overflow times

The approximation introduced in this section is based on reducing the size of the representation of the inter-overflow time process. As explained above, the size of the exact representation of this process is extremely large, making the analysis of even small networks infeasible. Here we introduce different methods to determine a reduced approximate representation of the overflow process. We first consider the case where the overflow process is approximated by a renewal process, assuming independent inter-overflow times. Next we allow the new process to be non-renewal by including information about the joint moments of successive inter-overflow times.

3.1 Renewal Approximation

3.1.1 3-moment match

When matching the inter-overflow time distribution, we make use of Phase-Type (PH) distributions [27, 32]. A PH distribution is defined as the absorption time in a CTMC with one transient class and one absorbing state. Let S be the sub-generator of the transient class, and let $\mathbf{s} = -S\mathbf{1}$ be the column vector of absorption rates at the transient states. The initial probability vector \mathbf{b} of the chain can be partitioned as $\mathbf{b} = [\tau_0 \ \boldsymbol{\tau}]$, with $\tau_0 = 1 - \boldsymbol{\tau}\mathbf{1}$ being the probability that the process starts in the absorbing state. Therefore a PH distribution is completely characterized by the parameters $(\boldsymbol{\tau}, S)$. From Equation (2) it is clear that E_0^i is the sub-generator of the PH representation of the inter-overflow times at site i . To represent the overflow process with a PH renewal process, we consider the stationary version of the inter-overflow times, with representation given by $(\boldsymbol{\gamma}_i, E_0^i)$. To define $\boldsymbol{\gamma}_i$ we first need the steady state probability vector of Q^i , i.e., the vector $\boldsymbol{\pi}^i$ such that $\boldsymbol{\pi}^i Q^i = \mathbf{0}$ and $\boldsymbol{\pi}^i \mathbf{1} = 1$. This vector can be partitioned as $\boldsymbol{\pi}^i = [\boldsymbol{\pi}_0^i, \boldsymbol{\pi}_1^i, \dots, \boldsymbol{\pi}_{C_i}^i]$, where $\boldsymbol{\pi}_j^i$ corresponds to the states with j busy servers in station i , for $j = 0, \dots, C_i$. Hence the steady state distribution of the phase of the arrival process after an overflow is given by

$$\boldsymbol{\beta}^i = \frac{\boldsymbol{\pi}_{C_i}^i D_+^i}{\boldsymbol{\pi}_{C_i}^i D_+^i \mathbf{1}}. \quad (4)$$

By partitioning $\boldsymbol{\gamma}_i$ in the same way as $\boldsymbol{\pi}^i$ and using the fact that overflows can only occur when the system is full, we obtain that

$$\boldsymbol{\gamma}^i = [\mathbf{0}, \dots, \mathbf{0}, \boldsymbol{\beta}^i].$$

Given the large number of phases in $(\boldsymbol{\gamma}_i, E_0^i)$, we consider moment-matching algorithms to obtain a PH representation with fewer phases. In particular, we make use of the algorithms described in [7, 33] to match the first three moments with an acyclic PH distribution with minimal number of phases. In particular, when the squared coefficient of variation (SCV) of the inter-overflow time distribution is greater than or equal to $\frac{1}{2}$, the resulting PH distribution has only two phases. If the SCV is greater than one, the method given in [37] can also be used to get a 3-parameter hyper-exponential representation, which is a particular case of the PH class with two phases. For the specific case of the inter-overflow time distribution, we

found that the SCV was always above one in our experiments. This is closely related to the high variability inherent in an overflow process where the jobs are only overflowed in a small set of the state space. Additionally, we observed that the two-phase representation was able to capture the third moment of the inter-overflow distribution in the numerical instances considered in Section 5, as well as in many other cases not presented here. However, this is not the case in general, since the set of moments to match $\{n_i, 1 \leq i \leq 3\}$ have to fulfill the following condition: if $n_1 > 0$ and the SCV is greater than 1, the third moment can be represented by an acyclic PH distribution [33] or a hyper-exponential distribution [37], both with two phases, if and only if $3n_2^2 \leq 2n_1n_3$. In case the set of moments cannot be represented by a PH distribution of order two, the method in [7] determines the minimum number of phases required to do it with an acyclic PH distribution.

Let (α_i, A_i) be the parameters of the reduced PH distribution obtained from the moment-matching method applied to the $\text{PH}(\gamma^i, E_0^i)$ representation of the stationary inter-overflow time distribution at station i . In this reduced representation the off-diagonal elements of the matrix A_i describe transition rates without arrivals, while the transition rates related to arrivals are included in the matrix $-A_i \mathbf{1} \alpha_i$. Therefore, the matrices $\{\bar{D}_0^i, \bar{D}_1^i, \dots, \bar{D}_N^i\}$ can be approximated as

$$\bar{D}_0^i = A_i, \quad \bar{D}_j^i = -A_i \mathbf{1} \alpha_i \frac{\pi_{C_i}^i D_j^i \mathbf{1}}{\pi_{C_i}^i (D_+^i - D_s^i) \mathbf{1}}, \quad j = 1, \dots, N, j \neq s, \quad (5)$$

and $\bar{D}_s^i = 0$. For $j \notin \{0, s\}$ the approximate \bar{D}_j^i is the product of the matrix of transition rates involving an overflow of any type and the probability that a particular overflow is of type j . When the overflowing job is coming from the station s (right after site i), it is dropped and none of those jobs are actually part of the overflow process. This reduction will be labeled PH(3) in the results, making explicit that it relies on a PH representation that matches the first three moments of the inter-overflow time distribution.

3.1.2 Matching higher moments

In order to analyze the effect of capturing more moments of the inter-overflow time distribution we rely on the class of Matrix-Exponential (ME) distributions. An ME random variable has a density function given by $f(x) = \alpha e^{Sx} \mathbf{s}$, for $x \geq 0$. Here α is a $1 \times m$ vector, S is a square matrix of size m , and \mathbf{s} is an $m \times 1$ vector (m is called the order of the representation) [6]. An ME representation can always be chosen such that $\mathbf{s} = -S \mathbf{1}$ and $\alpha \mathbf{1} = 1$ [1, 4]. Therefore, an ME distribution, as in the PH case, is characterized by the tuple (α, S) , but the parameters are less restricted than in the PH class. More specifically, the matrix S must be invertible and the density must be non-negative and integrate to 1. Even though the entries of α and S can be complex numbers, the ME class is equally broad if they are restricted to be real [3]. Therefore, the ME class can be seen as a generalization of the PH class, since the tuple (α, S) must satisfy some extra conditions to be a representation of a PH distribution: α must be non-negative and sub-stochastic, while S must be the sub-generator of a CTMC, as explained before. A discussion and further properties of ME distributions can be found in [1, 2, 3, 13, 17, 28].

As shown in [4], the traditional analysis of QBD processes can be extended to admit more general components, e.g., allowing ME instead of PH distributions. Therefore, we can extend the analysis of each site to allow for arrivals coming from a marked Rational Arrival Process (MRAP[N]) with parameters $\{D_0^i, D_1^i, \dots, D_N^i\}$, as a generalization of the MMAP[N] [23]. These matrices have real-valued entries, their sum $D^i = \sum_{j=0}^N D_j^i$ must satisfy $D^i \mathbf{1} = \mathbf{0}$ and

the real part of the dominant eigenvalue of D_0^i (resp. D^i) must be negative (resp. equal to zero). In this case the stationary inter-overflow time distribution has an ME representation given by (γ^i, E_0^i) as defined above. This result relies on the fact that the time between successive overflows can be represented as the time until absorption in a finite-state semi-Markov process with ME holding times, which is itself ME distributed [3]. In this case E_0^i is a real matrix with negative dominant eigenvalue and the vector γ^i is no longer a probability mass function, but a vector of weights of measures after an overflow [2, 4].

Relying on the ME class we can use the method proposed in [34, 35] to reduce the order of the representation (γ^i, E_0^i) by matching $2n - 1$ moments with an ME distribution of order n . This method is based on the algorithm for the partial realization problem proposed in [15]. Using the resulting representation and Equation (5), we get a smaller set of parameters to approximately represent the overflow process, matching not only three, but $2n - 1$ moments of the stationary inter-overflow time distribution. The results based on an ME representation of order n will be labeled ME($2n - 1$). Nevertheless, the algorithm in [35] does not assure that the kernel matrix obtained from a set of moments defines a distribution function, i.e., the function has the required moments, but it may be negative. To the best of our knowledge, a characterization of the moments representable by an ME distribution of arbitrary order is not available. To make use of these distributions we evaluate the density function at several points to numerically verify if it is non-negative. This test is done for every reduced representation (α_i, A_i) obtained from the matching algorithms.

From the description above it is clear that the size of the reduced ME representation can be defined *a priori* according to the number of moments to match. Another approach is to rely on the characterization of the minimal order of an ME distribution given in [17]. There the authors propose the use of a set of Hankel matrices to determine the minimal order of an ME distribution. One of the Hankel matrices is built from the moments of the distribution and its rank determines the minimal order of an ME distribution with the specified moments. If the minimal order is found to be n , then $2n - 1$ moments can be used as input for the method in [34, 35] to obtain an ME representation of the stationary inter-overflow time distribution. With this method the size of the ME representation could be made variable, however, the minimum order could be large, making the analysis of the next station a lengthy process. Nevertheless, determining the minimal order in advance is useful to avoid the computation of a redundant amount of moments when using the moment-matching method in [35], since the method returns an ME representation of order n as long as n is smaller than or equal to the minimal order.

3.2 Non-Renewal Approximation

The approximation defined above matches a predefined number of moments of the stationary inter-overflow time distribution. Nevertheless, we can also include information related to the joint moments of consecutive inter-overflow times by means of a Rational Arrival Process (RAP) with parameters $\{H_0^i, H_1^i\}$ using the approach proposed in [34]. In this process the inter-event times are ME distributed, the matrix H_0^i describes the evolution of the process between events and H_1^i contains the arrival intensities. In the method of [34] the inter-overflow times are first approximated with an ME distribution using the algorithm in [35], as described above. Based on that result and the joint moments of the inter-overflow times, the method computes the matrices $\{H_0^i, H_1^i\}$ that describe the reduced RAP. If the reduced process is of order n , it not only matches $2n - 1$ moments of the inter-overflow time distribution, but also $(n - 1)^2$ joint moments of successive inter-overflow times. These matrices are used to

approximate the overflow process as

$$\bar{D}_0^i = H_0^i, \quad \bar{D}_j^i = H_1^i \frac{\pi_{C_i}^i D_j^i \mathbf{1}}{\pi_{C_i}^i (D_+^i - D_s^i) \mathbf{1}}, \quad j = 1, \dots, N, j \neq s,$$

and $\bar{D}_s^i = 0$. The results obtained using a reduced RAP representation of order n are labeled $\text{RAP}(2n - 1)$.

A further step consists of using the method in [23] to directly approximate the matrices $\{\bar{D}_0^i, \bar{D}_1^i, \dots, \bar{D}_N^i\}$ as the parameters of an $\text{MRAP}[N]$. To do so the inter-overflow time is again approximated by an ME distribution. However, in this case the description includes the joint moments of successive inter-overflow times for each type of arrival. This makes the method able to determine not just a matrix H_1^i describing the arrival intensities of any type, but a set of matrices $\{H_1^i, \dots, H_N^i\}$ with the intensities for each type of arrival. These matrices, together with H_0^i , completely determine the approximate overflow process that is fed to the next station in the ring, i.e., $\bar{D}_j^i = H_j^i$, $j = 0, 1, \dots, N$, where $H_s^i = 0$ by construction. The label $\text{MRAP}(2n - 1)$ will be used to refer to the results obtained with an order- n $\text{MRAP}[N]$ representation. Notice that the two methods discussed in this section rely on the algorithm in [35] to compute an approximate ME representation of the inter-overflow time distribution. Therefore, the density function of this ME representation is evaluated at several points to test its non-negativity, in a similar way as with the ME renewal approximation.

Before turning to the second approximation method, it is important to emphasize that in this method the size of the reduced representation of the overflow process does not depend on the size of the arrival process representation. Actually, it only depends on the number of moments and joint moments of the inter-overflow time process to match. Therefore, this method can be used when the arrival of new jobs at each station is described by a MAP, as explained in Section 2.1, since this generalization has no effect on the size of the overflow representation.

4 ON-OFF approximation

This approximation aims at capturing the behavior of those periods where station i is not generating overflow jobs (OFF periods) and those where arriving jobs are sent to the next station in the ring (ON periods) separately. Furthermore, the reduction process is split into two steps: a first related to the reduction of the OFF period representation and a second to the reduction of the ON period representation. The first time each station is analyzed only the OFF period representation will be reduced. Thus, after analyzing all the stations once, the overflow process will include jobs from all the stations and those coming from the first station must be dropped. To eliminate those jobs we reduce the ON period representation by lumping the states related to the first station, i.e., the station generating the jobs that need to be eliminated. Hereafter the analysis of each station first reduces the representation of its OFF period, and then eliminates the jobs that must be dropped by reducing the ON period representation. The overflow process is again represented by an $\text{MMAP}[N]$. The first step computes the matrices $\{\hat{D}_0^i, \hat{D}_1^i, \dots, \hat{D}_N^i\}$ of an $\text{MMAP}[N]$ that still includes the jobs that must be dropped. The result of the second step is the set of matrices $\{\bar{D}_0^i, \bar{D}_1^i, \dots, \bar{D}_N^i\}$ which characterizes the approximating overflow process. The size of the approximate process grows the first time each station is analyzed, but remains fixed in the subsequent iterations. In fact, if the external arrivals at each station follow a Poisson process, the size of the approximated process depends linearly on the number of stations. Assuming a more general process would cause the size of the approximated process to grow exponentially, making it intractable except

for very small networks. Therefore, in the exposition to follow we assume that new incoming jobs at each station arrive according to a Poisson process.

To describe the behavior of the system during the OFF periods we consider, from the process with generator (1), those states where there is at least one idle server. The transient generator of these states is given by

$$E_{\text{OFF}}^i = \begin{bmatrix} D_0^i & D_+^i & 0 & \dots & 0 & 0 \\ \mu_i I & D_0^i - \mu_i I & D_+^i & \dots & 0 & 0 \\ 0 & 2\mu_i I & D_0^i - 2\mu_i I & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & D_0^i - (C_i - 2)\mu_i I & D_+^i \\ 0 & 0 & 0 & \dots & (C_i - 1)\mu_i I & D_0^i - (C_i - 1)\mu_i I \end{bmatrix}. \quad (6)$$

Since in steady state all the servers are busy with probability vector $\pi_{C_i}^i$, the distribution of the arrival phase after a service completion that causes the system to make a transition from the states $\{(C_i, l), 1 \leq l \leq m_i\}$ to the states $\{(C_i - 1, l), 1 \leq l \leq m_i\}$ is given by

$$\eta^i = \frac{\pi_{C_i}^i C_i \mu_i I}{\pi_{C_i}^i C_i \mu_i I \mathbf{1}} = \frac{\pi_{C_i}^i}{\pi_{C_i}^i \mathbf{1}}. \quad (7)$$

Thus, the stationary distribution of the duration of an OFF period can be described as a PH distribution with parameters $(\delta^i, E_{\text{OFF}}^i)$, where

$$\delta^i = [\mathbf{0} \quad \dots \quad \mathbf{0} \quad \eta^i].$$

Since the size of this representation is $m_i C_i$, we can reduce it using the methods in [7, 33, 37] to match the first three moments of the OFF period length distribution, in the same manner as described in Section 3.1.1 for the inter-overflow time distribution. Let (α^i, A^i) be the reduced PH representation of the duration of the OFF period in station i , then α_j^i is the probability of starting an OFF period in phase j and $\mathbf{a}^i = -A^i \mathbf{1}$ is the rate at which an ON period starts in each phase of the OFF period representation. In addition, higher moments can be captured using an ME representation as illustrated in Section 3.1.2 for the inter-overflow time distribution. The results obtained with an approximation based on matching n moments of the OFF period length distribution will be labeled ON-OFF(n). However, in the remainder of this section we assume a PH representation as its interpretation is more intuitive.

On the other hand, because of the exponential service times, the duration of an ON period is exponentially distributed with rate $C_i \mu_i$. Furthermore, the stationary distribution of the arrival phase when an ON periods begins is given by

$$\omega^i = \frac{\pi_{C_i-1}^i D_+}{\pi_{C_i-1}^i D_+ \mathbf{1}},$$

where $\pi_{C_i-1}^i$ is the stationary probability vector of having $C_i - 1$ busy servers in station i . Using this description we can connect the OFF and ON periods in a single MMAP[N] process with parameters $\{\hat{D}_0^i, \hat{D}_1^i, \dots, \hat{D}_N^i\}$ given by

$$\hat{D}_0^i = \begin{bmatrix} A^i & \mathbf{a}^i \omega^i \\ C_i \mu_i \mathbf{1} \alpha^i & D_0^i - C_i \mu_i I \end{bmatrix}, \quad \hat{D}_j^i = \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & D_j^i \end{bmatrix}, \quad j = 1, \dots, N. \quad (8)$$

The structure of the matrices $\{\hat{D}_j^i, 1 \leq j \leq N\}$ clearly shows that the description of the ON period remains unchanged. In contrast, the matrix \hat{D}_0^i contains the reduced representation of the OFF period. The reduced process resides in the OFF states according to matrix A^i , from which it can move to the ON states with rates \mathbf{a}^i and select a state in this set according to the vector $\boldsymbol{\omega}^i$. The transitions within the ON states are described by the matrices $\{D_j^i, 1 \leq j \leq N\}$ when an arrival is associated with the transition, and by the matrix D_0^i when no arrivals are generated by the transition. Finally, the process may move from any of the ON states to the OFF set with rate $C_i\mu_i$. When this happens, a new OFF state is selected according to the vector $\boldsymbol{\alpha}^i$.

Let n_i be the order of the representation $(\boldsymbol{\alpha}^i, A^i)$. Then, after analyzing station 1 for the first time (considering its own external Poisson traffic only), the approximate representation of the overflow process at this station is of size $n_1 + 1$. This process is combined with the Poisson process arriving at station 2 and, after reducing the OFF period representation at this station, the overflow process fed to station 3 is of order $n_2 + n_1 + 1$. Therefore, the size of the overflow process from station N to station 1 is $\sum_{k=1}^N n_k + 1$. Notice, the overflow process is built such that the first n_N states describe the OFF period in which no jobs are sent to station 1. In the next n_{N-1} states, the site only overflows jobs of type N since these states correspond to the case where the station N faces an ON period while the station $N - 1$ resides in an OFF period. Accordingly, in the next n_{N-2} states both jobs of type N and $N - 1$ are overflowing, and from this set of states the process can move to the first set of n_N states or to the second set of n_{N-1} (as well as to other states). As a result, the overflow process captures the actual behavior of the ring network where, for a type $N - 1$ job to be sent to station 1, both stations N and $N - 1$ have to be in their ON periods. Furthermore, if the site N moves to the OFF set, both jobs of type N and $N - 1$ stop overflowing to station 1.

To illustrate the transitions in the overflow process we consider a simple network made of three nodes. The approximate overflow process at the first station is characterized by the size- $(n_1 + 1)$ matrices

$$\hat{D}_0^1 = \left[\begin{array}{c|c} A^1 & \mathbf{a}^1 \\ \hline C_1\mu_1\boldsymbol{\alpha}^1 & -(\lambda_1 + C_1\mu_1) \end{array} \right], \quad \hat{D}_1^1 = \left[\begin{array}{c|c} 0 & 0 \\ \hline 0 & \lambda_1 \end{array} \right], \quad \hat{D}_2^1 = \hat{D}_3^1 = 0.$$

Notice, the vector $\boldsymbol{\omega}^1$ is not included explicitly in the definition of \hat{D}_0^1 since in this case the ON period is described by one state which is selected with probability one when a new ON period starts. At station two the arrival process is built by superposing the newly-generated and the overflow jobs, as shown in Equation (3). After reducing the OFF period in this station, the matrices that characterize the overflow process from station two to station three are

$$\hat{D}_0^2 = \left[\begin{array}{c|c} A^2 & \mathbf{a}^2\boldsymbol{\omega}^2 \\ \hline C_2\mu_2\mathbf{1}\boldsymbol{\alpha}^2 & \begin{array}{c|c} A^1 - (\lambda_2 + C_2\mu_2)I & \mathbf{a}^1 \\ \hline C_1\mu_1\boldsymbol{\alpha}^1 & -\sum_{j=1}^2 (\lambda_j + C_j\mu_j) \end{array} \end{array} \right],$$

$$\hat{D}_1^2 = \left[\begin{array}{c|c} 0 & 0 \\ \hline 0 & \begin{array}{c|c} 0 & 0 \\ \hline 0 & \lambda_1 \end{array} \end{array} \right], \quad \hat{D}_2^2 = \left[\begin{array}{c|c} 0 & 0 \\ \hline 0 & \begin{array}{c|c} \lambda_2 I & 0 \\ \hline 0 & \lambda_2 \end{array} \end{array} \right], \quad \hat{D}_3^2 = 0.$$

Finally, the reduction of the OFF period representation is applied to station three, which results in an overflow process described by an MMAP[3] of size $n_3 + n_2 + n_1 + 1$ with matrices

$$\hat{D}_0^3 = \left[\begin{array}{c|cc} A^3 & & \mathbf{a}^3 \boldsymbol{\omega}^3 \\ \hline & A^2 - (\lambda_3 + C_3 \mu_3) I & \mathbf{a}^2 \boldsymbol{\omega}^2 \\ C_3 \mu_3 \mathbf{1} \boldsymbol{\alpha}^3 & C_2 \mu_2 \mathbf{1} \boldsymbol{\alpha}^2 & \begin{array}{c} A^1 - \sum_{j=2}^3 (\lambda_j + C_j \mu_j) I \\ C_1 \mu_1 \boldsymbol{\alpha}^1 \\ - \sum_{j=1}^3 (\lambda_j + C_j \mu_j) \end{array} \end{array} \right],$$

$$\hat{D}_1^3 = \left[\begin{array}{c|cc} 0 & & 0 \\ \hline & 0 & 0 \\ 0 & 0 & 0 \\ & 0 & \lambda_1 \end{array} \right], \quad \hat{D}_2^3 = \left[\begin{array}{c|cc} 0 & & 0 \\ \hline & 0 & 0 \\ 0 & \lambda_2 I & 0 \\ & 0 & \lambda_2 \end{array} \right], \quad \hat{D}_3^3 = \left[\begin{array}{c|cc} 0 & & 0 \\ \hline & \lambda_3 I & 0 \\ 0 & 0 & \lambda_3 \\ & 0 & \lambda_3 \end{array} \right].$$

The matrix \hat{D}_0^3 shows how the process can move from an OFF state to any of the ON states, even those where jobs of all types are generated. It also reveals how the service rate $C_i \mu_i$ determines the transition rates from the states where station i is in an ON period to those where it enters an OFF period. It is clear that the last state is the only one where jobs of every type are generated, including those coming from the first station, which must be dropped. Also, the last $n_1 + 1$ states describe the ON period of station two, where jobs of type two are generated. This fact suggests that the overflow process may be reduced by combining the last $n_1 + 1$ states into one single state describing the ON period of station two.

In general, after analyzing all the stations once, the overflow process at station N still includes jobs of type one. Since the relevant streams for this overflow process are those of type $\{2, \dots, N\}$, the process could be reduced to have $\sum_{k=2}^N n_k + 1$ states, where the last $n_2 + 1$ states describe the OFF and ON periods of station two. Specifically, to reduce the representation of the overflow process and to eliminate the type-1 jobs, we consider lumping the last $n_1 + 1$ states. From [8], a CTMC is strongly lumpable with respect to some partition if, for every pair of sets \mathcal{A} and \mathcal{B} in the partition, the sum of the transition rates from a state in \mathcal{A} to the states in \mathcal{B} is the same for every state in \mathcal{A} . In fact, if we define the set \mathcal{A}_1 containing the last $n_1 + 1$ states of the process, and the partition $\mathcal{E} = \{1, 2, \dots, \sum_{k=2}^N n_k, \mathcal{A}_1\}$, then the underlying chain of the overflow process is strongly lumpable with respect to \mathcal{E} . This can be easily seen in the construction of the overflow process. The transitions from and to the set \mathcal{A}_1 are contained only in the matrix \hat{D}_0 , as can be seen from Equation (8). Since only the last $n_1 + 1$ states are lumped, the required condition for strong lumpability is automatically met for all the single-state sets in the partition. To see that the transition rates from any state $s \in \mathcal{A}_1$ to any state $t \notin \mathcal{A}_1$ are the same for every element in \mathcal{A}_1 , define \mathcal{A}_i as the set of states $\{\sum_{k=i+1}^N n_k + 1, \dots, \sum_{k=i}^N n_k\}$, i.e., the set of states describing the OFF period of station i , for $2 \leq i \leq N$. Then the transitions from \mathcal{A}_1 to \mathcal{A}_i occur with rates $C_i \mu_i \mathbf{1} \boldsymbol{\alpha}^i$, for $2 \leq i \leq N$, which do not depend on the specific state in \mathcal{A}_1 . Thus, it is clear that the transition rates from any state $s \in \mathcal{A}_1$ to $t \notin \mathcal{A}_1$ are the same for all $s \in \mathcal{A}_1$. Clearly, this argument applies when the reduced representation of the OFF period is a PH distribution, as in this case the underlying process is a CTMC. When the reduced representation is an ME distribution this result no longer applies and we are not aware of a similar result involving this more general process. However, we apply the same reduction of the ON period, lumping the last $n_1 + 1$ states of the process, to eliminate the type-1 jobs from the overflow process at station N . This implies an additional approximation, but the results in Section 5 show that matching more moments, using ME distributions, improve the performance of the ON-OFF

approximation.

Lumping the state space of the underlying CTMC (or the more general process based on ME distributions) of the overflow process is the key step to reduce the representation of the ON period and to determine the matrices $\{\bar{D}_0^N, \bar{D}_1^N, \dots, \bar{D}_N^N\}$ from the matrices $\{\hat{D}_0^N, \hat{D}_1^N, \dots, \hat{D}_N^N\}$. Let $\hat{D}_0^N(\mathcal{A}_j, \mathcal{A}_k)$ be the transition rates from the set \mathcal{A}_j to the set \mathcal{A}_k without arrivals. Also, let \mathcal{A}'_1 be the unitary set containing the last state of the overflow process created by lumping the states in \mathcal{A}_1 . Then, the sub-generator matrix of the overflow process at station N , partitioned according to $\{\mathcal{A}_N, \dots, \mathcal{A}_2, \mathcal{A}'_1\}$, is given by

$$\bar{D}_0^N = \left[\begin{array}{cccc|c} \hat{D}_0^N(\mathcal{A}_N, \mathcal{A}_N) & \hat{D}_0^N(\mathcal{A}_N, \mathcal{A}_{N-1}) & \cdots & \hat{D}_0^N(\mathcal{A}_N, \mathcal{A}_2) & \hat{D}_0^N(\mathcal{A}_N, \mathcal{A}_1)\mathbf{1} \\ \hat{D}_0^N(\mathcal{A}_{N-1}, \mathcal{A}_N) & \hat{D}_0^N(\mathcal{A}_{N-1}, \mathcal{A}_{N-1}) & \cdots & \hat{D}_0^N(\mathcal{A}_{N-1}, \mathcal{A}_2) & \hat{D}_0^N(\mathcal{A}_{N-1}, \mathcal{A}_1)\mathbf{1} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \hat{D}_0^N(\mathcal{A}_2, \mathcal{A}_N) & \hat{D}_0^N(\mathcal{A}_2, \mathcal{A}_{N-1}) & \cdots & \hat{D}_0^N(\mathcal{A}_2, \mathcal{A}_2) & \hat{D}_0^N(\mathcal{A}_2, \mathcal{A}_1)\mathbf{1} \\ \hline C_N \mu_N \alpha^N & C_{N-1} \mu_{N-1} \alpha^{N-1} & \cdots & C_2 \mu_2 \alpha^2 & -\sum_{k=2}^N (C_k \mu_k + \lambda_k) \end{array} \right].$$

We now partition the state space into two sets, one including the first $\sum_{k=2}^N n_k$ states and the other being \mathcal{A}'_1 . With this partition, the matrices $\{\bar{D}_1^N, \dots, \bar{D}_N^N\}$ can be expressed as

$$\bar{D}_1^N = 0, \quad \bar{D}_j^N = \begin{bmatrix} \bullet & 0 \\ 0 & \lambda_j \end{bmatrix}, \quad j = 2, \dots, N,$$

where \bullet denotes the transitions in \hat{D}_j^N from $\{\mathcal{A}_N, \dots, \mathcal{A}_2\}$ to itself. This step concludes the reduction of the overflow representation related to the ON period. Furthermore, this reduction must be repeated at each station from the second iteration onwards, since for the first iteration the reduction is only related to the OFF period. Thus, during any future iteration, the representation of the overflow process that arrives at station i is of order $\sum_{k=1, k \neq i}^N n_k + 1$.

Given that it is usually possible to match the first three moments of the OFF periods in each station with a PH distribution of order 2, the size of the overflow process representation is $2N - 1$, which increases linearly with the number of stations. In contrast to the approximation proposed in Section 3, the ON-OFF method is not able to efficiently include more general arrival processes at the stations because the size of the arrival process representation at each station would increase exponentially.

5 Performance Results and Comparisons

In this section we evaluate the performance of the approximations introduced in this paper by comparing their results against those obtained by simulation, for different values of the network parameters. We consider different values for the overall load, the number of sites in the network, the number of servers per site, and the squared coefficient of variation (SCV) of the arrival process at each site. The results of the approximations presented in Section 3 are labeled depending on the specific representation of the inter-overflow time process, which can be a PH or ME distribution, a RAP or a marked RAP (MRAP). In each case the number of matched moments of the inter-overflow time distribution is made explicit, e.g., ME(n) refers to an approximation based on an ME representation matching n moments. The results of the approximation introduced in Section 4 are labeled ON-OFF(n) when n moments of the OFF period length distribution are matched.

The main performance measures are the number of newly-generated jobs processed locally per unit of time (*local rate*), and the total traffic rate transmitted from one site to another

(*link traffic*). These measures have been chosen as they are of particular relevance when dimensioning the optical grid resources [10, 12]. Since these measures may be different for each station in the network, we present the maximum *relative* error found when comparing them against the simulation results. The widths of the 95% confidence intervals of the estimates obtained from simulation (computed with the batch means method) are always less than one percent of the respective mean. All the network configurations considered here are assumed to have an overall load between 75% and 95%. The load at each site is randomly chosen within the range $[x - 5\%, x + 5\%]$, where x is the overall load of the network. This setup is in line with the design of Grid networks, where one aims for loads between 80% and 90% at each site. In prior studies the number of sites ranges from less than ten to hundred, with typical values between twenty and fifty [9, 10, 11, 12]. Also, typical values for the number of servers per site are between twenty and fifty, although in some cases it can be larger than one hundred. We consider networks with 10, 20, and 40 stations, each having the same number of servers (20 or 40).

We start with a network made of ten nodes with 20 and 40 servers per station. The *relative* errors in the approximation of the local rate and link traffic are included in Tables 1 and 2, respectively. The approximate local rates are very close to those obtained with simulation, particularly for loads up to 90%. When the load is higher the approximations that also match the joint moments of successive inter-overflow times perform significantly better than those based on renewal processes. The ON-OFF method shows a competitive performance, with similar results to those of the other methods for loads under 90%, but better behavior for higher loads. Additionally, all the methods perform better when the number of servers increases from 20 to 40, especially for high loads. In relation to the link traffic (Table 2), the errors in the approximations are clearly larger than for the local rate. Although for loads from 75% to 85% the errors grow with the load, this is no longer the case if the load is further increased. Similarly, for the same range of loads, the errors are smaller for the system with 40 servers per station than for the one with 20. Again, this behavior does not hold for higher loads. For loads up to 85% the methods provide a reasonable approximation, especially the ones based on the RAP and RAPK representation of the overflow process.

Load	0.75		0.80		0.85		0.90		0.95	
# Servers	20	40	20	40	20	40	20	40	20	40
PH(3)	0.10	0.02	0.19	0.07	0.20	0.22	3.02	0.34	14.59	7.69
ME(5)	0.05	0.03	0.15	0.03	0.13	0.14	2.51	0.25	12.49	6.70
RAP(5)	0.03	0.03	0.09	0.03	0.19	0.09	2.91	0.22	12.15	7.21
RAPK(5)	0.03	0.03	0.09	0.03	0.19	0.09	2.83	0.22	11.34	7.04
ON-OFF(3)	0.11	0.02	0.29	0.08	0.45	0.29	1.64	0.67	9.66	5.01

Table 1: Maximum *relative* error (%) in local rate for $N = 10$ and $C = \{20, 40\}$

We now turn to the analysis of larger systems that consist of 20 and 40 stations, each with 40 servers. The approximation errors for the local rate are included in Table 3, while those for the link traffic are in Table 4. In this case we provide the results for an increasing number of matched moments of the distribution of both the inter-overflow time and the OFF period length. The results labeled RAP(K) correspond to both RAP and RAPK representations since both have very similar results. This similarity was also apparent in Tables 1 and 2, as well as in many other scenarios not presented here. The approximation of the local rate is again very close to the estimates obtained by simulation, with improved results when compared to the smaller network. In general, the error in approximating the local rate diminishes with

Load	0.75		0.80		0.85		0.90		0.95	
# Servers	20	40	20	40	20	40	20	40	20	40
PH(3)	-7.96	-3.74	-12.57	-7.30	-17.94	-13.76	-15.53	-25.39	10.62	-6.93
ME(5)	-3.94	-0.96	-8.26	-3.26	-13.47	-8.80	-10.51	-19.64	12.63	-1.70
RAP(5)	-3.06	-0.88	-6.87	-2.75	-11.23	-7.71	-6.29	-17.32	16.59	3.47
RAPK(5)	-3.06	-0.88	-6.87	-2.75	-11.25	-7.71	-6.68	-17.33	13.44	2.79
ON-OFF(3)	-5.17	-2.06	-9.58	-4.99	-14.91	-10.88	-13.24	-22.04	4.57	-6.60

Table 2: Maximum *relative* error (%) in link traffic for $N = 10$ and $C = \{20, 40\}$

the inclusion of higher moments, with the exception of high loads for the system with 40 stations. On the other hand, the errors in the link traffic approximation decrease consistently when matching more moments. This effect is more evident for the methods based on the inter-overflow time process than for the ON-OFF method. For all the methods, the accuracy of the approximations improves drastically when matching five instead of three moments, and the results are even better when matching up to seven moments, although the difference is not as large as in the first case. It is important to note that the approximation methods fail to include more than seven moments because the algorithm in [35] returns a kernel matrix that corresponds to a negative density. In many other cases we obtained a similar result when trying to include nine or more moments.

Load		0.75		0.80		0.85		0.90		0.95	
# Stations		20	40	20	40	20	40	20	40	20	40
ME	3	0.06	0.06	0.07	0.11	0.18	0.19	0.74	0.66	2.40	1.02
	5	0.05	0.05	0.04	0.06	0.11	0.12	0.64	0.56	1.68	1.63
	7	0.05	0.05	0.05	0.07	0.15	0.16	0.77	0.70	1.19	2.12
RAP(K)	3	0.06	0.06	0.06	0.10	0.16	0.17	0.63	0.56	2.97	0.40
	5	0.05	0.05	0.04	0.05	0.06	0.08	0.43	0.35	2.47	0.73
	7	0.05	0.05	0.04	0.05	0.10	0.11	0.56	0.49	1.86	1.32
ON-OFF	3	0.06	0.06	0.09	0.11	0.25	0.25	1.05	0.98	0.46	2.84
	5	0.05	0.04	0.06	0.08	0.19	0.19	0.94	0.87	0.47	2.82
	7	0.05	0.04	0.06	0.08	0.19	0.20	0.95	0.87	0.43	2.86

Table 3: Maximum *relative* error (%) in local rate for $N = \{20, 40\}$ and $C = 40$

When comparing the ON-OFF approach with the inter-overflow time method it is clear that neither outperforms the other. In particular, when the overall load is equal to 95% the ON-OFF method has a better performance in approximating the local rate for the scenario with 20 stations while it is worse for the one with 40 stations. For lower loads, both methods show a similar performance, with a slightly smaller error for the ME and RAP(K) methods. However, the link traffic results show a different behavior. As this measure includes a mixture of all the traffic in the network, the errors accumulate causing larger discrepancies between the approximate methods and the simulation results. These discrepancies tend to increase with the size and the overall load of the network, as each link carries more types of jobs when any of these parameters increases. In particular, the ON-OFF method shows a better performance than the other methods when matching three moments of the corresponding inter-event time distribution. When five moments are matched the RAP(K) method performs better than the

ON-OFF, which now has similar errors than the ME approach. Finally, all the methods based on the inter-overflow time process perform better than the ON-OFF approach when seven moments are matched. This reveals that the inclusion of higher moments has a greater effect on the methods based on the inter-overflow time process, and this effect is more apparent for the non-renewal representations.

Load		0.75		0.80		0.85		0.90		0.95	
# Stations		20	40	20	40	20	40	20	40	20	40
ME	3	-5.64	-6.03	-9.38	-9.80	-12.24	-12.30	-24.44	-21.91	-30.35	-37.76
	5	-2.07	-2.57	-4.48	-4.92	-7.18	-7.24	-18.73	-16.01	-25.58	-33.50
	7	-1.62	-2.12	-3.83	-4.28	-6.46	-6.53	-17.83	-15.08	-24.67	-32.69
RAP(K)	3	-5.55	-5.92	-9.10	-9.53	-11.82	-11.88	-23.56	-21.00	-28.70	-36.29
	5	-1.86	-2.29	-3.79	-4.24	-6.09	-6.15	-16.41	-13.61	-21.28	-29.66
	7	-1.37	-1.78	-2.98	-3.43	-5.10	-5.16	-14.93	-12.08	-19.29	-27.87
ON-OFF	3	-3.59	-4.06	-6.54	-6.97	-9.30	-9.36	-21.04	-18.39	-27.30	-35.00
	5	-1.74	-2.26	-4.11	-4.55	-6.88	-6.94	-18.58	-15.85	-25.69	-33.55
	7	-1.57	-2.10	-3.87	-4.32	-6.62	-6.68	-18.27	-15.53	-25.43	-33.30

Table 4: Maximum *relative* error (%) in link traffic for $N = \{20, 40\}$ and $C = 40$

In relation to the computation times, the inter-overflow time method has a clear advantage since the size of the approximate representation of the overflow process is independent of the number of stations. In the ON-OFF method this size grows linearly with the number of stations, which generates a proportional increase in the block size of the QBD that must be solved at each station. The relevance of the block size for the computation times comes from the fact that the time complexity of the algorithms to compute the stationary distribution and the moments of the distribution of the first-passage times to higher levels of the QBD is cubic in the block size. For the inter-overflow time method, the computation times were always below one minute in the scenarios considered here. The load also affects the computation times as more iterations are required for the traffic matrix to converge when the load is higher. For the scenarios shown here, usually less than five iterations were enough when the load was less than 80%, but this figure was between 20 and 30 when the load was 95%.

For the results considered thus far, the arrival process of newly-generated jobs at each station is assumed to be Poisson. Even though this is a usual assumption in prior studies on optical grid networks, more general processes can also be considered in order to capture the high variability of the arrival process [9, 36]. As mentioned in Section 2, the method based on the inter-overflow time distribution can be extended to allow for more general arrival processes without experiencing an exponential increase in the size of the overflow process representation. To consider the case of high variability in the arrival process, we assume that each station receives newly-generated jobs coming from a hyper-exponential renewal process with different arrival rates at each station, but the same SCV. Specifically, the cases where the SCV is equal to 5 and 20 are analyzed for a network with 20 stations and 40 servers per station. The approximation errors in local rate and traffic link are included in Tables 5 and 6, respectively. Again, the RAP and RAPK representations are presented together as their results are very similar. The first clear result is that the increase in the SCV causes larger errors in the approximation of both the local rate and the link traffic. Not only are the errors for SCV equal to 20 larger than those for SCV equal to 5, but in both cases the approximation is worse than under Poisson arrivals (Tables 3 and 4). With respect to the local rate, the

effect of matching higher moments becomes more evident for larger SCV. The approximations offer small errors for loads up to 90%, and for higher loads these are still below 10% (when matching seven moments of the inter-overflow time distribution). The results for the ME(7) approximation with an SCV equal to 20 are not included because the density function obtained with the algorithm in [35] was negative. For the link traffic, the effect of matching higher moments is evident, with a stronger effect for the non-renewal approximations. The errors tend to increase with the load, except for a load of 95%.

Load		0.75		0.80		0.85		0.90		0.95	
SCV		5	20	5	20	5	20	5	20	5	20
ME	3	0.14	0.76	0.40	1.09	0.33	1.25	1.59	4.63	10.73	15.22
	5	0.09	0.34	0.17	0.62	0.33	0.37	0.66	3.28	9.27	13.25
	7	0.14	0.59	0.28	1.11	0.59	1.21	0.59	0.94	8.24	-
RAP(K)	3	0.16	0.80	0.42	1.17	0.42	1.39	1.82	4.90	11.42	15.77
	5	0.08	0.39	0.20	0.49	0.17	0.66	0.98	3.87	10.21	14.51
	7	0.11	0.51	0.22	0.94	0.47	0.91	0.28	1.71	9.00	8.90

Table 5: Maximum *relative* error in local rate (%) for $N = 20$, $C = 40$ and $SCV = \{5, 20\}$

Load		0.75		0.80		0.85		0.90		0.95	
SCV		5	20	5	20	5	20	5	20	5	20
ME	3	-9.77	-20.91	-15.51	-29.90	-23.16	-38.03	-34.36	-42.86	-26.85	-34.42
	5	-5.79	-17.38	-10.46	-25.72	-17.87	-33.11	-28.87	-36.71	-20.77	-26.19
	7	-5.08	-16.16	-9.48	-24.32	-16.77	-31.53	-27.59	-34.48	-19.21	-
RAP(K)	3	-9.59	-20.69	-15.17	-29.61	-22.67	-37.67	-33.61	-42.41	-25.51	-32.62
	5	-5.28	-16.49	-9.48	-24.54	-16.49	-31.62	-26.71	-34.63	-16.74	-21.86
	7	-4.39	-14.84	-8.17	-22.54	-14.88	-29.23	-24.60	-31.25	-13.59	-16.20

Table 6: Maximum *relative* error in link traffic (%) for $N = 20$, $C = 40$ and $SCV = \{5, 20\}$

From the instances analyzed here, and others not included, we have found that a larger SCV of the arrival process and larger overall load deteriorate the approximation performance. The number of servers has an opposite effect, i.e., the approximation errors are smaller when the number of servers is larger. These observations are related to the impact that these parameters have on the spill probability (the probability that a job is sent to a remote station for service). Clearly, both the arrival process SCV and the load increase the spill probability, while, under the same load, a larger number of servers reduces the spill probability as more jobs can be handled locally. Figure 1 shows the errors in local rate and link traffic as a function of the spill probability for two specific methods: ON-OFF(3) and RAP(7). There we observe that the approximations behave better when the spill probability is small, typically under 0.3. For a larger spill probability the approximation errors become more pronounced, especially for the link traffic. Besides, in most of the cases the approximation methods underestimate the actual link traffic. Thus, the results of these methods can be regarded as a practical lower bound of the actual link traffic.

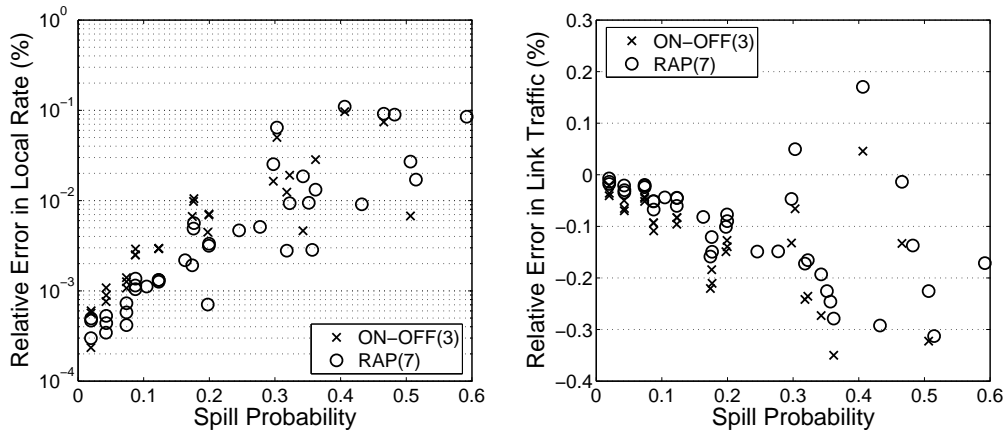


Figure 1: Maximum *relative* errors in local rate and link traffic as a function of the spill probability

A Computing the moments

A relevant issue for the approximation methods introduced in this paper is the computation of the moments of an inter-event time distribution. The purpose of this appendix is to provide an algorithm to compute these moments based on [14]. In both methods, the inter-event time distribution has a PH or ME representation characterized by a matrix with a QBD structure. For the method in Section 3, this matrix is shown in Equation (2), while for the ON-OFF method it is given by Equation (6). In both cases the inter-event time distribution can be seen as a first-passage time distribution to a higher level in a finite QBD (level k is the set of states $\{(k, l), 1 \leq l \leq m_i\}$). In the ON-OFF approximation the inter-event times correspond to the length of the OFF periods. An OFF period starts when the system moves from level C_i to level $C_i - 1$, where the initial phase in level $C_i - 1$ is selected according to $\boldsymbol{\eta}^i$, as defined in Equation (7). The OFF period ends as soon as the process reaches level C_i again, i.e., the first time it visits level C_i starting from level $C_i - 1$. A similar observation can be made for the approximation based on the overflow process. Since an overflow can only occur when the system is in level C_i , an inter-overflow time always starts in this level. In this case, however, there is no higher level than C_i , but we can define a fictitious level $C_i + 1$ of absorbing states that can only be accessed from level C_i with rates $D_+^i - D_s^i$. Then, the inter-overflow time can be seen as the first-passage time from level C_i to level $C_i + 1$, where the vector $\boldsymbol{\beta}^i$ in Equation (4) defines the initial phase in level C_i .

As both inter-event time distributions can be regarded as first-passage times to higher levels in a finite QBD, we can use the results in [14] to compute the moments of these distributions. In [14] the authors determine the generating function of the first-passage times to higher levels in a finite level-dependent QBD. Based on this, they derive an algorithm to compute the first two moments, but any higher moment can be computed from the generating function in a similar manner. Here we focus on the computation of N moments of the OFF-period length distribution, but this can be easily modified for the inter-overflow time distribution. Algorithm 1 presents the steps to compute these moments in an arbitrary station, dropping the super(sub)script i . For a better understanding of the algorithm we introduce a few definitions. As described in Section 2, the state of a station can be represented by the process $\{(N(t), J(t)), t \geq 0\}$ on the state space $\{(k, l), 0 \leq k \leq C, 1 \leq l \leq m\}$. Now let X_n be the first-passage time from level $n - 1$ to level n , for $n = 1, \dots, C$, i.e.,

$X_n = \inf\{t > 0 : N(t) = n | N(0) = n - 1\}$. Therefore, the length of the OFF-period is given by X_C . Let $U_n^{(k)}$ be the $m \times m$ matrix with entries

$$[U_n^{(k)}]_{ij} = E[X_n^k, J(X_n) = j | N(0) = n - 1, J(0) = i],$$

i.e, the k -th moment of the first-passage time from level $n - 1$ to level n , when the process starts in phase i and the first visit to level n occurs in phase j , for $1 \leq n \leq C$, $1 \leq i, j \leq m$ and $k \geq 1$. Also, let $u_n^{(k)}$ be the $m \times 1$ vector with entries

$$[u_n^{(k)}]_i = E[X_n^k | N(0) = n - 1, J(0) = i],$$

i.e, the k -th moment of the first-passage time from level $n - 1$ to level n , given that the process starts in phase i . Finally, let r_k be the k -th moment of the first-passage time distribution from level $C - 1$ to level C , i.e.,

$$r_k = E[X_C^k | N(0) = C - 1] = \boldsymbol{\eta} \mathbf{u}_C^{(k)},$$

where $\boldsymbol{\eta}$ is the probability distribution of the initial phase in level $C - 1$. Algorithm 1 computes r_k , for $1 \leq k \leq N$. The procedure starts by computing the moments of the first-passage times from level 0 to level 1, resulting in the vectors $\mathbf{u}_1^{(k)}$. The algorithm then computes the same quantities for levels 2 to C . To obtain the N vectors $\{\mathbf{u}_n^{(k)}, 1 \leq k \leq N\}$ at each step, the algorithm requires the $N - 1$ matrices $\{U_n^{(k)}, 1 \leq k \leq N - 1\}$ and the matrix K_n . The matrix $-K_n^{-1}$ is the generator of the process restricted to level $n - 1$ before the first visit to level n . Additionally, the $N - 1$ auxiliary matrices $\{T_n^{(k)}, 1 \leq k \leq N - 1\}$ are used to simplify the expressions.

Algorithm 1 Algorithm to compute N moments of the OFF period length distribution

```
 $K_1 \leftarrow -D_0^{-1}$ 
 $\mathbf{u}_1^{(1)} \leftarrow K_1 \mathbf{1}$ 
for  $i = 2$  to  $N$  do
   $\mathbf{u}_1^{(i)} \leftarrow iK_1 \mathbf{u}_1^{(i-1)}$ 
end for
 $U_1^{(1)} \leftarrow K_1^2 D_+$ 
for  $i = 2$  to  $N - 1$  do
   $U_1^{(i)} \leftarrow iK_1 U_1^{(i-1)}$ 
end for
for  $n = 2$  to  $C$  do
   $K_n \leftarrow -(D_0 - (n-1)\mu I + (n-1)\mu K_{n-1} D_+)^{-1}$ 
   $T_n^{(1)} \leftarrow K_n (I + (n-1)\mu U_{n-1}^{(1)})$ 
  for  $i = 2$  to  $N - 1$  do
     $T_n^{(i)} \leftarrow (n-1)\mu K_n U_{n-1}^{(i)}$ 
  end for
   $\mathbf{u}_n^{(1)} \leftarrow K_n (\mathbf{1} + (n-1)\mu \mathbf{u}_{n-1}^{(1)})$ 
  for  $i = 2$  to  $N$  do
     $\mathbf{u}_n^{(i)} \leftarrow (n-1)\mu K_n \mathbf{u}_{n-1}^{(i)} + \sum_{j=1}^{i-1} \binom{i}{j} T_n^{(i-j)} \mathbf{u}_n^{(j)}$ 
  end for
  for  $i = 1$  to  $N - 1$  do
     $U_n^{(i)} \leftarrow T_n^{(i)} K_n D_+ + \sum_{j=1}^{i-1} \binom{i}{j} T_n^{(i-j)} U_n^{(j)}$ 
  end for
end for
for  $k = 1$  to  $N$  do
   $r_k \leftarrow \boldsymbol{\eta} \mathbf{u}_C^{(k)}$ 
end for
```

References

- [1] S. Asmussen and M. Bladt. Renewal theory and queueing algorithms for matrix-exponential distributions. In S. Chakravarty and A. S. Alfa, editors, *Matrix-Analytic Methods in Stochastic Models*, pages 313–341. Marcel Dekker, New York, 1996.
- [2] S. Asmussen and M. Bladt. Point processes with finite-dimensional conditional probabilities. *Stochastic Processes and their Applications*, 82:127–142, 1999.
- [3] S. Asmussen and C. A. O’Cinneide. Matrix-exponential distributions. In S. Kotz, C. B. Read, and D. L. Banks, editors, *Encyclopedia of Statistical Sciences*, volume 2, pages 435–440. Wiley, New York, 1998.
- [4] N. G. Bean and B. F. Nielsen. Quasi-birth-and-death processes with rational arrival process components. Technical Report 2007-20, Informatics and Mathematical Modelling, Technical University of Denmark, DTU, 2007.
- [5] R. Bellman. *Introduction to Matrix Analysis*. McGraw-Hill, 1970.
- [6] M. Bladt and M. F. Neuts. Matrix-exponential distributions: calculus and interpretations via flows. *Stochastic Models*, 19:113–124, 2003.

- [7] A. Bobbio, A. Horváth, and M. Telek. Matching three moments with minimal acyclic phase type distributions. *Stochastic Models*, 21:303–326, 2005.
- [8] C. J. Burke and M. Rosenblatt. A Markovian function of a Markov chain. *The Annals of Mathematical Statistics*, 29:1112–1122, 1958.
- [9] K. Christodoulopoulos, M. Varvarigos, C. Develder, M. De Leenheer, and B. Dhoedt. Job demand models for optical grid research. In *Proc. 11th Conference on Optical Network Design and Modelling (ONDM)*, Athens, Greece, 2007.
- [10] M. De Leenheer, C. Develder, T. Stevens, B. Dhoedt, M. Pickavet, and P. Demeester. Design and control of optical grid networks (invited). In *Proc. 4th Int. Conf. on Broadband Networks (Broadnets 2007)*, Raleigh, NC, Sep. 2007.
- [11] M. De Leenheer, F. Farahmand, P. Thysebaert, B. Volckaert, F. De Turck, B. Dhoedt, P. Demeester, and J. Jue. Anycast routing in optical burst switched grid networks. In *Proc. 31st European Conference on Optical Communication (ECOC)*, Glasgow, Scotland, 2005.
- [12] C. Develder, M. De Leenheer, T. Stevens, B. Dhoedt, F. De Turck, and P. Demeester. Scheduling in optical grids: a dimensioning point of view. In *Proc. Conference on the Optical Internet - Australian Conference on Optical Fibre Technology (COIN-ACOFT)*, Melbourne, Australia, 2007.
- [13] M. Fackrell. Fitting with matrix-exponential distributions. *Stochastic Models*, 21:377–400, 2005.
- [14] D. P. Gaver, P. A. Jacobs, and G. Latouche. Finite birth-and-death models in randomly changing environments. *Advances in Applied Probability*, 16:715–731, 1984.
- [15] W. B. Gragg and A. Lindquist. On the partial realization problem. *Linear Algebra and its Applications*, 50:277–319, 1983.
- [16] Q. He and M. F. Neuts. Markov chains with marked transitions. *Stochastic Processes and their Applications*, 74:37–52, 1998.
- [17] Q. He and H. Zhang. On matrix exponential distributions. *Advances in Applied Probability*, 39:271–292, 2007.
- [18] H. Heffes. A class of data traffic processes - covariance function characterization and related queuing results. *Bell System Technical Journal*, 59:897–929, 1980.
- [19] A. Heindl. Decomposition of general queueing networks with MMPP inputs and customer losses. *Performance Evaluation*, 51:117–136, 2003.
- [20] A. Heindl, K. Mitchell, and A. van de Liefvoort. Correlation bounds for second-order MAPs with application to queueing network decomposition. *Performance Evaluation*, 63:553–577, 2006.
- [21] A. Heindl, Q. Zhang, and E. Smirni. ETAQA truncation models for the MAP/MAP/1 departure process. In *Proceedings of the First QEST: Quantitative Evaluation of Systems*, Enschede, The Netherlands, 2004.

- [22] A. Horváth, G. Horváth, and M. Telek. A joint moments based analysis of networks of MAP/MAP/1 queues. In *Proc. of the 5th International Conference on the Quantitative Evaluation of Systems (QEST)*, St Malo, France, Sept 2008.
- [23] A. Horváth, G. Horváth, and M. Telek. A traffic based decomposition of two-class queueing networks with priority service. *Computer Networks*, 53:1235 – 1248, 2009.
- [24] A. Kuczura. The interrupted Poisson process as an overflow process. *Bell System Technical Journal*, 52:437–448, 1973.
- [25] A. Kuczura and D. Bajaj. A method of moments for the analysis of a switched communication network’s performance. *IEEE Transactions on Communications*, COM-25:185–193, 1977.
- [26] P. J. Kuehn. Approximate analysis of general queueing networks by decomposition. *IEEE Transactions on Communications*, COM-27:113–126, 1979.
- [27] G. Latouche and V. Ramaswami. *Introduction to Matrix Analytic Methods in Stochastic Modeling*. ASA-SIAM Series on Statistics and Applied Probability. SIAM, Philadelphia, PA, 1999.
- [28] L. Lipsky. *Queueing Theory: a Linear Algebraic Approach*. Macmillan, New York, 1992.
- [29] D. Lucantoni, K. S. Meier-Hellstern, and M. F. Neuts. A single-server queue with server vacations and a class of non-renewal arrival processes. *Advances in Applied Probability*, 22:676–705, 1990.
- [30] J. Matsumoto and Y. Watanabe. Individual traffic characteristics of queueing systems with multiple Poisson and overflow inputs. *IEEE Transactions on Communications*, COM-33:1–9, 1985.
- [31] K. S. Meier-Hellstern. The analysis of a queue arising in overflow models. *IEEE Transactions on Communications*, 37:367–372, 1989.
- [32] M.F. Neuts. *Matrix-Geometric Solutions in Stochastic Models*. The John Hopkins University Press, Baltimore, 1981.
- [33] M. Telek and A. Heindl. Matching moments for acyclic discrete and continuous phase-type distributions of second order. *International Journal of Simulation Systems, Science & Technology*, 3:47–57, 2002.
- [34] M. Telek and G. Horváth. A minimal representation of Markov arrival processes and a moment matching method. *Performance Evaluation*, 64:1153–1168, 2007.
- [35] A. van de Liefvoort. The moment problem for continuous distributions. Technical report, University of Missouri, 1990.
- [36] B. Van Houdt, C. Develder, J. F. Pérez, M. Pickavet, and B. Dhoedt. Mean field calculation for optical grid dimensioning. To appear in *IEEE/OSA Journal of Optical Communications and Networking*.
- [37] W. Whitt. Approximating a point process by a renewal process, I: Two basic methods. *Operations Research*, 30:125–147, 1982.
- [38] R. I. Wilkinson. Theories for toll traffic engineering in the U.S.A. *Bell System Technical Journal*, 35:421–514, 1956.