

BOUNDING THE PARTITION FUNCTION OF BCMP MULTICLASS QUEUEING NETWORKS

J. Anselmi, P. Cremonesi

Politecnico di Milano, DEI

Milano, Italy

jonatha.anselmi@polimi.it

In this paper, we provide an inequality which bounds from above the partition function of multiclass, closed BCMP queueing networks. The inequality basically derives from the integral representation of the partition function and from the Holder's inequality. It essentially states that the partition function of a closed, multiclass network can be upper bounded by a product of partition functions related to singleclass networks. Hence, its computation is efficient even for very large models. The inequality is important from a theoretical point of view and provides a way to estimate the minimum amount of memory that exact solution algorithms implementations should allocate to avoid numerical instabilities.

Keywords: BCMP queueing networks, multiclass, partition function, normalizing constant, bound, Holder's inequality

1. INTRODUCTION

During the last decades, closed BCMP queueing networks [1] have been widely adopted to analytically evaluate the performance of computer and communication systems. Such networks extend the Gordon-Newell networks [5] to multiclass networks, i.e. networks with jobs of different types, more types of service disciplines and more general service times probability distributions. BCMP networks have the so called *product-form* solution which essentially means that the stationary probabilities of the states of the underlying Markov chain characterizing a BCMP queueing network can be expressed through the product of a number of simple terms related to network stations up to a multiplicative constant, called *partition function*, which normalizes product terms.

The problem of the efficient computation of such partition function is a well-known difficult problem that attracted the attention of many researchers during the last three decades and a number of works have been proposed, see e.g. [9, 8, 4]. However, in the general case they all exhibit numerical instabilities and this is mainly due to the fact that most of them are based on the partition function computation. We emphasize that even mean value approaches such as [8] are numerically unstable as long as Load-Dependent (LD) stations, i.e. stations with processing speeds variable with the number of jobs in their queues, are considered. These numerical instabilities strongly affect

the accuracy of the solution and eventually yield unfeasible solutions, e.g. negative throughputs, [8].

In this paper, we provide an inequality related to the partition function of closed, multiclass BCMP queueing network models. The inequality follows algebraically from the integral representation of the partition function [7] and the Holder's inequality. The inequality bounds from above the partition function of a multiclass network in terms of R partition functions related to singleclass models where R is the number of customers classes. It can be used to estimate the magnitude of the partition function itself as well as the precision size that current implementations of existing algorithms should allocate to avoid numerical instabilities.

2. MULTICLASS, CLOSED BCMP NETWORKS

2.1. Notation. We consider closed, multiclass BCMP queueing network models [1]. The network is composed of M stations and the jobs circulating in the network are partitioned into R classes. Stations are either load-dependent (LD) or load-independent (LI), i.e. their service speed can depend on the number of jobs in their queue or not. An Infinite Server (also known as pure delay) station is introduced and it is indexed by 0. If not otherwise specified, indices i and j will range from 0 to M indexing network stations.

The evolution of a BCMP network is markovian and we denote by $n_{ir} \geq 0$ the number of class- r jobs in station i observed in a given instant, by $\mathbf{n}_i = [n_{i1}, \dots, n_{iR}]$ the associated population vector in station i , by $n_i = n_{i1} + \dots + n_{iR}$ the total number of jobs in i , and by matrix $\vec{\mathbf{n}} = [\mathbf{n}_0, \mathbf{n}_1, \dots, \mathbf{n}_M]'$ an *aggregate* state of the underlying Markov chain of the queueing network.

We also denote by

- N_r , the (constant) number of class- r jobs circulating in the network,
- $\mathbf{N} = (N_1, N_2, \dots, N_R)$, the population vector,
- $N = N_1 + N_2 + \dots + N_R$, the total number of jobs without class distinction.
- $\pi(\vec{\mathbf{n}})$, the network stationary probabilities of being in state $\vec{\mathbf{n}}$,
- $\pi_i(\mathbf{n}_i)$, the stationary probability of having \mathbf{n}_i jobs in station i ,
- $\rho_{ir} \geq 0$, $i \geq 1$, is the mean *loading* of station i for class- r jobs (also known as relative utilization, or service demands) and for a closed BCMP network it represents the average time spent by a class- r job at station i during its full execution when using the network alone and visiting a reference station, say station 1, once.
- $\rho_{0r} \geq 0$, the mean class- r "think" time spent in the Infinite Server station,
- \mathbf{e}_r , the size- R orthogonal unit vector in direction r .

Let also $x_i : \mathbb{N} \setminus \{0\} \rightarrow \mathbb{R}^+$, $i \geq 1$, be an arbitrary positive function of the number of jobs which visit i . $x_i(n)$ represents the LD rate of service of i when there are n jobs in i *relative* to the service rate when $n = 1$, i.e. $x_i(1) = 1$. Analogously, let $y_{ir} : \mathbb{N} \setminus \{0\} \rightarrow \mathbb{R}^+$, $i \geq 1$, be the LD rate of service of class- r jobs in station i as function of the total number of jobs it contains *relative* to the class- r service rate of i when exactly one (class- r) job is present, i.e. $y_{ir}(1) = 1$. It is well-known that the

model discussed above with stations providing such types of load-dependencies satisfies the BCMP assumptions [1]. For simplicity, we say that station i provides a *type 1* load-dependence if its load-dependence is a function of the total number of jobs in i . Analogously, station i provides a *type 2* load-dependence if its load-dependence is a function of only the number of jobs in i of a specific class.

2.2. Product-form Solution. As long as closed BCMP networks with LI stations and a delay are considered, it is known that the stationary probabilities are given by the following product-form formula [1]

$$\pi(\vec{\mathbf{n}}) = G^{-1}(\mathbf{N}) \prod_{i=0}^M \pi_i(\mathbf{n}_i) \quad (1)$$

where

$$\pi_i(\mathbf{n}_i) = n_i! \prod_{r=1}^R \frac{\rho_{ir}^{n_{ir}}}{n_{ir}!}, i \geq 1, \quad \pi_0(\mathbf{n}_0) = \prod_{r=1}^R \frac{\rho_{0r}^{n_{0r}}}{n_{0r}!} \quad (2)$$

are the non-normalized probabilities of having \mathbf{n}_i jobs in station $i \geq 1$ and \mathbf{n}_0 jobs in the Infinite Server station and

$$G(\mathbf{N}) = \sum \prod_{i=0}^M \pi_i(\mathbf{n}_i) \quad (3)$$

is the partition function and the sum is taken over the state space defined by set

$$\mathbf{S}(\mathbf{N}) = \left\{ \vec{\mathbf{n}} : \sum_{i=0}^M n_{ir} = N_r, \forall r \right\} \quad (4)$$

(where each n_{ir} is a non-negative integer). For networks with LD stations, the formula in (2) becomes ($i \geq 1$)

$$\pi_i(\mathbf{n}_i) = \begin{cases} n_i! \prod_{r=1}^R \frac{\rho_{ir}^{n_{ir}}}{n_{ir}!} \cdot \frac{1}{\prod_{t=1}^{n_i} x_i(t)} & \text{for type-1 load-dependency,} \\ n_i! \prod_{r=1}^R \frac{\rho_{ir}^{n_{ir}}}{n_{ir}!} \frac{1}{\prod_{t=1}^{n_{ir}} y_{ir}(t)} & \text{for type-2 load-dependency} \end{cases} \quad (5)$$

and the partition function (3) must be modified accordingly to make all (5) sum to one.

It is evident that the partition function computation through the direct summation (3) is impractical from both a computational and numerical point of view.

3. INEQUALITY

Given a closed, multiclass BCMP network, we first make the following replacements

$$\rho_{ir} \leftarrow \begin{cases} \rho_{ir} / \inf_n x_i(n) & \text{if } i \text{ is type-1 load-dependent} \\ \rho_{ir} / \inf_n y_{ir}(n) & \text{if } i \text{ is type-2 load-dependent} \end{cases} \quad (6)$$

Within replacements (6), it is clear that the resulting queueing network is composed of only LI stations and the new value of the partition function bounds from above the original one.

The following integral representation of the partition function has been shown in [7]

$$G(\mathbf{N}) = \frac{1}{\prod_{r=1}^R N_r!} \int_{\mathfrak{R}^{+M}} \prod_{r=1}^R H(r, \mathbf{u})^{N_r} e^{-(u_1 + \dots + u_M)} d\mathbf{u} \quad (7)$$

where $H(r, \mathbf{u}) = \rho_{0r} + \tilde{\rho}_{1r}u_1 + \dots + \tilde{\rho}_{Mr}u_M$. The expression (7) can be rewritten as

$$G(\mathbf{N}) = \frac{1}{\prod_{r=1}^R N_r!} \int_{\mathfrak{R}^{+M}} \prod_{r=1}^R [H(r, \mathbf{u})^N e^{-(u_1 + \dots + u_M)}]^{\beta_r} d\mathbf{u} \quad (8)$$

where $\beta_r = N_r/N$. Given that, by definition, $\sum_r \beta_r = 1$, we can now apply the Holder's inequality. This yields

$$G(\mathbf{N}) \leq \prod_{r=1}^R \frac{1}{N_r!} \left[\int_{\mathfrak{R}^{+M}} H(r, \mathbf{u})^N e^{-(u_1 + \dots + u_M)} d\mathbf{u} \right]^{\beta_r}. \quad (9)$$

Now, multiplying and dividing by $N!^{\beta_r}$ each right-hand product term of (9), we obtain

$$G(\mathbf{N}) \leq \prod_{r=1}^R \frac{N!^{\beta_r}}{N_r!} \left[\frac{1}{N!} \int_{\mathfrak{R}^{+M}} H(r, \mathbf{u})^N e^{-(u_1 + \dots + u_M)} d\mathbf{u} \right]^{\beta_r} \quad (10)$$

and we note that the expression in the brackets can be interpreted as the integral representation of the partition function of a singleclass network populated by N class- r jobs only. Hence, (10) can be rewritten as

$$G(\mathbf{N}) \leq \prod_{r=1}^R \frac{[N!G(N\mathbf{e}_r)]^{\beta_r}}{N_r!} = \binom{N}{N_1, \dots, N_R} \prod_{r=1}^R G(N\mathbf{e}_r)^{\beta_r} \quad (11)$$

and the upper bound on the partition function of a closed, multiclass BCMP queueing network is provided.

We remark that the value of $G(N\mathbf{e}_r)$ can now be easily computed exploiting single-class algorithms [6, 3]. For instance, assuming $\rho_{0r} = 0$ and $\tilde{\rho}_{ir} \neq \tilde{\rho}_{jr}$ for all $i \neq j$, from the Koenigsberg's formula [6, 2]

$$G(N\mathbf{e}_r) = \sum_{i=1}^M \frac{\tilde{\rho}_{ir}^{N+M-1}}{\prod_{j=1, j \neq i}^M (\tilde{\rho}_{ir} - \tilde{\rho}_{jr})} \quad (12)$$

it follows

$$G(N\mathbf{e}_r) \approx \frac{(\max_{i \geq 1} \tilde{\rho}_{ir})^{N+M-1}}{\prod_{j=1, j \neq \arg \max_{i \geq 1} \tilde{\rho}_{ir}}^M (\max_{i \geq 1} \tilde{\rho}_{ir} - \tilde{\rho}_{jr})}. \quad (13)$$

which provides an estimate of the order of magnitude of $G(N\mathbf{e}_r)$ (exploiting the theory of residues, an analogous result can be obtained even in the case where $\tilde{\rho}_{ir} = \tilde{\rho}_{jr}$ for some i and j [2]). Thus, it follows that we can efficiently obtain an upper bound on the magnitude of $G(\mathbf{N})$.

REFERENCES

1. *F. Baskett, K. Chandy, R. Muntz, and F. Palacios*, Open, closed, and mixed networks of queues with different classes of customers. //Journal of the ACM 1975. V. 22.
2. *A. Bertozzi and J. McKenna*, Multidimensional residues, generating functions, and their application to queueing networks //SIAM Rev., 35 (1993), 1993. V. 22
3. *J. P. Buzen*, Computational algorithms for closed queueing networks with exponential servers //Communications of the ACM 1973
4. *G. Casale*, An efficient algorithm for the exact analysis of multiclass queueing networks with large population sizes, //Proc. of the ACM SIGMETRICS 2006.
5. *W. J. Gordon and G. Newell*, Closed queueing systems with exponential servers, //Operation Research 1967
6. *E. Koenigsberg*, Cyclic queues, //Operations Research Quarterly, 1958 V. 9
7. *J. McKenna, D. Mitra, and K. Ramakrishnan*, A class of closed markovian queueing networks: Integral representations, asymptotic expansions, generalizations., //Bell Syst. Tech. J. 1981
8. *M. Reiser*, Mean-value analysis and convolution method for queue-dependent servers in closed queueing networks, //Performance Evaluation, 1981
9. *M. Reiser and H. Kobayashi*, Queueing networks with multiple closed chains: Theory and computational algorithms, // IBM J. Res. Dev. 1975.