

A POPULATION-MIX DRIVEN APPROXIMATION FOR QUEUEING NETWORKS WITH FINITE CAPACITY REGIONS

J. Anselmi¹, G. Casale², P. Cremonesi¹

¹ Politecnico di Milano, Via Ponzio 34/5, I-20133 Milan, Italy

² Neptuny R&D, Via Durando 10-G, I-20158, Milan, Italy

{jonatha.anselmi, giuliano.casale, paolo.cremonesi}@polimi.it

In this paper we propose a novel approximate method for closed multiclass queueing networks with finite capacity regions and shared constraints. The approach is based on Norton’s theorem for queueing networks, in which each region is replaced by a Flow Equivalent Service Center (FESC). We propose a population-mix driven definition of FESCs service rates which provides improved accuracy with respect to existing methods. We solve the resulting non-product-form network with a new approximate variant of the convolution algorithm proposed in the paper. Comparison with simulation shows that the algorithm provides good accuracy.

Keywords: multiclass queueing network, finite capacity region, flow-equivalent server

1. INTRODUCTION

Among existing modeling techniques, queueing networks with “finite capacity regions” have largely proven to be effective in characterizing simultaneous resource possession in which a request can hold more resources simultaneously. Such queueing networks impose upper bounds on the number of jobs that can simultaneously reside in a set of stations [3, 6], and can be used to model application constraints. For instance, consider a web-based multi-level application having a pool of threads to process HTTP requests. The pool size is chosen in order to have a compromise between response time and number of rejected requests. When a new request arrives and finds no free thread, it waits in the web server queue. The web-based application can be described with a set of resources modeling the architecture (e.g., web, application and database servers) with a finite capacity region bounded by the size of the HTTP threads pool.

Inside the network, jobs are either in a thinking state, i.e., waiting at the terminals, or ready, i.e., visiting queues. Ready jobs can be either in a buffer queue waiting to enter into a finite capacity region according to a FCFS access rule, or active, i.e., circulating inside the region and competing for the resources. Finite capacity regions are characterized by two types of constraints: a dedicated constraint bounds the number of jobs in a region for a specified class; a shared constraint limits the number of jobs without class distinctions. Analytical methods for multiclass models with finite capacity

regions are mostly based on Norton's theorem for queueing networks [4] in which a region is replaced by a load-dependent station called FESC [6, 7].

Krzesinski et al. [5] propose an algorithm based on the homogeneity assumptions:

- 1) for each feasible network state, the class- r departure rate from region \mathcal{R} does not depend on the current population of the classes $s \neq r$ in \mathcal{R} , but instead depends only on their average population in \mathcal{R} over all states;
- 2) if a center does not contain class- r jobs, then an arriving class- r job immediately receives service regardless of the jobs of class $s \neq r$ in that station.

With these assumptions, the computational complexity needed to solve multiclass models can be dramatically reduced. Making an intensive use of the multiclass Linearizer technique (see, e.g., [3]), the algorithm groups together stations belonging to the same region and iteratively estimates the average number of jobs inside the regions.

A different approximate technique, proposed by Sauer in [7], defines FESC service rates that account for the finite capacity constraints. However, the way in which service rates are defined does not satisfy product-form assumptions, therefore the solution is computed analyzing the underlying Markov chain. Hence, the related computational complexity makes unfeasible the analysis of networks with multiple regions.

The contribution of the approximate method we propose is that it handles multiple finite capacity regions. We do not consider homogeneity assumptions, and rather draw concepts from the asymptotic theory of [1] for multiclass systems, replacing each region by a FESC with non-product-form service rates. Moreover, we solve the resulting network with an approximate variant of the well-known convolution algorithm instead of solving the underlying Markov chain which makes intractable the analysis.

The paper is organized as follows. In Section 2 we give notation and illustrate the approximate solution technique. In Section 3 we discuss characteristics of product-form service rates and propose the variant of the convolution algorithm which is able to handle networks with non-product-form service rates. Experimental results are further discussed in Section 4. Finally, Section 5 draws conclusions and outlines future work.

2. THE PROPOSED ALGORITHM

In this section we propose a new approximate method for multiclass closed queueing network models with finite capacity regions. We assume that if finite capacity constraints are removed, the model has product-form solution [3]. We also assume that queues have a single server and that the model includes only shared constraints.

2.1. Notation. Consider a multiclass closed queueing network model with R job classes, K service centers and M finite capacity regions. Each region m contains two or more service centers, and has a shared capacity constraint B_m . In other words, the total number of ready jobs in m never exceeds B_m . Each center belongs to a single region. If not otherwise stated, index k will range from 1 to K , r from 1 to R , and m from 1 to M .

We denote class- r population by N_r , the population vector by $\vec{N} = (N_1, \dots, N_R)$, and the total number of jobs by $N = \sum_r N_r$. We denote by $\vec{n} = (n_1, \dots, n_R)$ any population vector such that $\vec{0} < \vec{n} \leq \vec{N}$, and $n = \sum_r n_r$.

Other parameters are the mean number of class- r visits at center k , $V_{k,r}$, the mean class- r service time at center k , $S_{k,r}$, the mean class- r service rate when current population at center k is \vec{n} , $\mu_{k,r} \equiv \mu_{k,r}(\vec{n})$, and the mean class- r average think time, Z_r .

The performance measures of interest are the per-class throughputs, $X_r(\vec{n})$, and the aggregate throughput, $X(\vec{n}) = \sum_r X_r(\vec{n})$.

2.2. The Algorithm. We consider a heuristic that replaces service centers of finite capacity regions by FESCs. The general structure is shown in Algorithm 1.

Algorithm 1 General scheme of the approximate algorithm.

- 1: for each region m do
 - 2: Create a FESC analyzing m in isolation:
 - 3: a) Define service rates $\tilde{\mu}_{m,r}(\vec{n})$ for populations \vec{n} such that $n \leq B_m$
 - b) Approximate service rates $\tilde{\mu}_{m,r}(\vec{n})$ for populations \vec{n} such that $n > B_m$
 - 4: end for
 - 5: Replace all finite capacity regions with the FESCs
 - 6: Solve the resulting unconstrained network with an approximation of the load-dependent convolution algorithm.
-

From line 1 through 4, we replace each region with a suitable FESC: since the way in which we define FESCs service rates is the main innovation of our approach, a complete description is given below. In line 5, we build a non-product-form model which is solved (line 6) by an approximate variant of the convolution algorithm proposed later in Section 3.

2.3. FESC Construction. The main issue we address in this section is the definition of FESCs which replace regions and take into account finite capacity effects. To model such effects in single class networks, good results have been obtained using service rate saturation [6] which states that throughput cannot increase when the capacity constraint is active. In multiclass environments the approximation is more difficult to obtain. Since jobs of different classes have different service demands, the departure rate from a region depends on the mix of requests in m . As shown in [1], the population mix vector $\vec{\beta} = (\beta_1, \dots, \beta_R)$ where $\beta_r = \beta_r(\vec{N}) = N_r/N$ strongly affects the performance of a queueing network model. Indeed, whenever the number of jobs in a multiclass region overflows the capacity constraint, it is not sufficient to saturate the FESC service rates to a constant value as in single class models since the departure rate strongly depends on the population mix $\vec{\beta}_m(\vec{n}) = (\beta_{m1}(\vec{n}), \dots, \beta_{mR}(\vec{n}))$ of active jobs inside region m . Let \vec{n} be the vector of both waiting and active jobs of region m . We approximate the class- r mix of active jobs in region m as $\beta_{m,r}(\vec{n}) = n_r/n$. This is in general an approximation, since it considers also waiting jobs. Hence, FESCs service rates are defined as

$$\tilde{\mu}_{m,r}(\vec{n}) = \begin{cases} \tilde{X}_{m,r}(\vec{n}) & \{\vec{n} \mid n \leq B_m\} \\ \tilde{X}_{m,r}(\vec{s}_m) & \text{otherwise} \end{cases} \quad (1)$$

where $\tilde{X}_{m,r}$ denotes the class- r throughput of region m in isolation and $\vec{s}_m = (s_{m,1}, \dots, s_{m,R})$ is a population vector related to region m such that $s_{m,r} = s_{m,r}(\vec{n}) = \beta_r(\vec{n})B_m$. The exact

value of $\tilde{X}_{m,r}(\vec{n})$, for all population vectors $\{\vec{n} \mid n \leq B_m\}$, is computed by a single run of the MVA [3]. The term $\tilde{X}_{m,r}(\vec{s}_m)$ is the approximation which models finite capacity effects and encodes the fact that no more than B_m jobs can exist inside region m . Notice that $\sum_r s_{m,r} = B_m$. Since in general vector \vec{s}_m has non-integral elements, the computation of $\tilde{X}_{m,r}(\vec{s}_m)$, for each possible \vec{s}_m , requires non-trivial effort. Thus, for population vectors $\{\vec{n} \mid n > B_m\}$ we redefine $\tilde{\mu}_{m,r}(\vec{n})$ as

$$\tilde{\mu}_{m,r}(\vec{n}) = \begin{cases} \frac{\beta_r(\vec{n})B_m}{[s_{m,r}]} \tilde{X}_{m,r}([s_{m,1}], \dots, [s_{m,R}]) & [s_{m,r}] \neq 0 \\ \beta_r(\vec{n})B_m \tilde{X}_{m,r}([s_{m,1}], \dots, [s_{m,r-1}], 1, [s_{m,r+1}], \dots, [s_{m,R}]) & [s_{m,r}] = 0 \end{cases} \quad (2)$$

where $[\cdot]$ returns the nearest integer to \cdot , and $\frac{\beta_r(\vec{n})B_m}{[s_{m,r}]}$ and $\beta_r(\vec{n})B_m$ are scaling factors which represent the interpolation we use to approximate the actual throughput value $\tilde{X}_{m,r}(\vec{s}_m)$. With (2), we define FESC m service rates by simply reusing throughputs computed with the MVA. Notice that to improve the algorithm performance the approximation (2) can be performed on-line directly in the convolution algorithm.

3. PRODUCT-FORM APPROXIMATION OF SERVICE RATES

For product-form queueing networks the following steady state probability distribution holds (see, e.g., [3]): $P(\vec{n}_1, \dots, \vec{n}_K) = G^{-1} \prod_k F_k(\vec{n}_k)$, where G is the normalizing constant. From the discussions in [8], in our case it follows that for class-specific queue-dependent service rates $\mu_{k,r}(\vec{n})$, the following equation holds

$$F_k(\vec{n}) = \frac{1}{\mu_{k,r}(\vec{n})} F_k(\vec{n} - \vec{e}_r), \quad \vec{n} \geq \vec{e}_r \quad (3)$$

where r may be any of the job classes, $F_k(\vec{0}) = 1$ and \vec{e}_r is the unit vector in direction r . Let us define a path as a sequence of population vectors which can be spanned in the computation of $F_k(\vec{n})$ with (3). Since r is an arbitrary class, we can perform a recursion along a generic path from \vec{n} to $\vec{0}$ obtaining the same value of $F_k(\vec{n})$. Whenever (3) does not well define $F_k(\vec{n})$, the rate function $\mu_{k,r}(\vec{n})$ does not satisfy product form [8]. $F_k(\vec{n})$ is well defined if and only if the following condition holds true,

$$\mu_{k,i}(\vec{n})\mu_{k,j}(\vec{n} - \vec{e}_i) = \mu_{k,j}(\vec{n})\mu_{k,i}(\vec{n} - \vec{e}_j), \quad i \neq j, \quad \vec{n} - \vec{e}_i > \vec{e}_j, \quad \vec{n} - \vec{e}_j > \vec{e}_i \quad (4)$$

For stations with FCFS scheduling, (4) is a necessary condition to guarantee product-form. In our work, $\tilde{F}_m(\vec{n})$ is related to the FESC belonging to the resulting network representing region m . For population vectors $\{\vec{n} \mid n \leq B_m\}$, functions $\tilde{F}_m(\vec{n})$ are well defined since the way in which we define service rates guarantees product-form [4]. However, for population vectors $\{\vec{n} \mid n > B_m\}$, in general $\tilde{F}_m(\vec{n})$ is not well defined. This means that during the recursion the value of $\tilde{F}_m(\vec{n})$ depends, for such \vec{n} , on the particular path chosen from \vec{n} to $\vec{0}$. This can be easily verified considering a two-class, two-station model with population vector $\vec{n} = (2, 1)$, capacity $B = 2$ and loadings $D_{1,1} = 0.5$, $D_{2,1} = 2$, $D_{1,2} = D_{2,2} = 1$, where $D_{m,r} = V_{m,r}/\mu_{m,r}(\vec{e}_r)$.

In order to gain a path-independent value for $\tilde{F}_m(\vec{n})$ we compute the value of $\tilde{F}_m(\vec{n})$ by taking the average of all border paths of the state diagram, i.e. $\{\vec{n} - \vec{e}_{r_1}, \dots, \vec{n} - N_{r_1}\vec{e}_{r_1}, \vec{n} - N_{r_1}\vec{e}_{r_1} - \vec{e}_{r_2}, \dots, \vec{n} - N_{r_1}\vec{e}_{r_1} - N_{r_2}\vec{e}_{r_2}, \dots, e_{r_R}, \vec{0}\}$. Denote by $p \equiv p(\vec{n})$ the number of non-empty classes of vector \vec{n} . The set of all border paths has cardinality equal to $p!$. However, since in general queueing networks with finite capacity region with more than four classes are intractable and in practice are rarely used, cardinality of $\Phi(\vec{n})$, for $R \leq 4$, is small and does not affect algorithm performance. Our approximate variant of the convolution algorithm computes $\tilde{F}_m(\vec{n})$ for population vectors $\{\vec{n} \mid n > B_m\}$ as

$$\tilde{F}_m(\vec{n}) = \frac{1}{p!} \sum_{\varphi \in \Phi(\vec{n})} \tilde{F}_m^{\varphi}(\vec{n}), \quad \Phi(\vec{n}) = \{\varphi \mid \varphi \text{ is a border path from } \vec{n} \text{ to } \vec{0}\} \quad (5)$$

Using (5), we approximate the solution of the resulting network without solving the Markov chain which is needed to get exact solutions but makes the problem intractable.

4. EXPERIMENTAL RESULTS

In this section we show experimental results for the proposed approximation comparing our results with simulation results which have been obtained using the JMT simulation engine [2] choosing 99% level confidence intervals. We evaluate accuracy by measuring the percentage relative error $|X^S - X^A|/X^S \cdot 100\%$, where X^A and X^S are the overall throughputs obtained respectively with our approximation and simulation.

We validate the approximation on a class of models with exactly ten load-independent queues and a single delay. The delay does not belong to any finite capacity region. The number of regions M is randomly drawn from $\{1, 2, 3, 4\}$. Queues service rates $\mu_{k,r}(\vec{e}_r)$ range in $[0.01, 12]$ jobs per unit of time and visits $V_{k,r}$ in $[0.1, 5]$. Loadings $D_{k,r}$ are computed as $V_{k,r}/\mu_{k,r}(\vec{e}_r)$. Average think times Z_r range in $[0.01, 10]$. In order to limit the cost of simulations we decrease population sizes and proportionally scale capacity ranges when increasing the number of classes. For models with 2 (respectively 3 and 4) classes, we consider a maximum of 200 (100 and 60) jobs, a maximum of 100 (33 and 15) per-class jobs and capacities never exceed 40 (20 and 12) units. We generated 400 random models ensuring that $N > \max_{1 \leq m \leq M} B_m$. On an Intel Xeon 2.80GHz processor with hyperthreading technology, our method requires an average of 10 seconds while simulation takes more than 6 minutes to converge. The average percentage relative error computed is around 11%. We also compare our method against models published in the literature. Models parameters and related simulation results are discussed in [5]. In Tables 1(a) and 1(b) we show response times of our method (APP), simulation (SIM) and the method proposed in [5] (KT). Consider for instance the cases in which the capacity constraint strongly limits population inside the regions, i.e., when the approximation is highly stressed: we improve significantly the accuracy in all cases, except for case $\vec{N} = (20, 2)$, $B_1 = B_2 = 2$ where we are slightly less precise than KT. Consider also the case of a single region with $\vec{N} = (20, 2)$ and $B_1 = 2$: with respect to the average value of confidence interval, for class-1 response time we provide an error less than 3% instead of 28%. Also for the other cases, we conclude that in general our approximation returns more accurate results.

$N_1N_2B_1$	Class 1 Resp. time			Class 2 Resp. time		
	KT	SIM	APP	KT	SIM	APP
6	.73	(.71, .74)	0.73	4.71	(4.20, 5.02)	4.73
20 2 4	.83	(.72, .74)	0.75	4.65	(4.19, 4.97)	4.76
2	1.40	(.99,1.03)	1.04	4.44	(4.60, 5.55)	5.43
10	1.01	(.98,1.01)	1.01	7.34	(6.76, 8.03)	7.59
30 3 6	1.13	(1.01,1.04)	1.04	7.11	(6.17, 8.07)	7.63
2	2.50	(2.06,2.16)	2.32	6.09	(7.63, 9.33)	9.70
18	1.48	(1.40,1.44)	1.48	12.95	(10.54,12.46)	13.38
40 4 12	1.49	(1.46,1.52)	1.49	12.90	(12.80,15.26)	13.38
6	1.80	(1.57,1.62)	1.60	11.75	(11.68,13.91)	13.26

$N_1N_2B_1B_2$	Class 1 Resp. time			Class 2 Resp. time		
	KT	SIM	APP	KT	SIM	APP
3 3	.79	(.73, .75)	.74	4.51	(4.28,5.05)	4.75
20 2 2 2	.92	(.78, .80)	.78	4.64	(4.94,5.81)	4.88
1 1	.98	(.97,1.01)	1.06	4.10	(4.76,5.68)	5.94
5 5	1.04	(.97,1.00)	1.02	7.26	(6.84,8.06)	7.60
30 3 3 3	1.20	(1.06,1.09)	1.07	7.05	(6.23,7.43)	7.69
1 1	1.68	(1.79,1.85)	1.90	6.44	(7.87,9.67)	9.92
9 9	1.48	(1.46,1.50)	1.49	12.60	(12.82,15.15)	13.38
40 4 6 6	1.55	(1.42,1.47)	1.50	12.91	(11.53,13.61)	13.37
3 3	1.83	(1.59,1.64)	1.62	12.95	(11.19,13.48)	13.28

(a) First case: a single finite capacity region

(b) Second model: two finite capacity regions

Table 1: Comparison with published models

5. CONCLUSIONS

In this paper we propose an approximate method for multiclass closed queueing networks with finite capacity regions based on the construction of FESCs which replace regions and model finite capacity effects. We propose to solve the resulting non-product-form network with a variant of the convolution algorithm. A comparison with existing techniques on published models improves the accuracy and a comparison with simulation shows a small error. We leave as future work the extension to overlapped regions and open models in order to integrate the approximation in admission control schemas.

REFERENCES

1. G. Balbo and G. Serazzi, Asymptotic analysis of multiclass closed queueing networks: Multiple bottlenecks, *Perf. Eval.*, 30 (1997), pp. 115–152.
2. M. Bertoli, G. Casale, and G. Serazzi, Java modelling tools: an open source suite for queueing network modelling and workload analysis, in *Proc. of QEST 2006*, Riverside, US, Sep 2006, IEEE Press, pp. 119–120.
3. G. Bolch, S. Greiner, H. de Meer, and K. S. Trivedi, *Queueing Networks and Markov Chains*, Wiley-Interscience, 2005.
4. P. Kritzinger, S. V. Wyk, and A. Krzesinski, A generalization of norton’s theorem for multiclass queueing networks, *Perf. Eval.*, 2 (1982), pp. 98–107.
5. A. E. Krzesinski and P. Teunissen, Multiclass queueing networks with population constrained subnetworks., in *SIGMETRICS*, 1985, pp. 128–139.
6. E. D. Lazowska, J. Zahorjan, G. S. Graham, and K. C. Sevcik, *Quantitative System Performance*, Prentice-Hall, Englewood Cliffs, New Jersey, 1984.
7. C. H. Sauer, Approximate solution of queueing networks with simultaneous resource possession, *IBM Journal of Research and Development*, 25 (1981), pp. 894–903.

8. C. H. Sauer, Computational algorithms for state-dependent queueing networks, *ACM Trans. Comput. Syst.*, 1 (1983), pp. 67–92.