

June, 1977

72531
BNL-22915
AMD 758R1

77; 116800

Throughput in Locally Balanced Computer System Models*

Andrew S. Noetzel
Applied Mathematics Department
Brookhaven National Laboratory
Upton, NY 11973

The optimization of throughput in locally balanced queueing network models is investigated. A general result, useful in the design of computer system models, shows that throughput is a nondecreasing function of the number of customers contained in any subnetwork. Then processor allocation algorithms that maximize throughput are shown for the case where processing power can be switched between queues, as when several queues are served at a single multiprocessor system. The maximization of throughput is shown first in the case that processing power allocations to a queue depend on the queue state only, and then, in an extension of known locally balanced queue, the case in which processing power is allocated on the basis of an entire subnetwork state. The latter case provides a simple and optimum rule for processor allocations that maximize throughput in networks containing multiprocessor systems.

1. INTRODUCTION

The complete solution of queueing networks of arbitrary configuration is possible only for the class of networks that have been variously described as separable, locally balanced or admitting a product form solution [1,3,8]. Recently, this class of networks has been shown to include queues with general service time distributions under several different scheduling disciplines, and multiple classes of customers. As a result, the local balance model has found wide utility in the analysis of computer systems and networks [1,3]. It has been implemented as the basis of several interactive systems that yield quick analyses of computer systems and configurations [7,12]. Because of its unique tractability in general network configurations, this network model has also been studied as an approximation to networks whose queues do not meet the local balance requirement [3].

Some studies of locally balanced queueing models have assumed a limited variation of the processing rate of the queue with the number of customers contained therein. For example, in Gordon and Newell's classical work, a queue was considered to have a fixed number R of processors [6]. If the mean processing time of a single customer is $1/\mu$, then when the queue contains n customers, the mean processing rate of the queue is $n\mu$ for $n \leq R$, and $R\mu$ for $n > R$. However, it is not difficult to show that the local balance property will be retained if the number of processors in use at a queue is any general function $r(n)$ of the number n of customers at the queue. We investigate the implications of this and a more general form of processing power

*Work performed under the auspices of the ERDA.

NOTICE
This report was prepared as an account of work sponsored by the United States Government. Neither the United States nor the United States Energy Research and Development Administration, nor any of their employees, nor any of their contractors, subcontractors, or their employees, makes any warranty, express or implied, or assumes any legal liability or responsibility for the accuracy, completeness or usefulness of any information, apparatus, product or process disclosed, or represents that its use would not infringe privately owned rights.

MASTER

DISTRIBUTION OF THIS DOCUMENT IS UNLIMITED

Key

DISCLAIMER

This report was prepared as an account of work sponsored by an agency of the United States Government. Neither the United States Government nor any agency Thereof, nor any of their employees, makes any warranty, express or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States Government or any agency thereof. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States Government or any agency thereof.

DISCLAIMER

Portions of this document may be illegible in electronic image products. Images are produced from the best available original document.

allocation for the optimization of throughput.

Throughput is a central consideration in network models of computer systems. In open network models, it is fully specified by the assumed rates of the external sources. In closed networks, however, it is determined by the combined effects of congestion at the various queues. The question of processing power allocation to maximize throughput is important in the realistic case when processing power is limited, but can be traded off between queues.

Consider, for example, the several possible interpretations of the subnetwork of a closed queueing system, as shown in Figure 1. The queues might represent the sequential phases or steps of programs, with the parallel paths reflecting the differing requirements of identified special classes of programs, and the class of general user programs. The processing units would then be the individual computers at a computer center.

Alternatively, the queues may represent the phases of programs, either system functions or user-written routines, loaded into the executable memory of a multiprocessor system. Under either interpretation, the customers at the queues may be served by any of the processors. If time-sharing with a small time slice is used, it will be reasonable to speak of allocating a fractional number of processors to a queue. Throughout, we assume that the processing rate of a customer is linearly related to the processing power allocated him; the processing rate u on a single processor becomes rate ru when processing power r is allocated. The physical interpretations impose only the following constraints on the queueing model: since not more than n processors can be useful in serving a queue that contains n customers, $r(n) \leq n$; and not more than the fixed number R of processors may be in use by the subnetwork at any time.

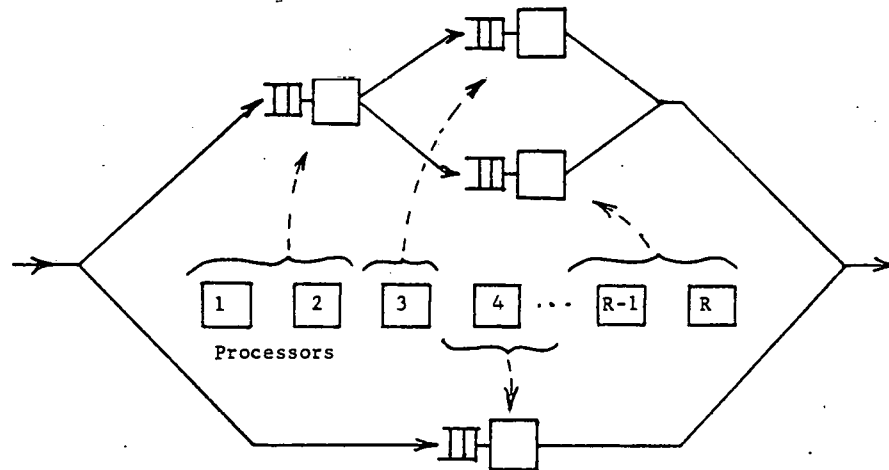


Figure 1: Processor Assignments in Subnetwork with Processing Power Tradeoffs

In order to establish the policies for maximization of throughput via processing-power tradeoffs, several general properties of throughput in locally balanced queueing networks must first be established. These general properties are discussed in Section 3.

The question of processing power allocation is examined in two fundamental cases. First, using the known model of locally balanced queues, queue-state dependent processing rates are considered. This is the case in which processing power $r_i(n_i)$ is allocated to queue i when it contains n_i customers. The processor allocation decisions can be made locally to the queue by the software processor in control, without considering the distribution $\{n_j: j \neq i\}$, of customers at the other queues of the subnetwork. It will be seen in Section 4 that optimization strategies under this form of processor allocation are severely constrained.

Then the case of subnetwork-state dependent processing rates is considered. The processing power allocations at queue i of the subnetwork take the form $r_i(n_1 \dots n_j \dots n_M)$, when there are n_j customers at queue j , for each of the M queues of the subnetwork. The strategy for processor allocation in this case is global to the subnetwork. Processing power allocations of this form require an extension of the known local balance queue model. The state probability solutions for cases of n -queue parallel and two-queue series subnetworks with this form of processing power allocation are shown in Section 5. The policy for allocating processing power for maximum throughput within the local balance constraint is demonstrated.

2. LOCALLY BALANCED NETWORKS AND NORTON'S THEOREM

The separable or locally balanced networks to be considered have fixed topology. After leaving queue i , a customer will go to queue j with fixed probability p_{ij} . Let P be the matrix of transition probabilities p_{ij} , $1 \leq i, j \leq M$, for a network of M queues. Let λ_i be the mean flow rate into queue i . If $L = \begin{matrix} \lambda_1 & \lambda_2 & \dots & \lambda_M \end{matrix}$ is a vector of relative flow rates or throughputs for the network, then $LP=L$. If the network is open, L is determined by the absolute input rate to the network. For closed networks, the above relation determines L to within a constant. Then, if the network contains N customers, the actual throughput at queue i is given by $\tau_i(N) = \lambda_i \frac{G(N-1)}{G(N)}$, where $G(n)$, $n = 1, \dots, N$, is the normalization factor computed

when the network contains n customers. $G(n)$ is computed by the convolution algorithm using the flow vector L [2,11,13]. Since the throughputs at the branches of the network remain fixed relative to each other as N and the processing rates at the queues are varied, the network throughput, or total processing rate, is optimized as throughput in any branch of the network is optimized.

The throughput characteristics of two queue networks are of considerable importance because of the reduction of arbitrary closed networks to an equivalent two queue network made possible by Norton's theorem for computer networks [4]. The behavior of a particular queue Q within a closed network η is the same as that of Q in a network with one other queue Q_e , which serves as the equivalent of all the queues of η except for Q (see Figure 2). The processing rate $U(n)$ of the equivalent queue, when it contains n customers, is determined by constructing network η' , which is the same as η except for a short circuit in place of Q . $U(n)$ is simply the mean throughput in the link replacing Q , when η' contains n customers. To study via the equivalent network, the effects of Q in η when η contains N customers, the rates $U(1), \dots, U(N)$ must be measured in the link in η' , when η' contains $1, \dots, N$ customers.

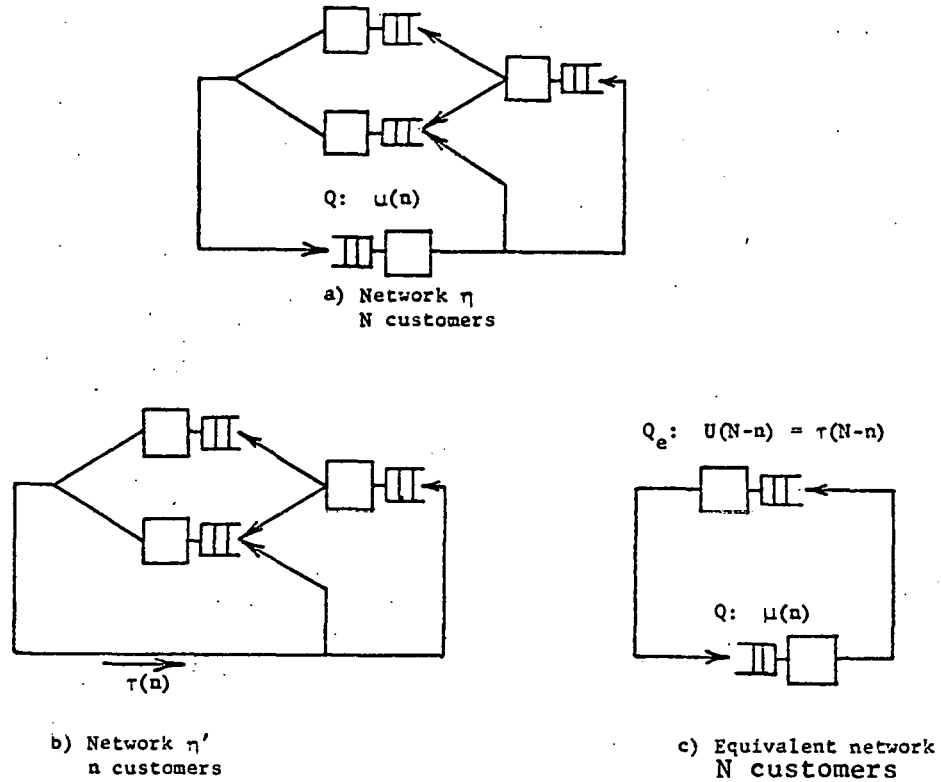


Figure 2: Nortons Equivalent Queue

It is possible to analyze in the same manner as Q the behavior of any subnetwork σ that can be isolated by a single pair of input and output terminals. Network η may be therefore reduced to a two queue network; an equivalent queue for σ , and an equivalent for the complement of σ in η . Study of throughput in the two queue network is straightforward, since the normalization constant $G(n)$ is expressed as a single convolution.

3. THROUGHPUT AS A FUNCTION OF LOAD

An important characteristic of separable queueing networks that can be determined with the aid of Norton's Theorem is that the throughput, or output rate, of any subnetwork is a nondecreasing function of the load, or the number of customers in the subnetwork. It is first shown that this holds for the equivalent two queue network.

Theorem 1. Let $u(m)$ and $U(n)$ be the processing rates of the queues in a closed locally balanced two queue network, when the queues contain m and n customers, respectively. Let $\tau(N)$ be the mean throughput of the closed network when it contains N customers. Then if $u(n+1) \geq u(n)$ and $U(n+1) \geq U(n)$ for all $0 < n \leq N$, then $\tau(N+1) \geq \tau(N)$.

The proof follows from the form of $\tau(N) = \frac{G(N-1)}{G(N)}$. The details are in Appendix A.

Through the use of Norton's theorem, the characteristic of nondecreasing throughput

with increasing load may now be proven for all locally balanced networks.

Theorem 2. Let $\tau(n)$ be the mean throughput in some branch of a closed locally balanced network that contains n customers. If the processing rate of each queue in the network is a nondecreasing function of the number of customers at the queue, then for all $n > 0$, $\tau(n+1) \geq \tau(n)$.

Proof. The proof is by induction on M , the number of queues in the network. Let $\tau_M(n)$ be the mean throughput in some branch b of a network of M queues that contains n customers. Let $u(n)$ be the processing rate of an arbitrary queue when it contains n customers, and let $u(n+1) \geq u(n)$. For $M = 1$, the network has only one queue, but may have a number of probabilistically selected branches from the queue back to itself. Let p be the probability that a customer takes branch b in returning to the queue. Then $\tau_1(n+1) = pu(n+1) \geq pu(n) = \tau_1(n)$.

Let $U(n)$ be the processing rate of the equivalent queue for an $M - 1$ queue network with respect to branch b . Then by Norton's Theorem and the inductive hypothesis, $U(n+1) = \tau_{M-1}(n+1) \geq \tau_{M-1}(n) = U(n)$. And if the equivalent queue is placed in series with a queue with processing rate $u(n+1) \geq u(n)$, $1 \leq n \leq N$, and there are N customers in this network, then the throughput of the closed two queue network is $\tau_M(N)$. But then $\tau_M(N+1) \geq \tau_M(N)$ by Theorem 1. This completes the proof.

Queues that have a fixed number of processors available meet the conditions of this theorem. But $u(n+1) \geq u(n)$ may not hold if processing power is traded off between queues.

4. THROUGHPUT WITH QUEUE-STATE DEPENDENT PROCESSING RATES

A queue may have the local balance property if each customer has an exponential processing time, or else if the queueing discipline is processor sharing or preemptive last-come-first-served. The total processing rate of the queue may be any function of the number of customers in the queue. We assume that $1/u$ is the mean processing time of one customer on one processor, and that the queue state is simply the number n of customers at the queue. If the number of processors assigned the queue is an arbitrary function $r(n)$ of the queue state, then $u(n) = r(n)u$. We will call processing rates of this form queue-state dependent.

There are several instances in which increased throughput can be obtained by dynamically switching processors from one queue to another. For example, two queues may represent two different processes or programs to be executed in one multiprocessor system. Then any or all of the processors may be available for either queue when they are not required for the other. However, with the developed model of a locally balanced queue (that is, with queue-state dependent processing rates), the analysis of processing power tradeoffs is severely constrained. When the network state changes by means of a customer moving from one queue to another, processing power exchanges may take place only between the two queues involved in the transition.

Throughput with processing power tradeoffs between a single queue and the remainder of the network can be analyzed in the network consisting of the selected queue and the Norton's equivalent queue. But first, the relationship of throughput to processing rates in a two queue network will be established when the rates at each of the queues are independent of each other. This is expressed in the following theorem.

Theorem 3. Let $u(m)$ and $U(n)$ be the processing rates at the queues of a locally balanced two queue network, when the queues contain m and n customers, respectively, and suppose all the $U(n)$ and the $u(m)$ are mutually independent for $0 < m, n \leq N$. Let $\tau(N)$ be the throughput when the network contains N customers. Then

- a) $\tau(N)$, as a function of $U(n)$, $0 < n \leq N$, has no extrema.
- b) $\tau(N)$ is a nondecreasing function of $U(n)$ if $\mu(i) \geq \mu(j)$ for all $j \leq N - n$ and $N - n < i \leq N$.
- c) If $\tau(N)$ is a nonincreasing function of $U(n)$, then it is a strictly increasing function of $\mu(N-n)$.

The theorem is proved by differentiating $\tau(N) = \frac{G(N-1)}{G(N)}$ with respect to $U(n)$, for $1 \leq n \leq N$. The details are in Appendix B.

As a simple example, consider optimizing a two-queue network, containing three customers in an environment where processing power tradeoffs between the queues might be possible. By Theorem 3b, both $\mu(3)$ and $U(3)$ can be maximized without decreasing throughput $\tau(3)$. Suppose $\mu(1) \leq \mu(2)$, $\mu(3)$. Then by 3b, $U(2)$ can be maximized without decreasing $\tau(3)$. Likewise, $U(1)$ can be maximized if $\mu(1)$, $\mu(2) \leq \mu(3)$. But suppose the condition $\mu(2) > \mu(3)$ holds. Then, it is possible that $\tau(3)$ will decrease as $U(1)$ increases. By 3c, $\tau(3)$ will then be a strictly increasing function of $\mu(2)$. Although independence of the rates was assumed in the theorem, switching processing power to minimize $U(1)$ and maximize $\mu(2)$ causes no conflict, and will be the optimal solution. But then, if the processing rate U varies linearly with the processing power, as in the assumption of queue-state dependency, $U(1)$ will be reduced to zero. The network will effectively contain two customers; which leads to a refutation of the $\mu(2) > \mu(3)$ assumption. Similar reasoning is employed in proof of the following theorems.

The optimum allocation of $R(N)$ available processors to a network that contains N customers will now be considered. Suppose $r_i(n_i)$ processors are allocated to queue i ; when it contains n_i customers. A possible allocation strategy may be to hold some processing power in reserve, that is, $\sum r_i(n_i) < R(N)$ for some network state $(n_1 \dots n_M)$; the reason for doing this may be to ensure that more processing power is available for $r_i(n_i+1)$, or even $r_i(n_i-1)$ if throughput might be more sensitive to those rates. But it can be shown that throughput cannot be optimized by such strategies.

Theorem 4. Let $R(N)$ be the processing power available to be allocated to the M queues of a closed locally balanced network that contains $N > 0$ customers. Let $r_i(n)$ be the processing power allocated queue i , when there are n customers at queue i , for $1 \leq i \leq M$. Let $\tau_M(N)$ be the mean throughput at some branch of the network.

- a) Then for maximum $\tau_M(N)$, the available processing power must always

$$\text{be fully utilized; that is, } \sum_{i=1}^M r_i(n_i) = R(N) \text{ for all } \sum_{i=1}^M n_i = N,$$

$$\text{and } r_i(0) = 0 \text{ for } 1 \leq i \leq M.$$

- b) And any processing power distribution meeting the above constraints provides the maximum throughput, and when $\tau_M(N)$ is maximum,

$$\tau_M(N) = kR(N), \text{ where } k \text{ is a constant.}$$

The proof is again by induction on the number of queues in the network. The M queue network is reduced to a two queue network by Norton's Theorem. Then

Theorem 3 is used to show that for maximum throughput, each processing rate of the two queue network must be maximized. The details of the proof are in Appendix C.

Therefore, for maximum throughput with queue-state dependent rate assignments, the processing power added to one queue must be exactly that taken from another whenever there is a transition between the queues. But for closed networks with more than two queues, it can be seen that this determines the processing power uniquely.

Theorem 5. Consider any network of $M > 2$ locally balanced queues, containing N customers, in which any fraction of the processing power of $R(N)$ processors may be assigned to any queue. If $r_i(n)$ is the processing power assigned queue i when it contains n customers, for $1 \leq i \leq M$, then throughput is maximized by the assignment $r_i(n) = \frac{n}{N}R(N)$.

Proof. From Theorem 4, for maximum throughput, one must consider processing power allocations such that

$$\sum_{i=1}^M r_i(n_i) = R(N) \text{ for each network state } (n_1 \dots n_M). \text{ When there are } N - k$$

customers at queue one, the processing rate at queue one remains constant and independent of the distribution of the remaining k customers. Therefore, for $1 < i \leq M$, $r_i(k) = R(N) - r_1(N-k)$. Similarly, when there are $N - k$ customers at queue two, for $2 < i \leq M$, $r_i(k) = R(N) - r_2(N-k)$. Hence, $r_i(k) = r_j(k)$ for $1 \leq i, j \leq M$, as long as $M > 2$. Consider the case $k = 2$. With $N - 2$ customers at queue one, the remaining two customers may both be at queue i , or may be at queues i and j , $1 < i < j \leq M$, while the processing rate at queue one remains constant. Hence, $r_1(2) = r_i(1) + r_j(1) = 2r_1(1)$. Similarly, considering $k = 3, 4, \dots, N$ it is seen that $r_i(k) = kr_1(1)$. And, since $r_i(N) = Nr_1(1)$, $r_i(k) = \frac{k}{N}r_i(N)$ for $1 \leq i \leq M$ and $1 \leq k \leq N$. The proof is complete.

It is seen that there is a simple rule for optimally allocating processing power if queue-state dependencies are assumed. But, it applies only if all processors may be freely switched to any queue in the network. This limits the applicability of the model to networks of logical processes entirely within a single multiprocessor system, in which any processor may execute any logical function. But in such case there is little scheduling difficulty. In a model of a computer network, the capability of processing power tradeoffs must be limited to subnetworks representing processes that can be served by compatible devices located at the same processing center.

But in this case, it can be seen that queue-state dependencies cannot yield an efficient solution. Maximum throughput requires full utilization of the available processors. Since the number of customers in a subnetwork does not remain constant, either processing power must be held in reserve when the number of customers in the subnetwork is less than N , or else there is the possibility of an arrival to the subnetwork when all of the processors are busy. If the arriving customer joins a queue that is idle, a processor must be taken from another subnetwork queue, even though there is no change in the number of customers at that queue. This violates the assumption of queue-state rate dependency.

5. THROUGHPUT WITH SUBNETWORK-STATE DEPENDENT PROCESSING RATES

Dynamic processing power allocations in a subnetwork can be useful to optimize throughput only if they are subnetwork-state dependent; in this case, the processing power at queue i of an M -queue subnetwork is determined by the function $r_i(n_1 \dots n_M)$ of the subnetwork state. We assume as before that a total processing

power $R(N)$ is available to a subnetwork when it contains N customers, and that the subnetwork retains the local balance property. Consider, for example, the parallel and series subnetworks of Figure 3. For either case, when there are m and n customers at queues one and two, respectively, the processing rates are $u_1(m,n)$ and $u_2(m,n)$, respectively. Each of these two-queue subnetworks will have the local balance property if and only if

$$\frac{u_1(m,n-1)}{u_1(m,n)} = \frac{u_2(m-1,n)}{u_2(m,n)}, \text{ for all } m, n > 0. \quad (5-1)$$

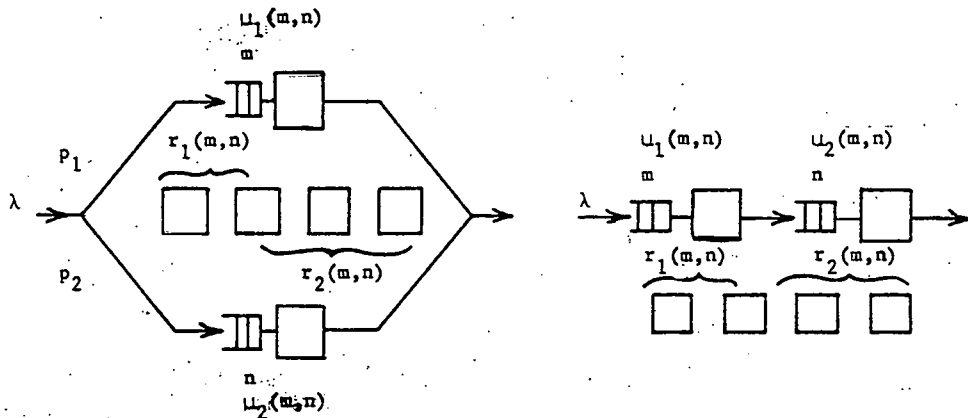


Figure 3: Parallel and Series Subnetworks

A generalization for more than two queues is straightforward for the parallel case, but not for the series subnetwork case. The state probabilities for either subnetwork are expressed by

$$P(m,n) = \frac{\lambda_1^m}{\prod_{i=1}^m u_1(i,n)} \cdot \frac{\lambda_2^n}{\prod_{i=1}^n u_2(0,i)} P(0,0), \quad (5-2)$$

where $\lambda_1 = p_1\lambda$, $\lambda_2 = p_2\lambda$ for the parallel case, and $\lambda_1 = \lambda_2 = \lambda$ for the series case.

According to Norton's Theorem, the processing rate of a subnetwork that contains n customers is the same as the throughput of the subnetwork when it stands alone as a closed network containing n customers. This is because the distribution of the n customers in the subnetwork is the same in either case. For both the series and parallel subnetworks, then, maximization of network throughput apparently entails maximization of each processing rate $\tau(1), \tau(2) \dots \tau(N)$, which are the throughputs of the two queue network (i.e., the isolated subnetwork) containing $n = 1, 2, \dots, N$ customers. By Theorem 4, these rates may be achieved by any processing power distribution such that no processing power is wasted; i.e., for maximum $\tau(n)$, $r_1(n,0) = r_2(0,n) = R(n)$. But notice that if the local balance condition (5-1) holds, $\tau(n)$ cannot be maximized independently of $\tau(n-1)$, for $1 < n \leq N$. The optimization procedure is, therefore, not as straightforward as indicated above; the details can be found in [10]. The result is that the optimal processing rate

assignments are consistent with full processing power utilization, and are only further constrained by (5-1). For maximum network throughput, when processing-power tradeoffs in either subnetwork are possible,

$$r_1(m,n) = \frac{m}{N} R(N), \text{ and } r_2(m,n) = \frac{n}{N} R(N),$$

where m and n are the number of customers at queues one and two, respectively, and $N = m + n$.

The state probabilities for the subnetwork with this processing power assignment are

$$P(m,n) = \binom{m+n}{n} \frac{1}{\alpha(m+n)} \rho_1^m \rho_2^n P(0,0),$$

where $\rho_i = \frac{\lambda_i}{\mu_i}$, and $\alpha(n) = \prod_{i=1}^n R(i)$. If R processors are available to the subnetwork,

$$\alpha(n) = \begin{cases} n! & , \text{ for } n \leq R \\ n! R^{n-R} & , \text{ for } n > R. \end{cases}$$

The state probability of this form will be a component of the product form solution of any locally balanced network containing the subnetwork.

The result, that maximum throughput is achieved by sharing the processing power proportionately based on the load at each queue, can be generalized for the case of any fixed number of queues in parallel. This will serve as a local balance model of a multiprocessor system handling parallel queues in a computer network.

ACKNOWLEDGEMENT

The author would like to thank the session chairman, Ken Sevcik, and the referees, particularly Tom Giammo, for their helpful suggestions.

REFERENCES

1. Baskett, F., K. M. Chandy, R. R. Muntz and F. Palacios-Gomez, "Open, Closed and Mixed Networks of Queues with Different Classes of Customers", JACM 22 2 (1975), pp. 248-260.
2. Buzen, J. P., "Computational Algorithms for Closed Queueing Networks with Exponential Servers" CACM 16, September 1973.
3. Chandy, K. M., "The Analysis and Solutions for General Queueing Networks", Proc. Sixth Annual Princeton Conf. on Information Sciences and Systems, Princeton Univ., Princeton, N. J., (March 1972), pp. 219-224.
4. Chandy, K. M., U. Herzog, and L. Woo, "Parametric Analysis of Queueing Network Models" IBM Journal of Research and Development, 19,1 (January 1975), pp. 36-42.
5. Chandy, K. M., U. Herzog, and L. Woo, "Approximate Analysis of General Queueing Networks" IBM Journal of Research and Development, 19, 1 (January 1975).

- 11367
6. Gordon, W. J. and G. F. Newell, "Closed Queueing Systems with Exponential Servers", Oper. Res. 15, 2 (1967), pp. 252-267.
 7. Keller, T. W., "ASO Manual" Dept. of Computer Sciences Report TR-27, University of Texas, Austin, Texas, 1973.
 8. Kleinrock, L., Queueing Systems, Vol. 1: Theory. John Wiley and Sons, N. Y., 1975.
 9. Noetzel, A. S., "Analysis of Discrete-Time Queues for Random Output Processes and Network Solutions" Proc. 1977 Conf. on Information Sciences and Systems, The Johns Hopkins University, March 1977.
 10. Noetzel, A. S., "Product-Form Queueing Networks with Processing-Rate and Arrival-Rate Tradeoffs" Technical Report 53, Department of Computer Sciences, The University of Texas at Austin, December 1975.
 11. Reiser, M., and H. Kobayashi, "Recursive Algorithms for General Queueing Networks with Exponential Servers" IBM Res. Report. RC-4254, March 1973.
 12. Reiser, M., "QNET4 User's Guide". IBM Research Report RA71, June 1975.
 13. Reiser, M. and H. Kobayashi, "On the Convolution Algorithm for Separable Queueing Networks" Proceedings, 1976 International Symposium on Computer Performance Modeling, Measurement, and Evaluation, Harvard University, Cambridge, Mass.

Appendix A: Proof of Theorem 1.

The throughput $\tau(N)$ of the locally balanced two queue network can be expressed as the ratio of the normalization constants with $N - 1$ and N customers in the network. Hence, the theorem is proved if

$$\tau(N+1) = \frac{G(N)}{G(N+1)} \geq \frac{G(N-1)}{G(N)} = \tau(N), \quad (A-1)$$

where

$$G(n) = \sum_{i=0}^n X(n-i)Y(i), \text{ and } X(n) = \prod_{i=1}^n \frac{1}{u(i)} \text{ and } Y(n) = \prod_{i=1}^n \frac{1}{U(i)}$$

for $n = 0, 1 \dots N$.

The inequality (A-1) can be written

$$G^2(N) \geq G(N+1)G(N-1) \quad (A-2)$$

The terms of $G^2(N)$ contain the factors $Y(i)Y(j)$ for $0 \leq i, j \leq N$. The terms of $G(N+1)G(N-1)$ contain the factors $Y(i)Y(j)$ for $0 \leq i \leq N - 1, 0 \leq j \leq N + 1$. In both cases $0 \leq i + j \leq 2N$. The inequality (A-2) is demonstrated by grouping the terms into $2N + 1$ inequalities. Inequality $k, 0 \leq k \leq 2N$, will have all the terms with factors $Y(i)Y(j)$ such that $i + j = k$.

First, consider the case $0 \leq k < N$. Collecting terms from (A-2),

$$\sum_{i=0}^k X(N-i)Y(i)X(N-k+i)Y(k-i) \geq \sum_{i=0}^k X(N+1-i)Y(i)X(N-1-k+i)Y(k-i). \quad (A-3)$$

Grouping coefficients of $Y(i)Y(k-i)$,

$$\sum_{i=0}^k [X(N-i)X(N-k+i) - X(N+1-i)X(N-1-k+i)]Y(i)Y(k-i) \geq 0. \quad (A-4)$$

This sum can be rewritten as two summations; first for index $i = 0$ to $\lfloor \frac{k}{2} \rfloor$, and then for $i = k - \lfloor \frac{k}{2} \rfloor + 1$ to k . If k is even, these two ranges obviously cover the range 0 to k . If k is odd, the term for $i = \lfloor \frac{k}{2} \rfloor + 1$ is missing. But the term in the summation for this value of i is zero. Hence, rewriting (A-4) in two summations, and letting $j = k + 1 - i$ replace i as the index of the second summation, one has

$$\begin{aligned} & \sum_{i=0}^{\lfloor \frac{k}{2} \rfloor} [X(N-i)X(N-k+i) - X(N+1-i)X(N-1-k+i)]Y(i)Y(k-i) \\ & + \sum_{j=1}^{\lfloor \frac{k}{2} \rfloor} [X(N-k-1+j)X(N+1-j) - X(N-k+j)X(N-j)]Y(k+1-j)Y(j-1) \geq 0. \quad (A-5) \end{aligned}$$

Note that for each $X(m)X(n)$ in (A-5), $m + n = 2N - k$. Also, note that for any m, n with $m \geq n$ and any $j \leq n$,

$$\begin{aligned} & X(m)X(n) - X(m+j)X(n-j) \\ & = X(m)X(n-j) \left[\prod_{i=n-j+1}^n \frac{1}{u(i)} - \prod_{i=m+1}^{m+j} \frac{1}{u(i)} \right] \geq 0, \quad (A-6) \end{aligned}$$

since all of the indices i , and hence rates $u(i)$ in the second product are greater than those of the first product. Therefore, the products $X(m)X(n)$ for all $m + n = 2N - k$ are ordered inversely as $|m-n|$, or directly as $\min(m, n)$. If $i = \min(m, n)$ let $\bar{X}(i) = X(m)X(n)$ and let $\bar{Y}(i) = Y(m)Y(n)$. Then the inequality (A-5) can be rewritten by selecting the smaller index of each product, as follows:

$$\sum_{i=0}^{\lfloor \frac{k}{2} \rfloor} (\bar{X}(N-k+i) - \bar{X}(N-k+i-1))\bar{Y}(i) + \sum_{j=1}^{\lfloor \frac{k}{2} \rfloor} (\bar{X}(N-k+j-1) - \bar{X}(N-k+j))\bar{Y}(j-1) \geq 0. \quad (A-7)$$

The summations of (A-7) can be combined and terms rearranged to obtain

$$(\bar{X}(N-k) - \bar{X}(N-k-1))\bar{Y}(0) + \sum_{i=1}^{\lfloor \frac{k}{2} \rfloor} [\bar{X}(N-k+i) - \bar{X}(N-k+i-1)][\bar{Y}(i) - \bar{Y}(i-1)] \geq 0. \quad (A-8)$$

It is seen that each factor of every term of the summation is nonnegative. Hence, the inequality is demonstrated.

Now consider the case $k = N$. In collecting all terms of (A-2) with factors $Y(i)Y(j)$ where $i + j = N$, $G(N-1)$ contributes terms with factors $Y(i)$ for $0 \leq i \leq N - 1$. Hence, $G(N+1)$ contributes terms with factors $Y(j)$ for $1 \leq j \leq N$. The inequality corresponding to (A-3) is

$$\sum_{i=0}^N X(N-i)Y(i)X(i)Y(N-i) \geq \sum_{i=0}^{N-1} X(i+1)Y(N-i)X(N-1-i)Y(i). \quad (A-9)$$

Collecting terms, and then adjusting the index of the summation to range from 1 to N , this inequality is written as follows:

$$X(0)X(N)Y(0)Y(N) + \sum_{i=1}^N [X(N+1-i)X(i-1) - X(i)X(N-i)]Y(N+1-i)Y(i-1) \geq 0. \quad (A-10)$$

The summation can be expressed as two summations, first with index $i = 1$ to $\lfloor \frac{N}{2} \rfloor$, then with $i = N - \lfloor \frac{N}{2} \rfloor + 1$ to N , noting that if N is odd, the term for $i = \lfloor \frac{N}{2} \rfloor + 1$ disappears. Then rewriting the second summation with index $j = N + 1 - i$, one obtains

$$X(0)X(N)Y(0)Y(N) + \sum_{i=1}^{\lfloor \frac{N}{2} \rfloor} [X(N+1-i)X(i-1) - X(i)X(N-i)]Y(N+1-i)Y(i-1) + \sum_{j=1}^{\lfloor \frac{N}{2} \rfloor} [X(j)X(N-j) - X(N+1-j)X(j-1)]Y(j)Y(N-j) \geq 0. \quad (A-11)$$

Then, if $i = \min(m, n)$, let $\bar{X}(i) = X(m)X(n)$, and $\bar{Y}(i) = Y(m)Y(n)$.

The summations of (A-11) can be combined and terms rearranged to obtain

$$\bar{X}(0)\bar{Y}(0) + \sum_{i=1}^{\lfloor \frac{N}{2} \rfloor} [\bar{X}(i) - \bar{X}(i-1)][\bar{Y}(i) - \bar{Y}(i-1)] \geq 0. \quad (A-12)$$

Since each term in the summation is (A-12) nonnegative, the inequality is demonstrated.

Last, the case for $N < k \leq 2N$ must be considered. But from the symmetry of X and Y in the definition of the function G , inequality k of A-2 is exactly inequality $2N - k$ with the roles of X and Y interchanged. Hence, it has been demonstrated in the first case. The theorem is proved.

Appendix B: Proof of Theorem 3

Let $Y(n) = \prod_{i=1}^n \frac{1}{U(i)}$ and $X(n) = \prod_{i=1}^n \frac{1}{u(i)}$ for $0 < n \leq N$. Let

$G(N) = \sum_{i=0}^N Y(i)X(N-i)$ be the normalization constant for the state probabilities of

the two queue network with N customers. Then for all N > 0,

$$\tau(N) = \frac{G(N-1)}{G(N)}. \tag{B-1}$$

For $0 < n \leq N$ let $G_n^+(N)$ be all of the terms of $G(N)$ that have the factor $U^{-1}(n)$,

$$G_n^+(N) = \sum_{i=n}^N Y(i)X(N-i), \tag{B-2}$$

and $G_n^-(N) = G(N) - G_n^+(N)$.

Then, differentiating $G(N)$ with respect to $U(n)$,

$$\frac{\partial G(N)}{\partial U(n)} = -\frac{1}{U(n)} G_n^+(N),$$

and differentiating $\tau(N)$ with respect to $U(n)$,

$$\begin{aligned} \frac{\partial \tau(N)}{\partial U(n)} &= \left[G(N) \frac{\partial G(N-1)}{\partial U(n)} - G(N-1) \frac{\partial G(N)}{\partial U(n)} \right] \frac{1}{G^2(N)} \\ &= \left[-(G_n^-(N) + G_n^+(N))G_n^+(N-1)U(n)^{-1} \right. \\ &\quad \left. + (G_n^-(N-1) + G_n^+(N-1))G_n^+(N)U(n)^{-1} \right] \frac{1}{G^2(N)} \\ &= \left[-G_n^-(N)G_n^+(N-1) + G_n^-(N-1)G_n^+(N) \right] \frac{1}{U(n)G^2(N)} \end{aligned} \tag{B-3}$$

The derivative will be nonnegative if

$$G_n^-(N-1)G_n^+(N) \geq G_n^-(N)G_n^+(N-1). \tag{B-4}$$

But each term of this inequality has exactly one factor $U^{-1}(n)$. Hence, it may be cancelled out of the inequality. As a function of $U(n)$, $\tau(N)$ is, therefore, either always increasing, always decreasing, or is constant. This proves part a) of the theorem.

The terms on the right of the inequality (B-4) have factors $Y(k)Y(i)$, $0 \leq k < n$, $n \leq i < N$, and the terms on the left have factors $Y(k)Y(i)$, $0 \leq k < n$, $n \leq i \leq N$. The coefficient of each $Y(k)Y(i)$ that appears on the right is $X(N-k)X(N-1-i)$, and the coefficient of that term on the left is $X(N-1-k)X(N-i)$. Hence, if $X(N-1-k)X(N-i) \geq X(N-k)X(N-1-i)$, the inequality holds. But since $i > k$, $X(N-1-i)X(N-1-k)$ can be factored out of this inequality, leaving $\frac{1}{U(N-i)} \geq \frac{1}{U(N-k)}$.

Therefore, if $u(N-k) \geq u(N-i)$ for all $0 < k < n \leq i \leq N$, $\tau(N)$ is a nondecreasing function of $U(n)$. This proves part b) of the theorem.

In particular, it should be noted that if $u(i) \geq u(j)$ for all $i > j$, which is the usual case, then $\tau(N)$ is a nondecreasing function of $U(n)$, for all n.

Now let $H_m^+(N)$ be the sum of all of the terms of $G(N)$ that have the factor $U^{-1}(m)$.

$$H_m^+(N) = \sum_{i=m}^N Y(N-i)X(i) \quad (B-5)$$

and

$$H_m^-(N) = G(N) - H_m^+(N).$$

Then, from the definition of $H_m^+(N)$, the following relationships are noted:

$$H_{N-n}^+(N) = G_{n+1}^-(N) \quad (B-6)$$

and

$$H_{N-n}^+(N-1) = G_n^-(N-1). \quad (B-7)$$

By the steps leading to (B-4) the condition for $\frac{\partial \tau(N)}{\partial u(N-n)} > 0$ is determined to be

$$H_{N-n}^-(N-1)H_{N-n}^+(N) > H_{N-n}^-(N)H_{N-n}^+(N-1). \quad (B-8)$$

Using (B-6) and (B-7) this can be expressed as

$$G_n^+(N-1)G_{n+1}^-(N) > G_{n+1}^+(N)G_n^-(N-1). \quad (B-9)$$

And then $G_{n+1}^+(N)$ can be related to $G_n^+(N)$,

$$G_n^+(N-1)[G_n^-(N) + Y(N-n)X(n)] > [G_n^+(N) - Y(N-n)X(n)]G_n^-(N-1), \quad (B-10)$$

which can be written

$$[G_n^+(N-1)G_n^-(N) - G_n^+(N)G_n^-(N-1)] + Y(N-n)X(n)G(N-1) > 0 \quad (B-11)$$

The inequality (B-11) must be satisfied if the term in brackets is nonnegative. But, this term expresses the condition (B-4); it will be nonnegative if $\tau(N)$ is a nonincreasing function of $U(n)$. This proves part c) of the theorem.

Appendix C: Proof of Theorem 4

The proof is by induction on M . First, consider any network with $M = 1$. Let p be the probability that a customer leaving the queue uses a particular branch b in returning to the queue. Let the processing rate at the queue be u when a single processor is assigned the queue. If $r(N)$ processors are assigned the queue when the network contains N customers, the throughput in branch b is $\tau_1(N) = pr(N)u$. The throughput is maximum when $r(N) = R(N)$. Then $\tau_1(N) = puR(N)$, which satisfies the theorem.

Suppose the theorem holds for all networks of $M - 1$ queues. Let $\tau_{M-1}(n)$ be the throughput at some branch b of an $M - 1$ queue network when there are n customers in the network. Then by Norton's Theorem, the $M - 1$ queue network may be represented by an equivalent queue with respect to branch b . If $U(n)$ is the processing rate of the equivalent queue when it contains n customers, then $U(n) = \tau_{M-1}(n)$. Let $R_1(n)$ be the processing power available to the $M - 1$ queue network. By the inductive hypothesis, the maximum throughput at branch b is $\tau_{M-1}(n) = UR_1(n)$, where U is a constant, and is achieved when processing power $R_1(n)$ is fully utilized.

This is also the maximum processing rate of the equivalent queue when it contains n customers and has available processing power $R_1(n)$.

Let $\tau_M(N)$ be the throughput in the two queue network consisting of the equivalent queue in series with queue M ; this will be the same as the throughput in branch b of the $M - 1$ queue network with queue M inserted in branch b . Assuming the processing rate at queue M is $\mu(n) = r(n)u$ when it contains n customers and is using processing power $r(n)$, the throughput of the equivalent two queue network is determined as follows.

Let $G(n)$, $X(n)$ and $Y(n)$ be as defined in (A-1).

$$\text{Then } \tau_M(N) = \frac{G(N-1)}{G(N)} \quad (C-1)$$

Let $R(N)$ be the processing power available to the M queue network when it contains N customers. And let $\rho(n) = \frac{r(n)}{R(N)}$ be the optimum fraction of the available processing power to be used by queue M when it contains $n \leq N$ customers, in order to maximize $\tau_M(N)$. Then $u(n) = \rho(n)R(N)u$ are the processing rates at queue M that maximize $\tau_M(N)$.

Let $\bar{\rho}(n) = 1 - \rho(n)$. Then, for maximum $\tau_M(N)$, processing power $R_1(n) = \bar{\rho}(n)R(N)$ is available to be allocated to the equivalent queue when it contains $N - n$ customers.

Examining (C-1) shows that for maximum $\tau_M(N)$, $U(N)$ and $u(N)$ are to be maximized.

Clearly, all available processing power is used for these rates, so that

$$\rho(N) = \bar{\rho}(0) = 1.$$

Note that $\rho(n) > 0$ for all $n > 0$ may be assumed. For if this is not the case, let m be the largest integer for which $\rho(m) = 0$. Then at least m customers will always be at queue M . Let $N' = N - m$ and $\mu'(n) = \mu(n+m)$. Maximization of $\tau_M(N)$ is accomplished in this case by considering only the rates $U(N'-n)$ and $u'(n)$, for $n \leq N'$. The result will be the same as maximization with $\rho(n) > 0$ for $n \leq N$, if it is shown that processing power is fully utilized when there are m customers in queue M . But note that $\rho(m) = 0$ only if $\tau_M(N)$ is a nonincreasing function of $u'(0)$. By Theorem 3C, then, $\tau_M(N)$ is an increasing function of $U(N')$, and hence, $\bar{\rho}(m) = 1$.

With the rates $\rho(n)R(N)$ for queue M fixed at the values required for maximum $\tau_M(N)$, the rates $U(n)$ may be selected within the range $0 \leq U(n) \leq \bar{U}_0(N-n)R(N)$ to maximize $\tau_M(N)$.

Suppose $\tau_M(N)$ is not an increasing function of $U(n)$, for some $n < N$. Then, by Theorem 3C, it must be an increasing function of $u(N-n)$. Then $\rho(N-m) = 1$, and, therefore, $U(m) = 0$. If m is the largest integer for which $\tau_M(N)$ is not an increasing function of $U(m)$, then there will never be less than m customers at the equivalent queue. Hence, letting $N' = N - m$ and $U'(n) = U(n+m)$, only rates $U'(n)$ for $0 \leq n \leq N'$ must be considered in maximizing $\tau_M(N)$. And this maximization will yield the same result as maximization with $\tau_M(N)$ an increasing function of $U(n)$ for all $n > 0$. Therefore, the maximum value $U(n) = \bar{U}_0(N-n)R(N)$ must be chosen for $U(n)$, $n \leq N$ in order to maximize $\tau_M(N)$. This proves part a) of the theorem. Then each term $X(j)Y(k)$ of $G(N-1)$, where $j + k = N - 1$, can be written

$$X(j)Y(k) = \mu^j U^k R^{j+k} (N) \prod_{i=1}^j \rho(i) \prod_{i=1}^k \rho(N-1) \quad (C-2)$$

The terms of $G(N)$ have the same form, but $j + k = N$. Each term $X(j)Y(k)$ of $G(N)$ with $j, k > 0$ contains the product

$$\begin{aligned} \frac{1}{\rho(j)} \cdot \frac{1}{\rho(j)} &= \frac{1}{\rho(j)} + \frac{1}{\rho(j)}. \text{ Hence, it can be written} \\ X(j)Y(k) &= \frac{1}{UR(N)} \left[\mu^{j-1} U^k R^{N-1} (N) \prod_{i=1}^{j-1} \rho(i) \prod_{i=1}^k \rho(N-1) \right]^{-1} \\ &+ \frac{1}{UR(N)} \left[\mu^j U^{k-1} R^{N-1} (N) \prod_{i=1}^j \rho(i) \prod_{i=1}^{k-1} \rho(N-1) \right]^{-1} \\ &= \frac{1}{UR(N)} X(j-1)Y(k) + \frac{1}{UR(N)} X(j)Y(k-1). \end{aligned} \quad (C-3)$$

$G(N)$ also includes the terms

$$X(N)Y(0) = \left[\mu^N R^N (N) \prod_{i=1}^{N-1} \rho(i) \right]^{-1} = \frac{1}{UR(N)} X(N-1)Y(0) \quad (C-4a)$$

and

$$X(0)Y(N) = \left[\mu^N R^N (N) \prod_{i=1}^{N-1} \rho(N-1) \right]^{-1} = \frac{1}{UR(N)} X(0)Y(N-1). \quad (C-4b)$$

$G(N)$ may then be expressed as follows:

$$\begin{aligned} G(N) &= X(N)Y(0) + \sum_{\substack{j+k=N \\ j,k>0}} X(j)Y(k) + X(0)Y(N) \\ &= \frac{1}{UR(N)} X(N-1)Y(0) + \sum_{\substack{j+k=N \\ j,k>0}} \frac{1}{UR(N)} X(j-1)Y(k) \\ &+ \frac{1}{UR(N)} X(0)Y(N-1) + \sum_{\substack{j+k=N \\ j,k>0}} \frac{1}{UR(N)} X(j)Y(k-1) \\ &= \frac{1}{UR(N)} \sum_{j+k=N-1} X(j)Y(k) + \frac{1}{UR(N)} \sum_{j+k=N-1} X(j)Y(k) \\ &= \left(\frac{1}{U} + \frac{1}{U} \right) \frac{1}{R(N)} G(N-1). \end{aligned} \quad (C-5)$$

$$\text{Therefore, when } \tau_M(N) \text{ is maximum, } \tau_M(N) = \frac{G(N-1)}{G(N)} = \left(\frac{1}{U} + \frac{1}{U} \right)^{-1} R(N). \quad (C-6)$$

Then part b) of the theorem is proved.