ON HOMOGENEITY AND

ON-LINE=OFF-LINE BEHAVIOR

IN M/G/1 QUEUEING SYSTEMS

by

Raymond M. Bryant

ON HOMOGENEITY AND ON-LINE=OFF-LINE BEHAVIOR

IN M/G/1 QUEUEING SYSTEMS

Raymond M. Bryant
Computer Sciences Department
and
Madison Academic Computing Center
University of Wisconsin-Madison

ABSTRACT

Operational analysis replaces certain classical queueing theory assumptions with the conditions of "homogenous service times" and "on-line=off-line behavior." It has been conjectured that these conditions hold as $t \to \infty$ only if the service times are exponentially distributed. In this paper, we show that this is correct for M/G/1 queueing systems. We also state dual results for inter-arrival times in G/M/1. Finally, we consider the relationship between the operational quantities $S(n)$, $n=1, 2, \ldots$ and the mean service time in M/G/1. This relationship is shown to depend on the form of the service time distribution.

# ON HOMOGENEITY AND ON-LINE=OFF-LINE BEHAVIOR

# IN M/G/1 QUEUEING SYSTEMS

## 1. INTRODUCTION

Operational Analysis [2; 3; 4] is a non-classical approach to analysis of queueing systems in which the system parameters $\lambda(n)$ and $\mu(n)$ are replaced by observed quantities $I(n)$ and $S(n)$ respectively. Assumptions about arrival and service time distributions are replaced by conditions on $S(n)$ and $I(n)$. Two of the key conditions are "homogenous service times" which states that $S(n)$ is constant in n and "on-line=off-line behavior" which states that the $S(n)$ do not depend on the system arrival rate.

Whenever a new idea like this appears, it is natural to explore its relation to the existing theory. This paper examines the relationship between operational and classical concepts by considering the limiting values (as $t\rightarrow\infty$) of $I(n)$ and $S(n)$ for the sample paths of an M/G/1 queueing system. The primary results are that on-line=off-line behavior and homogenous service times occur in M/G/1 if and only if the service times are exponentially distributed. Dual results for the G/M/1 queue are stated. It is also shown that open, feed-forward networks of single-server queues with Poisson external arrivals can have product form solutions with load independent behavior if and only if all the service times are exponential. Finally, exact values

for  S(n)  in  M/G/1  queueing  systems  are  derived  and  their
dependence  on  the  mean  service  time  is depicted for several
standard service time distributions.  This discussion shows  that
to  estimate  how  an  observed value of S(n) would change if the
server's rate were increased or decreased, one must  specify  the
service time distribution.

In Section 2 we describe the notation of the paper and  give
definitions  of "homogeneous service times" and "on-line=off-line
behavior."  Section 3  discusses  what  it  means  for  an  M/G/1
queueing  system  to  have  these  operational  properties;  this
section also contains the  main  results  of  the  paper.   These
results  are  used  to  provide  a method of calculating S(n) for
arbitrary service times  in  an  M/G/1  queue.   Graphs  of  these
values versus mean service time are then given in Section 4.


## 2.   NOTATION AND DEFINITIONS


Throughout this paper, whenever we are considering an  M/G/1
queueing  system,  we  will assume that the system is stable, has
arrival rate $\lambda$ and service distribution B(t).  We  let  $\bar{x}$  denote
the  mean  service  time  and  $\mu=1 / \bar{x}$.  We will let $\rho$ denote the
system utilization and p(n) denote the stationary probability  of
finding  n  customers in system.  For a G/M/1 queueing system, we
let A(t) denote the inter-arrival time distribution, $\bar{a}$ denote the
mean inter-arrival time, and $\mu$ denote the system service rate.

We will use a superscript * to indicate the Laplace-Stieljes
transform; for example, $B^*(s)$ is the transform  of  B(t)  and  is

defined as:

$$B^*(s) = \int_{-\infty}^{+\infty} e^{-st} \, dB(t)$$

Where necessary to distinguish real numbers from real valued random variables, we will use an underline to indicate the random variable.

For any particular realization of an M/G/1 queueing system, we define the sample path $\omega(t)$ as the right continuous function giving the number of jobs in system versus time. We assume that $\omega(0) = 0$ for all sample paths, and that $\omega$ is a sample point in some probability space $\Omega$.

We now give some definitions from operational analysis. Most of this material is contained in [3], however we prefer a notation more similar to that of [4]. To emphasize the fact that these quantities depend on values observed during a finite time interval [0,t), we will modify the notation of [4] to explicitly include the parameter t.

We begin by defining the "basic operational measures" of a queueing system during [0,t):

A(n,t)    is the number of customers who arrive in [0,t) to find exactly n customers already in system.

C(n,t)    is the number of customers who left the system during [0,t) when there were exactly n customers in system.

T(n,t)    is the amount of time during [0,t) when there were exactly n customers in system.

Given these quantities, we then may define the following

"operational performance measures."    (We follow the convention of [4] and leave undefined any quantity with a zero denominator.):

S(n,t)    =T(n,t)/C(n,t) is the mean service time between job departures during [0,t) given n jobs in system.

I(n,t)    =T(n,t)/A(n,t) is the mean inter-arrival time during [0,t) given n jobs in system.

P(n,t)    =T(n,t)/t is the proportion of time there were n jobs in system during [0,t).

We note that in operational analysis, I(n,t) and S(n,t) serve the roles of conditional arrival and service rates in classical queueing theory (see [3]).

We will be primarily interested in the asymptotic values of S(n,t), I(n,t), and P(n,t) as t-->∞, assuming that these values can be defined in a reasonable way. We will indicate this limiting value (assuming it exists) by dropping the parameter t. Thus:

$$S(n) = \lim_{t \to \infty} S(n,t).$$

We will refer to S(n) and I(n) as the (asymptotic) service and arrival functions, respectively.

Finally, we wish to define certain operational terms so that they can be conveniently referred to in the sequel:

Definition 2.1: A queueing system is said to have homogenous arrivals during [0,t) if I(n,t) is constant in n.

Definition 2.2: A queueing system is said to have homogenous service times during [0,t) if S(n,t) is constant in n.

Definition 2.3: A queueing system is said to satisfy on-line=off-line behavior during [0,t) if it has homogenous service times during [0,t) and the common value of S(n,t) is equal to the average customer service time.

We point out that in operational analysis, the homogenous service time condition is the counterpart of the assumption of exponential service times in classical queueing analysis [4].

We now discuss the meanings of these conditions with regard to an M/G/1 queueing system.

## 3.  OPERATIONAL ANALYSIS AND M/G/1 QUEUEING SYSTEMS

The operational performance measures defined in the last section are calculated from observations of a system during a particular time interval [0,t). In the context of an M/G/1 queueing system, we would say that they have been defined for a particular sample path, $\omega_0$. Thus, we have defined what it means to say that "$\omega_0$ has homogenous service times during [0,t)" but we have yet to define what it means to say that "an M/G/1 queueing system has homogenous service times." It is the purpose of this section to define the latter phrase in what we believe is a natural way and to explore the consequences of such a definition.

For any sample path $\omega$ in $\Omega$ , let A(n,t,$\omega$), C(n,t,$\omega$), S(n,t,$\omega$), and I(n,t,$\omega$) be the values of A(n,t), C(n,t), etc. associated with $\omega$ during [0,t). Let A(n,t), C(n,t), etc. denote

the random variables thus defined on $\Omega$. Let $E_a(n)$ be the event that an arrival occurs to find n jobs already in system, and let $E_d(n)$ be the event that a departure occurs when there were n jobs in system. Finally, if E is a recurrent event, let m(E) denote the mean recurrence time of the event. Then we note that for any stable M/G/1 queueing system:

(1)    With probability one, $\underline{T}(n,t)/t \longrightarrow p(n)$ as $t \longrightarrow \infty$.

(2)    Since the embedded Markov Chain defined at departure instants is irreducible and positive recurrent, it follows that $m(E_d(n)) < \infty$, for all $n > 0$. Furthermore, since the probability of two or more arrivals in [t,t+h) is o(h), it follows that $0 < m(E_d(n))$.

(3)    The recurrence times of $E_d(n)$ are asymptotically independent random variables, since recurrence times during distinct busy cycles must be independent. Therefore, by an elementary result of renewal theory [8, p. 36]:

$$\lim_{t \longrightarrow \infty} \underline{C}(n,t)/t = 1/m(E_d(n))$$

with probability one.

(4)    $\underline{S}(n) = \lim_{t \longrightarrow \infty} (\underline{T}(n,t)/t)/(\underline{C}(n,t)/t)$.

We have thus shown:

Theorem 3.1: The limiting random variables $\underline{S}(n)$ are constant with probability one and $S(n) = p(n) \, m(E_d(n))$. $\square$

To get a similar statement for $\underline{I}(n)$, we need the following Lemma, which we will find useful later in this section:

Lemma 3.2: In any stable M/G/l queueing system,

$$\underset{t \to \infty}{\text{Lim}} \; \frac{A(n-1,t)}{t} = \underset{t \to \infty}{\text{Lim}} \; \frac{C(n,t)}{t},$$

with probability one, for all $n \geq 1$.

Proof: Let $\{t_i\}$ be the starting instances of the busy cycles

of the queue. Clearly $t_i \to \infty$ as $i \to \infty$ and $t_i < \infty$ for all i,

both statements with probability one. Similarly,

$A(n-1,t_i) = C(n,t_i)$ with probability one, since the number of

up-crossings of level n-1 must be the same as the number of

down-crossings level n at the start of each busy cycle. (Note

that the arrival at time $t_i$ is not counted in $A(0,t_i)$ since

$A(0,t_i)$ is the number of arrivals which found the system empty

during $[0,t_i)$.) Finally, we note that $A(n-1,t,\omega) \geq C(n,t,\omega) \geq$

$A(n-1,t,\omega)-1$ for all t and all sample paths $\omega$. Thus with

probability one

$$\underset{t \to \infty}{\text{Lim}} \; \frac{A(n-1,t)}{t} = \underset{t \to \infty}{\text{Lim}} \; \frac{C(n,t)}{t}. \quad \square$$

Therefore, $E_a(n)$ is a recurrent event whenever $E_d(n+1)$ is,

and we have

Theorem 3.3: The limiting random variables $I(n)$ are constant

with probability one and $I(n) = p(n) \; m(E_a(n))$. $\square$

Since $I(n)$ and $S(n)$ are almost everywhere constant, we will

drop the distinction between these random variables and their

values.

With these facts in mind, is seems natural to suggest the following definitions:

Definition 3.4: An M/G/1 queueing system will be said to have homogenous arrival rates if and only if $I(n)$ is constant in n.

Definition 3.5: An M/G/1 queueing system will be said to have homogenous service times if and only if $S(n)$ is constant in n.

Definition 3.6: An M/G/1 queueing system will be said to satisfy on-line=off-line behavior if and only if $S(n) = \bar{x}$ for all n.

Now we wish to determine what types of M/G/1 queueing systems satisfy these definitions. We begin with a basic Lemma:

Lemma 3.7: In any stable M/G/1 queueing system:

$$(3.1) \quad S(1) = \frac{1}{\lambda} \left[ \frac{1}{B^*(\lambda)} - 1 \right] .$$

Proof: $E_d(1)$ occurs if and only if the system becomes idle. Thus $m(E_d(1))$ is the mean busy cycle length. Now the Laplace transform of the busy period distribution, $G^*(s)$, is known to satisfy the functional equation:

$$G^*(s) = B^*[s + \lambda - \lambda G^*(s)]$$

(see, for example, [6, p. 212]). From this equation it is easy to determine the mean busy period length, and upon adding the mean idle time we obtain the mean busy cycle length:

$$1 / \lambda + \bar{x} / (1 - \rho).$$

Also, we know $p(1) = Q'(0)$, where $Q(z)$ is the generating

function of p(n). Thus p(1) can be found from the Pollaczek-Khinchin transform equation (see, for example, [6, p. 194]):

$$(3.2) \quad Q(z) = B^*(\lambda - \lambda z) \frac{(1 - \rho)(1-z)}{B^*(\lambda - \lambda z) - z}.$$

Calculating p(1) from equation (3.2) and using Theorem 3.1 shows that S(1) has the indicated form. □

We observe that on-line=off-line behavior means that the S(n)'s cannot depend on $\lambda$. This observation is the basis for the following theorem.

Theorem 3.8: Suppose $B^*(s)$ is analytic for $0 < Re(s) < \mu$. Then the M/G/1 queueing system satisfies on-line=off-line behavior if and only if B(t) is exponential.

Proof: (i) If B(t) is exponential then the result is straightforward.

(ii) Suppose that the system satisfies on-line=off-line behavior. Then, in particular, S(1) does not depend on $\lambda$. Solving equation (3.1) for $B^*(\lambda)$ we get

$$B^*(\lambda) = \frac{1}{1 + \lambda S(1)}, \quad 0 < \lambda < \mu.$$

But we have thus determined $B^*(s)$ on a set with limit point, and hence determined $B^*(s)$ throughout its region of analyticity. Therefore $B(t) = 1 - \exp(t S(1))$. □

If B(t) is non-exponential, the dependence of S(1) on $\lambda$ can be quite pronounced. In Figure 3.1 we have plotted S(1) versus $\lambda$ for some typical service distributions. (All distributions have $\bar{x}=1.0$.) The horizontal line at S(1)=1.0 represents the S(1) versus $\lambda$ curve for exponential service times. We see that those distributions with coefficients of variation less than one have S(1) versus $\lambda$ curves which lie above the exponential case; and that those distributions with coefficients of variation greater than one have S(1) versus $\lambda$ curves which lie below the exponential case.

We have shown that if the S(n)'s are constant in n and do not depend on $\lambda$, then the service time must be exponential. But is it possible that the S(n)'s are constant in n, but all of them depend on $\lambda$ in the same way? If so, we would still have homogenous service times but not on-line=off-line behavior. To study this situation, we begin with:

Lemma 3.9: In M/G/1, $I(n)=\lambda^{-1}$ for all n.

Proof: Because the arrival process is Poisson, $m(E_a(a)) = p(n) \lambda$. The result then follows from Theorem 3.3. □


Lemma 3.10: In M/G/1, p(n) = p(n-1) S(n) / I(n-1), n $\geq$ 1.

Proof: By Lemma 3.2 we know that $m(E_a(n-1)) = m(E_d(n))$. Using this fact, combining the formulas for S(n) and I(n) from Theorems 3.1 and 3.3, and solving for p(n) gives the desired result. □


We can now show that homogeneity is essentially equivalent

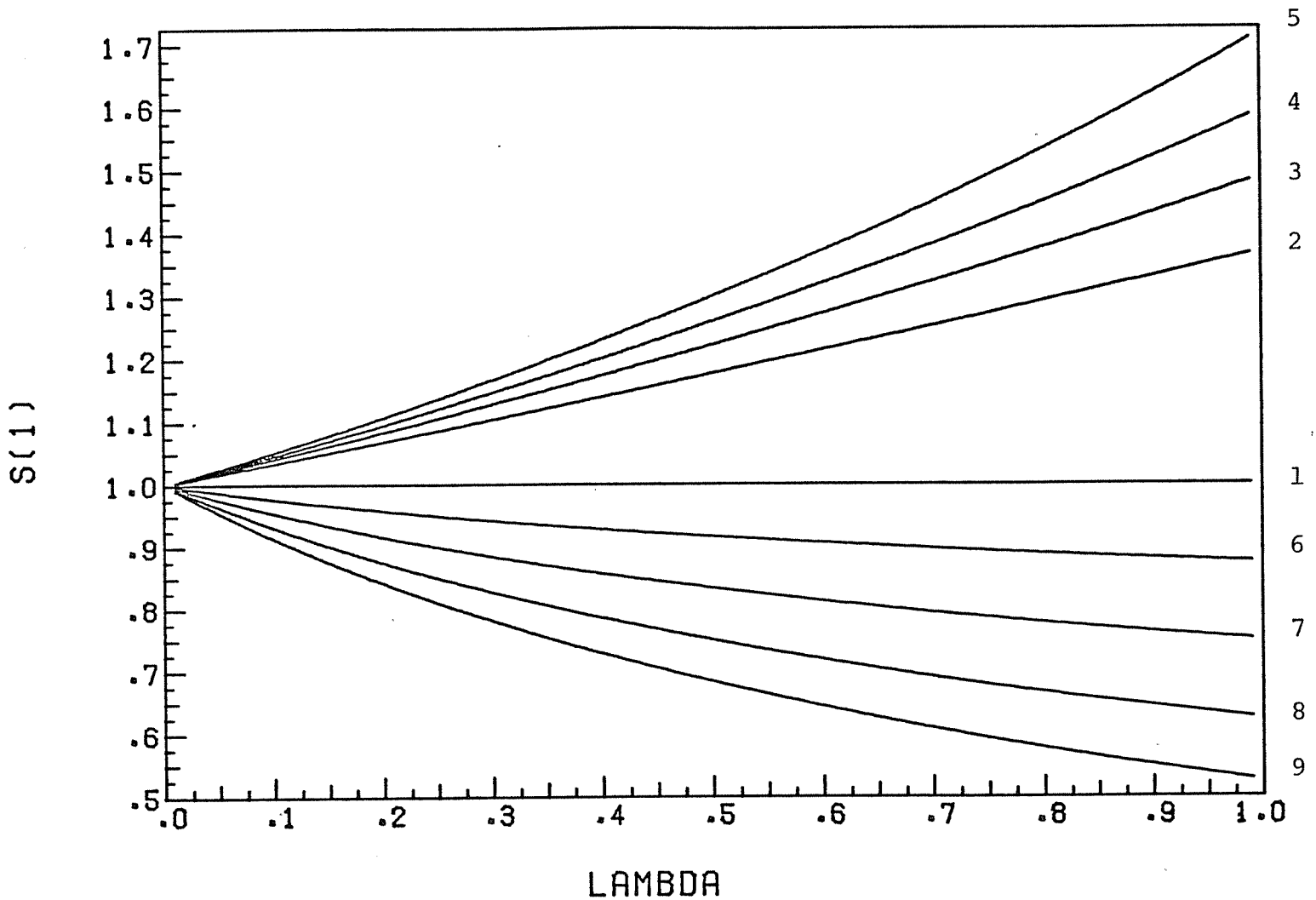# SERVICE FUNCTION S(1) VERSUS ARRIVAL RATE



Figure 3.1

| Curve Number | Distribution |
|---|---|
| 1 | Exponential |
| 2 | Erlang-3 |
| 3 | Erlang-5 |
| 4 | Erlang-10 |
| 5 | Constant |
| 6 | Hyperexponential, $\alpha=0.5$, $CV_2^2=1.5$ |
| 7 | Hyperexponential, $\alpha=0.5$, $CV_2^2=2.0$ |
| 8 | Hyperexponential, $\alpha=0.5$, $CV_2^2=2.5$ |
| 9 | Hyperexponential, $\alpha=0.5$, $CV^2=2.9$ |

to on-line=off-line behavior:

Theorem 3.11: An M/G/1 queueing system has homogeneous service times if and only if the service distribution is exponential.

Proof: (i) As before, if the service times are exponential the result is straightforward.

(ii) Let S denote the common value of S(n). From the Lemmas it follows that $p(n) = \lambda S p(n-1)$, $n > 1$. We know $p(0) = 1 - \rho$ in any M/G/1 queue. Thus

$$p(n) = (1 - \rho) (\lambda S)^n, n \geq 0.$$

But $\sum_n p(n) = 1$ implies that

$$\frac{1}{1 - \lambda S} = \frac{1}{1 - \rho}.$$

Hence $\rho = \lambda S$ or $S = 1 / \mu$. Thus S cannot depend on $\lambda$.

To finish the proof without the explicit analyticity condition of Theorem 3.8, we let P(z) be the generating function of p(n). Equating P(z) and Q(z) from equation (3.2) gives us:

$$\frac{(1 - \rho)}{(1 - \rho z)} = B^*(\lambda - \lambda z) \frac{(1 - \rho)(1-z)}{B^*(\lambda - \lambda z) - z}.$$

Solving for $B^*(\lambda - \lambda z)$ and substituting $s = \lambda - \lambda z$ yields:

$$B^*(s) = \frac{1}{1 + S s}. \quad \square$$

Before turning to the case of G/M/1, we note the following limited result for queueing networks:

Theorem 3.12: Consider any open, stable, feed-forward network of single-server queues, and assume that external arrivals to the network are Poisson and independent of the network state. Then the network stationary probability distribution p($\underline{n}$) is of the product form

$$(3.3) \quad p(\underline{n}) = G \prod_{i=1}^{N} \rho_i^{n_i}$$

if and only if all of the service times are exponential.

Proof: (i) If the service times are all exponential, the result is due to Jackson [5].

(ii) Pick any queue with only external arrivals. By the form of equation (3.3) this queue is an M/G/1 queue with homogenous service times. Hence, by Theorem 3.8, the service time at this queue is exponential. By Burke's Theorem [1], we then know that the departure process from this queue is Poisson. Repeating this argument at all queues with only external arrivals shows that all queues in the network have Poisson arrivals. Hence all queues in the network have Poisson arrivals and homogenous service times. Therefore all of the service times must be exponential. $\square$

We note that in the case of closed networks, there are several examples known of queueing systems which satisfy a type of product form, but which have decidedly non-exponential service times [4]. In such networks, the arrival process at each node is definitely not Poisson due to dependencies among inter-arrival

times.  We will not consider closed networks here.

We state without proof some dual results for G/M/1  queueing

systems:

Theorem 3.13: If $A^*(s)$ is analytic in  re(s)  >  $1 / \bar{a}$,  and

I(n)  does  not  depend  on  $\mu$,  then  A(t)  is  the  exponential

distribution.  □


Theorem 3.14: A G/M/1 queueing system has homogenous arrival

times (in the  sense  of  Definition 3.4)  if  and  only  if  the

inter-arrival time distribution is exponential.  □


## 4.  S(n) EVALUATION IN M/G/1

Let  us  consider  the  following  performance  prediction

problem:

> Random arrivals from a very large population are served
> one at a time, in FCFS order.  The system appears to be
> stable, but still is very heavily loaded.  What  would
> the  mean number of jobs in system be if a server twice
> as fast as the current one  were  installed?  Observed
> values of I(n,t) and S(n,t) are available.

To solve this  problem  using  operational  analysis,  one  would

adjust  the  S(n,t)  values to represent the service function for

the faster server and leave the I(n,t)  values  unchanged.   From

the  new  values  of S(n,t), the current values of I(n,t) and the

relation P(n,t) = P(n-1) S(n,t)/I(n-1,t), one could estimate  new

values of P(n,t) and hence new values of the mean number of jobs in system.

The difficult part of the problem is estimating the new values for S(n,t). One obvious estimate is to let the new values be one-half of the old values. However, simulation evidence indicates that this is not always the most accurate approach. Instead the exact way that the S(n,t) depend on $\bar{x}$ varies with the service time distribution at the server. To illustrate this dependence, we will consider the relationship between S(n) and $\bar{x}$ in an M/G/1 queueing system.

To evaluate S(n) we recall Lemma 3.10, which relates S(n), I(n), and p(n). Since in M/G/1, I(n)=$\lambda$, it follows that the p(n) determine the S(n). This observation (originally due to Buzen [3]), allows the calculation of S(n) from the service time distribution. To do so we can use a power series expansion of the Pollaczek-Khinchin transform formula (equation 3.2) to calculate p(n). While this process is algebraically involved, suitable tools exist to assist in the calculation and evaluation of the derivatives of Q(z). (The calculations in this paper used the FORMAL system [7], a FORMAC like system developed at the University of Maryland.) Thus for a particular service distribution, one can determine p(n) and hence S(n) as a function of the distribution parameters, at least for a few values of n.

For values of n > 10, the expressions for p(n) become too complex for FORMAL to handle. Even if such large expressions could be generated, round off error would probably make any evaluated results meaningless. As we shall see below, values of

$S(n)$ for $n > 5$ are usually not needed.

We believe it is unlikely that significantly simpler formulas for $S(n)$ will ever be found. It is clear that if simple closed form expressions for $S(n)$ were known, then simple closed form expressions for $p(n)$ could easily be constructed. Since no known formulas for the latter exist, it seems unlikely that any will be found for $S(n)$. (However, see [3] for an alternative approach.) We now return to our discussion of the relationship between $S(n)$ and $\bar{x}$.

We have already given a formula for $S(1)$ in M/G/1, and we begin our discussion by considering a graph of $S(1)$ versus $\bar{x}$ for some typical service distributions. (See Figure 4.1.) Throughout this discussion, we are considering a stable M/G/1 queueing system with $\lambda=1.0$. In this figure we have included the $S(1)$ versus $\bar{x}$ curve for exponential service times as a comparison. We see that of the service distributions we have considered, the distributions with squared coefficients of variation ($CV^2$) greater than 1 have $S(1)$ versus $\bar{x}$ curves which lie below the exponential case, and those distributions with $CV^2 < 1$ have curves which lie above the exponential case. Thus to estimate new values of $S(1,t)$ in our performance prediction problem, we should more than halve the observed $S(1,t)$ values when the service distribution has $CV^2 > 1$, and less than halve the observed $S(1,t)$ values when $CV^2 < 1$.

Figures 4.2, 4.3, and 4.4 give $S(n)$ versus $\bar{x}$ graphs for higher values of n. In each graph the $S(n)$ values for a specific distribution have been plotted. (As before, we have included the

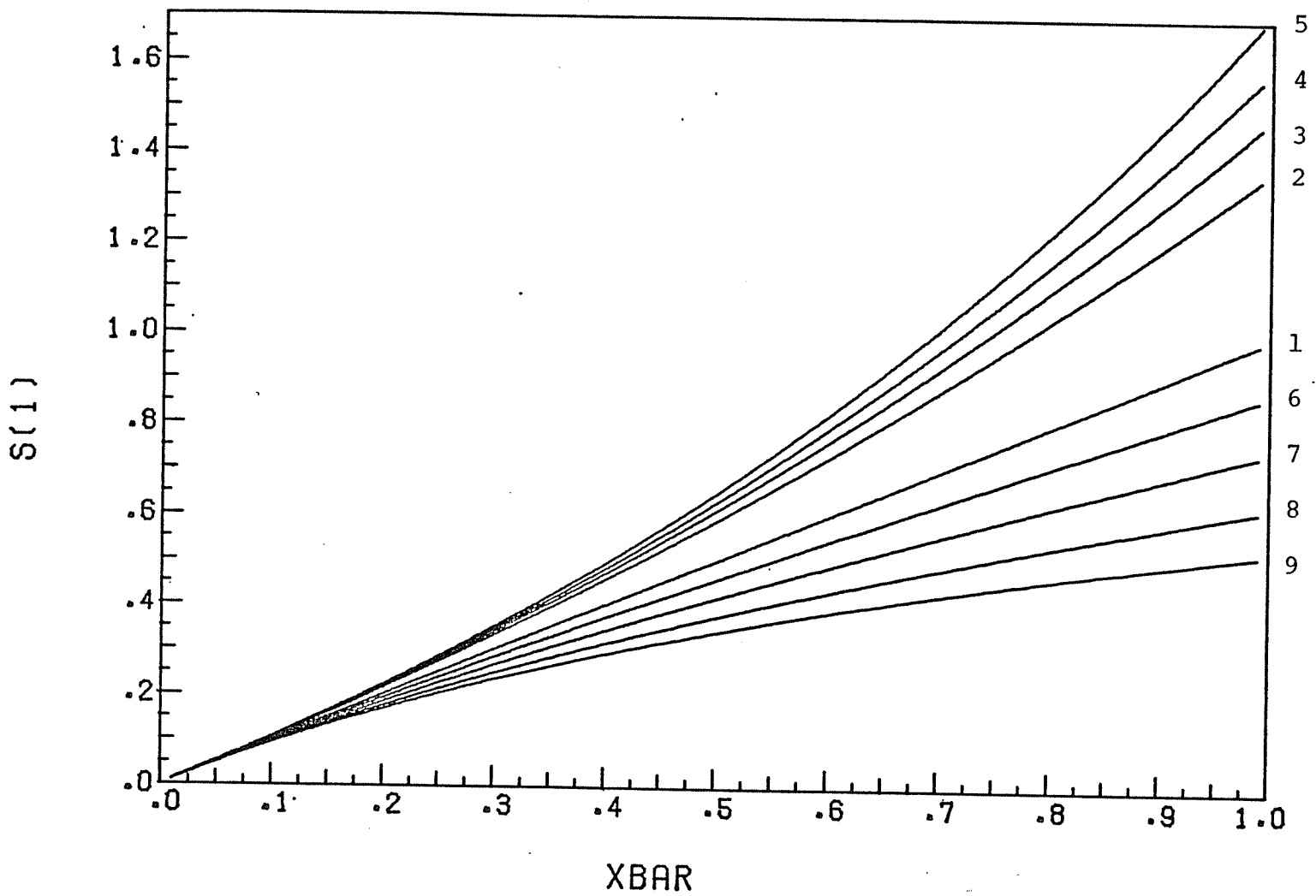# SERVICE FUNCTION S(1) VERSUS MEAN SERVICE TIME



Figure 4.1

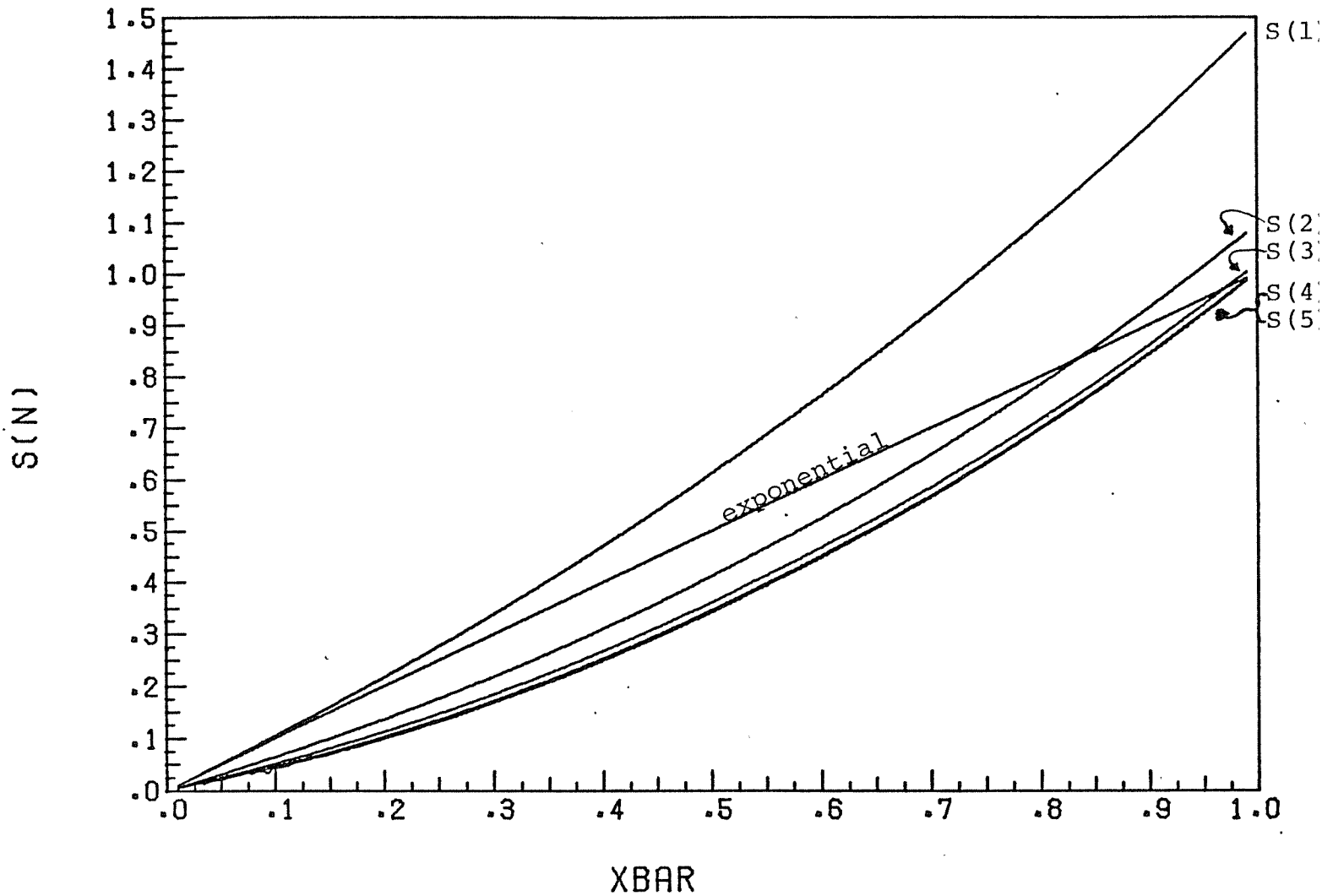| Curve Number | Distribution |
|---|---|
| 1 | Exponential |
| 2 | Erlang-3 |
| 3 | Erlang-5 |
| 4 | Erlang-10 |
| 5 | Constant |
| 6 | Hyperexponential, $\alpha=0.5$, $CV_2^2=1.5$ |
| 7 | Hyperexponential, $\alpha=0.5$, $CV_2^2=2.0$ |
| 8 | Hyperexponential, $\alpha=0.5$, $CV_2^2=2.5$ |
| 9 | Hyperexponential, $\alpha=0.5$, $CV_2^2=2.9$ |

SERVICE FUNCTIONS S(N) FOR M/G/1 WITH ERLANG-R SERVICE TIMES



Figure 4.2

Values of service functions S(1), . . . ,S(5) versus
mean service time for Erlang-5 distribution.

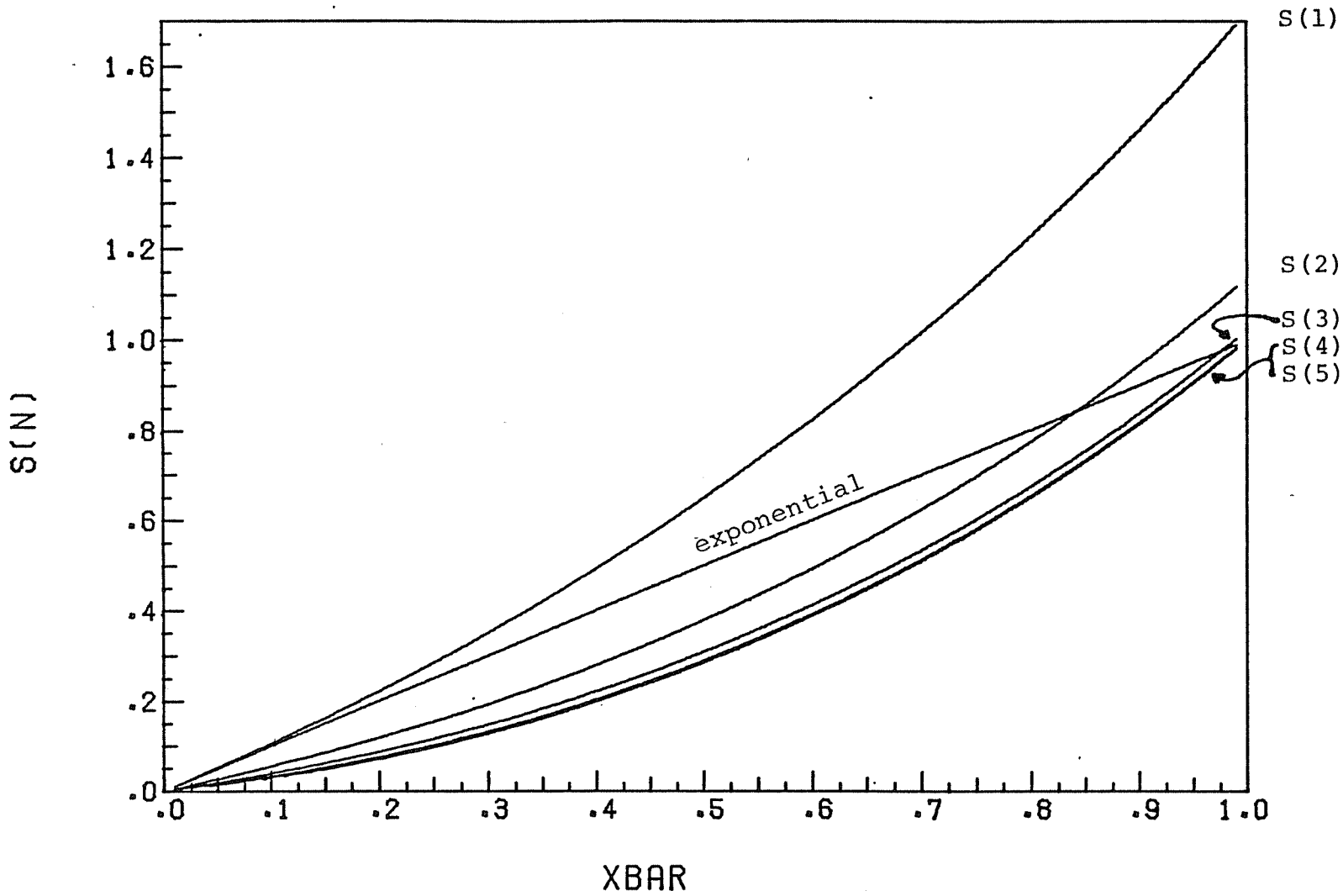# SERVICE FUNCTIONS S(N) FOR M/G/1 WITH CONSTANT SERVICE TIMES



Figure 4.3

Values of service functions S(1),. . .,S(5) versus mean
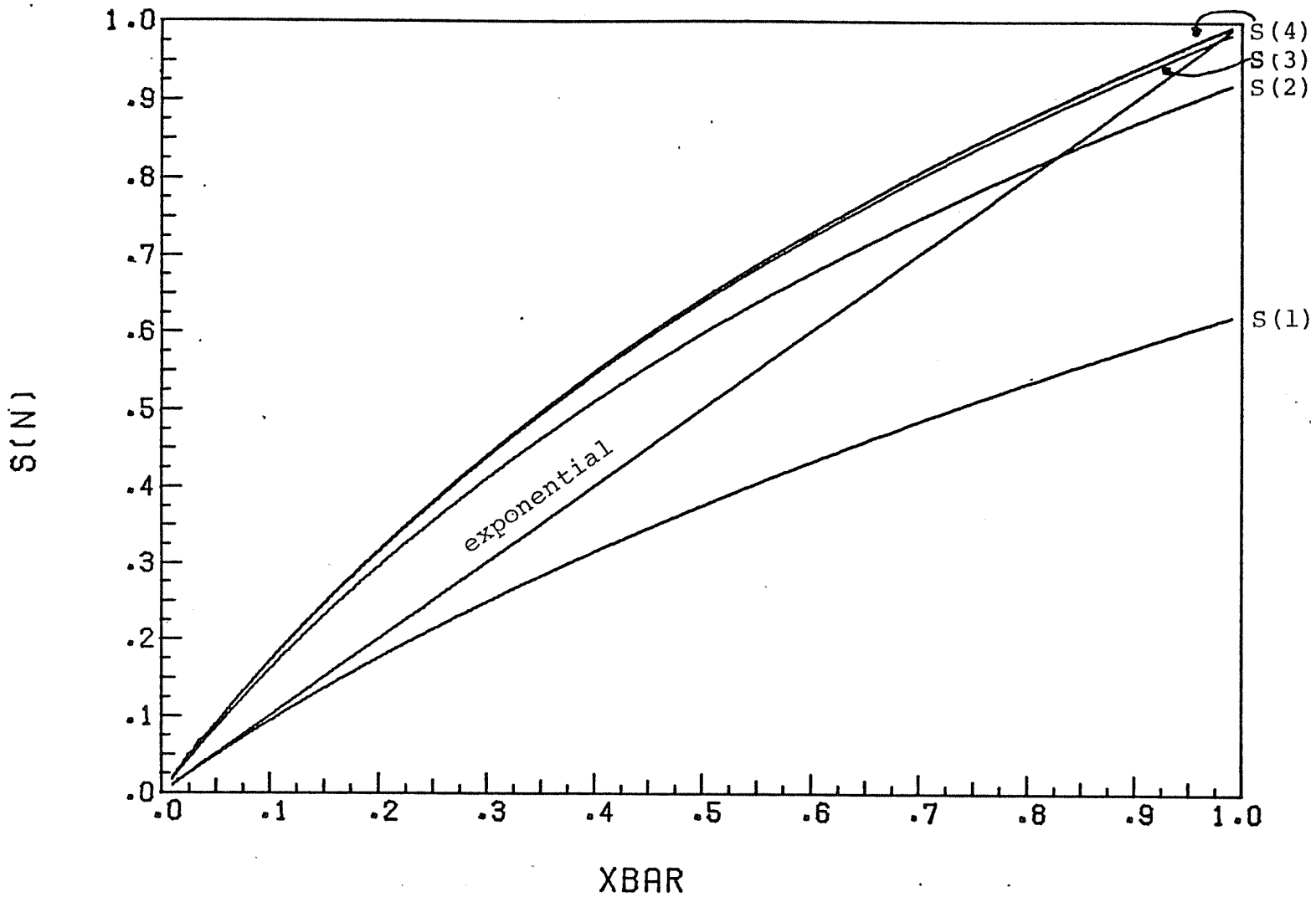service time for constant service times

SERVICE FUNCTIONS S(N) FOR M/G/1 WITH HYPEREXPONENTIAL SERVICE    TIMES



Figure 4.4

Values of service functions S(1), . . .,S(5) versus
mean service time for two stage hyperexponential
distribution, $\alpha$=0.5, $CV^2$=2.5.

$S(n)$ versus $\bar{x}$ curve for the exponential case as a reference.) The first observation about these graphs is that $S(n)$ rapidly approaches a limiting value as n increases. Apparently, the tail of the distribution of number in system is approximately geometric for large values of n. A second observation is that the limiting $S(n)$ versus $\bar{x}$ curve always lies on the other side of the exponential case curve from the $S(1)$ curve. Thus we cannot extend the statement of the last paragraph to higher values of n. Exactly how to estimate these values of $S(n,t)$ depends on the current value of $\bar{x}$ and the other parameters of the service distribution. Third, from these examples it appears that if $CV^2$ < 1, then the $S(n)$ versus $\bar{x}$ curve is convex upward; if $CV^2$ > 1 then the curve is convex downward.

In brief, the performance problem we have posed does not appear to be solveable without making additional assumptions about the distribution of service time. With such assumptions, graphs like Figure 4.2 could be drawn, and new values of $S(n,t)$ could be determined from observed values in a methodical way. Thus instead of assuming that the service time distribution is of a certain form, fitting the observed service times to that form, and then applying the M/G/1 solution with G replaced by the fitted distribution, we could instead estimate the new values of $S(n,t)$ and solve directly for $P(n,t)$. The relative merits of these two approaches have yet to be explored.

## 5. CONCLUDING REMARKS

The subject of this paper has been an amalgam of two almost antithetical disciplines: operational and stochastic analysis of queueing systems. It has not been the intent of this paper to criticise or praise either of these approaches. Each is useful in its own way, when applied to certain types of problems, or when used to best advantage by certain groups of people.

The point of this paper is that there is not that much difference between the two approaches, as the equivalence proofs of Section 3 have shown. Indeed, one of the most useful results of this paper might be the graphs of Section 4; results which come from an interplay of both operational and stochastic concepts.

BIBLIOGRAPHY

[1]    Burke, P. J., "The output of a queueing system,"
       Operations Research, 4, 699-704 (1956).

[2]    Buzen, J. P., "Fundamental operational laws of computer
       system performance," Acta Informatica, 7, 2 (1976),
       pp. 167-182.

[3]    Buzen, J. P., "Operational analysis: an alternative to
       stochastic modeling," Proceedings of the International
       Concference on the Performance of Computer Installations,
       D. Ferrari (ed.), North-Holland Publishing Co., Amsterdam,
       The Netherlands (1978), pp. 175-194.

[4]    Denning, P. J. and J. P. Buzen, "The operational analysis
       of queueing network models," ACM Computing Surveys, 10, 3
       (September 1978), pp. 226-261.

[5]    Jackson, J. R., "Networks of waiting lines," Operations
       Research, 5, 518-521 (1957).

[6]    Kleinrock, L.  Queueing Systems Volume I: Theory, John
       Wiley and Sons, New York (1975).

[7]    Mesztenyi, C. K., "FORMAL - a formula manipulation
       language," Computer note CN-1.1, University of Maryland
       Computer Science Center, College Park, Maryland (October
       1971).

[8]    Ross, S. M., Applied Probability Models with Optimization
       Applications, Holden-Day, San Francisco (1970).