

ON WAITING FOR SIMULTANEOUS ACCESS TO TWO RESOURCES

Michael L. Honig and Teunis J. Ott
Bell Communications Research
Morristown, New Jersey 07960

ABSTRACT

This paper studies the $M/G/1$ queue where a special (test) customer can get service only if he has simultaneous access to the server and a second resource. All other customers only need access to the server. The second resource becomes available after an exponentially distributed amount of time. The ordinary customers are served according to the FIFO discipline. The test customer has the freedom to leave his place in the queue at any time and join the end of the queue. If he reaches the server before the second resource becomes available, he then must return to the back of the queue.

We derive the waiting time distribution of the test customer given that he always maintains his position in the queue until he reaches the server. A number of conditions are given under which this "move-along" policy is optimal, i.e., minimizes the test customer's mean delay until service. These conditions depend on the amount of information and freedom of action available to the test customer.

1. PROBLEM DESCRIPTION, AND OUTLINE

The problem studied in this paper is derived from the scheduling problems which occur in a service system with multiple resources, where a customer can get service only when all the resources it needs are simultaneously available. The prototype multi-resource service system can be thought of as a multi-processor computer system or a database system with locking mechanisms for integrity protection. The simplified problem considered here is described in [Gopinath, 1984], and has become known as the "Waiting for Godot" problem. This simplification still captures some of the effects of waiting for simultaneous availability of multiple resources, and, as is seen in this paper, is reasonably tractable.

In the model studied there is one special customer, called the test customer, who is waiting for simultaneous availability of two resources: the server and an "extraneous" resource. All other customers only need the server. The other customers arrive according to a Poisson process with intensity λ and have service times (with the server) which are i.i.d. random variables with distribution function $F(\cdot)$, Laplace Stieltjes Transform (LST) $\phi(\cdot)$ and expected value m . For these other customers the service discipline is First In, First Out (FIFO). The time until the extraneous resource becomes available is an independent, exponentially distributed random variable with parameter α (expected value α^{-1}). At the time the extraneous resource becomes available we say that the event E occurs, and once E has occurred the extraneous resource remains constantly available.

At any time the test customer has two options. He can either maintain his place in the queue, or he can voluntarily leave his place in the queue and move to the end of the queue. Whenever he reaches the server after event E has occurred, service starts. If the test customer reaches the server before event E has occurred, however, he must go back to the end of the queue. An interesting question, which will be partially answered, is whether and why it is ever profitable for the test customer to use the option of moving to the end of the queue without being forced to do so.

For a further study of waiting for simultaneous access to multiple resources it will be necessary to consider situations where there is competition for the external resources, or where many other customers are also waiting for their own (possibly different and independent) resources, or both.

As stated before, the test customer has the freedom to give up his place in the queue and go back to the end of the line even when he is not facing the server. The "move-along policy" is the policy where the test customer never uses that option. With that policy he maintains his place in the queue until he reaches the server. At that point either service starts (if in the meantime event E has occurred), or he goes back to the end of the line.

Our first results, stated in Section 2, are expressions for the distribution (in fact its LST) and the expected value of the time until service starts for the test customer, given that he uses the move-along policy, and given that at time $t=0$ there is a random variable X representing the total amount of work in front of him in the queue, and a random variable Y representing the total amount of work behind him in the queue. These expressions of course are in terms of the joint distribution of X and Y . Other results, also stated in Section 2, describe under what circumstances the move-along policy is better than competing policies. The competing policies depend on the information and degree of freedom available to the test customer. The test customer always knows λ and $F(\cdot)$. A policy is said to be optimal if it minimizes the expected value of the test customer's delay until service, starting from an arbitrary state. The situations considered are:

1.1 Complete Information and Freedom

In these policies the test customer always knows the (remaining) service times for all customers in the system. He therefore exactly knows the total amount of work in front of him in the queue (t_f) and the total amount of work behind him in the queue (t_b). At any time he can decide (based on t_f and t_b) to give up his place and join the end of the line. It will be proved that as long as the move along policy is the best among all such "complete information and freedom" policies.

1.2 Partial Information and Complete Freedom

In these policies the test customer knows, at any time, the (remaining) service times of the customers in front of him, but only the number of customers behind him. As a result he knows the total amount t_f of work in front of him and the total number j of customers behind him. At any point in time he can decide (based on t_f and j) to give up his place and join the end of the queue. It will be proved that if then the move-along policy is the best among all such "partial information and complete freedom" policies.

1.3 Minimal Information and Complete Freedom

In these policies the test customer only knows the numbers n_i and n_j of customers in front of him, respectively behind him, and for the customer currently being served he also knows the elapsed service time t_{au} . At any point in time he can decide (based on n_i , n_j , and t_{au}) to give up his place and join the end of the queue. It is clear that if (1.2) holds then the move along policy is the best of all "minimal information and complete freedom" policies. It is possible (because of the smaller amount of information available) to replace (1.2) with a weaker (larger) upper bound for λ (see (2.26)).

1.4 Minimal Information and Limited Freedom

In these policies the test customer only knows the numbers n_i and n_j of customers in front of him and behind him, and is allowed to leave his place and go to the end of the queue only at service completion epochs. This situation is called the "discrete-epoch" situation. It will be proved that if then the move along policy is the best among all discrete-epoch policies.

It is shown in [Li, 1987] that the move-along policy is the best among all discrete-epoch policies if $\lambda \leq 1 / \bar{m} \phi(\alpha)$ where \bar{m} is the expected service time of a customer given that event E did not occur during (or before) his service. Since the two results are equivalent, although the proof given by Li is different from that given in Section 3.

In [Honig, 1987] it is shown for deterministic service that there exists a threshold $\lambda_{\text{sub } 0}$, which depends on α and m , such that for $\lambda > \lambda_{\text{sub } 0}$, the move along policy is not optimal. In Section 3 it will be proved that in the case of a general service time distribution there exists a threshold $\lambda_{\text{sub } 0}^*$ such that for $\lambda > \lambda_{\text{sub } 0}^*$ the move along policy is not the optimal discrete-epoch policy. This result is easily explained by the following simple intuitive argument. Assume that λ is quite large (e.g., 200), $m \approx 1$, and that $\alpha \approx 0.1$, so that the expected time until E occurs is on the order of 10 service times. Suppose that initially there is one person ahead of the test customer. While the test customer is waiting for the customer ahead of him to finish service, new arrivals are rapidly joining the queue behind him. Consequently, if the test customer chooses to maintain his position until he reaches the server, he will most likely have to wait in back of all of the (approximately 200) new arrivals. Alternatively, if the test customer decides to

join the back of the queue after, say, the first 100 new arrivals, which most likely occurs during the service epoch, he will almost certainly reduce his delay.

The previous results suggest the following conjecture:

Conjecture: Given α and $F(\cdot)$ there exist critical levels λ_k^* , $1 \leq k \leq 4$, $k=1$ for "complete information, complete freedom", $k=2$ for "partial information, complete freedom", $k=3$ for "minimal information and complete freedom", and $k=4$ for "minimal information and limited freedom", with such that in situation k ($k=1,2,3,4$) the move-along policy is the best of all situation k policies if and only if $\lambda_k^* \geq 1 / \phi(\alpha)$. From (1.3) it is apparent that $\lambda_4^* \geq 1 / \phi(\alpha)$.

In Section 2 we give, without proofs, the main results of this paper. The proofs of these results are given in Section 3. Section 4 discusses a related problem, where the test customer can decide to wait outside the queue before joining the end of the line. Finally, some specific service distributions are considered in Section 5, and Section 6 discusses some other related problems.

2. NOTATION AND THE MAIN RESULTS

At time zero the test customer has a total amount of work X in front of him and a total amount of work Y behind him in the queue. The joint distribution of X and Y is given by $F(x,y)$ and the marginal distributions of X and Y are denoted as $F_X(x)$ and $F_Y(y)$.

Given any distribution $R(\cdot)$, the LST of R is denoted by $\psi_R(\cdot)$. In particular,

The distribution of T , the time until the test customer starts service, of course depends on $F(\cdot, \cdot)$ and on the policy used. $\psi_T(\cdot)$ denotes the LST of T given that the move along policy is used:

If X and Y are independent we denote this as $\psi_T(s)$. By a (hopefully not confusing) abuse of notation we define

and

(where $F^{*i}(\cdot)$ denotes the i -fold convolution of $F(\cdot)$). Finally, we define:

(where $\eta_{i,0}(s)$ is defined as in (2.6)), and we write $\eta_{G,H}(s)$ for $\eta_{G,H}(s)$ when $\Pr\{Y=0\}=1$ (no customers are behind the test customer).

2.1 Waiting Time Distribution and Move-Along Mean Delay

Theorems 1 and 2, which follow, give explicit expressions for $\eta_i(s)$ and $\eta_P(s)$ in terms of $\lambda, \mu, \alpha, \phi(s)$, and $\psi_P(s_1, s_2)$, and will be proved in Section 3. The basic idea of the proofs is as follows: If there is a deterministic amount t_f of work in front of the test customer, and a random amount Y of work behind him, then the test customer first waits for an amount of time t_f . If by that time event E has occurred, $T \leq t_f$. If not, at time t_f the test customer is at the end of the queue with an amount of work in front of him equal to Y plus the service times of all customers who arrived in the time interval $[0, t_f]$. Averaging over the distribution of t_f expresses $\eta_{G,H}(s)$ in terms of the sequence $\{\eta_{\{H^{*k}\}}(s)\}_{k=0}^{\infty}$.

Choosing $Y=0$ and $G \sim F^{*i}$ expresses $\eta_i(s)$ in terms of $\{\eta_{\{k\}}(s)\}_{k=0}^{\infty}$, and makes it possible to compute $\eta_i(s)$ and thus prove Theorem 1. Theorem 2 is then proved by repeated use of the same idea. Section 3 not only contains the proofs of Theorems 1 and 2, but also a number of intermediate results such as expressions for $\eta_{G,H}(s)$ and $\eta_{\{t_f, t_b\}}(s)$. Some readers may prefer to read Section 3 before reading the remainder of this section.

The results in this section are expressed in terms of the sequences $\{y_k(s)\}$, $\{y_{tk}(s)\}$, $\{x_k(s)\}$, and $\{x_{tk}(s)\}$, which are defined as:

For $s \geq 0$ it is easily shown that

and

where the 'less than or equal to' signs are equalities if and only if $s=0$. The sequences $y_k(s)$ and $x_k(s)$ are shown graphically in Figure 1. $y_{\infty}(s)$ and $x_{\infty}(s)$ satisfy

where $\beta(s)$ is the LST of the lengths of the busy periods in the $M/G/1$ queue with $\lambda, F(\cdot), \phi(\cdot)$. In section 3 a number of results related to (2.9) and (2.10) will be derived which show that the infinite

series in Theorems 1, 2, and 3 below converge uniformly for $\text{Re}(s) \geq 0$.

Theorem 1:

where

Theorem 2:

where $x_{-1}(s) = 0$.

The expected value of the time until the test customer starts service, assuming the move-along policy is adopted and that the joint distribution of work in front of and in back of the observer is given by (2.1), is denoted as t_{gh} . In analogy with the notation introduced before, t_{gh} denotes the expected delay when $P(x,y) = G(x)H(y)$, t_{fb} denotes the expected delay given that $X = t_f$ and $Y = t_b$, t_{bj} denotes the expected delay when $G(t) = F^{*i}(t)$ and $H(t) = F^{*j}(t)$, and t_{bi} denotes the expected delay when $G(t) = F^{*i}(t)$ and $\Pr\{Y = 0\} = 1$. Taking the derivative of the expression in (2.13) gives the next Theorem.

Theorem 3:

where

is the mean delay given that $X = Y = 0$. X_{ba} is the expected value of X , and $y_k = y_k(0)$, where $y_k(s)$ is defined in (2.8).

If $G(t) = F^{*i}(t)$ and $H(t) = F^{*j}(t)$, i.e., there are i customers ahead of the test customer, none of whom have received any service yet, and j customers in back of the test customer, (2.15) becomes

where $x_k = x_k(0)$, and $x_k(s)$ is defined in (2.8).

2.2 Conditions for Move-Along Optimality

A policy is a sequence of actions which the test customer may take, and in general each action

depends on the entire history of states visited. The only allowable action the test customer may take is to give up his current position in the queue, and move to the back of the queue. A policy is said to be optimal if it minimizes (over the entire class of allowable policies) the test customer's expected delay until the start of service, given some arbitrary initial state.

Let $D(X, Y; \pi)$ denote the expected time until the start of service for the test customer given that initially the amounts of work in front of him and behind him are X , respectively Y , and given that consistently policy π is used. Let π_{MA} denote the move along policy. The maximum principle from dynamic programming suggests that π_{MA} is optimal if and only if

for all (nonnegative) random variables X and Y . This is a consequence of well known results in Markov Decision theory. Intuitively, it can be seen as follows: suppose there is a time L such that for $t > L$ the move-along policy will be used. The problem is to find the optimal policy for $t \leq L$. But this now is a finite horizon dynamic programming problem, and (2.18) implies that the move-along policy is always optimal. By choosing L sufficiently large, the probability that $T > L$, where T is the time until the test customer starts service, can be made arbitrarily small. This implies that the "end effect" of what happens after time L becomes irrelevant, and that the move-along policy is optimal. This argument can be made rigorous by observing that in the worst case ($\rho > 1$), the expected amount of work in the system at time t grows linearly with t , while the probability that the extraneous resource is not yet available at time t is $e^{-\alpha t}$, and $\lim_{L \rightarrow \infty} \int_L^{\infty} \int_0^L e^{-\alpha t} dt = 0$ for any constant c .

If $X = tf$ and $Y = tbb$, then (2.18) becomes

For the case $G(t) = F^{*i}(t)$ and $H(t) = F^{*j}(t)$, the condition (2.18) becomes

for all positive integers i and j . It is therefore of interest to study $t^{b_i + b_j} e^{-\lambda t}$ as a function of λ , α , and m . The next corollary is obtained from Theorem 3.

Corollary 1:

Averaging over tf and tbb for the case where X and Y are independent gives

If $G(t) = F^{\sup *i}(t)$ and $H(t) = F^{\sup *j}(t)$, (2.22) becomes

The expressions (2.21)-(2.23) are used to prove the following four Theorems, which imply the results stated in Sections 1.1 - 1.4. The condition on λ in each case ensures that every term in the corresponding sum in (2.21)-(2.23) is positive. The following theorems therefore give sufficient, but not necessary, conditions for the move-along policy to be optimal.

Theorem 4: If $\lambda \leq 1$, the condition $b \leq b + b$ holds for all $t, \tau \geq 0$.

Theorem 4 applies to the "complete information and freedom" situation and establishes the statement made in subsection 1.1.

Theorem 5: If $\lambda \leq \alpha / \{1 - \phi(\alpha)\}$, the condition $b \leq b + G^*H$ holds for all G, H such that $X = t$ and $H(t) = F^{\sup *j}(t)$, for any $t \geq 0, j \geq 0$.

Theorem 5 establishes the statement made in subsection (1.2). To prove the statement in subsection (1.3) some more notation is needed. Let

(where superscript "+" denotes limit from the right and superscript "-" denotes limit from the left), i.e., $F_{\tau}(t)$ is the probability distribution of the remaining service time given that the customer has been in service τ time units, and let

Suppose that initially there are i customers ahead of the test customer, and j customers behind the test customer, and that the elapsed time since the customer at the front of the queue started service is τ . For this case $\phi_i(s) = \phi^{\sup i-1}(s)$ and $\phi_{ih}(s) = \phi^{\sup j}(s)$. The next theorem gives a weaker condition on λ than that given in Theorem 5.

Theorem 6. If there are initially i customers in front of the test customer, and j customers in back of the test customer, then the condition $b \leq b + G^*H$, where $\phi_{ih}(s) = \phi^{\sup j}(s)$, $\phi_{ig}(s) = \phi^{\sup i-1}(s)$ and $\phi_{it}(s)$ is defined in (2.25), holds for all i, j , and τ if

This condition on λ is weaker than the condition stated in Theorem 5 since $\{1 - \phi(\alpha)\}^i$ over $\{\sum_{i=1}^n \phi(\alpha)\}^i \geq \alpha$ for any distribution $G(t)$ over $[0, \infty)$. If τ , the elapsed time since the customer at the front of the queue started service, is taken to be zero, then $\phi(s) = \phi(\sup_i(s))$, and the upper bound on λ in (2.26) can be evaluated to give:

Corollary 2: If $\lambda \leq 1 / \sum_{i=1}^n \phi(\alpha)$, the condition $t_{ij} < t_{i+j}$ holds for all positive i and j .

The move-along policy is therefore optimal for the discrete-epoch problem if $\lambda \leq 1 / \sum_{i=1}^n \phi(\alpha)$.

It will be shown in the next section that

so that Theorems 4, 5, 6, and Corollary 2 give progressively weaker conditions on λ corresponding to less information or freedom available to the test customer.

We also have the following conjecture about the expression in (2.26):

Conjecture:

Namely, we believe that for $\tau \geq 0$ fixed the expression $\{1 - \phi(\alpha)\}^i / \sum_{i=1}^n \phi(\alpha)$ is increasing in i . We have not succeeded in proving this. Neither have we succeeded in proving the even stronger statement (which probably is not always true) that $\{1 - \phi(\alpha)\}^i / \sum_{i=1}^n \phi(\alpha)$ is increasing if $G(\cdot)$ is stochastically increasing.

To show that the move-along policy is not optimal for a given λ , m , and α , it suffices to find a particular i and j such that $t_{ij} > t_{i+j}$. As an example, if $i=j=1$, then (2.23) becomes

If $\lambda > 1 / \sum_{i=1}^n \phi(\alpha)$, then the first term in the sum in (2.29) will be negative.

However, it is not true that all of the remaining terms become negative for large enough λ . In particular,

so that $\lambda \leq 1$ for large enough λ , and $\lambda \leq \lambda \sum_{i=1}^n \phi(\alpha) \dots$

$> \lambda \sum_{i=1}^n \phi(\alpha)$ (see Lemma 2 in Section 3). It therefore is conceivable that the sum (2.29) is

positive for all λ . Nevertheless, the following theorem states that in fact for any α , the move-along

policy is not optimal if λ is large enough.

Theorem 7: Given any α and β there exist two numbers λ_0 and λ_1 such that $t_{i+j} < t_{ij}$ for $\lambda > \lambda_0$ and $t_{i+j} > t_{ij}$ for $\lambda < \lambda_1$.

The previous results suggest that for any of the situations considered, there exists a threshold, λ_0 , such that the move-along policy is optimal *if and only if* $\lambda \leq \lambda_0$. To prove this result one must show that if for some $\lambda = \lambda'$, $t_{i+j} \leq t_{ij} + b$ for all positive t_i and t_j , then it must also be true for all $\lambda < \lambda'$. (Alternatively, one could show that if for specific t_i and t_j , $t_{i+j} \geq t_{ij} + b$ for $\lambda = \lambda'$, then $t_{i+j} \geq t_{ij} + b$ for any $\lambda > \lambda'$.) This appears to be difficult, and it is as yet undetermined whether or not this is true.

3. PROOFS

The sequences $y_k(s)$, $y_{tk}(s)$, $x_k(s)$, $x_{tk}(s)$ are based on the map f_s defined by

In particular,

where $f_s^{(k)}$ is the k times iterated map.

If $\text{Re}(s) \geq 0$, then f_s maps the half plane $\text{Re}(z) \geq 0$ into the half plane $\text{Re}(z) \geq \alpha + \text{Re}(s)$. If $\rho = \lambda m < 1$ then, for $\text{Re}(s) \geq 0$, f_s is a contraction map on $\text{Re}(z) \geq 0$:

where the fact that $|\text{arg}(z)| \leq |\text{arg}[\text{Re}(z)]|$, which is decreasing in $\text{Re}(z)$, has been used.

Suppose now that $\rho = \lambda m \geq 1$. Since for real $s \geq 0$, $\phi(s)$ is decreasing in s , $y_k(s)$ and $y_{tk}(s)$ in (2.8) are increasing in s , and $x_k(s)$, $x_{tk}(s)$ are decreasing in s . Hence, for $\text{Re}(s) \geq 0$, f_s maps the half plane $\text{Re}(z) \geq y_k(0)$ into the half plane $\text{Re}(z) \geq y_{k+1}(0)$. Since (see Figure 1) $\lambda |\text{arg}[\text{Re}(z)]| < 1$, there exist a ρ^* , $0 < \rho^* < 1$, and a k_0 , with the property that if $k \geq k_0$ then

for all $s \geq 0$, z_1, z_2 with $\text{Re}(s) \geq 0$, $\text{Re}(z_i) \geq 0$. For ρ^* we could choose

and for k_0 we choose

(3.6) uses the fact that for $s \geq 0$, $Y_k(s)$ is increasing in both k and s (see Figure 1).

As a straightforward application of (3.2)-(3.4) we can obtain relations such as

where ρ , ρ^* , k_0 , do not depend on s . The series in Theorems 1, 2, and 3 therefore converge uniformly on compact subsets of $\{s \in \mathbb{R} : s \geq 0\}$. In the remainder of this section we will mostly disregard convergence issues.

The following lemma is needed to prove Theorems 1 and 2.

Lemma 1.

Proof: Suppose that initially the amount of work in front of the test customer is exactly t . If by the time this work has been done by the server event E has occurred, then $T = t$. Otherwise, the test customer first waits time t , and then becomes the last customer in a queue with an amount of work in front of him equal to Y , the initial amount of work behind him, plus the work required by customers who arrived in the meantime. (3.8) is a formal statement of this observation, and allows the amount of work in front of the test customer to be random, as long as it is independent of the amount of work behind the test customer. \square

Remark. By defining

we can rewrite (3.8) as

This result makes it easy to prove Theorems 1 and 2. Theorem 1 is obtained by choosing $Y = 0$ and $G = F$ for $x \geq 0$ and $G = 0$ for $x < 0$.

Proof of Theorem 1: Assume that at time zero there is no work behind the test customer, and there are i customers in front of him, none of whom have received any service yet, so that $Y = 0$ and $G = F$ for $x \geq 0$. From (3.10) we have for $i \geq 1$:

In addition to (3.11) we also have the boundary condition

Namely, if the test customer is alone in the system, then the first event to occur is either E , or the arrival of an

ordinary customer.

To prove Theorem 1 we must show that (2.11), (2.12) are the unique solution to (3.11), (3.12). First we substitute (3.12) into (3.11), and obtain for $i \geq 1$,

This can be rewritten as

where

>From (3.9) and (3.15) we see that

and for $i \geq 1, \operatorname{Re}(s) \geq 0$,

Hence, the solution to (3.14) is unique. Moreover, the solution to (3.14) can be obtained by choosing $\eta_{i \sup}(s)$ arbitrarily, and then iterating the contraction map (in $\|\cdot\|_{\sup}$)

Clearly, $\lim_{n \rightarrow \infty} \eta_{i \sup}(n)(s) = \eta_{i \sup}(s)$. This is one of the ways (2.11), (2.12) can be derived. Our original derivation was a form of "clever trying". Since the solution is available, it suffices to verify that (2.11), (2.12) indeed form a solution to (3.11), (3.12). It is easily seen that for any distribution G on $[0, \infty)$, and with $p_{k \sup}(G)(s)$ defined as in (3.9), and $x_m(s)$ and $y_m(s)$ as in (2.8):

It now is easy to verify that (2.11), (2.12) indeed is a solution to (3.11) - (3.13). This completes the proof of Theorem 1. \square

Proof of Theorem 2. The proof of Theorem 2 consists of the following steps: (1) Use Theorem 1 and (3.9), (3.10) to obtain an expression for $\eta_G(s)$. (2) Use this expression for $\eta_G(s)$ and (3.9), (3.10) to get an expression for $\eta_{G,H}(s)$ for general G, H . This gives $\eta_{\{tf, tbb\}}(s)$ as a special case, which immediately gives $\eta_P(s)$.

>From (3.9), (3.10), and Theorem 1,

which is the promised expression for $\eta(s)$. Substituting for $\eta_{H^*F \supset k}(s)$ in (3.10), and noting that $\psi_{H^*F \supset k}(s) \sim \phi_{H^*F \supset k}(s)$, gives

For the case $\Pr\{X \sim 1\} = 1$ and $\Pr\{Y \sim 1\} = 1$, $\phi_{H^*F \supset k}(s) \sim e^{-s}$, and $\psi_{H^*F \supset k}(s) \sim e^{-s}$, and (3.21) specializes to

where $x_{-1}(s) = 0$. Averaging over t and b gives (2.13). \square

Proof of Theorem 3: To compute the derivative of the expression (2.13), it is necessary to compute

where $y_k \sim y_k(0)$. From (2.13),

Using the fact that $y_k(0) \sim y_k(0) \sim y_k$,

Combining (2.14), (3.24), and (3.25), and noting that $\phi_{H^*F \supset k}(0) \sim X$, gives (2.15).

From (2.12),

\square

As a side remark, we outline an alternative derivation of (2.22). In analogy with the recurrence relation (2.14), the following recurrence relation can be derived for the mean delay $t_{H^*F \supset k}$,

where $t_{H^*F \supset k}$ is the mean delay until service given that the amount of work in front of the test

customer is the random variable Y plus k service times and $p_k \supset (G) \sim p_k \supset (G) (0)$. If $Y \sim 0$, then

where $t_k \sim t_{\{F \supset *k\}}$ and $p_k \supset (i) \sim p_k \supset \{(F \supset \{star i\})\}$. The boundary condition, in analogy with (3.12), is

The contraction mapping technique, which can be used to obtain the solution to (3.9), (3.10), can also be applied to (3.28), (3.29), thereby giving explicit expressions for t_{bij} and t_{bgh} .

Before proving Theorems 4-7, some basic properties of the sequences y_k and x_k , which follow directly from the discussion at the beginning of this section, are stated. The sequences y_k and x_k are illustrated graphically in Figure 1.

Lemma 2. For real $s \geq 0$, the sequence $y_k(s)$, defined by (2.8), increases monotonically and converges to $y_{\infty}(s)$, and the sequences $x_k(s)$ and $|\phi(y_k(s))|$ each decrease monotonically and converge to $x_{\infty}(s)$ and $|\phi'(y_{\infty}(s))|$, respectively. Also, $y_{\infty}(s) < s + \lambda \mu$.

Proof of Corollary 2: From (2.15), for the case $\phi(s_1, s_2) = e^{-s_1 t_f + s_2 t_b}$,

so that

□

The next lemma proves Theorem 4 by showing that if $\lambda \mu \leq 1$, then all terms in the series (2.21) are nonnegative.

Lemma 3: The function

is greater than or equal to zero for $t_f \geq 0$, $t_b \geq 0$, and $\lambda |\phi(y_k)| \leq 1$.

Proof: Observe that

Also,

Lemma 2 states that $y_k \geq y_{k+1}$, therefore the derivative (3.34) is positive, which implies that $f(y_k, y_{k+1})$ is positive, if $\lambda \leq \phi(y_k)$. Every term in the sum (3.31) is guaranteed to be positive if $\lambda \leq \phi(y_0)$. \square

Notice that if $\lambda = 1$, then the sum (3.31) will be strictly positive, since x_k is strictly decreasing with k . It therefore seems likely that there exists a threshold $\lambda_0 > 1/m$, such that if $\lambda \leq \lambda_0$ the sum (3.31) is positive for all positive t and b .

Proof of Theorem 5: For the case $\phi(s) = e^{-st}$ and $\psi(s) = \phi^j(s)$, (2.22) becomes

where

Now,

and

which is positive for $k \geq 1$ if $1 - \lambda \phi(\alpha) \geq 0$. For $k=0$,

which is greater than or equal to zero if

It is easily shown that $\{e^{-\alpha t}\}^j$ over $t \geq \alpha$, with equality as t approaches zero. Furthermore, it can be shown that the function j over $\{1 - \lambda \phi^j\}$, where $0 \leq \lambda \leq 1$ and $j \geq 1$ is an integer, increases with j , so that if

then $f(y_k, y_{k+1}) \geq 0$ for all $k \geq 0$. Referring to Figure 2, it is clear that

which proves the result. \square

Proof of Theorem 6: Substituting $\psi(s) = \phi^j(s)$ into (2.22) gives

where

As in the preceding proofs, it is easily shown that for $k \geq 1$, $\partial f / \partial y_k \leq 0$, and hence $f(y_k, y_{k+1}) \geq 0$, if $\lambda \leq \phi(\alpha)$. For $k=0$,

which is greater than or equal to zero if

for all possible $\phi_i(s)$ and $\phi_j \geq 1$. Substituting $\phi_i(s) = \phi_{i-1}(s) \phi_i(s)$, and noting that the right side of (3.45) is minimized by setting $\phi_j = 1$, and is equal to $1 / \sum_{i=1}^n \phi_i(\alpha)$ for $\phi_i = 1$ and $\tau = 0$, gives Theorem 6. \square

Proof of Corollary 2: For this case $\phi_i(s) = \phi_{i-1}(s)$, and we show that

increases with ϕ_i . Assume that this is false. Then

for some ϕ_i . This implies that

The left side assumes its maximum value, however, when $\phi_i(\alpha) = 1$, therefore this cannot be true.

Consequently,

and Corollary 2 follows from Theorem 6. \square

We remark that Corollary 2 can also be proved directly from (2.23). In particular, it is easily shown that

for $\phi_i \geq 1$, $\phi_j \geq 1$, and $\lambda \sum_{i=1}^n \phi_i(\alpha) \leq 1$.

Proof of Theorem 7: From (2.23) and Lemma 2

assuming λ is large enough so that $\lambda \sum_{i=1}^n \phi_i(\alpha) < 1$. For fixed α it can be shown that the last term goes to zero faster than $O(1/\lambda)$. Therefore, for large λ ,

which can be negative only if

or

Since $\phi_i < 1$, ϕ_i can be selected large enough so that (3.51) is true for any ϕ_j , and if λ is greater than some threshold λ_0 , then from (3.51), $t_{i+j} - t_{ij} < 0$. \square

4. A MODIFIED PROBLEM

So far we have assumed that the test customer can always give up his place in the queue, and move to

the back of the queue. It has been shown that if λ is large enough, using this option will decrease the test customer's expected delay until service. Suppose, however, that the test customer cannot move to the back of the queue once he is in the queue, but upon reaching the server before E has occurred, he can choose to wait outside the queue any amount of time before rejoining the back of the queue. The amount of time the test customer waits is determined according to some policy, i.e., it may be determined by observing the length of the queue. Initially, then, the test customer may wait before joining the queue, but once in the queue he must stay in the queue until he reaches the server. This version of the problem was in fact the original version [Gopinath (1984)], and will be referred to as Problem 2 (P2). The problem considered so far will be referred to as Problem 1 (P1).

Lemma 4. Given any λ , α , and m , if the move-along policy is optimal for P1, then it is also optimal for P2.

Proof: Any allowable policy for P2 can be effectively duplicated by a policy for P1 (but not vice versa). Therefore the optimal policy for P1 must perform at least as well as the optimal policy for P2. \square

Theorems 4-6 and Corollary 2 therefore also apply to P2. Because any policy for P1 cannot in general be duplicated by a policy for P2, the converse to Lemma 4 may not be true. That is, if the move-along policy is not optimal for P1, it is unknown whether or not this implies that the move-along policy is not optimal for P2. The following Theorem states the analogous result for P2 as was stated in Theorem 6.

Theorem 8: For P2, given any α and m , there exists a λ_0 such that if $\lambda \geq \lambda_0$, the move-along policy is not optimal.

Proof: Assume that initially there are i customers in the queue, and that the test customer must decide to either join the queue immediately, or wait until either there are i' customers in the queue, or until E occurs, whichever occurs first. If the test customer chooses to wait, the mean delay until service is

where $p_{\{i,i'\}}$ is the probability that E occurs before the queue length becomes i' , D is the mean delay until service given that E occurs before the queue length becomes i' , and $\tau_{\{i,i'\}}(\lambda)$ is the mean time it takes to go from a queue length of i to i' (not including the test customer). Since D must be less than the mean delay given that E occurs after the queue length becomes i' ,

waiting outside the queue reduces the mean delay if

for some i . Clearly,

Using (2.17) and Lemma 2 gives the lower and upper bounds

and

The infinite series can be summed provided that $\lambda < 1$. From (2.30) this is always true if λ is large enough. Consequently, (4.5) becomes

The condition (4.2) is satisfied if $\tau_{i,i} > \tau_{i,i}(\lambda)$, or

For fixed α, m, i , and $i > 1$,

Consequently, for large λ the left hand side of (4.7) becomes

which is negative for large enough λ if

Since the function $f(u) = \sum_{i=1}^u$ decreases with i when $u < 1$ and i is large enough, (4.9) is true for large enough i . Consequently, for any α, m , and initial number of customers i , if λ is large enough, the observer can reduce his delay until service by waiting for the number of customers in the queue to increase to i . Therefore if λ is greater than some threshold value λ_0 , then the condition (4.2) is true for some i and $i > 1$, and the move-along policy is not optimal. \square

5. EXAMPLES

We conclude with two specific examples, namely the $M/D/1$ and $M/M/1$ queues.

Deterministic Service Time (M/D/1)

For this case

and

To compute $\tau_{i,i}$, it is necessary to compute the sequence

For the case $\phi_i(s) = \sum_{i=1}^i \phi_i(s) = e^{-s(i-\tau)}$, and $\phi_i(s) = e^{-s i}$, (2.22) becomes

where

>From Theorem 5, assuming that the test customer can move to the back of the queue at any time, the move-along policy is optimal if

and from Corollary 2, the move-along policy is optimal for the discrete-epoch problem if

Exponential Service Time (M/M/1)

For this case

and

and

For the case $\phi_i(s) = \tilde{\phi}^{sup i-1}(s) \phi_i(s) = \tilde{\phi}^{sup i}(s)$, and $\phi_j(s) = \tilde{\phi}^{sup j}(s)$, the expression for β_{gh} becomes

where

In the case of exponentially distributed service times the situations described in subsections (1.3) (minimal information and complete freedom) and (1.4) (minimal information and limited freedom) become identical. From Theorem 6 or Corollary 2 the move-along policy is guaranteed to be optimal if

>From Theorem 5, if the test customer knows the service times of the customers ahead of him, then the move-along policy is optimal if

6. UNANSWERED QUESTIONS

Assuming the conjecture stated in Section 1 is true, then Theorems 4 through 6 give lower bounds on the critical levels λ_k^* . Further improvements to these bounds have not yet been obtained.

The following Theorem applies to the case where the queueing system has a "finite capacity" $C < \infty$, where incoming customers are blocked and disappear if there are already C customers in the system (not including the test customer).

Theorem 9: For an $M/M/1$ queue with $C \leq 4$, the condition $t_{i+j} \leq t_i + t_j$ is always true, independent of λ , m , and α .

Proof: For finite C , a finite set of linear difference equations for the mean delay t_{ij} can be written down by inspection from which the theorem is easily verified. Details are omitted. \square

It is not known if Theorem 9 is true for any $C > 4$. For a very large capacity queueing system, Theorem 7 must apply. Consequently, there must exist a threshold λ_{under} (not necessarily finite), which is a function of C , such that the move-along policy is optimal if $\lambda \leq \lambda_{\text{under}}$. How does λ_{under} behave as C goes to infinity?

Perhaps the most interesting generalization of the problem studied here is the case where other customers in the queue are also waiting for events to occur before they can be served. For example, all customers may be waiting for independent events, each of which occurs after an exponentially distributed amount of time, and each may decide to follow the move-along policy. Is the move-along policy optimal for a particular test customer?

Acknowledgment

The authors thank A. Descloux for his careful review of the original manuscript.

REFERENCES

1. B. Gopinath, "Waiting for Godot", oral presentation at the *1984 Workshop on Open Problems in Communications and Computation*, Morristown, New Jersey.
2. M. L. Honig, "On Waiting for Simultaneous Access to Two Resources - Deterministic Service Distribution", *IEEE Transactions on Automatic Control*, Vol. AC-32, No. 11, pp. 1022-1025.
3. S. Y. R. Li, "Optimal Control of Premature Queueing," *IEEE Transactions on Automatic Control*, Vol. 33, No. 4, April 1988.