

INEQUALITIES AND AN APPROXIMATION FORMULA FOR THE MEAN DELAY TIME IN TANDEM QUEUEING SYSTEMS

HIROTAKA SAKASEGAWA AND GENJI YAMAZAKI

(Received July 13, 1976; revised Sept. 28, 1976)

1. Introduction

This paper deals with multi-stage tandem queueing systems (TQ for short) where each stage consists of a single server. An infinite queue is allowed to form before the first stage, whereas only a finite queue is allowed before all of the other stages. This assumption brings to a result of the blocking of servers besides the last one. In a study of the TQ with blocking effect it is difficult to analyze the system because it is impossible to regard its servers independent and treat them separately. Therefore, an analytical study for this system has not been done so extensively. Hunt [3] and others have analyzed a maximum utilization and a queue length of two- and three-stage TQ, in which an interarrival time and each service time distributions are Markovian type. Suzuki [11] derived steady state distributions of a queue size and a sojourn time in two-stage TQ with no intermediate queue, where the input stream forms a Poisson process. Often, these assumptions are not considered practical. Moreover, these analyses have been essentially based on characteristics of the system and so it is difficult to extend results to more complex systems.

Alternatively, there are some tendencies for proving general theorems of a qualitative character, valid under quite general conditions, but little study has been done under this heading in this area except that there are important theorems of Suzuki [12] and Hildebrand [2] on the ergodicity of the system.

In a practical viewpoint, another fruitful approach is to establish inequalities and approximation formulae concerning such variables as a delay time and a queue length, which are valid for a wide range of input and service mechanisms. For a single server queueing system (SQ for short) there exist such useful formulae including pioneering studies by Kingman ([5], [6]).

The present paper focuses on such inequalities and approximation formulae for a multi-stage TQ. Our approach developed in this paper

is to reduce it to a SQ equivalent in some sense to the TQ. First, begin with comparing the system with a certain SQ, some inequalities are derived for a mean delay time and a mean waiting time in two-stage TQ by using a mean waiting time in the SQ (Section 3). Second, as an extension of a case of two-stage TQ we try to reduce a multi-stage TQ to a certain SQ and derive an inequality for a mean waiting time in the TQ by using a mean waiting time in the SQ. We also derive upper and lower bounds for a capacity of the TQ (Section 4). From these inequalities some approximation formulae are set up. The validity of these formulae proposed in this paper is established by comparing with known analytical results whenever they are available and with simulation results otherwise (Section 5).

2. Description of the system and some notations

In this section we sum up the description of the system and some notations used in this paper. First, K -stage TQ where each stage consists of a single server is considered as follows: There are K service facilities (or servers for short, numbered 1 through K) arranged in tandem. Customers arrive at the system independently and queue up for a service by the first server (server 1). Each customer receives the service by the server 1, next the service by the server 2 and so forth until the service by the last server (server K). The service discipline is first-come-first-service at each stage. The queue before the first server may be allowed to grow unlimitedly, but on and after the second server, only a fixed finite number of customers are permitted to wait. If a queue before the k th server is full when another customer completes his service by the $(k-1)$ th server, this customer stays at the $(k-1)$ th stage and blocks further service at the stage until the service by the k th server is finished. There are no customer defections at any point.

When we consider such TQ, we can identify it with a TQ where none of customers are permitted to wait between servers. This fact, which is assured by considering a suitable number of servers with a service time identically zero, is due to Avi-Itzhak and Yadin [1]. Accordingly, from now on we assume without loss of generality that an intermediate queue is not allowed. For such a system, a notation of $GI/G_1 \rightarrow G_2 \rightarrow \dots \rightarrow G_K$ is employed where G_k is a service time distribution function (d.f.) by the server k .

Let $\{A_n\}$ and $\{S_{k,n}\}$ be a sequences of mutually independent identically distributed (i.i.d.) random variables (r.v.'s) defined on a probability space (Ω, \mathcal{B}, P) . A_n represents an interarrival time between $(n-1)$ th customer (C_{n-1}) and C_n , and $S_{k,n}$ represents a service time of C_n by the

server k . By using these r.v.'s we define other r.v.'s as follows:

$$(2.1) \quad a_n = A_1 + A_2 + \cdots + A_n \quad (\text{time epoch of the } n\text{th arrival to the system}),$$

$$(2.2) \quad T_{k,n} = \begin{cases} \max(T_{k+1,n-1}, T_{k-1,n} + S_{k,n}) & \text{for } k=1, 2, \dots, K-1 \\ T_{k-1,n} + S_{k,n} & \text{for } k=K \end{cases}$$

(time epoch at which C_n leaves the server k),

$$(2.3) \quad T_{0,n} = \max(T_{1,n-1}, a_n) \quad (\text{time epoch at which } C_n \text{ enters his first service by the server 1}),$$

$$(2.4) \quad W_n^0 = T_{0,n} - a_n \quad (\text{waiting time of } C_n \text{ in front of the first server}),$$

$$(2.5) \quad W_n = T_{K-1,n} - a_n \quad (\text{sojourn time of } C_n \text{ until a service by the last server begins since his arrival to the system})$$

and

$$(2.6) \quad B_{k,n} = T_{k,n} - T_{k-1,n} - S_{k,n} \quad (\text{blocking time of } C_n \text{ in the } k\text{th stage}).$$

We often use an r.v., say X , without a subscript n , say X_n (for example, W instead of W_n) which indicates an r.v. with a limiting d.f. of X_n . Let $X \vee Y$ and $X \wedge Y$ denote a maximum and a minimum of two r.v.'s X and Y , respectively. Further, we write $X \subset Y$ if X is stochastically smaller than Y , i.e. $\Pr(X > x) \leq \Pr(Y > x)$ for any x and $X \stackrel{d}{\sim} Y$ if a d.f. of X is identical with that of Y , i.e. $\Pr(X \leq x) = \Pr(Y \leq x)$ for any x .

3. Inequalities for 2-stage tandem queueing systems

Consider a 2-stage TQ with no intermediate queue in which a service time d.f. at each stage and an interarrival time d.f. are arbitrary. In [12], Suzuki proved that a sequence of a waiting time d.f. of n th customer in this system converges to an 'honest' d.f. as $n \rightarrow \infty$ iff $E A > E(S_1 \vee S_2)$. Namely, the TQ is identical with the SQ with a d.f. of $S_1 \vee S_2$ as the service time d.f., at least, on the equilibrium condition of the system. We denote such SQ by the notation $GI/\tilde{G}/1$. Though these systems have the same equilibrium condition, it cannot truly be regarded their nature as the same. However, if we can find some relationship between the TQ and the SQ, $GI/\tilde{G}/1$, it would be possible to apply precise results in the single server queueing theory to the TQ. The

main objective of this section is to find some relations between the $GI/G_1 \rightarrow G_2$ queue and the $GI/\tilde{G}/1$ queue. From now on, the $GI/\tilde{G}/1$ queue is called as a reduced single server queueing system (RSQ for short) of the $GI/G_1 \rightarrow G_2$ queue. Let X_n and \tilde{W}_n be a service time and a waiting time of the n th customer in the RSQ, respectively. We assume that $\{X_n\}$ is a sequence of i.i.d. r.v.'s with a d.f. of $S_1 \vee S_2$.

By comparing the TQ with its RSQ, we can obtain the following inequalities in the steady state.

THEOREM 3.1.

$$(3.1) \quad E W^0 \leq E \tilde{W} \leq E W \leq E \tilde{W} + E X.$$

PROOF. A proof of $E W^0 \leq E \tilde{W}$ is given in Section 4 for more general case (see Theorem 4.1). We prove $E \tilde{W} \leq E W$ and $E W \leq E \tilde{W} + E X$ here. By using (2.1), (2.2), (2.3) and (2.5), we have a recursive relation for W_n in the 2-stage TQ as follows.

$$(3.2) \quad \begin{aligned} W_{n+1} &= T_{1,n+1} - a_{n+1} \\ &= S_{1,n+1} \vee (T_{1,n} + S_{1,n+1} - a_{n+1}) \vee (T_{1,n} + S_{2,n} - a_{n+1}) \\ &= S_{1,n+1} \vee (S_{1,n+1} \vee S_{2,n} - A_{n+1} + T_{1,n} - a_n) \\ &\equiv S_{1,n+1} \vee (U_n + W_n) \end{aligned}$$

where $U_n = S_{1,n+1} \vee S_{2,n} - A_{n+1}$. If we assume that the system starts from scratch so that $W_1 = S_{1,1}$, we have as the solution of (3.2),

$$(3.3) \quad \begin{aligned} W_n &= S_{1,n} \vee (U_{n-1} + S_{1,n-1}) \vee (U_{n-1} + U_{n-2} + S_{1,n-2}) \vee \cdots \\ &\quad \vee (U_{n-1} + \cdots + U_1 + S_{1,1}). \end{aligned}$$

On the other hand, according to the famous relationship by Lindley [8], we have for the RSQ

$$(3.4) \quad \tilde{W}_{n+1} = 0 \vee (\tilde{U}_n + \tilde{W}_n)$$

where $\tilde{U}_n = X_n - A_{n+1}$. If we assume that $\tilde{W}_1 = 0$, we have as the solution of (3.4)

$$(3.5) \quad \tilde{W}_n = 0 \vee \tilde{U}_{n-1} \vee (\tilde{U}_{n-1} + \tilde{U}_{n-2}) \vee \cdots \vee (\tilde{U}_{n-1} + \cdots + \tilde{U}_1).$$

Since \tilde{U}_n has the same distribution as U_n , and $S_{1,i} \geq 0$, we have from (3.3) and (3.5)

$$(3.6) \quad \tilde{W}_n \subset W_n.$$

Now, replacing $S_{1,i}$ in each term of the right-hand side of (3.3) by $S_{1,i} \vee S_{2,i-1}$, we have

$$\begin{aligned}
 (3.7) \quad W_n &\subset (S_{1,n} \vee S_{2,n-1}) \vee (U_{n-1} + S_{1,n-1} \vee S_{2,n-2}) \vee \dots \\
 &\quad \vee (U_{n-1} + \dots + U_1 + S_{1,1} \vee S_{2,0}) \\
 &= S_{1,n} \vee S_{2,n-1} + 0 \vee U_{n-1}^* \vee (U_{n-1}^* + U_{n-2}^*) \vee \dots \vee (U_{n-1}^* + \dots + U_1^*)
 \end{aligned}$$

where $U_i^* = S_{1,i} \vee S_{2,i-1} - A_{i+1}$. Since $U_i^* \stackrel{d}{\sim} \tilde{U}_i$ and $S_{1,n} \vee S_{2,n-1} \stackrel{d}{\sim} X_n$, we have

$$(3.8) \quad W_n \subset X_n + \tilde{W}_n.$$

From (3.6) and (3.8), (3.1) follows.

Q.E.D.

From (3.1), the following inequalities can be obtained for the mean delay time $E(W^0 + B_1)$ ($= E W - E S_1$) in the TQ.

$$(3.9) \quad E \tilde{W} - E S_1 \leq E(W^0 + B_1) \leq E \tilde{W} + E X - E S_1.$$

If $\rho = (E X / E A) < 1$ so that $E S_1 / E A < 1$ and $(E X - E S_1) / E A < 1$, we have

$$(3.10) \quad |L_q - \tilde{L}_q| < 1$$

where $L_q = E(W^0 + B_1) / E A$ and $\tilde{L}_q = E \tilde{W} / E A$. If $E A$ is slightly greater than $E(S_1 \vee S_2)$, i.e. in heavy traffic situation, when \tilde{L}_q becomes much larger than unity, so (3.10) implies that

$$(3.11) \quad L_q \simeq \tilde{L}_q \quad \text{or} \quad E(W^0 + B_1) \simeq E \tilde{W}$$

in the sense that its relative error becomes small.

Furthermore, the following upper bound for $E(W^0 + B_1)$ can be obtained.

THEOREM 3.2. For a $GI/G_1 \rightarrow G_2$ queue, if $E(S_1 \vee S_2)^2$, $E A^2 < \infty$ and $E U < 0$, then

$$(3.12) \quad E(W^0 + B_1) \leq \frac{\text{Var}(U)}{2 E(-U)} + (E X - E S_1) \wedge \left(E S_1 \cdot \frac{\rho}{1 - \rho} \right)$$

where $\rho = E X / E A$.

PROOF. First we show that

$$(3.13) \quad E(W^0 + B_1) \leq \frac{\text{Var}(U)}{2 E(-U)} + E S_1 \cdot \frac{\rho}{1 - \rho}.$$

Now, (3.2) implies that, if S_1 and U are independent of W and have the common distribution of $S_{1,n}$ and U_n , respectively, then

$$(3.14) \quad W \stackrel{d}{\sim} S_1 \vee (U + W).$$

Hence,

$$(3.15) \quad E W = E(S_1 \vee (U + W)) \quad \text{and} \quad E W^2 = E(S_1 \vee (U + W))^2.$$

Let Z_n be a r.v. defined by $-\{S_{1,n} \wedge (S_{1,n} \vee S_{2,n-1} - A_n + W_{n-1})\}$. Then

the following functional relationships hold :

$$(3.16) \quad \begin{aligned} & S_1 \vee (U+W) - Z = S_1 + U + W \\ & \text{and} \\ & (S_1 \vee (U+W))^2 + Z^2 = S_1^2 + (U+W)^2. \end{aligned}$$

And so,

$$(3.17) \quad \begin{aligned} & E(S_1 \vee (U+W)) - E Z = E S_1 + E U + E W \\ & \text{and} \\ & E(S_1 \vee (U+W))^2 + E Z^2 = E S_1^2 + E(U+W)^2. \end{aligned}$$

By using (3.15), (3.17) can be simplified as

$$(3.18) \quad -E Z = E S_1 + E U \quad \text{and} \quad E Z^2 = E S_1^2 + E U^2 + 2 \cdot E W \cdot E U.$$

From (3.18), we obtain

$$(3.19) \quad \text{Var}(Z) = \text{Var}(S_1) + \text{Var}(U) + 2(E W - E S_1) E U,$$

so that

$$(3.20) \quad E(W^0 + B_1) = E W - E S_1 = \frac{\text{Var}(U) + \text{Var}(S_1) - \text{Var}(Z)}{2 E(-U)}.$$

Well, let R_n be a r.v. defined as

$$(3.21) \quad R_n \equiv Z_n + S_{1,n} = 0 \vee (A_n + 0 \wedge (S_{1,n} - S_{2,n-1}) - W_n).$$

In the steady state, we have

$$(3.22) \quad E R = E(-U)$$

and

$$(3.23) \quad \text{Var}(Z) - \text{Var}(S_1) = \text{Var}(R) - 2 \text{Cov}(R, S_1).$$

On the other hand, from (3.21)

$$(3.24) \quad R_n \subset A_n \quad \text{or} \quad S_{1,n} \cdot R_n \subset S_{1,n} \cdot A_n.$$

Since $S_{1,n}$ is independent of A_n , we have

$$(3.25) \quad E(S_1 \cdot R) \leq E S_1 \cdot E A.$$

Using (3.22), (3.23) and (3.25), we can rewrite (3.20) as

$$(3.26) \quad \begin{aligned} E(W^0 + B_1) &= \frac{\text{Var}(U) - \text{Var}(R) + 2 \text{Cov}(R, S_1)}{2 E(-U)} \\ &\leq \frac{\text{Var}(U)}{2 E(-U)} + \frac{E(S_1 \cdot R) - E S_1 \cdot E R}{E(-U)} \end{aligned}$$

$$\begin{aligned} &\leq \frac{\text{Var}(U)}{2E(-U)} + \frac{ES_1 \cdot EA - E(-U) \cdot ES_1}{E(-U)} \\ &= \frac{\text{Var}(U)}{2E(-U)} + ES_1 \cdot \frac{EX}{E(-U)} \end{aligned}$$

which is equivalent to (3.13). On the other hand, the following inequality holds for the mean waiting time in the RSQ (Kingman [5]):

$$(3.27) \quad E\tilde{W} \leq \frac{\text{Var}(A) + \text{Var}(X)}{2(EA - EX)} = \frac{\text{Var}(\tilde{U})}{2E(-\tilde{U})} \quad \left(= \frac{\text{Var}(U)}{2E(-U)} \right).$$

From (3.9) and (3.27), we can obtain

$$\begin{aligned} (3.28) \quad E(W^0 + B_1) &= EW - ES_1 \leq E\tilde{W} + EX - ES_1 \\ &\leq \frac{\text{Var}(U)}{2E(-U)} + EX - ES_1. \end{aligned}$$

From (3.26) and (3.28), (3.12) follows.

Q.E.D.

We can derive further inequalities than (3.1) and (3.12) for some special cases. Now let begin with inequalities for a mean delay time in a $1/\lambda (=EA)$ -MRLB/ $G_1 \rightarrow G_2$ and DFR/ $G_1 \rightarrow G_2$ queues (that is, the class of $GI/G_1 \rightarrow G_2$ queues whose arrival process has $1/\lambda$ -MRLB or DFR property). Before deriving these, we mention about γ -MRLB and DFR properties (Marshall [9]). Let $F(x)$ be a d.f. defined on a positive real axis. We define $F(x)$ to have γ -MRLB (Mean Residual Life bounded by γ from Below) property when

$$\int_x^\infty (1-F(y))dy / (1-F(x)) \geq \gamma \quad \text{for all } x \geq 0 \text{ s.t. } F(x) \neq 1.$$

And we define $F(x)$ to have DFR (Decreasing Failure Rate) property when

$$(F(x+\Delta) - F(x)) / (1-F(x))$$

is decreasing in x s.t. $F(x) \neq 1$ for any $\Delta > 0$. Then we have the following

THEOREM 3.3. *For a $1/\lambda$ -MRLB/ $G_1 \rightarrow G_2$ queue,*

$$(3.29) \quad EW \leq \frac{\text{Var}(U)}{2E(-U)} - \frac{E(-U)}{2}.$$

For a DFR/ $G_1 \rightarrow G_2$ queue,

$$(3.30) \quad EW \leq \frac{\text{Var}(U)}{2E(-U)} - \frac{EA(C_a^2 - \rho)}{2}$$

where C_a is a coefficient of variation (c.v.) of an interarrival time d.f.

PROOF. Marshall [9] derived that

$$(3.31) \quad E \tilde{W} \leq \frac{\text{Var}(U)}{2 E(-U)} - \frac{E A(1+\rho)}{2}$$

for a $1/\lambda\text{-MRLB}/\tilde{G}/1$ queue and

$$(3.32) \quad E \tilde{W} \leq \frac{\text{Var}(U)}{2 E(-U)} - \frac{E A(C_a^2 + \rho)}{2}$$

for a $\text{DFR}/\tilde{G}/1$ queue. From (3.1) and (3.31) we can obtain

$$(3.33) \quad \begin{aligned} E W &\leq E \tilde{W} + E X \leq \frac{\text{Var}(U)}{2 E(-U)} - \frac{E A(1+\rho)}{2} + E X \\ &= \frac{\text{Var}(U)}{2 E(-U)} - \frac{E(-U)}{2} \end{aligned}$$

for the $1/\lambda\text{-MRLB}/G_1 \rightarrow G_2$ queue. Similarly, from (3.1) and (3.32) we can obtain (3.30). Q.E.D.

Because a Poisson arrival process has $1/\lambda\text{-MRLB}$ property, (3.29) holds for $M/G_1 \rightarrow G_2$ queue, too. Especially in an $M/M \rightarrow M$ queue, we have the following lower bound for $E(W^0 + B_1)$.

PROPOSITION 3.1. For an $M/M \rightarrow M$ queue,

$$(3.34) \quad E(W^0 + B_1) \geq E \tilde{W}.$$

PROOF. Let $F(x) = 1 - e^{-\lambda x}$ be an interarrival time d.f. and $G_i(x) = 1 - e^{-\mu_i x}$ be a service time d.f. by the i th server ($i=1, 2$). Using Pollaczek-Khinchin's formula, a mean queue length of the RSQ, \tilde{L}_q , becomes

$$(3.35) \quad \tilde{L}_q = \lambda E \tilde{W} = \frac{1 + C_s^2}{2} \frac{\rho^2}{1 - \rho}$$

where C_s is a c.v. of $S_1 \vee S_2$. Insert

$$C_s^2 = \frac{\text{Var}(S_1 \vee S_2)}{\{E(S_1 \vee S_2)\}^2} = \frac{1 + 2\alpha - \alpha^2 + 2\alpha^3 + \alpha^4}{(1 + \alpha + \alpha^2)^2} \quad \text{and} \quad \rho = \rho_1 \frac{1 + \alpha + \alpha^2}{\alpha + \alpha^2}$$

where $\alpha = \mu_2/\mu_1$ and $\rho_i = \lambda/\mu_i$ ($i=1, 2$) into (3.35), we have

$$(3.36) \quad \tilde{L}_q = \frac{\rho_1^2(1 + 2\alpha + \alpha^2 + 2\alpha^3 + \alpha^4)}{\alpha(1 + \alpha)D}$$

where $D = \alpha + \alpha^2 - (1 + \alpha + \alpha^2)\rho_1$. On the other hand, Kishi [7] gave an explicit formula of L_q ($= \lambda E(W^0 + B_1)$) as follows:

$$(3.37) \quad L_q = \frac{H(\rho_1^2 + \rho_2^2) - \rho_1^2 \rho_2^2 / (\rho_1 + \rho_2)}{H(H - \rho_1 - \rho_2)}$$

where $H = 1 + (\rho_1 \rho_2 / (\rho_1 + \rho_2))$. After tedious computations, L_q is rewritten as follows:

$$(3.38) \quad L_q = \frac{\rho_1^2(1+\alpha)(1+\alpha+\alpha^2+\alpha^3+(1-\alpha+\alpha^2)\rho_1)}{\alpha(1+\alpha+\rho_1)D}.$$

Now, subtracting \tilde{L}_q from L_q , we obtain

$$(3.39) \quad L_q - \tilde{L}_q = \frac{\rho_1^2}{(1+\alpha)(1+\alpha+\rho_1)}$$

which is clearly positive for all positive values of ρ_1 and α . Q.E.D.

An important point to note about the derivation of this proposition is that the difference between L_q and \tilde{L}_q for $M/M \rightarrow M$ queue is very small (order of ρ_1^2) for all over the range $0 < \rho < 1$, i.e. $L_q \simeq \tilde{L}_q$ (or, equally, $E(W^0 + B_1) \simeq E\tilde{W}$) for $0 < \rho < 1$. For some special cases, it will be shown below that $E(W^0 + B_1)$ equals to $E\tilde{W}$. These results, together with (3.10), suggest that (3.11) holds for considerably wide subclass of $GI/G_1 \rightarrow G_2$ queues not only in heavy traffic situation but for all over the range $0 < \rho < 1$, which will be assured numerically in Section 5.

For the TQ where a service time of at least one server is constant, we have the following

THEOREM 3.4. *For a $GI/D \rightarrow G$ queue,*

$$(3.40) \quad E(W^0 + B_1) = E\tilde{W} \leq \frac{\text{Var}(U)}{2E(-U)}.$$

For a $GI/G \rightarrow D$ queue,

$$(3.41) \quad E(W^0 + B_1) \geq E\tilde{W}.$$

PROOF. We prove (3.40) first. Let C be a constant service time by the first server (that is, $S_{1,n} = C$ for all n), then from (3.2)

$$W_{n+1} = C \vee (U_n + W_n)$$

where $U_n = C \vee S_{2,n} - A_{n+1}$. Using this equality

$$(3.42) \quad W_{n+1}^0 + B_{1,n+1} = W_{n+1} - C = 0 \vee (U_n + W_n - C) = 0 \vee (U_n + W_n^0 + B_{1,n})$$

where U_n is independent of $W_n^0 + B_{1,n}$. On the other hand, (3.4) holds for \tilde{W}_{n+1} in the RSQ. Since \tilde{U}_n has the same d.f. as U_n , $W_n^0 + B_{1,n}$ has the same d.f. as \tilde{W}_n , i.e.

$$(3.43) \quad W_n^0 + B_{1,n} \stackrel{d}{\sim} W_n.$$

From (3.43) and (3.27), we can obtain (3.40).

Now, let $(W_n^0 + B_{1,n})^*$ be a delay time of C_n in a dual system of $GI/D \rightarrow G$ queue (denoted by $GI/G \rightarrow D$) which is obtained by interchanging two service facilities. The following result has been derived by Kawashima [4]:

$$(3.44) \quad W_n^0 + B_{1,n} \subset (W_n^0 + B_{1,n})^*.$$

It is clear that the RSQ for the $GI/D \rightarrow G$ queue is identical with that for the $GI/G \rightarrow D$ queue, accordingly, from (3.43) and (3.44) we can obtain (3.41). Q.E.D.

If a service time by the first server is always longer than that by the second server or vice versa, the following relations hold.

PROPOSITION 3.2. If $\Pr(S_1 \geq S_2) = 1$, then

$$(3.45) \quad E(W^0 + B_1) = E\tilde{W} \leq \frac{\text{Var}(U)}{2E(-U)}.$$

If $\Pr(S_1 \leq S_2) = 1$, then

$$(3.46) \quad E(W^0 + B_1) \geq E\tilde{W}.$$

PROOF. Let $W_n^0 + B_{1,n}$ be a delay time of C_n in a $GI/G_1 \rightarrow G_2$ queue where $\Pr(S_1 \geq S_2) = 1$ and $(W_n^0 + B_{1,n})^*$ be a delay time of C_n in a dual system of the $GI/G_1 \rightarrow G_2$ queue. Since the first server is never blocked in the $GI/G_1 \rightarrow G_2$ queue, we have

$$(3.47) \quad W_n^0 + B_{1,n} \stackrel{d}{\sim} \tilde{W}_n$$

and (3.45) follows. On the other hand, (3.44) is also true in this situation (Kawashima [4]). Using (3.44) and (3.47), we get (3.46). Q.E.D.

4. Inequalities for K -stage tandem queueing systems ($K \geq 3$)

For any $K (\geq 3)$ -stage TQ, it is quite difficult to treat a waiting time, a delay time, etc., analytically, and the authors have very little information on them. Same as the preceding section, it seems useful to reduce the TQ to an equivalent SQ in some sense and to approximate a mean delay time in the TQ by a mean waiting time in the SQ. An upper bound for a mean waiting time before the first server EW^0 in the TQ (i.e. the mean delay time minus the mean blocking time) is given by the following theorem. The theorem shows that EW^0 is always smaller than the mean waiting time in an SQ ($E\tilde{W}$) if we take

a d.f. of $S_1 \vee S_2 \vee \dots \vee S_K$ as a service time d.f. of the SQ.

THEOREM 4.1. For K -stage TQ and the SQ mentioned above,

$$(4.1) \quad W_n^0 \subset \tilde{W}_n \quad \text{for any } n.$$

Hence,

$$(4.2) \quad E W^0 \leq E \tilde{W}.$$

PROOF. Using (2.3), (2.4) and (2.6), we can obtain

$$(4.3) \quad W_n^0 = 0 \vee (W_{n-1}^0 + S_{1,n-1} + B_{1,n-1} - A_n) \quad \text{for any } n \text{ and any } \omega \in \Omega.$$

We define a r.v. U_n as $S_{1,n-1} + B_{1,n-1} - A_n$. Similarly for the SQ

$$(4.4) \quad \tilde{W}_n = 0 \vee (\tilde{W}_{n-1} + X_{n-1} - A_n)$$

where X_n is a service time of C_n in the SQ. Here we also define \tilde{U}_n as $X_{n-1} - A_n$. If both systems start from scratch, i.e. $W_0^0 = \tilde{W}_0 = 0$, the following two equalities are derived as the solutions of (4.3) and (4.4).

$$(4.5) \quad W_n^0 = 0 \vee U_n \vee (U_n + U_{n-1}) \vee \dots \vee (U_n + \dots + U_1)$$

and

$$(4.6) \quad \tilde{W}_n = 0 \vee \tilde{U}_n \vee (\tilde{U}_n + \tilde{U}_{n-1}) \vee \dots \vee (\tilde{U}_n + \dots + \tilde{U}_1).$$

Now, the following inequality holds (Hildebrand [2]):

$$(4.7) \quad S_{1,n-1} + B_{1,n-1} \leq S_{1,n-1} \vee S_{2,n-2} \vee \dots \vee S_{K,n-K} (\equiv Y_{n-1}).$$

So that

$$(4.8) \quad W_n^0 \leq 0 \vee U_n^* \vee (U_n^* + U_{n-1}^*) \vee \dots \vee (U_n^* + \dots + U_1^*)$$

for any n and any $\omega \in \Omega$

where $U_i^* = Y_{i-1} - A_i$. The right-hand side of (4.8) is identical with \tilde{W}_n in distribution because Y_i has the same distribution as X_i . Accordingly, (4.1) is concluded, so that

$$(4.9) \quad E W_n^0 \leq E \tilde{W}_n.$$

Since both sides of (4.9) have their limiting values $E W^0$ and $E \tilde{W}$, respectively, (4.2) follows. Q.E.D.

It is still an open problem to evaluate a mean delay time in the TQ by a closed form. Practically, however, it is necessary to estimate the mean delay time, though the estimation is slightly rough. Now we proceed to search more effective SQ than that used in the Theorem 4.1 in order to extend the results derived in Section 3. In 2-stage TQ,

since $EA > E(S_1 \vee S_2)$ is an equilibrium condition, we considered as its RSQ the SQ whose service time d.f. is that of $S_1 \vee S_2$. In doing so, we derived some closed relations between the mean delay time in the TQ and the mean waiting time in the RSQ. For any K -stage TQ also, if there exists a d.f. whose mean is connected with an equilibrium condition in the TQ, we may expect that an SQ with this d.f. as its service time d.f. plays an important role in estimating the mean delay time in the TQ. From now on, this SQ is called by the RSQ same as the case of 2-stage TQ. In order to derive this d.f. we consider K -stage TQ where an infinite customers are always queueing up in front of the first stage, and let $B'_{1,n}$ be a blocking time of the C_n in the first stage in this imaginary system. Hildebrand [2] proved that the sequence $\{S_{1,n} + B'_{1,n}\}$ converges in distribution and that there exists an 'honest' equilibrium d.f. of a waiting time if $\lim_{n \rightarrow \infty} E(S_{1,n} + B'_{1,n}) < EA$. Let $F(\cdot)$ be a limiting d.f. of $S_{1,n} + B'_{1,n}$, then this $F(\cdot)$ is just what we want and it becomes the service time d.f. of the RSQ. We note in parentheses that in 2-stage TQ, a d.f. of $S_1 \vee S_2$ becomes the $F(\cdot)$.

In order to approximate the mean delay time in K -stage TQ by the mean waiting time in the RSQ with $F(\cdot)$ as the service time d.f., at least the first few moments of $F(\cdot)$ are needed. Unfortunately, the existence is the only matter which is known to us about $F(\cdot)$. Accordingly, we require to derive approximation formulae for these parameters. The following theorem concerns the first moment of the $F(\cdot)$.

THEOREM 4.2. *For K -stage TQ ($K \geq 3$), the following inequalities hold:*

$$(4.10) \quad E\{S_1 \vee S_2 \vee (S_3 + 0 \wedge (S'_2 - S'_1))\} \\ \vee E\{S_K \vee S_{K-1} \vee (S_{K-2} + 0 \wedge (S'_{K-1} - S'_K))\} \\ \leq \lim_{n \rightarrow \infty} E(S_{1,n} + B'_{1,n}) \leq E(S_1 \vee S_2 \vee \dots \vee S_K)$$

where $\{S_1, S'_1, S_{1,n}\}$, $\{S_2, S'_2\}$, $\{S_{K-1}, S'_{K-1}\}$ and $\{S_K, S'_K\}$ are i.i.d. r.v.'s, respectively.

PROOF. In this proof, we treat only the imaginary system described above. Then regarding all a_n equals to zero, from (2.1), (2.2) and (2.3) we have

$$(4.11) \quad T_{1,n} - T_{0,n} = S_{1,n} \vee S_{2,n-1} \vee (S_{3,n-2} + T_{2,n-2} - T_{1,n-1}) \vee \dots \\ \vee (S_{K,n-K+1} + T_{K-1,n-K+1} - T_{1,n-1}) \\ \geq S_{1,n} \vee S_{2,n-1} \vee (S_{3,n-2} + T_{2,n-2} - T_{1,n-1}) \\ \text{for any } n \text{ and any } \omega \in \Omega.$$

Furthermore,

$$\begin{aligned}
 (4.12) \quad T_{2,n-2} - T_{1,n-1} &= T_{2,n-2} - (T_{1,n-2} + S_{1,n-1}) \vee T_{2,n-2} \\
 &= 0 \wedge (T_{2,n-2} - T_{1,n-2} - S_{1,n-1}) \geq 0 \wedge (S_{2,n-2} - S_{1,n-1}) .
 \end{aligned}$$

Hence,

$$(4.13) \quad T_{1,n} - T_{0,n} \geq S_{1,n} \vee S_{2,n-1} \vee (S_{3,n-2} + 0 \wedge (S_{2,n-2} - S_{1,n-1})) .$$

By taking a limit of an expectation of both sides, we can obtain

$$\begin{aligned}
 (4.14) \quad \lim_{n \rightarrow \infty} E(S_{1,n} + B'_{1,n}) &= \lim_{n \rightarrow \infty} E(T_{1,n} - T_{0,n}) \\
 &\geq E\{S_1 \vee S_2 \vee (S_3 + 0 \wedge (S'_2 - S'_1))\} .
 \end{aligned}$$

By using the duality of the TQ derived by the authors [13], we can obtain

$$(4.15) \quad \lim_{n \rightarrow \infty} E(S_{1,n} + B'_{1,n}) \geq E\{S_K \vee S_{K-1} \vee (S_{K-2} + 0 \wedge (S'_{K-1} - S'_K))\} .$$

From (4.14) and (4.15), the first inequality in (4.10) follows. On the other hand, by taking a limit of an expectation of both sides in (4.7) we can obtain the second inequality in (4.10). Q.E.D.

In the following section, we will derive an approximation formula for the mean of $F(\cdot)$ based on Theorem 4.2. We know, for the present, nothing about the second moment of $F(\cdot)$, however, we will give an approximation formula for it based on much numerical experiments. Using those results, we will show that the mean delay time in the TQ is well approximated by the mean waiting time in the RSQ in the following section.

5. Numerical examples

For a given tandem queueing system like as the one treated in this paper, the only way to know such characteristics as a mean queue length, a mean waiting time, etc., is the method of simulation experiments, because of the lack of analytical results. Simulation experiments are, however, much time-consuming and the estimated value is only assured to be within such and such a confidence interval. Consequently, it is very convenient if such characteristics can be estimated, though approximately, without using simulation experiments. In this section we propose some approximation formulae which are based on results given in the previous sections and show with many numerical examples that they are practically useful for a wide range of TQ's. In the following, we treat a mean queue length $L_q = \lambda E(W^0 + B)$ instead of $E(W^0 + B)$ because of its dimensionless property. In practical situation, balanced-service-time systems are recognized to be most important, so

a mean service time of each server is set to unity.

In Section 3, we showed that the mean delay time in the 2-stage TQ is approximately the same as the mean waiting time in its RSQ in heavy traffic situation. Now we examine the following approximation formula for many systems not only in heavy traffic but also in light traffic situations.

$$(5.1) \quad L_q \equiv \lambda E(W^0 + B_1) \simeq \tilde{L}_q \equiv \lambda E\tilde{W}.$$

5.1. 2-stage TQ with Poisson arrival process

For $M/M \rightarrow M$, $M/M \rightarrow D$ and $M/D \rightarrow M$ queueing systems, tables of

the mean queue length for various system parameters are published ([14]). On the other hand, the mean queue length of the RSQ for each system can be calculated using Polaczek-Khinchin's formula, i.e.,

$$(5.2) \quad L_q = \frac{1 + C_s^2}{2} \frac{\rho^2}{1 - \rho}.$$

The mean queue length of the TQ and that of the RSQ, together with the upper bound of the mean queue length of the RSQ (Kingman [5]),

$$(5.3) \quad \frac{\lambda \text{Var}(U)}{E(-U)} = \frac{\rho^{-2} + C_s^2}{2} \frac{\rho^2}{1 - \rho}$$

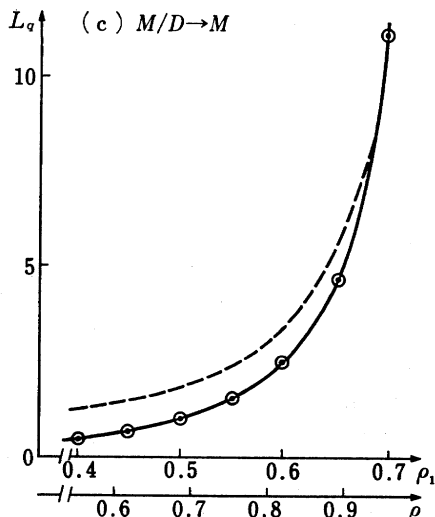
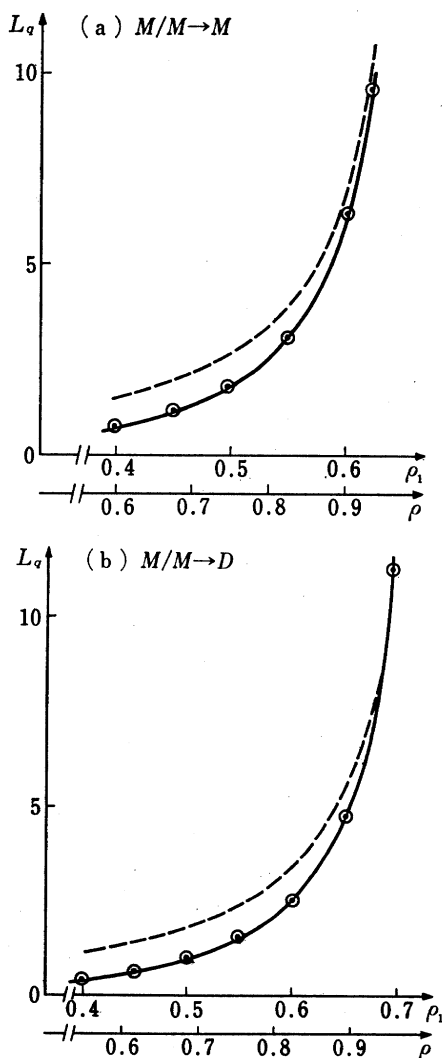


Fig. 5.1. Mean queue length in $M/M \rightarrow M$, $M/M \rightarrow D$ and $M/D \rightarrow M$ queueing systems.

○... L_q ; —... \tilde{L}_q ; - - - $\lambda \text{Var}(U)/2E(-U)$

are illustrated in Fig. 5.1. According to these figures, L_q 's are well-approximated by \tilde{L}_q 's not only in heavy traffic situation but for any value of $\rho < 1$.

When service time d.f.'s are arbitrary, the exact formula for the mean queue length is not known, so we resort to a simulation experiment to know it. In a single simulation run, we start the system from scratch and average whole data of 5000 customers. Let $s_{1,n}$ and $s_{2,n}$ be a service time of C_n at the first and the second stage, respectively, a_n be an interarrival time between C_{n-1} and C_n and w_n be a delay time of C_n . Then, $\bar{\rho}$ and \bar{L}_q are evaluated as follows:

$$\bar{\rho} = \frac{\bar{s}}{\bar{a}}$$

where $\bar{a} = (1/N) \sum a_n$ and $\bar{s} = \mu(1/2N) \sum (s_{1,n} + s_{2,n})$, $\mu = E(S_1 \vee S_2)$ and

$$\bar{L}_q = \frac{1}{\bar{a}} \frac{1}{N} \sum w_n.$$

Then a point $(\bar{\rho}, \bar{L}_q)$ is illustrated in a ρ - L_q plane. The mean queue length of the RSQ, \tilde{L}_q , can be obtained from (5.2). Among many simulated systems, two of them are illustrated in Fig. 5.2 (a) and (b). On the same plane, the upper bound of the mean queue length of the RSQ which is given by (5.3) is put together. In each figure, α denotes $(C_a^2 + C_i^2)/2$ with c.v.'s C_a and C_i of the RSQ. According to these figures, the approximation formula (5.1) seems to be successful for all $\rho < 1$.

5.2. 2-stage TQ other than 5.1

The mean queue length for $GI/G_1 \rightarrow G_2$ cannot be expressed in a closed form in general, so a simulation experiment is carried out to know it. The mean queue length for its RSQ is also unknown in general, but there is a well-going approximation formula for the mean queue length for a single server queueing system has been proposed by one of the authors [10], then we use this formula (5.4).

$$(5.4) \quad \tilde{L}_q \simeq \frac{C_a^2 + C_i^2}{2} \frac{\rho^2}{1 - \rho}.$$

Several computer runs were made for typical queueing systems in which an interarrival time and service time d.f.'s are Erlangian type. Fig. 5.2 (c)~(f) are four examples of our simulation experiments. According to them, our attempt to reduce the TQ to the RSQ has a good chance of success.

In this case, most points lie beneath the dotted line, and so it seems to be correct, though we cannot prove analytically, that $\text{Var}(U)/2 E(-U)$ is an upper bound of $E(W^0 + B_1)$.

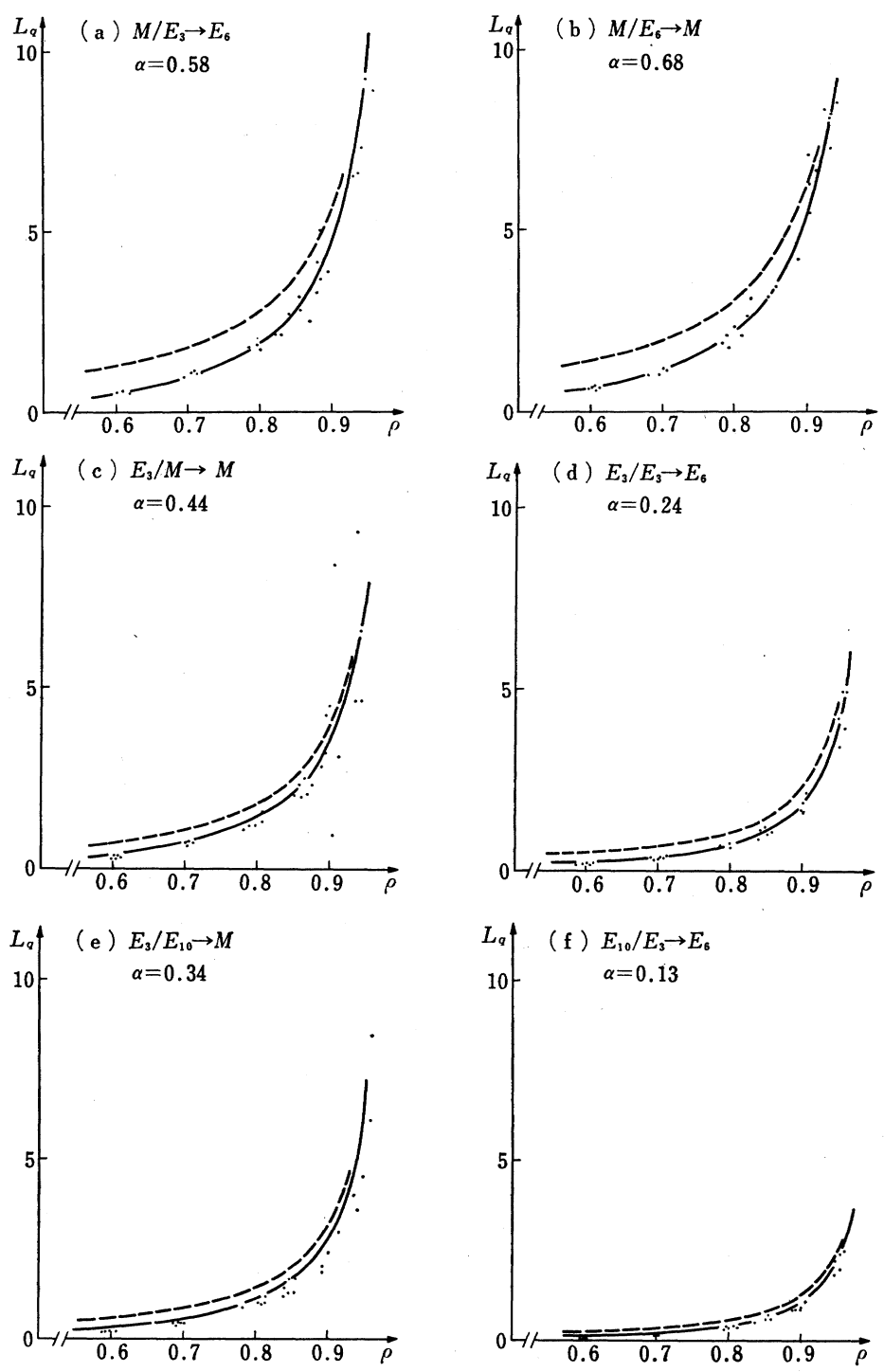


Fig. 5.2. Simulation results for 2-stage tandem queueing systems.

•... L_q ; —... \tilde{L}_q ; - - - $\lambda \text{Var}(U)/2 E(-U)$

5.3. K -stage TQ ($K \geq 3$)

A mean queue length cannot also be calculated analytically in these cases and simulation experiments are designed for Erlangian type TQ's. On the other hand, a mean queue length of the RSQ is approximately evaluated from (5.4) with a knowledge of a c.v. of a service time d.f. of the RSQ. Following the preceding section, the service time d.f. of the RSQ is a limiting d.f. of $S_{1,n} + B_{1,n}$, say $F(\cdot)$, but unfortunately, this $F(\cdot)$ is only known to exist. Now we here use some approximation formulae to evaluate first two moments and a c.v. of $F(\cdot)$.

For the first moment of $F(\cdot)$, say μ , we showed inequalities about μ in the preceding section, Theorem 4.2. With our much numerical experiments, the following formula is useful to evaluate μ approximately (see Table 5.1).

$$(5.5) \quad \mu \approx \frac{1}{4} \{ 2 E(S_1 \vee \dots \vee S_K) \\ + E(S_1 \vee S_2 \vee S_3) \\ + E(S_K \vee S_{K-1} \vee S_{K-2}) \} \equiv \tilde{\mu}.$$

For the second moment of $F(\cdot)$, we have no analytical discussion but we conclude with our much numerical experiments that a variance of an

Table 5.1. Test of an approximation formula (5.3) for $\infty/E_{k_1} \rightarrow \dots \rightarrow E_{k_K}$

(k_1, \dots, k_K)	μ	$\tilde{\mu}$
(1, 1, 1, 1)	1.943	1.958
(3, 3, 3, 3)	1.549	1.565
(1, 3, 6, 10)	1.580	1.569
(1, 1, 1, 1, 1)	2.059	2.058
(3, 3, 3, 3, 3)	1.607	1.616
(1, 2, 3, 6, 10)	1.678	1.659
(3, 3, 1, 3, 10)	1.686	1.698
(1, 1, 1, 1, 1, 1)	2.141	2.142
(1, 1, 1, 1, 1, 1, 1)	2.208	2.213

Table 5.2. Comparison between $Y = S_1 + B_1$ and $X = \max(S_1, \dots, S_K)$ for $\infty/E_{k_1} \rightarrow E_{k_2} \rightarrow \dots \rightarrow E_{k_K}$

(k_1, \dots, k_K)	Var(Y)	Var(X)	c.v. ² (Y)	$\frac{\text{Var}(X)}{\tilde{\mu}^2}$	α	
					Y	X
(1, 1, 1)	1.38	1.36	0.44	0.41	0.72	0.70
(3, 3, 3)	0.30	0.34	0.14	0.15	0.57	0.57
(3, 1, 10)	0.58	0.64	0.25	0.26	0.62	0.63
(10, 1, 3)	0.64	0.64	0.27	0.26	0.64	0.63
(1, 1, 1, 1)	1.49	1.42	0.39	0.37	0.70	0.69
(3, 3, 3, 3)	0.35	0.33	0.15	0.14	0.57	0.57
(1, 3, 6, 10)	0.67	0.58	0.27	0.24	0.63	0.62
(10, 6, 3, 1)	0.59	0.58	0.24	0.24	0.62	0.62
(1, 1, 1, 1, 1)	1.50	1.46	0.35	0.35	0.68	0.68
(3, 3, 3, 3, 3)	0.33	0.33	0.13	0.13	0.56	0.56
(1, 2, 3, 6, 10)	0.76	0.62	0.27	0.23	0.63	0.60
(10, 6, 3, 2, 1)	0.57	0.62	0.20	0.23	0.60	0.60

* For 3-stage TQ, $\tilde{\mu} = E(S_1 \vee S_2 \vee S_3)$.

** $\alpha = \frac{C_a^2 + C_i^2}{2}$ in (5.4) where $C_a^2 = 1$ and $C_i^2 = \text{c.v.}^2(Y)$ or $\frac{\text{Var}(X)}{\tilde{\mu}^2}$ in this table.

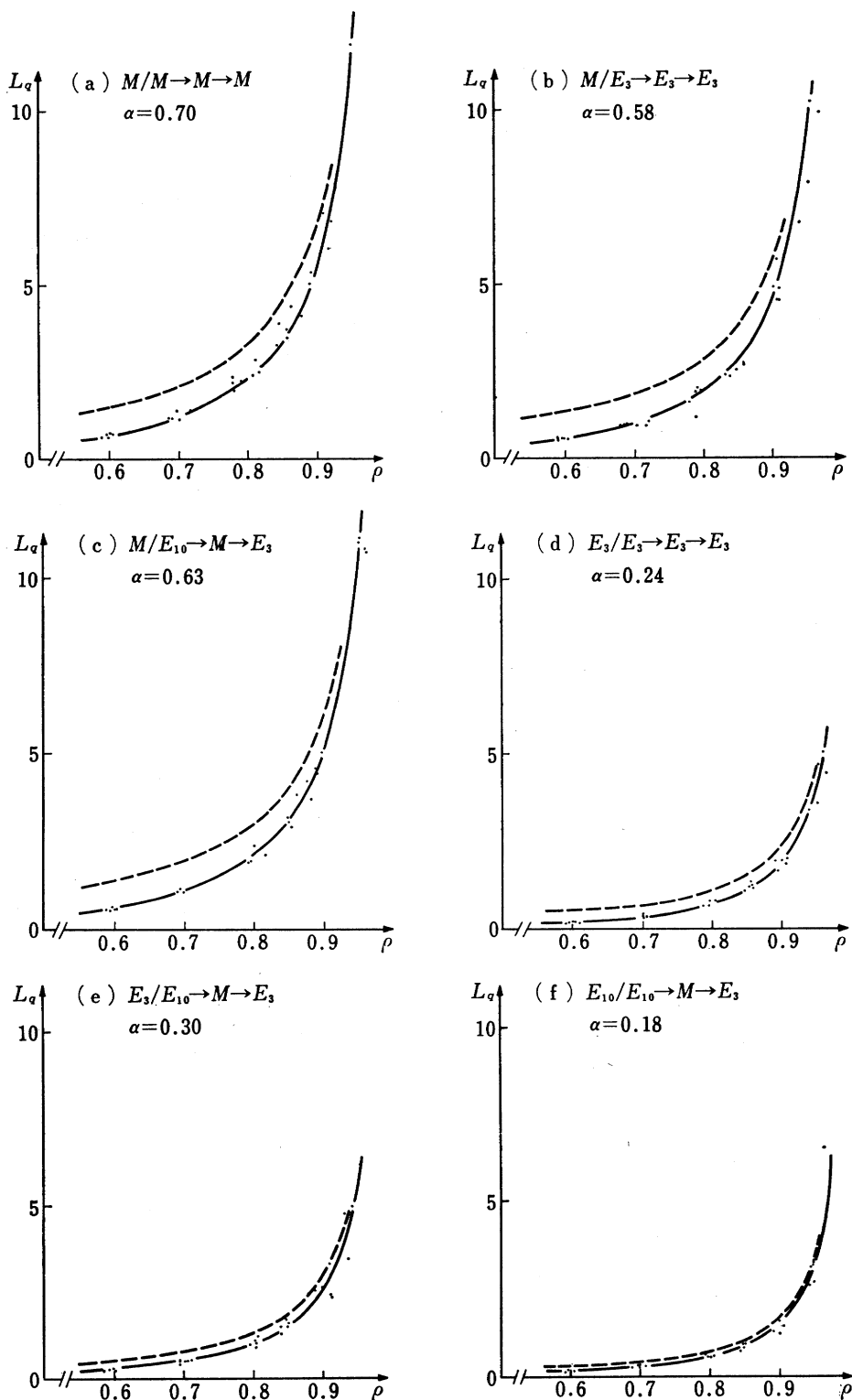


Fig. 5.3. Simulation results for 3-stage tandem queueing systems.

$\bullet \dots L_q$; $\text{---} \dots \bar{L}_q$; $\text{---} \dots \lambda \text{Var}(U)/2E(-U)$

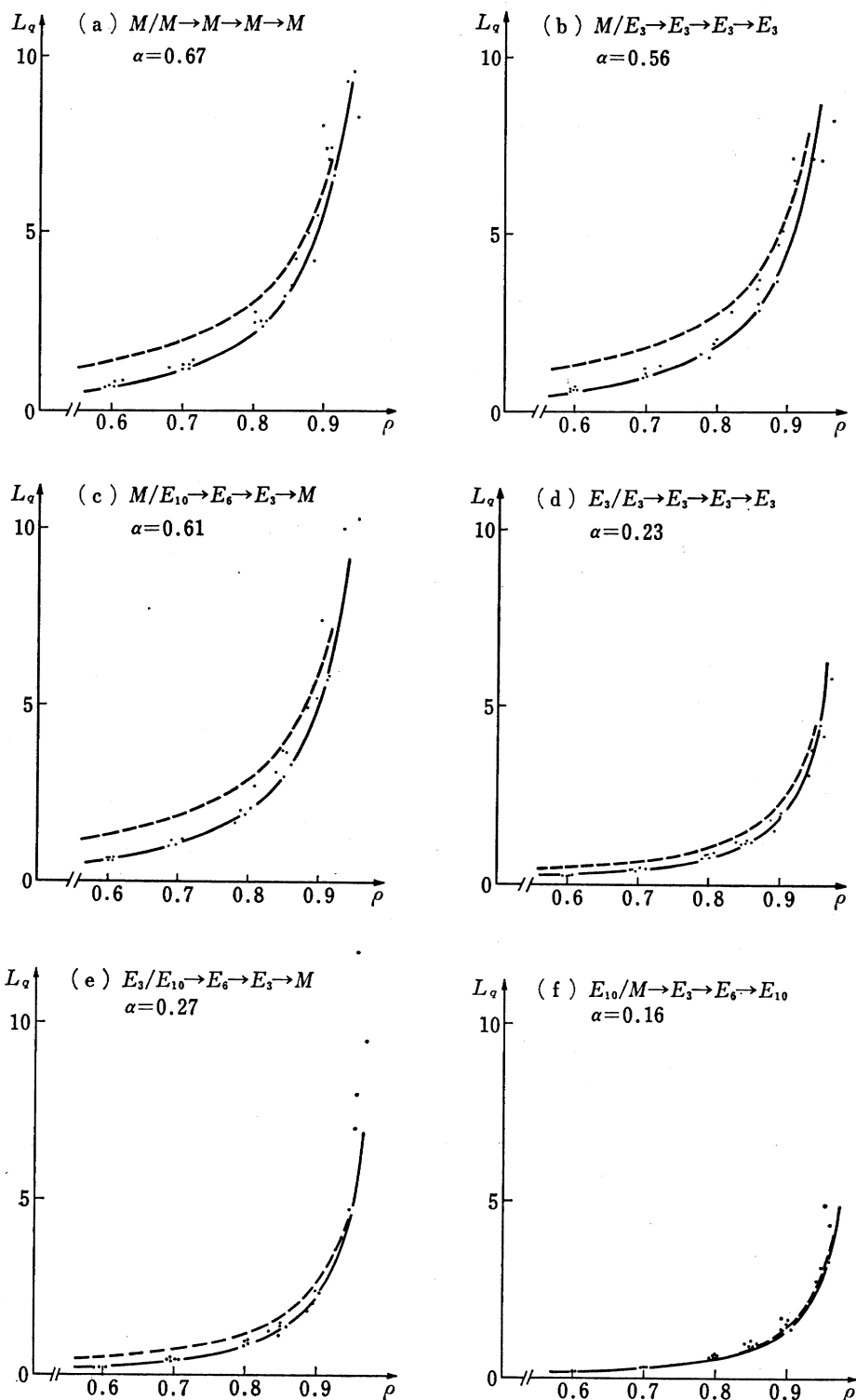


Fig. 5.4. Simulation results for 4-stage tandem queueing systems.

 $\bullet \dots L_q$; $\text{—} \dots \bar{L}_q$; $\text{---} \dots \lambda \text{Var}(U)/2E(-U)$

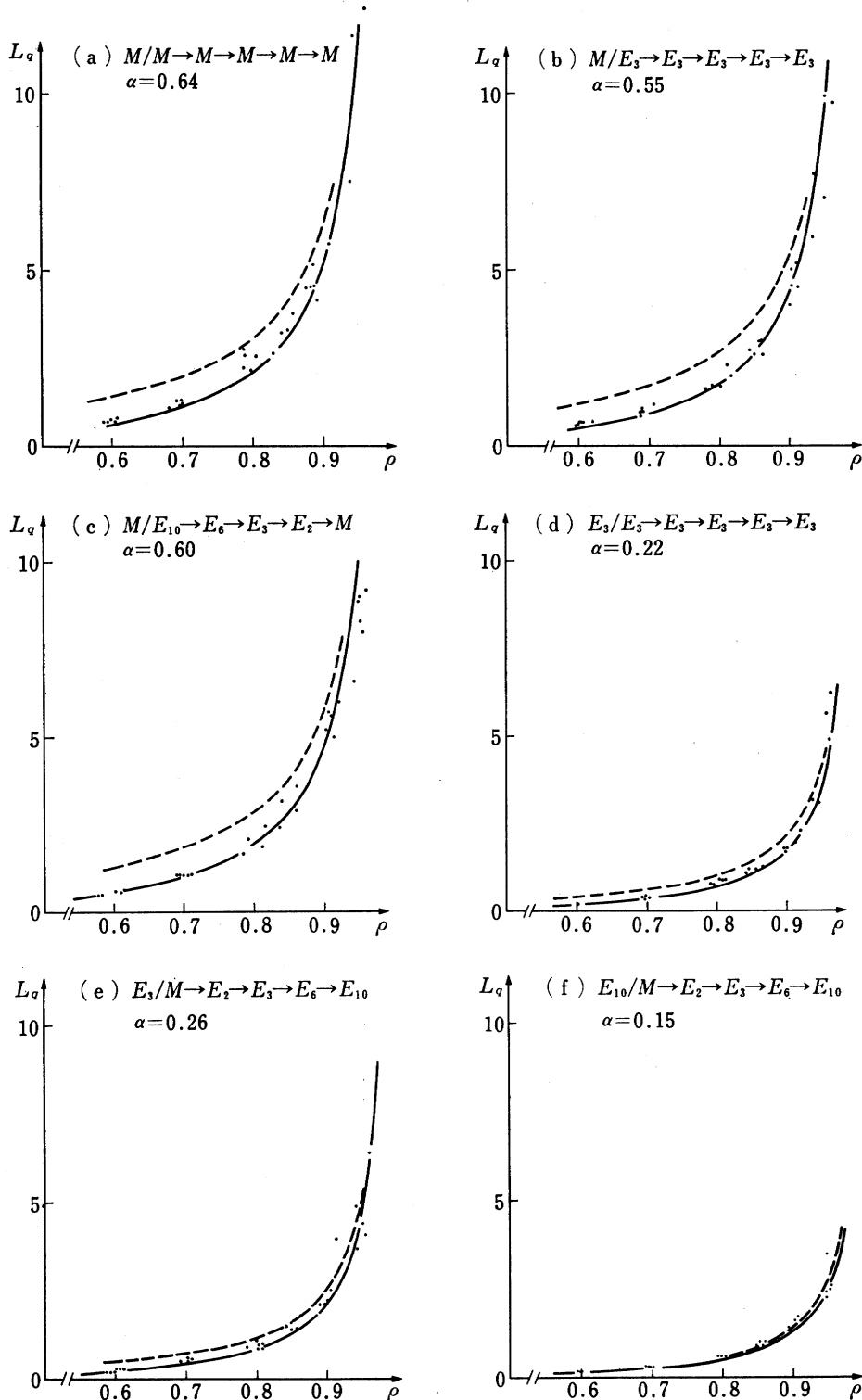


Fig. 5.5. Simulation results for 5-stage tandem queueing systems.

$\bullet \dots L_q$; $\text{---} \dots \bar{L}_q$; $\text{---} \dots \lambda \text{Var}(U)/2 E(-U)$

r.v. $S_1 \vee S_2 \vee \dots \vee S_K$ approximates that of $F(\cdot)$ with practical accuracy (see Table 5.2). Finally, the c.v. of $F(\cdot)$ is calculated by the above two values. According to our simulation experiments, these formulae are not so bad in many cases. The mean queue length for the TQ are examined against this approximated values of the mean queue length for the RSQ. Some parts of our simulation experiments are shown in Fig 5.3, Fig 5.4 and Fig 5.5. Method of generation of each point in these figures is the same as Section 5.2 above. According to these figures the mean queue length for the TQ is well approximated by the mean queue length for the RSQ in many cases.

Same as Fig. 5.2, a curve of $\lambda \text{Var}(U)/2 E(-U)$ is superposed to each figure. If C_a is large ((a)~(e) in Fig 5.3, Fig 5.4 and Fig 5.5), $\text{Var}(U)/2 E(-U)$ seems to be an upper bound of $E(W^0 + B_1)$ but if C_a is relatively small (Fig. 5.3 (f), Fig. 5.4 (f) and Fig. 5.5 (f)), this proposition is doubtful.

Well, we examined applicability of the approximation formula (5.1) with many queueing systems but for most of them the right-hand side of (5.1) could not be calculated exactly but evaluated approximately by using (5.4). So, we may propose the next approximation formula without using the mean waiting time in the RSQ rather than (5.1).

$$(5.6) \quad L_q \simeq \frac{C_a^2 + C_s^2}{2} \frac{\rho^2}{1 - \rho}.$$

Let us describe in detail. In our approximation formula we use $\text{Var}(S_1 \vee \dots \vee S_K)$ instead of $\text{Var}(S_1 + B_1)$ and $\tilde{\mu}$ instead of μ . This means that if d.f.'s of service time are not changed, arrangement of servers is no effect on the delay time. According to our simulation experiments, however, a mean delay time in a system, say Q_1 , seems to be slightly longer than that in a service-order-reversed system, say Q_2 , when a c.v. of the first stage service time d.f. in Q_1 is greater than that in Q_2 . On the other hand, under the assumption on Q_1 and Q_2 above, $\text{Var}(S_1 + B_1)$ in Q_1 seems to be slightly greater than $\text{Var}(S_1 + B_1)$ in Q_2 in our simulation experiments. Hence, a c.v. of $S_1 + B_1$ is not the same in both systems (see Table 5.2). So, if we can calculate a c.v. of $S_1 + B_1$ exactly, the approximation formula (5.1) or (5.6) will become more accurate using this value.

6. Concluding remarks

We have shown that the approximation formula (5.6) is practically useful for any tandem queueing systems. But there remain further studies to take the place of a simulation method. They are the study of the determination of the range of subclass in which such and such

a property holds, the study of the property of $F(\cdot)$ appeared in Section 5.3, and so on.

In conclusion we give two conjectures, which we cannot prove analytically at yet, in evidence of much simulation experiments.

CONJECTURE 1*. For 2-stage tandem queueing system,

$$E \tilde{W} \leq E(W^0 + B_1) \leq \frac{\text{Var}(U)}{2E(-U)}.$$

CONJECTURE 2. For K -stage tandem queueing system,

$$|L_q - \tilde{L}_q| < 1$$

where \tilde{L}_q is a mean queue length of the RSQ with $F(\cdot)$ (c.f. Section 5.3) as the service time d.f.

Acknowledgement

Authors wish to thank Dr. R. Kato and Mr. T. Kawashima for their valuable discussions. They also wish to thank the referee for his pertinent criticism, to which they are deeply indebted in refining the first version of this paper.

THE INSTITUTE OF STATISTICAL MATHEMATICS
KOGAKUIN UNIVERSITY

REFERENCES

- [1] Avi-Itzhak, B. and Yadin, M. (1965). A sequence of two servers with no intermediate queue, *Management science*, 11, 553-564.
- [2] Hildebrand, D. K. (1967). Stability of finite queue, tandem server system, *J. Appl. Prob.*, 4, 571-583.
- [3] Hunt, G. C. (1956). Sequential arrays of waiting lines, *Operat. Res. Quart.*, 4, 674-683.
- [4] Kawashima, Y. (1976). Reverse ordering of services in tandem queues, *Memoirs of Defense Academy*, 15, 151-159.
- [5] Kingman, J. F. C. (1962). On queues in heavy traffic, *J. R. Statist. Soc.*, B-24, 383-392.
- [6] Kingman, J. F. C. (1970). Inequalities in the theory of queues, *J. R. Statist. Soc.*, B-32, 102-110.
- [7] Kishi, T. (1960). Queues with parallel phases (in Japanese), *Keiei-Kagaku*, 3, 156-165.
- [8] Lindley, D. V. (1952). The theory of queues with a single server, *Proc. Camb. Phil. Soc.*, 48, 277-289.
- [9] Marshall, K. T. (1968). Some inequalities in queuing, *Operat. Res. Quart.*, 16, 651-665.
- [10] Sakasegawa, H. (1976). An approximation formula $L_q \simeq \alpha \cdot \rho^2 / (1 - \rho)$, *Ann. Inst. Statist. Math.*, 29, 67-75.
- [11] Suzuki, T. (1964a). On a tandem queue with blocking, *J. Operat. Res. Soc. Japan*, 6, 137-157.
- [12] Suzuki, T. (1964b). Ergodicity of a tandem queue with blocking, *J. Operat. Res. Soc. Japan*, 7, 48-75.
- [13] Yamazaki, G. and Sakasegawa, H. (1975). Properties of duality in tandem queueing systems, *Ann. Inst. Statist. Math.*, 27, 201-212.
- [14] *Queueing Tables* (in Japanese), (1970), Iwanami.

* Recently, the first inequality of Conjecture 1 was proved to be correct, which will appear in near future.