

# TOTAL WAITING TIME DISTRIBUTION FUNCTION AND THE FATE OF A CUSTOMER IN A SYSTEM WITH TWO QUEUES IN SERIES

Wolfgang Krämer  
University of Stuttgart  
Stuttgart, Federal Republic of Germany

## ABSTRACT

The study of traffic flow in modern computer and communication networks and various other technical systems leads in many cases to systems or subsystems with queues arranged in series or tandem.

This paper is concerned especially with a system of two unlimited queues in series, when Poisson traffic is offered to the first stage and the service times in the two stages are independent of each other and negative exponentially distributed. The first stage includes one service-unit, while the second stage is allowed to be a multiserver queuing system.

In this paper, by pursuing a particular customer (call or request) at his walk through the whole system, the fate of a customer in the second stage is determined as a function of all possible components of fate in the first stage. Though the flow times (waiting plus service times) of the same customer in the successive stages are independent, other values of fate are not independent (waiting times or the numbers of customers met upon arriving at the single stages).

Since the fate of a customer in the second stage is independent of the concrete value ( $>0$ ) of his waiting time in the first stage, it is possible to determine the total waiting time distribution function by convolution of special terms.

Finally, also such customers are investigated which have met a known number of predecessors in the first stage upon their arrival. Considering all possible paths in a RANDOM WALK diagram, the queue-length met upon arrival in the second (single server) stage is determined, as well as the expected waiting time or flow time.

Beitrag des Instituts für  
Nachrichtenvermittlung und  
Datenverarbeitung der Universität  
Stuttgart zum 7th International  
Teletraffic Congress Stockholm  
vom 13.-20. Juni 1973

## 1. INTRODUCTION

### 1.1 DESCRIPTION OF THE SYSTEM

The system dealt with consists of two unlimited queues arranged in series, where the input process to the first stage is a Poisson process with mean arrival rate  $\lambda$ . The arriving customers, calls or requests, shortly referred to as calls, first are served by a single server and then by one server of the second stage, which is allowed to be a multiserver system (fig. 1.1).

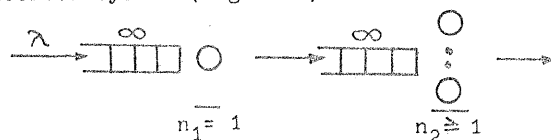


Fig.1.1 The system

The service or holding times  $T_{Hi}$  ( $i=1,2$ ) of a call in both stages are independent of each other and negative exponentially distributed with distribution functions (d.f.)

$$P(T_{Hi} \leq t) \stackrel{\text{def}}{=} H_i(t) = 1 - e^{-\epsilon_i t} \quad (1.1)$$

and means

$$E(T_{Hi}) = \frac{1}{\epsilon_i} = h_i \quad (1.2)$$

The traffics offered  $A_i$  are defined by

$$A_i = \frac{\lambda}{\epsilon_i} = \lambda \cdot h_i \quad (1.3)$$

and the utilizations

$$S_i = \frac{A_i}{n_i} = \frac{\lambda}{\mu_i} \quad \text{with } \mu_i = n_i \cdot \epsilon_i \quad (1.4)$$

Both stages are assumed to be in statistical equilibrium, so that

$$\lambda < \min(\mu_1, \mu_2) \quad (1.5)$$

Considering a general call at its walk through the whole system, the following time diagram is obtained:

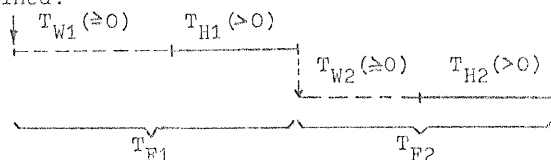


Fig.1.2 General time diagram

$T_{Wi}$  is the random waiting time not including service  $T_{Hi}$ , and  $T_{Pi}$  the random flow time in stage  $i$ . The queue disciplines in both stages are, as long as no d.f.'s of the waiting times are concerned, arbitrary, otherwise in order of arrival (FIFO).

### 1.2 KNOWN RESULTS

It is well known, that the output of the first stage is a Poisson process with mean output rate  $\lambda$  (BURKE [1,3], and others), so that each single stage may be computed completely according to the formulae for the M/M/n queueing system, cf. e.g. SYSKI [11].

JACKSON [5] proved the state probabilities of the single stages of such a system at the same time to be independent of each other. So it is possible to quote directly the state probabilities of the whole system.

Considering a particular call at its walk through the whole system, a further group of traffic characteristics may be obtained which will be regarded in this paper in more detail.

Fig. 1.3 represents a graphical survey of some relations between various fate values (random variables) of a certain arbitrary call in such a two stage system ( $n_1 > 1$  included).

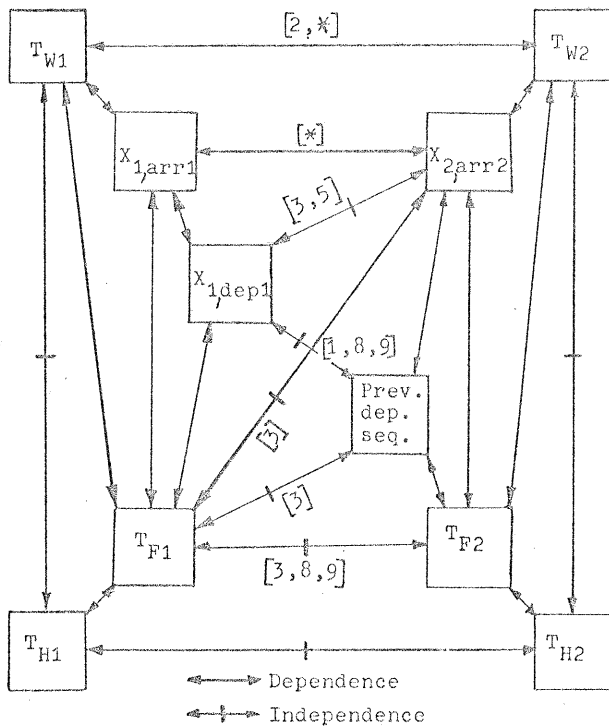


Fig. 1.3 Relations between fate values

Besides  $T_{Wi}$ ,  $T_{Hi}$ ,  $T_{Fi}$  and the number of calls met upon arrival in stage  $i$  ( $X_{i, arr i}$ ), also the number of calls left behind in the first stage ( $X_{1, dep 1}$ ) and the departure sequence of the first stage previous to the departure of the considered call ( $v, dep. seq.$ ) are involved.

The arrows connecting two values showing whether (separately considered) they are dependent or not, are labeled with references to the literature ( $[*]$  means basic parts of this paper, see 1.3).

Using the concept of reversibility of a Markov chain, the independence of the flow times was shown by REICH [9,10] for single server and by BURKE [3] for multiserver systems. (Further results from [3] see 2.1).

NELSON [7] derived by convolution an expression for the distribution function of the total waiting time in a (more generally structured) system, assuming independence of the waiting times in the single stages.

BURKE [2] proved that the waiting times of a call in the investigated system with  $n_2=1$  are dependent. Using the theorem of JACKSON and the virtual delay in stage 2, he showed by explicit calculation that the probabilities of waiting  $W_i = P(T_{Wi} > 0)$  in the first and second stage are not independent. Though the traffic offered to stage 2 is pure chance traffic, the future fate of a call (number  $X_{2, arr 2}$  of calls met in stage 2 upon arrival there and/or waiting time  $T_{W2}$ ) is not independent of its previous fate in stage 1.

### 1.3 TREATED PROBLEMS

The aim of this paper is to investigate and to throw more light upon the dependencies between these various fate values of a call in the two successive stages and to extend known results. Therefore, several test-calls with different and more or less known fate in the first stage are considered.

First, there are calls with special assumptions about the waiting time in the first stage (no waiting, waiting in the first stage of unknown duration, known waiting time ( $>0$ )).

Determining the number of calls met by these test-calls upon arrival in the second stage, the influences of these values to the further fate (waiting or flow time) are obtained. Admitting both calls with known and unknown service time in the first stage, also the influence of the service time in the first stage is shown (chapter 2).

The knowledge of all these dependencies is the precondition for the determination of total fates (total waiting time d. f. with total probability of waiting and mean total waiting time of the waiting calls). Since in a preliminary chapter it is shown, that the additional assumption of a concrete waiting time ( $>0$ ) in the first stage has no further effect (called 'limited dependency'), it is possible to determine the d. f. of the total waiting time by convolution of special terms in chapter 3.

In the last chapter, a further group of test-calls is considered, namely such calls which have found a certain number of predecessors in the first stage. Considering all possible paths in a RANDOM WALK diagram, the number of calls met in a second (single server) stage is determined. It is shown that for the queue lengths no such 'limited dependency' is valid as for the waiting times.

### 2. TEST CALLS WITH GIVEN WAITING TIME IN STAGE 1

#### 2.1 OUTPUT PROCESS OF M/M/n DURING CONCRETE WAITING TIME

The aim of this preliminary investigation is to obtain statements about the behaviour of a stationary single stage M/M/n system with FIFO during a certain waiting time  $T_W = t_0$  ( $>0$ ) of a test-call. Let  $p_W(j, t_0) = P(X_{arr} = j | T_W = t_0)$  be the probability that this test-call has met  $j$  ( $\geq n$ ) calls upon its arrival in the whole stage. Applying the theorem of BAYES it holds:

$$P(X_{arr} = j | T_W \in [t_0, t_0 + dt]) = \frac{P(X_{arr} = j)}{P(T_W \in [t_0, t_0 + dt])} \cdot P(T_W \in [t_0, t_0 + dt] | X_{arr} = j) \quad (2.1)$$

Inserting known expressions into the right side and making the limit transition  $dt \rightarrow 0$ , it is obtained

$$p_W(j, t_0) = e^{-\lambda t_0} \cdot \frac{(\lambda t_0)^{j-n}}{(j-n)!} \quad t_0 > 0, j \geq n \quad (2.2)$$

It is obvious that this expression (which is independent of  $\epsilon$ ) is identical with the probability that  $j-n=z$  calls arrive during  $t_0$ :

The probability that a call with concrete waiting time  $T_W = t_0$  ( $>0$ ) has met  $z$  ( $\geq 0$ ) calls in the waiting storage, is the same as the probability that the same call upon leaving the waiting storage (to begin service) leaves  $z$  ( $\geq 0$ ) calls behind.

If a call with concrete waiting time  $t_0$  ( $>0$ ) has met  $j$  ( $\geq n$ ) calls upon arrival in the whole stage (this occurs with probability  $p_W(j, t_0)$ ), exactly  $j-n+1$  calls must be served until its service begins. Since the departure of the last one of them coincides with the end of  $t_0$ , exactly  $j-n$  calls leave the system during  $t_0$ . If we forget the value of  $j$ , nevertheless the Poisson distribution (2.2) must be fulfilled for reasons of stationarity. This implies that the output intervals during  $t_0$  are negative exponentially distributed with mean  $1/\lambda$  (Poisson).

Remark:

More detailed investigations in this direction were performed by BURKE [3]. Considering so-called 'partial delays' (partial waiting times for calls who find at least  $j-1$  calls in the waiting storage to leave waiting place number  $x$ ,  $x > 0$ ), he proved that the conditional d.f. of a partial delay given a call has met  $j$  calls in the waiting storage is the same as under the condition that  $j$  calls arrived during this partial delay. Hereof a lemma was derived which states, that a partial delay given that  $j$  calls arrived during this partial delay, is independent of the previous departure sequence.

Finally, it may be noted that in case of  $n_1=1$  with the very same method used here for the waiting times, an alternative proof of the independence between the fate in the second stage and the flow time in the first stage can be obtained.

## 2.2 GENERAL WAY OF CALCULATION

The observation of the system starts at time  $T$ , when the fixed service time  $T_{H1}=t_1$  of a test-call begins.

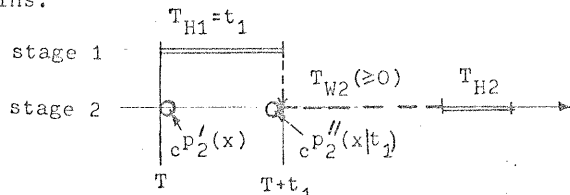


Fig. 2.1 Time diagram

From the state probabilities  $p_2'(x)$  at time  $T^+$  (called starting probabilities), which are independent of the required service time  $t_1$ , the state probabilities  $cP_2''(x|t_1)$  at time  $(T+t_1)^+$  (called meeting probabilities) are determined. (Prefix  $c$  means that the probability is a conditional probability related to a special call, where  $c=0$  refers to a call with  $T_{W1}=0$  and  $c=1$  to a call with  $T_{W1}>0$  of known or unknown duration.)

The values received are valid for calls with a certain required service time in the first stage. By integration over all possible service times also results are obtained for calls with unknown service time  $T_{H1}$ .

## 2.3 STARTING PROBABILITIES IN STAGE 2

Using the theorem of JACKSON [5], it is obvious that the state probabilities  $p_2'(x)$  at the arrival of a nonwaiting call in stage 1 are independent and according to the absolute state probabilities:

$$p_2'(x) = p_2(x) \quad (2.3)$$

where

$$p_2(x) = \begin{cases} p_2(0) \cdot \frac{A_2^x}{x!} & 0 \leq x \leq n_2 \\ p_2(0) \cdot \frac{A_2^x}{n_2! n_2^{x-n_2}} & x \geq n_2-1 \end{cases} \quad (2.4)$$

with

$$p_2(0) \cdot \frac{A_2^{n_2}}{n_2!} \cdot \frac{1}{1 - \frac{A_2}{n_2}} = E_{2,n_2}(A_2) = W_2 \quad (2.5)$$

which is the absolute probability of waiting in the second stage according to the second formula of ERLANG, cf. e.g. SYSKI [1].

For calls which have to wait in the first stage the following situation prevails: at time  $T-t_0$  the considered test-call with waiting time  $t_0 (>0)$  arrives in stage 1. Due to JACKSON's theorem the state probabilities of the second stage at time  $T-t_0$  are identical with the absolute values according to (2.4). During the subsequent time  $t_0$  the input process of the second stage is Poisson (shown in 2.1), so the state probabilities of the second stage at time  $T$  will also be distributed according to the absolute

values. So

$$p_2'(x) = \begin{cases} 0 & \text{for } x = 0 \\ p_2(x-1) & \text{for } x > 0 \end{cases} \quad (2.6)$$

It is obvious that (2.6) is also valid when the concrete value of  $T_{W1}$  is unknown, or when a concrete service time  $T_{H1}=t_1$  is preassigned to a test-call.

## 2.4 MEETING PROBABILITIES IN STAGE 2

Let  $p(i, t_1)$  be the probability that  $i$  calls leave the second stage during the service time  $T_{H1}=t_1$  of a test-call. As long as stage 2 is fully occupied ( $X_2 \geq n_2$ ), the rate of the whole stage to serve one call will be  $\mu_2 = n_2 \epsilon_2$  (cf. fig. 2.1).

So the meeting probabilities are

$$cP_2''(x|t_1) = \sum_{i=0}^{\infty} p(i, t_1) \cdot cP_2'(x+i) \quad \begin{matrix} x \geq n_2 \\ c=0,1 \end{matrix} \quad (2.7)$$

where

$$p(i, t_1) = e^{-\mu_2 t_1} \cdot \frac{(\mu_2 t_1)^i}{i!} \quad i \geq 0 \quad (2.8)$$

is the probability of  $i$  Poisson events during  $t_1$ .

By integration we will get the meeting probabilities for corresponding test-calls with unknown service time  $T_{H1}$ :

$$cP_2''(x) = \int_0^{\infty} cP_2''(x|t_1) dH_1(t_1) \quad c=0,1 \quad (2.9)$$

For test-calls which have not waited in stage 1 it is obtained from (2.7)

$$P(X_{2, \text{arr}} = x | T_{W1}=0, T_{H1}=t_1) \stackrel{\text{def}}{=} cP_2''(x|t_1) = p_2(x) \cdot e^{-(\mu_2 - \lambda)t_1} \quad x \geq n_2 \quad (2.10)$$

If the assumption of a certain service time is dropped, we receive by (2.9) with (1.1)

$$P(X_{2, \text{arr}} = x | T_{W1}=0) \stackrel{\text{def}}{=} oP_2''(x) = \frac{\epsilon_1}{\epsilon_1 + \mu_2 - \lambda} \cdot p_2(x) \quad x \geq n_2 \quad (2.11)$$

Similarly, we obtain the meeting probabilities for calls with  $T_{W1}>0$  of known or unknown duration ( $c=1$ ):

$$1P_2''(x|t_1) = p_2(x-1) \cdot e^{-(\mu_2 - \lambda)t_1} \quad x \geq n_2 \quad (2.12)$$

$$1P_2''(x) = \frac{\epsilon_1}{\epsilon_1 + \mu_2 - \lambda} \cdot p_2(x-1) \quad (2.13)$$

Considering a call of which only the service time in the first stage is known, we get from (2.10) and (2.12) by weighting summation

$$P(X_{2, \text{arr}} = x | T_{H1}=t_1) \stackrel{\text{def}}{=} p_2''(x|t_1) = (1-s_1 + \frac{s_1}{s_2}) \cdot p_2(x) \cdot e^{-(\mu_2 - \lambda)t_1} \quad x \geq n_2 \quad (2.14)$$

## 2.5 CONDITIONAL PROBABILITIES OF WAITING

Summing up the meeting probabilities for all numbers of calls met in the waiting storage, the probability of waiting in stage 2 is obtained for the various test-calls:

$$cP(T_{W2} > 0 | T_{H1}=t_1) = \sum_{x=n_2}^{\infty} cP_2''(x|t_1) \quad (2.15)$$

If service time is unknown:

$$cW_2 = \sum_{x=n_2}^{\infty} cP_2''(x) \quad c=0,1 \quad (2.16)$$

(2.10) and (2.12) yield with (2.15) after some intermediate calculations

$$P(T_{W2} > 0 | T_{W1}=0, T_{H1}=t_1) \stackrel{\text{def}}{=} oP(T_{W2} > 0 | T_{H1}=t_1) = E_{2,n_2}(A_2) e^{-(\mu_2 - \lambda)t_1} \quad (2.17)$$

and

$$1P(T_{W2} > 0 | T_{H1}=t_1) = \frac{n_2}{\lambda} \cdot E_{2,n_2}(A_2) e^{-(\mu_2 - \lambda)t_1} \quad (2.18)$$

Considering test-calls with unspecified service times, the following formulae are obtained:

$$P(T_{W2} > 0 | T_{W1}=0) \stackrel{\text{def}}{=} oW_2 = oF \cdot E_{2,n_2}(A_2)$$

$$\text{with } oF = \frac{\epsilon_1}{\epsilon_1 + \mu_2 - \lambda} = \frac{s_2}{s_1 + s_2 - s_1 s_2} \quad (<1) \quad (2.19)$$

$$P(T_{W2} > 0 | T_{W1}>0) \stackrel{\text{def}}{=} 1W_2 = 1F \cdot E_{2,n_2}(A_2)$$

$$\text{with } 1F = \frac{\mu_2 \epsilon_1}{\lambda(\epsilon_1 + \mu_2 - \lambda)} = \frac{1}{s_1 + s_2 - s_1 s_2} \quad (>1) \quad (2.20)$$

Comparing the two probabilities of waiting, we get the simple relation

$${}_0W_2 = \frac{A_2}{n_2} \cdot {}_1W_2 \quad (2.21)$$

Especially, for  $n_2=1$  (2.19) and (2.20) simplify to the same results derived by BURKE [2]. In fig. 2.2 these conditional probabilities of waiting are plotted versus the utilization  $S_2 = A_2/n_2$  of stage 2 with  $S_2/S_1$  as parameter.

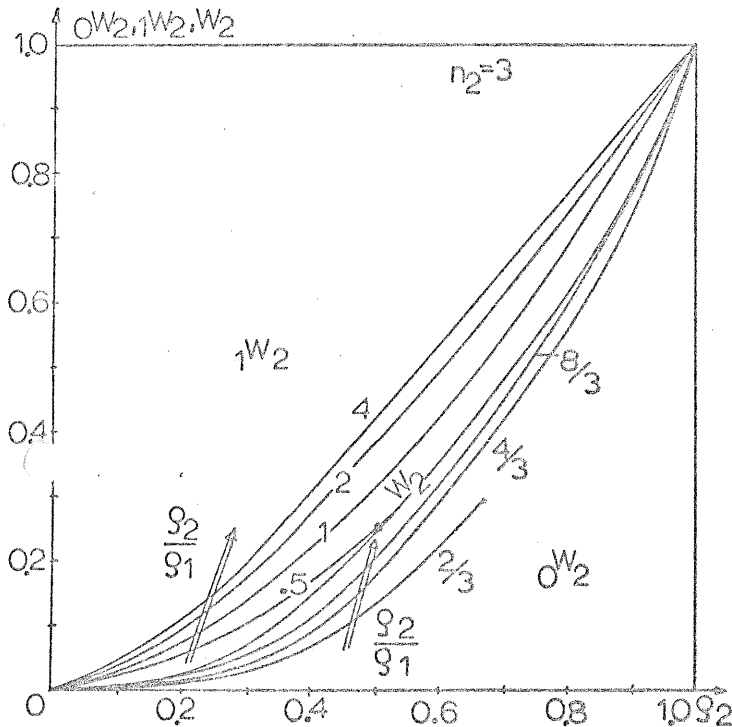


Fig. 2.2 Conditional probabilities of waiting

All curves for  ${}_0W_2$  are situated below  $W_2 = E_{2,n_2}(A_2)$ , all curves for  ${}_1W_2$  above.

## 2.6 CONDITIONAL WAITING TIME DISTRIBUTION FUNCTIONS IN STAGE 2

If it is further assumed that a considered test-call must wait in the second stage, for each test-call the same conditional meeting probabilities are obtained as for all waiting calls in stage 2 together:

$$P(X_{2,arr 2} = x | T_{W2} > 0) = \left(1 - \frac{A_2}{n_2}\right) \cdot \left(\frac{A_2}{n_2}\right)^{x-n_2} \quad x \geq n_2 \quad (2.22)$$

Therefore the following is true:

As soon as a call must wait in the second stage, the number of calls met and therefore (arbitrary queue discipline allowed) the further fate (waiting time, flow time) is independent of the previous fate in stage 1 (waiting and/or service time).

(With the results of chapter 4 it can be shown that this is only true, if no special number of calls met in the first stage  $X_{1,arr 1}(>0)$  is assumed.)

If, for example, the queue discipline in stage 2 is FIFO, each of these test calls waiting in the second stage has the same (complementary) conditional d.f. as all calls together:

$$P(T_{W2} > t | T_{W2} > 0) \stackrel{\text{def}}{=} W_{2W}(>t) = e^{-(\mu_2 - \lambda)t} \quad (2.23)$$

ropping the condition that the test-call has to wait in the second stage, it is obtained by (2.23) with (2.17) to (2.20)

$${}_cP(T_{W2} > t | T_{H1} = t_1) = \frac{E_{2,n_2}(A_2)}{\left(\frac{A_2}{n_2}\right)^c} \cdot e^{-(\mu_2 - \lambda)(t_1 + t)} \quad c=0,1; \quad (2.24)$$

$${}_cW_2(>t) = {}_cF \cdot E_{2,n_2}(A_2) \cdot e^{-(\mu_2 - \lambda)t} \quad (2.25)$$

Formula (2.25) for unknown service time in stage 1 yields for  $n_2=1$  the same expressions given by BURKE [2].

From (2.24) it is possible to determine the fate of a call with a certain required service time in the first stage by weighting summation:

$$P(T_{W2} > t | T_{H1} = t_1) = P(T_{W2} > t, T_{W1} = 0 | T_{H1} = t_1) + P(T_{W2} > t, T_{W1} > 0 | T_{H1} = t_1) \\ = \left(1 - S_1 + \frac{S_1}{S_2}\right) \cdot E_{2,n_2}(A_2) \cdot e^{-(\mu_2 - \lambda)(t_1 + t)} \quad (2.26)$$

## 3. DETERMINATION OF TOTAL FATES

The independence of the flow times in the two stages proved in the literature causes the total flow time d.f. to be a simple convolution of the two single stage flow time d.f.'s, cf. 1.2. Determining the total waiting time d.f., the dependencies of the waiting times in the single stages must be taken into account. This is now possible, because in the previous chapter these dependencies not only have been shown to exist but also completely determined, for known and for unknown service time as well.

### 3.1 TOTAL PROBABILITY OF WAITING

The total probability of waiting  $W$  is the probability that a call has to wait somewhere in the system.

$$W = P(T_W > 0) = P(T_{W1} > 0) + P(T_{W1} = 0) \cdot P(T_{W2} > 0 | T_{W1} = 0) \\ = P(T_{W1} > 0) + P(T_{W2} > 0) - P(T_{W1} > 0, T_{W2} > 0) \quad (3.1)$$

Using (2.19) or (2.20) the following form for  $W$  can be obtained

$$W = W_1 + W_2 - \frac{W_1 \cdot W_2}{S_1 + S_2 - S_1 S_2} \quad \text{where} \quad W_1 = S_1 = \frac{\lambda}{\epsilon_1}, \quad W_2 = E_{2,n_2}(A_2) \quad S_2 = \frac{A_2}{n_2} = \frac{\lambda}{n_2 \epsilon_2} \quad (3.2)$$

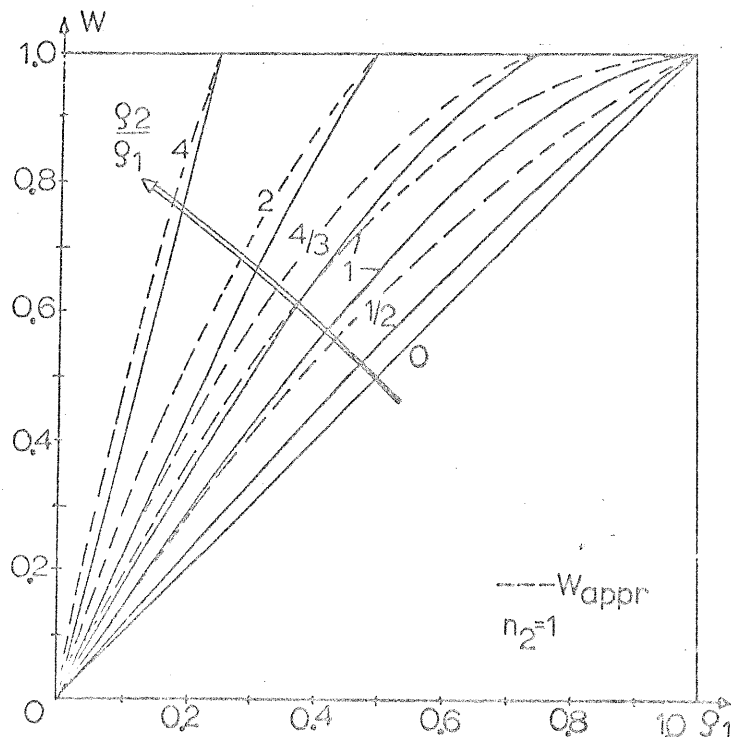


Fig. 3.1 Total probability of waiting

In fig. 3.1 the total probability of waiting is shown for  $n_2 = 1$ . For comparison, the total probability of waiting assuming independence is also depicted:

$$W_{\text{appr}} = W_1 + W_2 - W_1 \cdot W_2 \quad (3.3)$$

The total probability of waiting always is less than that value obtained by assuming independence of the waiting times, because preferably such calls wait in the second stage, which have already waited in the first stage.

An analogous relation to (3.1) yields with  $P(T_{W2} > 0 | T_{W1} = 0, T_{H1} = t_1)$  from (2.17) a formula describing the dependence of  $W$  from the service time in the first stage:

$$P(T_W > 0 | T_{H1} = t_1) = W_1 + (1 - W_1) \cdot W_2 \cdot e^{-(\mu_2 - \lambda)t_1} \quad (3.4)$$

### 3.2 MEAN TOTAL WAITING TIME OF THE WAITING CALLS

The mean total waiting time related to all calls is (independent of the sequence of the stages) the sum of the mean waiting time in the single stages related to all calls, even if the waiting times are dependent:

$$E(T_W) = E(T_{W1}) + E(T_{W2}) \quad (3.5)$$

So the mean total waiting time of the waiting calls is

$$t_W = \frac{E(T_W)}{W} \quad (3.6)$$

If, for example, 2 single server stages with  $h_1 = h_2 = h$  are in series ( $A_1 = A_2 = A$ ) it holds

$$t_W = \frac{2-A}{(1-A)(3-2A)} \cdot 2h \quad (3.7)$$

### 3.3 DISTRIBUTION FUNCTION OF THE TOTAL WAITING TIME

In a 2-stage system it generally holds

$$P(T_W > t) = P(T_{W2} = 0) \cdot P(T_{W1} > t | T_{W2} = 0) + P(T_{W1} = 0) \cdot P(T_{W2} > t | T_{W1} = 0) + \\ + P(T_{W1} > 0, T_{W2} > 0) \cdot P(T_{W1} + T_{W2} > t | T_{W1} > 0, T_{W2} > 0) \quad (3.8)$$

In this weighted summation of (complementary) conditional d.f.'s all weighting probabilities are known and  $P(T_{W2} > t | T_{W1} = 0)$  was derived in (2.25). The last d.f. is concerned with calls waiting in both stages. Because it was shown that in the considered system the random values of  $T_{W1}$  and  $T_{W2}$  are independent of each other if and only if both are  $> 0$ , convolution only is allowed under these additional conditions:

$$P(T_{W1} + T_{W2} > t | T_{W1} > 0, T_{W2} > 0) = P(T_{W1} > t | T_{W1} > 0) \cdot P(T_{W2} > t | T_{W2} > 0) \\ = e^{-(\epsilon_1 - \lambda)t} \cdot e^{-(\mu_2 - \lambda)t} \quad (3.9)$$

By the irrelevance of the concrete value of  $T_{W1} (> 0)$  it can be shown that

$$P(T_{W1} > t | T_{W2} = 0) = P(T_{W1} > 0 | T_{W2} = 0) \cdot P(T_{W1} > t | T_{W1} > 0) \quad (3.10)$$

So, finally, the d.f. of the total waiting time can be written as:

$$P(T_W > t) \stackrel{\text{def}}{=} W(t) = \begin{cases} \left\{ \frac{\lambda}{\epsilon_1} - \frac{\epsilon_1 - \lambda}{\epsilon_1 + \mu_2 - \lambda} \cdot \frac{\mu_2}{\epsilon_1 - \mu_2} \cdot W_2 \right\} \cdot e^{-(\epsilon_1 - \lambda)t} + \\ + \frac{\epsilon_1 - \lambda}{\epsilon_1 + \mu_2 - \lambda} \cdot \frac{\epsilon_1}{\epsilon_1 - \mu_2} \cdot W_2 \cdot e^{-(\mu_2 - \lambda)t} & \text{for } \epsilon_1 \neq \mu_2 \\ \left\{ \frac{\lambda}{\mu} + \frac{\mu - \lambda}{2\mu - \lambda} \cdot W_2 (1 + \mu t) \right\} \cdot e^{-(\mu - \lambda)t} & \text{for } \epsilon_1 = \mu_2 = \mu \end{cases} \quad (3.11)$$

$$(3.12)$$

In case of  $n_2 = 1$  a symmetrical expression in  $\epsilon_1$  and  $\epsilon_2$  can be achieved:

$$W(t) = \frac{\lambda}{\epsilon_1} \cdot \frac{\epsilon_2}{\epsilon_2 - \epsilon_1} \cdot \frac{\epsilon_2 - \lambda}{\epsilon_1 + \epsilon_2 - \lambda} \cdot e^{-(\epsilon_1 - \lambda)t} \\ + \frac{\lambda}{\epsilon_2} \cdot \frac{\epsilon_1}{\epsilon_1 - \epsilon_2} \cdot \frac{\epsilon_1 - \lambda}{\epsilon_1 + \epsilon_2 - \lambda} \cdot e^{-(\epsilon_2 - \lambda)t} \quad \text{for } \epsilon_1 \neq \epsilon_2 \quad (3.13)$$

If 2 single server stages are in series, the total waiting time d.f. is independent of the sequence of the two stages, i.e. the two stages are interchangeable in relation to the total waiting time. This term was used by REICH [10] for systems with sequence-independent total flow time distribution function.

It is obvious that the total waiting time d.f. for calls with a preassigned service time in the first stage is obtained by the very same procedure.

### 3.4 CORRELATION AND ERROR CONSIDERATIONS

Assuming independence of the waiting times, the approximate total waiting time d.f. simply is obtained by convolution:

$$W(t)_{\text{appr}} = P(T_{W1} > t) * P(T_{W2} > t) \quad (3.14)$$

As a simple example, for 2 equal single server stages ( $A_1 = A_2 = A, h_1 = h_2 = h$ ) it holds

$$W(t)_{\text{appr}} = A \left\{ 2 - A + A(1 - A) \cdot \frac{t}{h} \right\} \cdot e^{-(1-A) \cdot \frac{t}{h}} \quad (3.15)$$

The (complementary) d.f. according (3.12) is

$$W(t) = \frac{A}{2-A} \left\{ 3 - 2A + (1-A) \cdot \frac{t}{h} \right\} \cdot e^{-(1-A) \cdot \frac{t}{h}} \quad (3.16)$$

In fig. 3.2 these two d.f.s, which have the same expectation  $E(T_W)$ , are depicted.

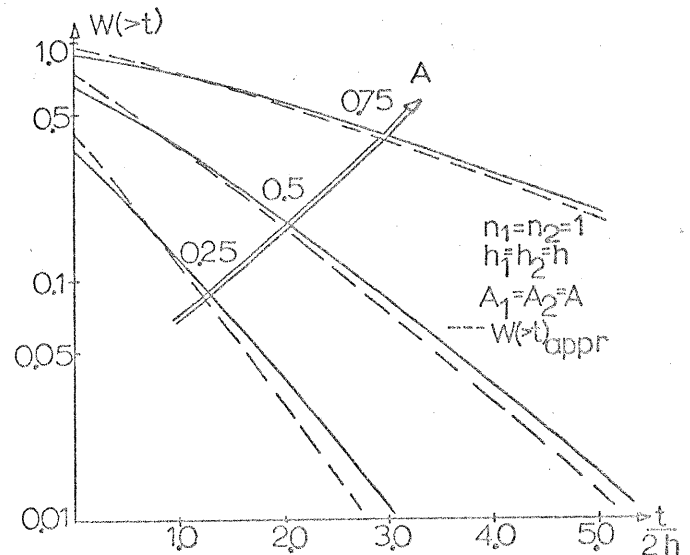


Fig. 3.2 Total waiting time distribution function

It is easily seen, that the exact curves yield, due to the dependence of the waiting times, a greater variance than the approximated ones.

We have

$$\text{var}(T_W) = \text{var}(T_{W1}) + \text{var}(T_{W2}) + 2\text{cov}(T_{W1}, T_{W2}) \\ \text{with } \text{cov}(T_{W1}, T_{W2}) = E(T_{W1} \cdot T_{W2}) - E(T_{W1}) \cdot E(T_{W2}) \\ = \frac{\lambda}{(\epsilon_1 - \lambda)(\epsilon_2 - \lambda)} \left\{ \frac{1}{\epsilon_1 + \epsilon_2 - \lambda} - \frac{\lambda}{\epsilon_1 \cdot \epsilon_2} \right\} \quad (3.17)$$

The correlation coefficient

$$r(T_{W1}, T_{W2}) = \frac{\text{cov}(T_{W1}, T_{W2})}{\sqrt{\text{var}(T_{W1})} \cdot \sqrt{\text{var}(T_{W2})}}$$

which is in case of  $\text{var}(T_{W1}) = \text{var}(T_{W2})$  identical with the relative increase of variance, when (positive) covariance is taken into account, is shown in fig. 3.3 for  $n_2 = 1$ . For comparison, a simple example for  $n_2 > 1$  is also depicted.

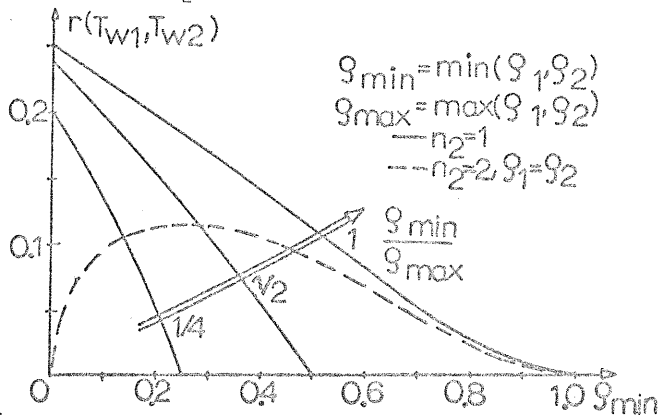


Fig. 3.3 Correlation coefficients

Determining the relative errors  $(W_{\text{appr}} - W)/W$ , we will get in principle the same traffic dependence as for the correlation coefficients. So only some calculated maximum values of this relative error as a function of the number  $n_2$  of servers in the second stage are given here:

$n_2$	1	2	3	6	10
possible error	33%	13%	9%	6%	4%

For  $n_2 > 1$  these errors can be obtained, when the utilizations are such, that the probabilities of waiting in both stages are nearly equal and in the range of 0.2 to 0.3.

#### 4. TEST-CALLS WITH GIVEN STARTING POSITION IN STAGE 1 (RANDOM WALK)

In this chapter the conditional meeting probabilities

$$P(X_{2, \text{arr}} = x | X_{1, \text{arr}} = x_1) \stackrel{\text{def}}{=} p_2(x | x_1)$$

are considered, that is to say, the dependence of the number of calls (queue lengths) met upon arrivals of the same call in the two stages. In order to obtain explicit results, the case of two single server stages is considered. (It is clear that  $X_{1, \text{arr}} = 0$  is identical with  $T_{W1} = 0$ , a condition considered in chapter 2.)

##### 4.1 GENERAL WAY OF CALCULATION

The walk of a call through the system may be described by a sequence of flow-states, the call is engaged with (path in a RANDOM WALK diagram). Since the fate of the considered call is not influenced by succeeding calls (FIFO), the RANDOM WALK diagram is a directed graph without loops (fig. 4.1).

A general flow-state  $[i_1, i_2]$  related to a considered call is defined such, that  $i_1$  calls are in the first and  $i_2$  in the second stage, where succeeding calls are irrelevant, but including the call in question. If this flow-state exists, the next event is either the ending of service in the first or second stage with probabilities

$$p_1 = \frac{\epsilon_1}{\epsilon_1 + \epsilon_2}; \quad p_2 = \frac{\epsilon_2}{\epsilon_1 + \epsilon_2} \quad (4.1)$$

(a stage is empty, double arrows show that the single-step probability is equal to 1 (reflecting barrier)).

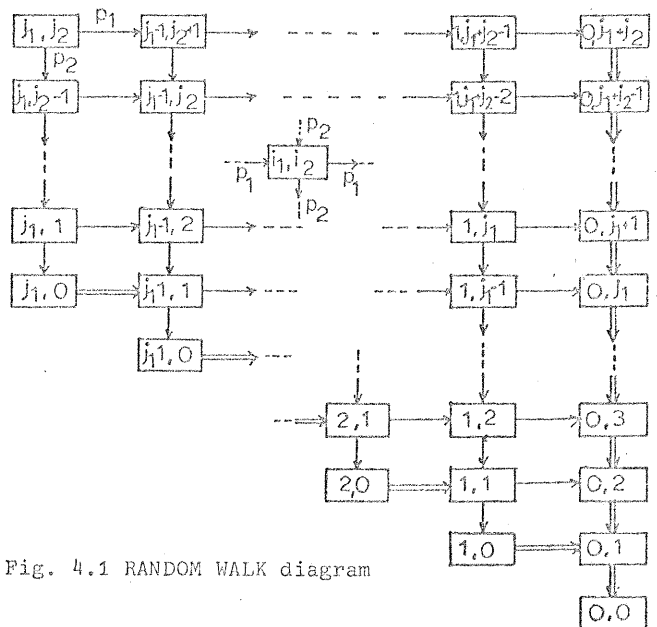


Fig. 4.1 RANDOM WALK diagram

Each call starts its walk at the starting flow-state  $[j_1, j_2]$  and moves on a certain path with certain probability to the absorbing flow-state  $[0, 0]$  where the call leaves the system. It is not difficult to quote the d.f. of the time a call spends in a flow-state; then it is obvious, how  $T_{W1}$ ,  $T_{H1}$ ,  $T_{W2}$  and  $T_{H2}$  are reflected in this diagram.

Since this RANDOM WALK diagram now contains all necessary information, it would be possible in principle to make also time considerations concerning  $T_{W1}$  as in chapter 2, but the method used there proved to be more effective for this case.

Let be:

- $p_F(i_1, i_2)$  the flow-state probability, that the walk of a call from  $[j_1, j_2]$  to  $[0, 0]$  touches  $[i_1, i_2]$  ( $j_1, j_2$  implicitly understood)
- $d_i$  the number of ' $p_i$ -transitions' of a path ( $p_i$ -distance),  $i = 1, 2$
- $d_0$  the number of transitions of a path without alternative.

Then the probability of a certain path is equal to

$$p_{\text{path}} = 1^{d_0} \cdot p_1^{d_1} \cdot p_2^{d_2} \quad (4.2)$$

and the flow-state probability is equal to the sum of probabilities of all paths leading from  $[j_1, j_2]$  to the considered flow-state.

If a call uses the transition from  $[1, x]$  to  $[0, x+1]$ , it meets exactly  $x$  calls upon arrival in the second stage. Thus it is sufficient to calculate the flow state probabilities  $p_F(1, x)$  ( $x = 0 \dots j_1 + j_2 - 1$ ), from which directly the conditional meeting probabilities for a certain starting pattern  $\{x_1, x_2\}$

$$P(X_{2, \text{arr}} = x | X_{1, \text{arr}} = x_1, X_{2, \text{arr}} = x_2) \stackrel{\text{def}}{=} p_2(x | x_1, x_2)$$

are obtained. Finally a weighted summation over  $x_2$  yields  $p_2(x | x_1)$ .

##### 4.2 FLOW-STATE PROBABILITIES

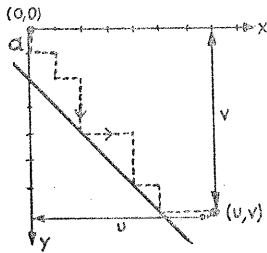
For a fixed starting flow-state  $[j_1, j_2]$  the probability of a path to a fixed flow-state  $[1, x]$  depends on its number  $d_i$  of flow-states with idle time for stage 2. For  $0 \leq x < j_1 + j_2 - 1$  there is no

path which touches such an idle state. Therefore, all possible paths from a fixed starting position have the same probability and the number of paths is simply to quote. For  $0 < x \leq j_1 - 1$  there are paths having  $d_0 = 0, 1, \dots, d_{\text{omax}}$  idle states. If  $k(d_0, x)$  is the number of paths from flow-state  $[j_1, j_2]$  ( $j_1 > 1, j_2 \geq 0$ ) to a flow state  $[1, x]$  ( $x > 0$ ), which touches  $d_0$  idle states, so it holds:

$$p_F(1, x) = \sum_{d_0=0}^{d_{\text{omax}}} k(d_0, x) \cdot p_1^{d_1} \cdot p_2^{d_2} \quad (4.3)$$

with  $d_1 = j_1 - 1 - d_0, d_2 = j_1 - 1 + j_2 - x, d_{\text{omax}} = j_1 - x$ .

To determine  $k(d_0, x)$ , the following definition and lemma is made (fig. 4.2):



Let  $\Psi(d_0, u, v, a)$  be the number of paths leading from  $(0,0)$  to  $(u,v)$ , which touch but not cross the line  $y = x + a$  exactly  $d_0$  times.  
( $a \geq 0, u > 0, 0 \leq v < u + a$ )

Fig. 4.2 Directed path

LEMMA:

$$\Psi(d_0, u, v, a) = \begin{cases} \binom{u+v}{u} - \binom{u+v}{u+a} & \text{for } d_0 = 0 \end{cases} \quad (4.4.1)$$

$$\Psi(d_0, u, v, a) = \begin{cases} \binom{u+v-d_0}{u+a-1} - \binom{u+v-d_0}{u+a} & \text{for } 0 < d_0 \leq v - a + 1 \end{cases} \quad (4.4.2)$$

PROOF: For  $a > 0, u > 0, 0 \leq v < u + a$ , (4.4.1) is identical with lemma 2 of MILCH and WAGGONER [6] and may be proved, also for  $a \geq 0, u \geq 0, 0 \leq v \leq u + a$ , with the reflection principle, cf. e.g. FELLER [4]. (4.4.2) is for the same conditions part of corollary 2 in [6] and has been proved by the so-called telescope principle. Because of  $\Psi(d_0, u, v, 0) = \Psi(d_0 - 1, u - 1, v, 1)$  for  $d_0 > 0$ ,  $a = 0$  is also allowed. Then the validity for  $v = 0$  is obvious.

Going back to the random walk diagram with

$$k(d_0, x) = \Psi(d_0, j_1 - 1, j_1 - 1 + j_2 - x, j_2) \quad (4.5)$$

an explicit general formula for the flow-state probabilities  $p_F(1, x)$  can be derived.

#### 4.3 CONDITIONAL MEETING PROBABILITIES

If a call starts its walk with flow-state  $[j_1, j_2]$ , upon its arrival  $x_1 = j_1 - 1$  respectively  $x_2 = j_2$  calls already have been in the system. So with

$$p_2(x | x_1, x_2) = p_F(1, x) \cdot p_1 \quad \text{for } x > 0 \quad (4.6)$$

(4.3)-(4.5) yield

$$p_2(x | x_1, x_2) = p_2^{x_1 + x_2 - x} \left\{ p_1^{x_1 + 1} \binom{2x_1 + x_2 - x}{x_1} + p_1^x \sum_{i=0}^{x_1 - x} p_1^i \left\{ \binom{x_1 + x_2 - 1 + i}{x_1 + x_2 - 1} - p_1 \binom{x_1 + x_2 + i}{x_1 + x_2} \right\} \right\} \quad (4.7)$$

If  $x > x_1$  the sum should be replaced by 0.

With  $\binom{-1}{-1} \stackrel{\text{def}}{=} 1$ , it can be shown that (4.7) generally holds for  $x_1, x_2 \geq 0, 0 \leq x \leq x_1 + x_2$ .

REMARK: These meeting probabilities for a certain starting pattern  $\{x_1, x_2\}$  may also be interpreted as state probabilities of a single stage M/M/1 with arrival rate  $\lambda_1$  and service rate  $\mu_2$  ( $> \text{or } \leq \lambda_1$ ) upon the arrival of call number  $x_1 + 1$ , when at the begin of the arrival process exactly  $x_2$  calls have been in the system. (Investigation of 'time'-dependence where 'time' is represented by the ordinal number of the arriving call.)

This kind of consideration was made by TAKÁCS [12], who derived with the theory of homogeneous Markov-chains an extensive expression for the bivariate generating function of these higher transition probabilities. For a rather simplifying special case of this function an explicit result was stated, which yields transferred to this problem

$$p_2(x | x_1, 0) = \left( \frac{p_1}{p_2} \right)^x \cdot \sum_{i=x}^{x_1} \frac{i}{2x_1 - i} \left( \frac{p_1}{x_1} \right)^{x_1 - i} \cdot p_1^{x_1} \cdot p_2^{x_1} \quad (4.8)$$

By complete induction, accordance between (4.8) and (4.7) with  $x_2 = 0$  can be shown.

By weighting summation over  $x_2$  the wanted conditional meeting probabilities are obtained:

$$p_2(x | x_1) = \sum_{x_2=0}^{\infty} P(X_{2, \text{arr}} = x_2 | X_{1, \text{arr}} = x_1) \cdot p_2(x | x_1, x_2) \quad (4.9)$$

So the general formula is

$$p_2(x | x_1) = (1 - A_2) \cdot p_2^{x_1 - x} \cdot \left\{ [A_2 p_2]^{j_2} \left\{ p_1^{x_1 + 1} \binom{2x_1 + j_2 - x}{x_1} + p_1^x \sum_{i=0}^{x_1 - x} p_1^i \left\{ \binom{x_1 + j_2 - 1 + i}{x_1 + j_2 - 1} - p_1 \binom{x_1 + j_2 + i}{x_1 + j_2} \right\} \right\} \right\} \quad (4.10)$$

The sum over  $i$  should be replaced by 0 if  $x > x_1$ .

In fig. 4.3 these meeting probabilities are plotted and compared with the absolute meeting probability  $p_2(x)$ . In this example, where the second stage is slower than the first stage ( $\lambda_2 < \lambda_1$ ), it is vividly shown, how momentary variations (traffic peaks) in the first stage are continued later on in the second stage.

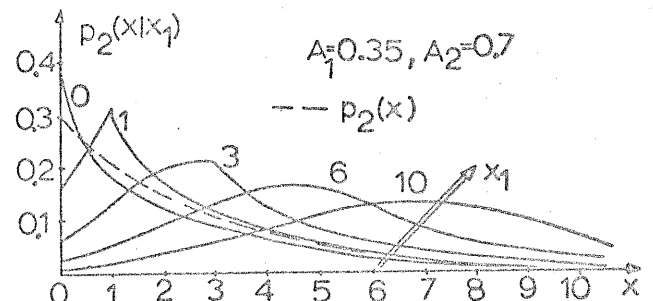


Fig. 4.3 Conditional meeting probabilities

It is clear, that in case of  $\lambda_2 > \lambda_1$  the influence of  $x_1$  to  $p_2(x | x_1)$  is weaker and that  $p_2(x | x_1)$  must tend for  $x_1 \rightarrow \infty$  to the absolute state probabilities of a M/M/1 system with arrival rate  $\lambda_1$  and service rate  $\mu_2$ . If  $x > x_1$ , no idle situation is possible for stage 2 and (4.10) can be simplified to

$$p_2(x | x_1) = (1 - A_2) \cdot A_2^{x_1 + 1} \cdot \left( \frac{1}{A_1 + A_2 - A_1 A_2} \right)^{x_1 + 1} \quad (4.11)$$

which is for  $x_1 = 0$  identical with  $p_2^{(1)}(x)$  according (2.11), derived in chapter 2. Equation (4.10) shows that, for instance, the probability of waiting in the second stage is the higher, the greater the number  $x_1$  of calls met in the first stage. This means that here no such 'limited dependency' is valid as obtained for the waiting times.

#### 4.4 FURTHER FATE VALUES

The number of calls met in the second stage is completely sufficient for the determination of the further fate in stage 2, as there are the conditional waiting and flow time distributions. E.g. the mean waiting time and the mean flow time can

be simply deduced from the conditional expectation

$$E(x|x_1) \stackrel{\text{def}}{=} E(x_{2, \text{arr}2} | x_{1, \text{arr}1} = x_1) = \sum_{x=0}^{\infty} x \cdot p_2(x|x_1) \quad (4.12)$$

which are shown in figs. 4.4 and 4.5.

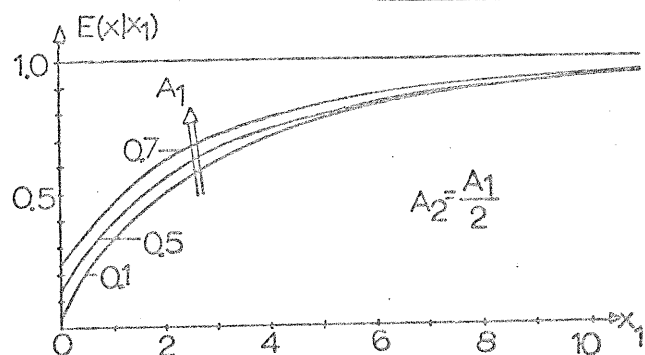


Fig. 4.4 Conditional expectations in stage 2 ( $\epsilon_2 > \epsilon_1$ )

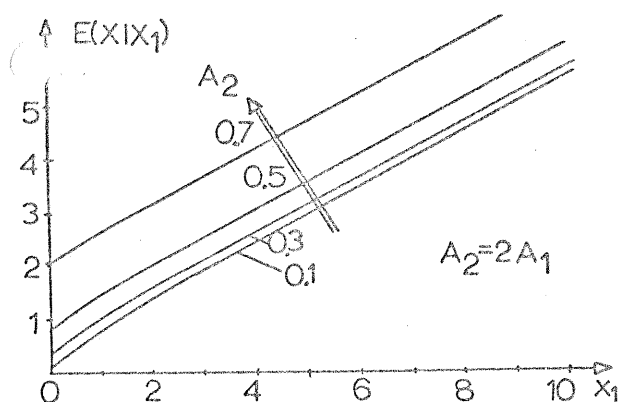


Fig. 4.5 Conditional expectations in stage 2 ( $\epsilon_2 < \epsilon_1$ )

By these curves the following behaviours (obtained by plausible separate arguments) are confirmed:  
If  $\epsilon_2 > \epsilon_1$ , according to a quasistationary behaviour

$$\lim_{x \rightarrow \infty} E(x|x_1) = \frac{A_2'}{1-A_2'} \text{ with } A_2' = \frac{\epsilon_1}{\epsilon_2} \quad (4.13)$$

whereas in case of  $\epsilon_2 < \epsilon_1$  the expectation value nearly linearly increases with gradient

$$\frac{E(x|x_1 + \Delta x_1) - E(x|x_1)}{\Delta x_1} \sim 1 - \frac{\epsilon_2}{\epsilon_1} \quad (4.14)$$

since the probability that the second stage is idle tends to 0 for sufficiently high  $\epsilon_1$  and  $x_1$ .

Considering total fates as in chapter 3, besides the total probability of waiting and the inherent mean total waiting time of the waiting calls, the total number of predecessors in stage 1 and 2

$$P(x_{1, \text{arr}1} + x_{2, \text{arr}2} = x) = \sum_{j=0}^x P(x_{1, \text{arr}1} = x-j) p_2(j|x_j) \quad (4.15)$$

further can be determined, which is equivalent to the number of phases a total waiting time is composed of.

## CONCLUSION

For a system with two queues in series the fate of a call in the second stage was determined as a function of all components of the previous fate in stage 1, existing dependencies were illuminated and calculated. Therefore, it is possible to give a more detailed fate prediction of special calls (e.g. with certain required service time).

It was shown that in relation to the waiting times the dependency is limited, so that as practical result the total waiting time d.f. could be determined.

The error made by assuming independence was shown, which cannot always be accepted in this system. Moreover, the dependencies including the correlation coefficient may be considered to give orientation values for corresponding systems, if service time d.f. has no memoryless property as investigated; also it is hoped that these results may give more insight into these problems if the structure is more complicated.

The results given for the total waiting time d.f. are valid for systems with  $n_1=1$  and  $n_2 \geq 1$ . In case of  $n_1 > 1$  the fate in the second stage depends on the concrete value of the waiting time in the first stage, because of the possibility of overtaking in the first stage.

Finally, it may be noted that it is possible to extend the results of chapters 2 and 3 to a system with several parallel queuing systems in the second stage (tandem queues with group selection).

The author wants to express his thanks to the Deutsche Forschungsgemeinschaft (German Research Society) which supported this work, being part of a research project.

## REFERENCES

- [1] BURKE, P.J., The Output of a Queuing System. J.Op.Res.4 (1956), 699-704.
- [2] BURKE, P.J., The Dependence of Delays in Tandem Queues. Ann.Math.Stat.35 (1964), 874-875.
- [3] BURKE, P.J., The Output Process of a Stationary M/M/s Queuing System. Ann.Math.Stat.39 (1968), 1144-1152.
- [4] FELLER, W., An Introduction to Probability Theory and Its Applications. Vol. 1, 2nd ed. John Wiley, New York, 1957.
- [5] JACKSON, R.R.P., Random Queuing Processes with Phase-Type Service. J.R.S.S. Ser.B.18 (1956), 129-132.
- [6] MILCH, P.R., WAGGONER, M.H., A Random Walk Approach to a Shutdown Queuing System. SIAM J. Appl.Math. 19 (1970) 1, 103-115.
- [7] NELSON, R.T., Waiting Time Distributions for Application to a Series of Service Centers. J.Op.Res.6 (1958), 856-862.
- [8] REICH, E., Waiting Times when Queues are in Tandem. Ann.Math.Stat.28 (1957), 768-773.
- [9] REICH, E., Note on Queues in Tandem. Ann.Math. Stat.34 (1963), 338-341.
- [10] REICH, E., Departure Processes. Proceedings of the Symposium on Congestion Theory (1964), North Carolina, 439-457.
- [11] SYSKI, R., Introduction to Congestion Theory in Telephone Systems. Oliver and Boyd, London, 1960.
- [12] TAKÁCS, L., Introduction to the Theory of Queues. Oxford University Press. New York, 1962.