## Deadline-Oriented Servicing: Waiting-Time Distributions

#### B. WALKE AND W. ROSENBOHM

Abstract—An infinite queue single server model is considered where requests arrive from independent Poisson streams and demand service according to arbitrary distribution functions which may be different for different requests. Associated with each request is an urgency numer which, together with request's time of arrival, defines a deadline for beginning its service. This relative urgency discipline has at its two limiting cases the first-come first-serve and head of the line discipline. In [1] the mean waiting time is computed approximately and close bounds are derived there. Here we present simulation results, derive close approximations for the tails of the waiting-time distribution functions and compare them to those of the two limiting cases.

Index Terms-Computer performance evaluation, queuing analysis.

#### I. INTRODUCTION AND DEFINITION OF MODEL

In real-time computer-control systems it is important that incoming service requests from a running process be completed in time, i.e., within a given time limit. In this short note we will consider the case of requests arriving at random (not predictable) times. Given the randomness it is impossible

Manuscript received December 6, 1978; revised August 25, 1979. This work was partially supported by the 3rd EDP Program of the Federal Government of West Germany under Contract DV 0812248.

The authors are with AEG-TELEFUNKEN, Research Institute of Ulm, Ulm, West Germany.

to guarantee that all finite deadlines for beginning or ending these requested services can be met. The best that can be done is to guarantee a high probability for meeting such deadlines.

Let us consider a computer model with unlimited waiting space for incoming requests which may possess various properties (to be denoted in the following as "types"). Further, suppose these type-i requests  $(1 \le i \le N)$  arrive at a rate  $\lambda_i$  from a Poisson process, have an arbitrary service-time  $b_i$  according to a distribution function  $F_i(t) = P(b_i \le t)$  with a finite second moment, and that each request has a deadline  $t_i$  for beginning based on its arrival time  $T_i$  and its urgency  $\omega_i$ . Without any loss of generality we order the types i of requests such that

$$0 \leq \omega_1 < \omega_2 < \dots < \omega_i < \dots < \omega_N \tag{1}$$

where N is the total number of types and  $\omega_i$  is the urgency of a type-i request. The smaller the value of  $\omega_i$ , the higher the urgency.

Fig. 1 shows a model of our system. We will be concerned almost exclusively with the case of requests which, once servicing has begun, may not be interrupted (nonpreemptive priority). After a request has been serviced, the request among those waiting to be processed which has the highest dynamic priority  $q_i(t)$  is chosen for servicing. The dynamic priority at time t is given by

$$q_i(t) = \omega_i - t + T_i, \tag{2}$$

where  $T_i$  is the arrival time and  $\omega_i$  the urgency. If we define the waiting time  $w_i(t)$  at time t as

$$w_i(t) = t - T_i \tag{3}$$

0098-5589/80/0500-0304\$00.75 © 1980 IEEE

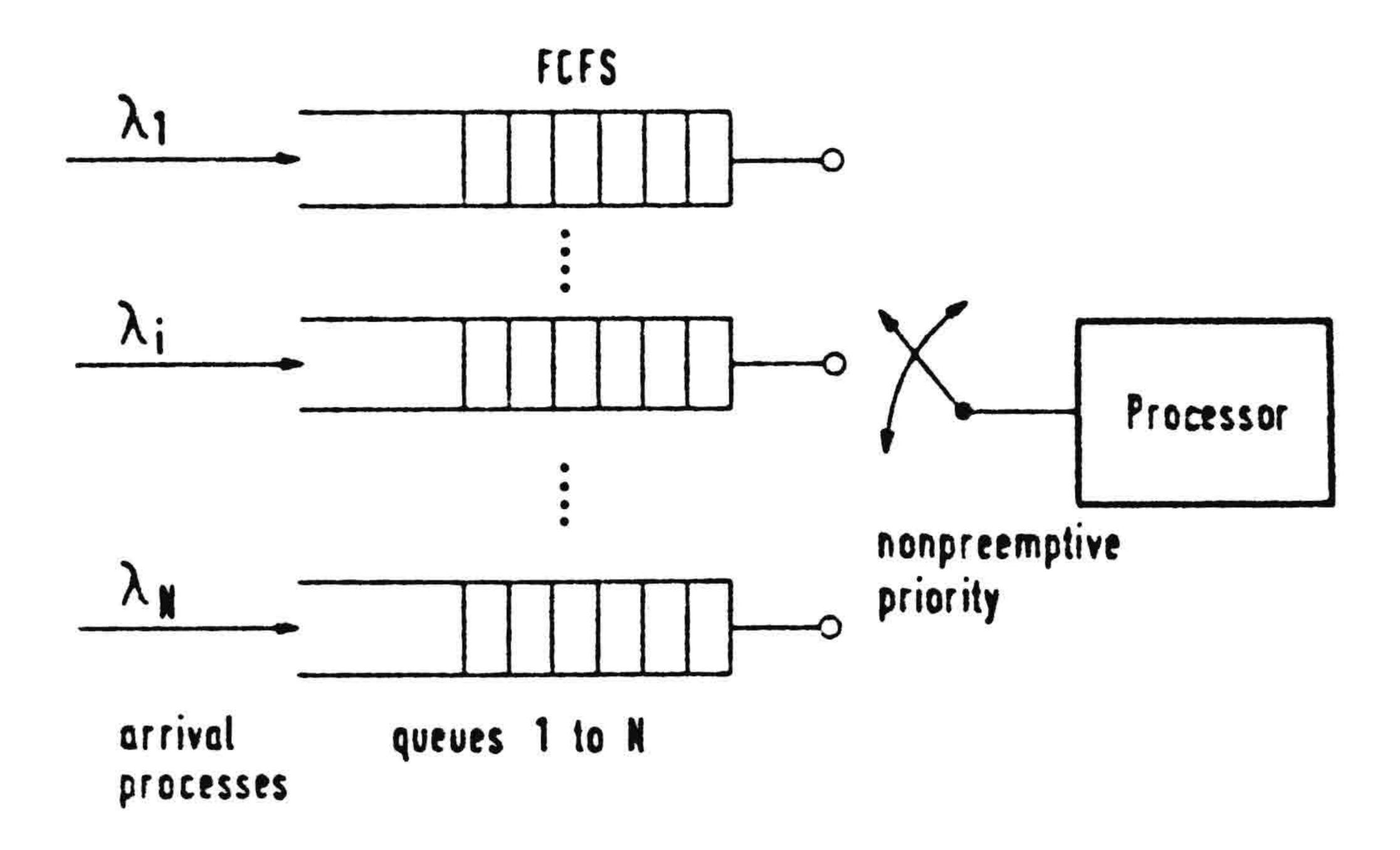


Fig. 1. Model of a real-time computer-control system,  $\lambda_i$  = arrival rate.

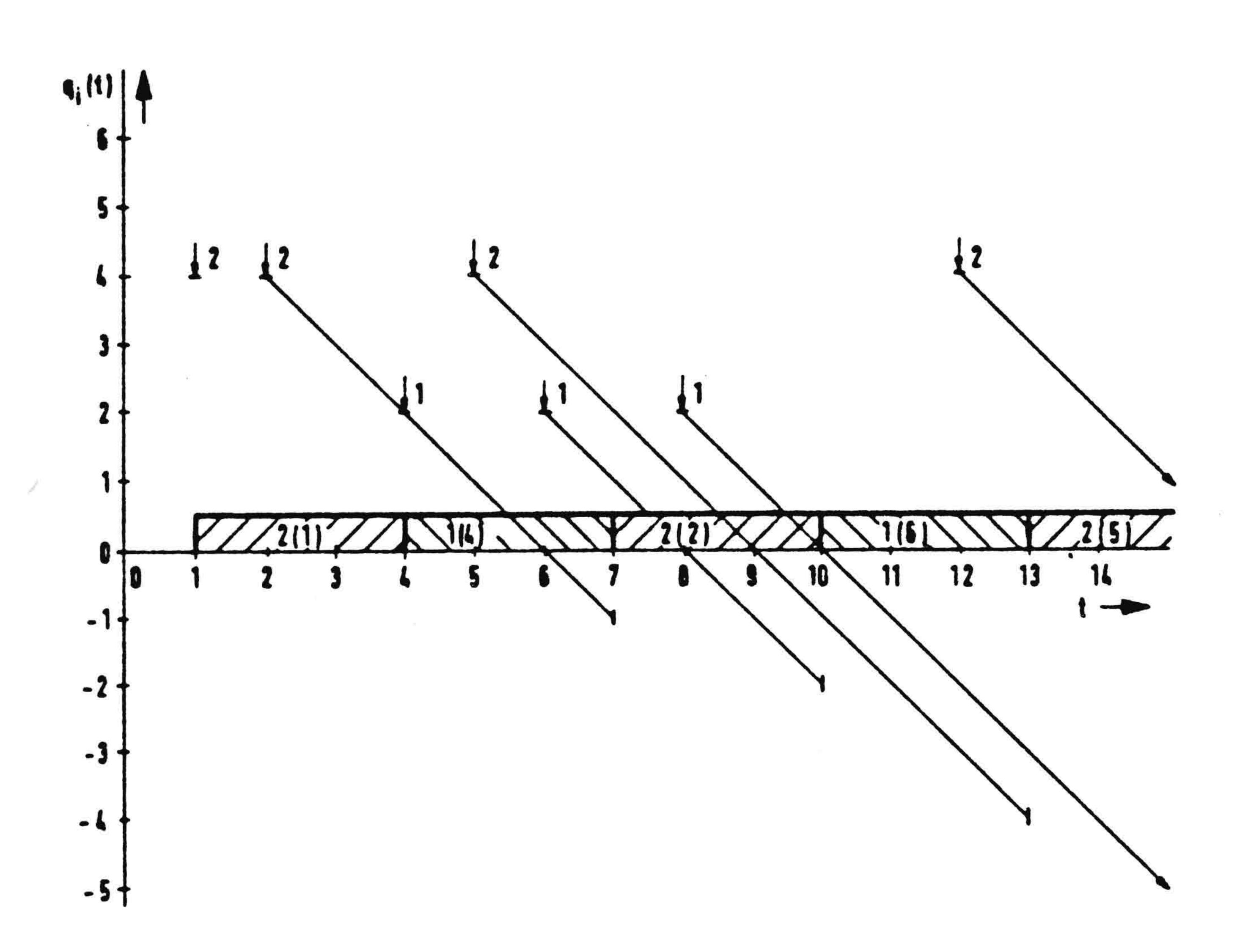


Fig. 2. Example showing the dynamic priority  $q_i(t)$  over time t(N = 2, $\omega_1 = 2$ ,  $\omega_2 = 4$ , equal and constant service times). Times of arrival of type-i requests are marked with an indexed arrow  $\downarrow i$ .

then the current priority is

$$q_i(t) = \omega_i - w_i(t) \tag{4}$$

and depends linearly only on the current waiting time. As usual, small values of  $q_i(t)$  indicate a high priority. It is interesting to note that in this model the priorities of all outstanding requests, regardless of type, grow uniformly. Requests of the same type are processed in the order they arrive [firstcome first-serve (FCFS)]. Since requests are serviced according to individual urgencies and arrival times, we speak of a relative-urgency (RU) discipline. One possible interpretation is that a request's deadline is reached if its waiting time  $w_i(t)$ equals ω<sub>i</sub>.

Fig. 2 shows an example with N = 2 types of requests with urgencies  $\omega_1 = 2$ ,  $\omega_2 = 4$ , and constant and equal service times  $\beta_1 = \beta_2 = 3$ . The ordinate shows the actual dynamic priority  $q_i(t)$  and the abscissa shows the time t. At its arrival time  $t = T_i$  each request has a priority  $q_i(0) = \omega_i$  because  $w_i(t) = 0$ . The priority of a waiting request increases linearly [the value of  $q_i(t)$  decreases] with its waiting time  $w_i(t)$ . Service requests arrive at times t = 1, 2, 4, 5, 6, 8, 12. Their individual priority functions are shown. The server occupation by a service request is shown by means of a hatched bar with intervals denoted by  $i(T_i)$ , the type i, and arrival time  $T_i$  of an individual request.

At t = 1 a busy period starts. The next decision is needed at t = 4 when service of the first request is finished. At this time both the type-2 request which arrived at t = 2 and the type-1 when  $\rho > 1$ . However, we will not consider this case here.

request which just has arrived have the same priority. Such ties are broken by servicing the most urgent type, i.e., type 1. At t = 7 the next decision is needed. At this time the type-2 request which arrived at t = 2 will be elected, although a type-1 request is waiting also. This follows from a comparison of the actual priorities at t = 7. Apparently less urgent requests are preferred if they have waited long enough.

In addition to the nonpreemptive discipline (RU-NPRE) in Section IV we will consider a few examples of a preemptive discipline (RU-PRE). The following abbreviations are used throughout:

d.f.  

$$W_{i}(t) = P(w_{i} \leq t)$$

$$W_{i}^{(r)}, W_{i}^{(1)} = W_{i}$$

$$\lambda = \sum_{j=1}^{N} \lambda_{j}$$

$$W(t) = \lambda^{-1} \sum_{i=1}^{N} \lambda_{i} W_{i}(t)$$

$$F_{i}(t) = P(b_{i} \leq t)$$

$$\beta_{i}^{(r)}, \beta_{i}^{(1)} = \beta_{i}$$

$$\rho_{i} = \lambda_{i} \beta_{i}$$

$$\rho \leq i = \sum_{j=1}^{i} \rho_{j}$$

$$\rho = \rho \leq N$$

$$F(t) = \lambda^{-1} \sum_{i=1}^{N} \lambda_{i} F_{i}(t)$$

$$\beta_{i}^{(r)} = \lambda^{-1} \sum_{i=1}^{N} \lambda_{i} \beta_{i}^{(r)}$$

$$C^{2} = \beta^{(2)}/\beta^{2} - 1$$

$$\overline{\omega} = \rho^{-1} \sum_{i=1}^{N} \rho_{i} \omega_{i}$$

$$\operatorname{var}(\omega) = \sum_{i=1}^{N} (\rho_{i}/\rho) \omega_{i}^{2} - \overline{\omega}^{2}$$

$$M/A/1/B$$

 $w_i, b_i$ 

and service time of a type-i request; distribution function; waiting-time d.f. of a type-i request; rth moment of  $W_i(t)$ ; total arrival rate of all requests together; common waiting time d.f. of all requests together; service-time d.f.; service moment time d.f.; offered traffic of type-i requests; offered traffic of types 1 to i; total offered traffic; common service-time d.f. of all requests together; rth moment of F(t); squared coefficient of variance of F(t); mean weighted urgencies; variance of weighted urgencies; model with Markovian (= Poisson) arrival process, a distinct service-time d.f. (A = D) deterministic, A =M Markovian, A = G general), one server, and service discipline B; head of the line (static

nonpreemptive priority)

discipline.

random variables: waiting

HOL

The RU discipline uses dynamic priorities based on waiting times. The problem described at the beginning of this paper nowadays is usually treated with a service discipline which recognizes only static priorities. Arriving requests are assigned a static priority i according to their type i. Among the waiting requests the one with the highest priority (lowest i) which has waited the longest of all requests of the same type i (FCFS) is processed first. Here, too, the model in Fig. 1 applies, except that there is no urgency  $\omega_i$ .

If one is interested in the probability of missing deadlines, then one should know the waiting-time d.f.  $W_i(t) = P(w_i \le t)$ . From this the probability that the waiting time wi is not greater than a given time t can be obtained.

We will consider the model in Fig. 1 only in a state of equilibrium, which occurs only when the total offered traffic  $\rho < 1$ . In disciplines with static priorities, a stationary equilibrium is achieved for requests with high priorities, sometimes even

### II. WAITING-TIME D.F.'s IN THE RU AND HOL DISCIPLINES; EXAMPLES

Let us assume N = 4 different types of requests, each with its arrival rate  $\lambda_i$ , urgency  $\omega_i$ , and service-time d.f.  $F_i(t)$ . Since the waiting-time d.f. cannot be calculated, we will consider two examples which were simulated. We used an M/D/1 model with all requests having a constant service time.

Example 1:  $\rho_i = \rho/N$ ,  $\beta = 2.39$ , C = 0.81,  $\overline{\omega} = 30$ , var  $(\omega) =$ 125

i	1	2	3	4
$\beta_i$	1	3	5	7
$\omega_i$	15	25	35	45

Example 2: corresponds to Example 1 but with inverse order of the mean service time  $\beta_i$  with regard to the type i.

	1	2	3	4
$eta_i$ $\omega_i$	7 15	<b>5 25</b>	3 5	1 45

Note that our findings stated in the rest of this paper are also valid for quite different load parameters of our model. This can be seen from [4] where results of M/G/1 and M/M/1models are presented.

From earlier investigations, cf. [1] and [2], it is known that the smaller the variance of the urgency numbers is, the less the mean waiting times w; of the RU-NPRE discipline differ from is fulfilled, then the probabilities of meeting the deadlines the mean waiting time

$$W_{\text{FCFS}} = 1/2\lambda\beta^{(2)}/(1-\rho) \tag{5}$$

in the first-come first-serve discipline. On the other hand if the variance var  $(\omega)$  is very large, the RU-NPRE discipline produces mean waiting times approximating those of the head of the line (HOL) discipline, a discipline with static nonpreemptive priorities

$$W_{i \text{ HOL}} = \lambda \beta^{(2)} / [2(1 - \rho_{\leq i})(1 - \rho_{\leq i})]. \tag{6}$$

extreme cases are known and their first two moments can easily be calculated.

Assume  $p(w_i)$  to be the equilibrium probability that a typei request must wait no longer than wi units of time for service. For  $0 \le f \le 1$ , define

$$Q_i(f) = \inf \{ w_i | p(w_i) \ge f \}. \tag{7}$$

This is the waiting time corresponding to fractile f in the cumulative waiting-time d.f. of type-i requests.  $Q_i(f)$  usually having a different variance var  $(\omega)$ . If var  $(\omega)$  is small the reis called the f-quantile.

#### A. Observations from a Simulation Model

From a simulation model we have made three observations.

First observation: the simulation shows that the waitingtime d.f.'s, when plotted semilogarithmically, asymptotically approach parallel straight lines regardless of the example and the individual parameters. This seems to indicate that for every M/G/1/RU-NPRE model we have

$$\lim_{f \to 1} [Q_i(f) - Q_j(f)] = \omega_i - \omega_j, \quad (i, j = 1, 2, \dots N). \quad (8)$$

cally approaches the difference in their urgencies  $\omega_i$  and  $\omega_j$  within the accuracy obtainable by simulation. For f < 1, this distance is smaller than the difference  $\omega_i - \omega_i$ .

Second observation: the waiting-time d.f.'s in the RU discipline only depend on the differences of the urgency numbers, i.e., for both parameter sets  $\{\omega_1, \omega_2, \dots, \omega_N\}$  and  $\{\omega_1 + r, \omega_2 + r, \cdots, \omega_N + r\}$  (r real number) the same d.f.'s appear.

Third observation: the waiting-time d.f. in the FCFS discipline runs parallel to the tails of the waiting-time d.f.'s in the RU discipline. This and the two other observations seem to be valid for every M/G/1/RU model. They are the basis for an approximate analytical computation of the waiting-time d.f.'s in the RU discipline.

An example: Consider an example with  $p = W_i(t) = p = W_i(t) = P$ 0.985 and  $\rho = 0.75$ . From Fig. 3 it can be seen that the pquantile of a type-i request, for which the abbreviation  $Q_i(p)$ just has been introduced, equals the urgency number  $\omega_i$ . Apparently for any other probability  $f(f \ge p)$  the following approximation holds:

$$Q_i(f) - Q_i(f) = w_i - w_i \approx \omega_i - \omega_i. \tag{9}$$

Observe that  $f = P(w_i \leq Q_i(f)) = P(w_j \leq Q_j(f))$ .

Very large waiting times of different type requests differ by the constant difference in their urgencies. Smaller waiting times differ by less. The property formulated in (8) implies that for very large waiting times  $w_i$  the probability  $P(w_i \leq t)$ of meeting deadlines is practically independent of the type of request, i.e., all requests are handled with equal fairness.

If the condition

$$\min_{i} \{\omega_i\} >> \max_{i} \{\beta_i\}, \quad (1 \leq i \leq N),$$

 $P(w_i \leq \omega_i)$  are about the same, regardless of type.

#### (5) B. Dependency on the Total Offered Traffic

Fig. 3 shows simulation results of the waiting-time d.f.'s for Example 1. Two different total offered traffics  $\rho = 0.75$ and 0.85 are assumed. If the total offered traffic is increased, the probabilities for meeting a given set of waiting times  $\{\omega_1, \omega_2, \omega_3, \omega_4\}$  are decreased. Our observations, mentioned above, only can be made for such offered traffics for which, in fact, the RU discipline does not degenerate to the HOL disci-Laplace-Stieltjes transforms of the waiting-time d.f.'s for these pline. Such degeneration appears if the probability of more than one request being waiting is small enough to prevent a nonurgent request to become (during its wait) more urgent than an urgent request which arrived later.

#### C. Dependency of the Waiting-Time d.f.'s on the Variance of the Urgencies

From Fig. 4 waiting-time d.f.'s in the FCFS and RU disciplines can be seen. We used two sets of urgency numbers sulting dotted curves are very much closer to the waitingtime d.f. in the FCFS discipline than is the case for a four times larger variance (dashed curves).

### D. Dependency on the Attachment of a Set of Mean Service Times to a Set of Urgency Numbers

Assume the model in Fig. 1 in the HOL discipline and define the common-mean waiting time W, cf. (6),

$$W = \lambda^{-1} \sum_{i=1}^{N} \lambda_i W_{i\text{HOL}}.$$
 (10)

As the probability f approaches 1, the distance between the The common-mean waiting time results from weighting and waiting-time d.f.'s for type-i and type-j requests asymptoti- summing over the individual mean waiting times of priority

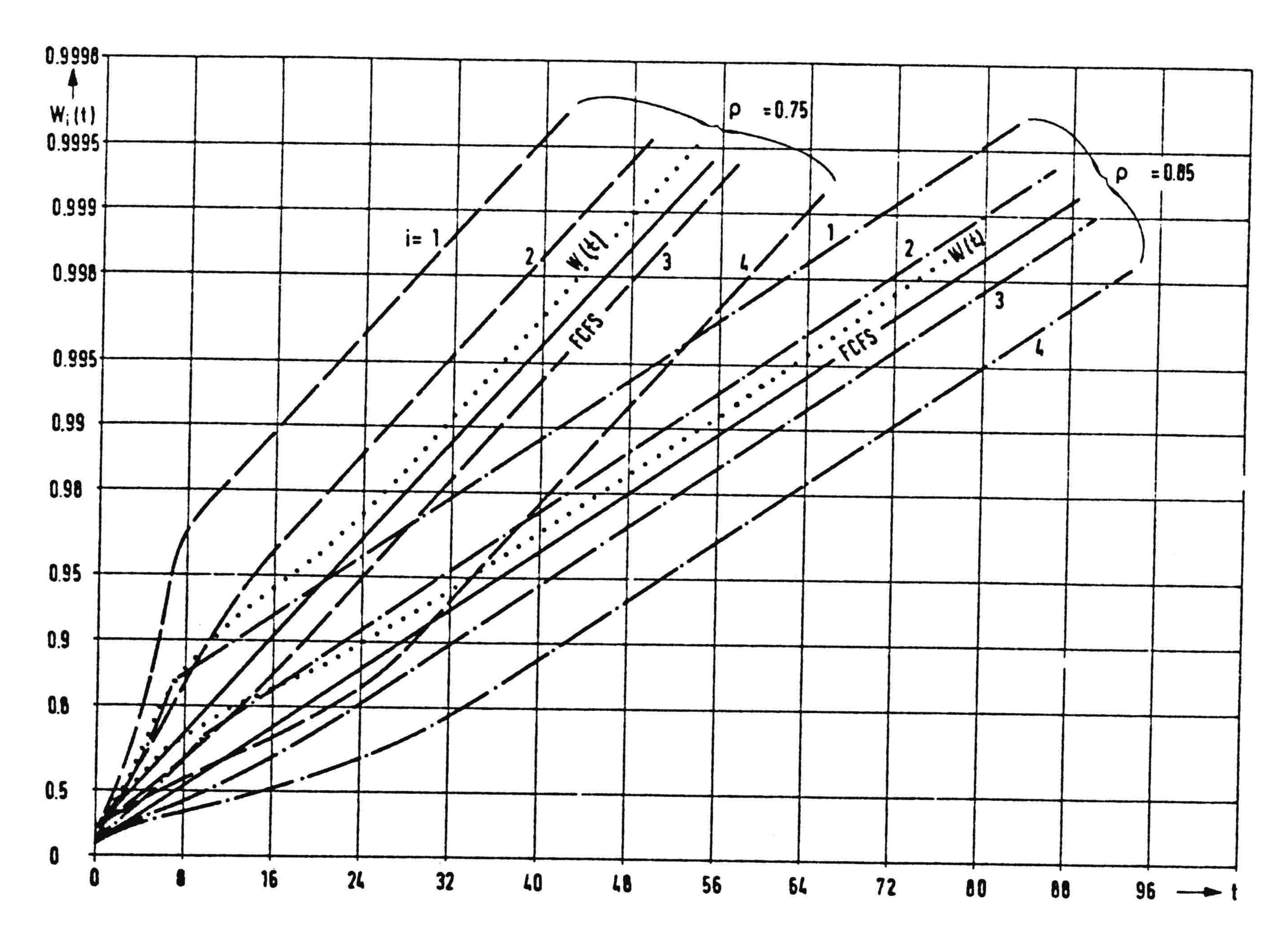


Fig. 3. Simulated waiting-time d.f.'s in the FCFS and RU disciplines, cf. Example 1. The type i and total offered traffic ρ are parameters. The solid and dashed curves result from the FCFS and RU disciplines, respectively. The dotted curve shows the common waiting-time d.f. W(t).

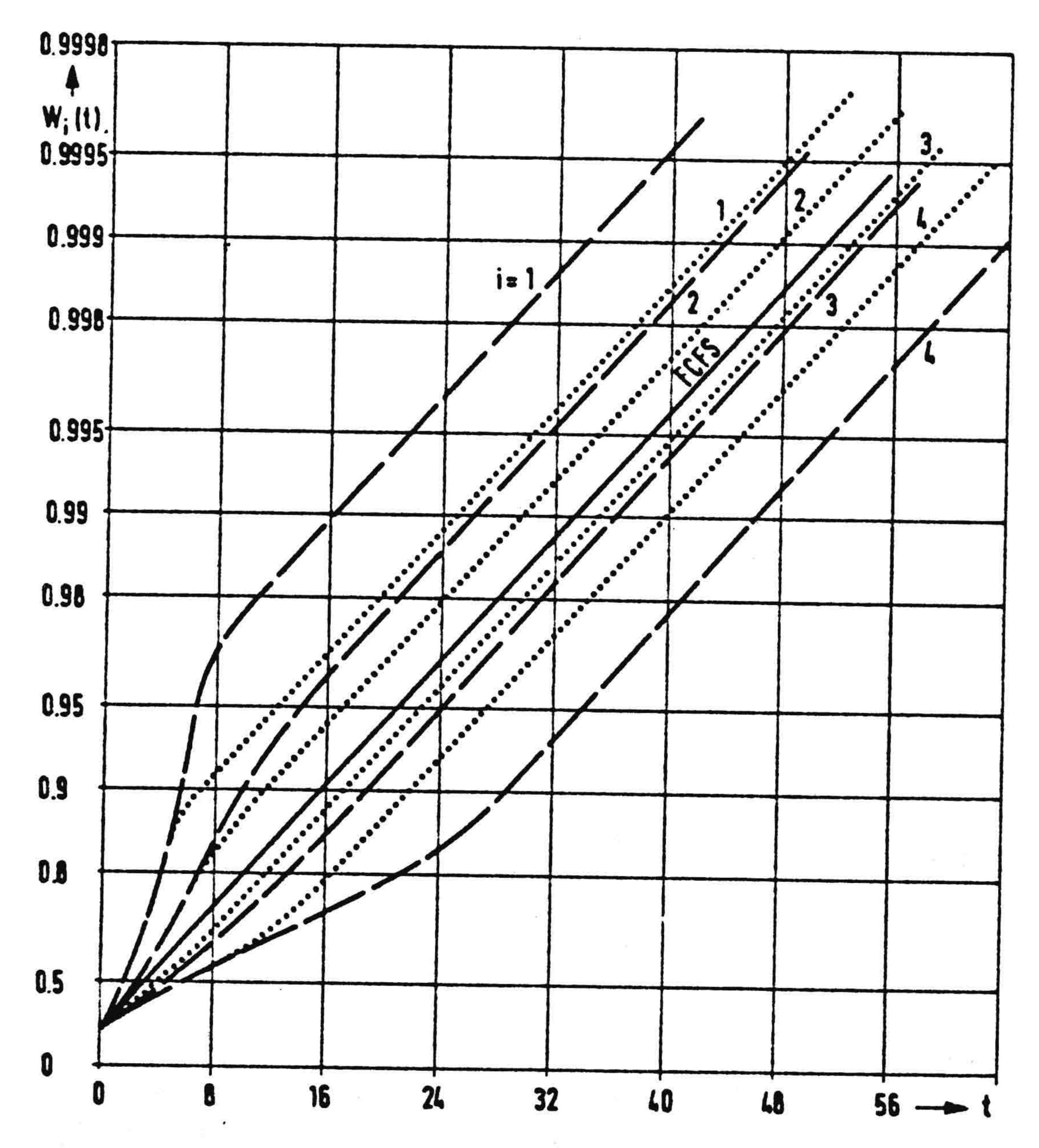


Fig. 4. Waiting-time d.f.'s in the FCFS and RU disciplines, cf. Example 1. The dashed and dotted curves appear for a large and small variance of the urgency numbers, respectively. The total offered traffic is

levels i. W equals W<sub>FCFS</sub>, cf. (5), only if all mean service times  $\beta_i$  are equal. This is not the case in Examples 1 and 2. From [3] it is known that W can be minimized in a static priority model by satisfying the condition

$$\beta_1 \leqslant \beta_2 \leqslant \cdots \leqslant \beta_i \leqslant \cdots \leqslant \beta_N,$$
 (11)

where 1 is the highest priority. The kind of service-time d.f.'s themselves do not matter; only (11) must be satisfied. Since W and W<sub>FCFS</sub> are the means of the corresponding waiting-time d.f.'s  $P(w \le t)$  and  $P(w_{FCFS} \le t)$ , respectively, then for the It follows that depending on how well or badly the relation

case of unequal  $\beta_i$ 's, if (11) is satisfied, the waiting-time d.f.  $P(w \le t)$  must be better in some sense than  $P(w_{FCFS} \le t)$ .

The opposite should be expected if the worst possible priority assignment is chosen, namely  $\beta_i \geqslant \beta_{i+1}$ . The waiting-time d.f.  $P(w \le t)$  simply is a weighted combination of the waiting-(11) time d.f.'s  $P(w_{iHOL} \le t)$  of the static priority model using

$$w = \lambda^{-1} \sum_{i=1}^{N} \lambda_i w_i \text{HOL}.$$

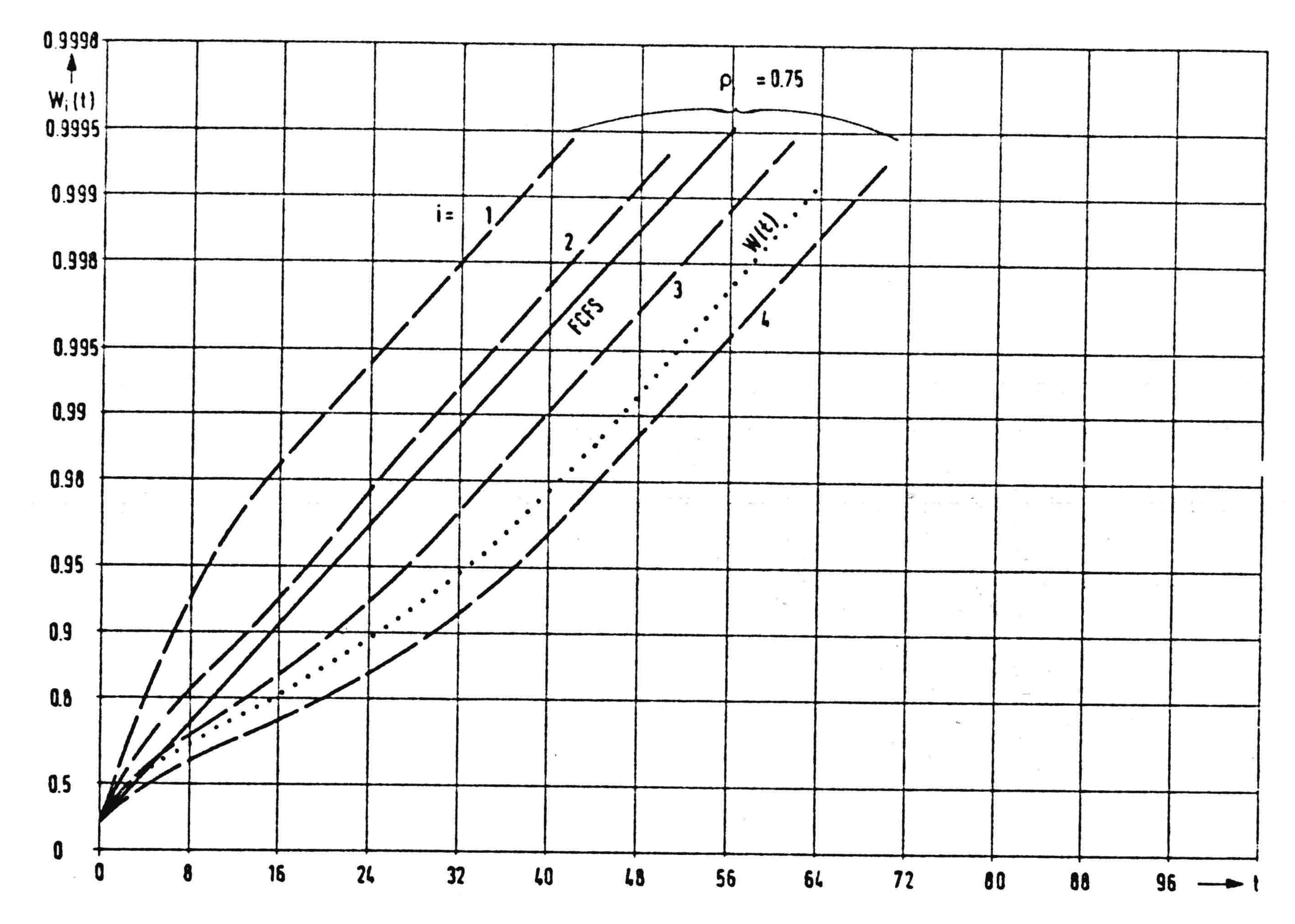


Fig. 5. Waiting-time d.f.'s in the FCFS and RU disciplines, cf. Example 2. The type i and total offered traffic  $\rho$  are parameters. The solid and dashed curves result from the FCFS and RU disciplines, respectively. The dotted curve shows the common waiting-time d.f. W(t).

in (11) is observed, the d.f.  $P(w_{FCFS} \le t)$ —which is independent of priority assignment—has a worse or better position in the set of curves  $P(w_{i \text{ HOL}} \le t)$ , resulting.

The assignment of urgency numbers to customers with different mean service times  $\beta_i$  has the same effect as the assignment of static priorities: small urgencies correspond to small priority numbers (high priorities) while large urgencies correspond to large priority numbers. In the next step define the common waiting-time d.f. in the RU discipline.

$$W(t) = \lambda^{-1} \sum_{i=1}^{N} \lambda_i P(w_i \leq t).$$

Now it can be argued that assigning small urgency numbers to customers with small mean service times  $\beta_i$  and large  $\omega_i$ 's to customers with large  $\beta_i$ 's results in a bad position of  $P(w_{FCFS} \leq t)$  compared to W(t). Just this urgency number assignment minimizes the common-mean waiting time in the RU discipline and optimizes the common waiting-time d.f. W(t). An inverse urgency number assignment results in the least favorable set of curves  $P(w_i \leq t)$  compared to  $P(w_{FCFS} \leq t)$ .

The optimal urgency number assignment just introduced is used in Example 1, cf. Fig. 3. It can be seen that  $P(w_{FCFS} \le t)$  (the curve indexed FCFS) has a bad position compared to the set of curves  $W_i(t)$ . This corresponds to a common waiting-time d.f. W(t) being preferable to  $P(w_{FCFS}) \le t$ ).

The opposite is true for Example 2: there, the smaller  $\beta_i$  is, the larger the urgency  $\omega_i$ . In Fig. 5, consequently,  $P(w_{FCFS} \le t)$  has a good position compared to the set of curves  $W_i(t)$ .

### E. The RU and HOL Disciplines Compared

The RU-NPRE discipline may be thought of as an alternative to the HOL discipline. From [1] and [2] it is known that the greater the variance of the urgency numbers, the more the mean waiting time  $W_i$  in the RU discipline differs from  $W_{FCFS}$ , and the more  $W_{iHOL}$  is approximated. Now, depending on the set of  $\omega_i$ 's, the question arises how the waiting-time d.f.'s in the RU and HOL disciplines will differ. First we recall that only the variance of the urgency numbers plays a role. Our observation is that the waiting-time d.f.'s in the HOL and RU

disciplines are very similar for "small" waiting times, cf. Fig. 6. But what a small waiting time is heavily depends on the total offered traffic  $\rho$  and the variance of the  $\omega_i$ 's.

One substantial difference between the HOL and the RU disciplines can be seen by comparing the individual probabilities of type-i requests in the same discipline: While the RU discipline meets all deadlines of different requests given by the set of the  $\omega_i$ 's [or a comparable set, cf. the second observation following (8)] with nearly the same probability, the corresponding probabilities  $P(w_i \leq \omega_i)$  may differ by factors in the HOL discipline. In the RU discipline all requests are handled on an equal basis while the opposite is true for the HOL discipline.

By choosing the urgency numbers, it is possible in the RU discipline to relate the waiting-time d.f.'s to one another in a prescribed way which is impossible in the HOL discipline.

#### III. APPROXIMATION OF THE WAITING-TIME D.F.'s

#### A. Previously Known Results

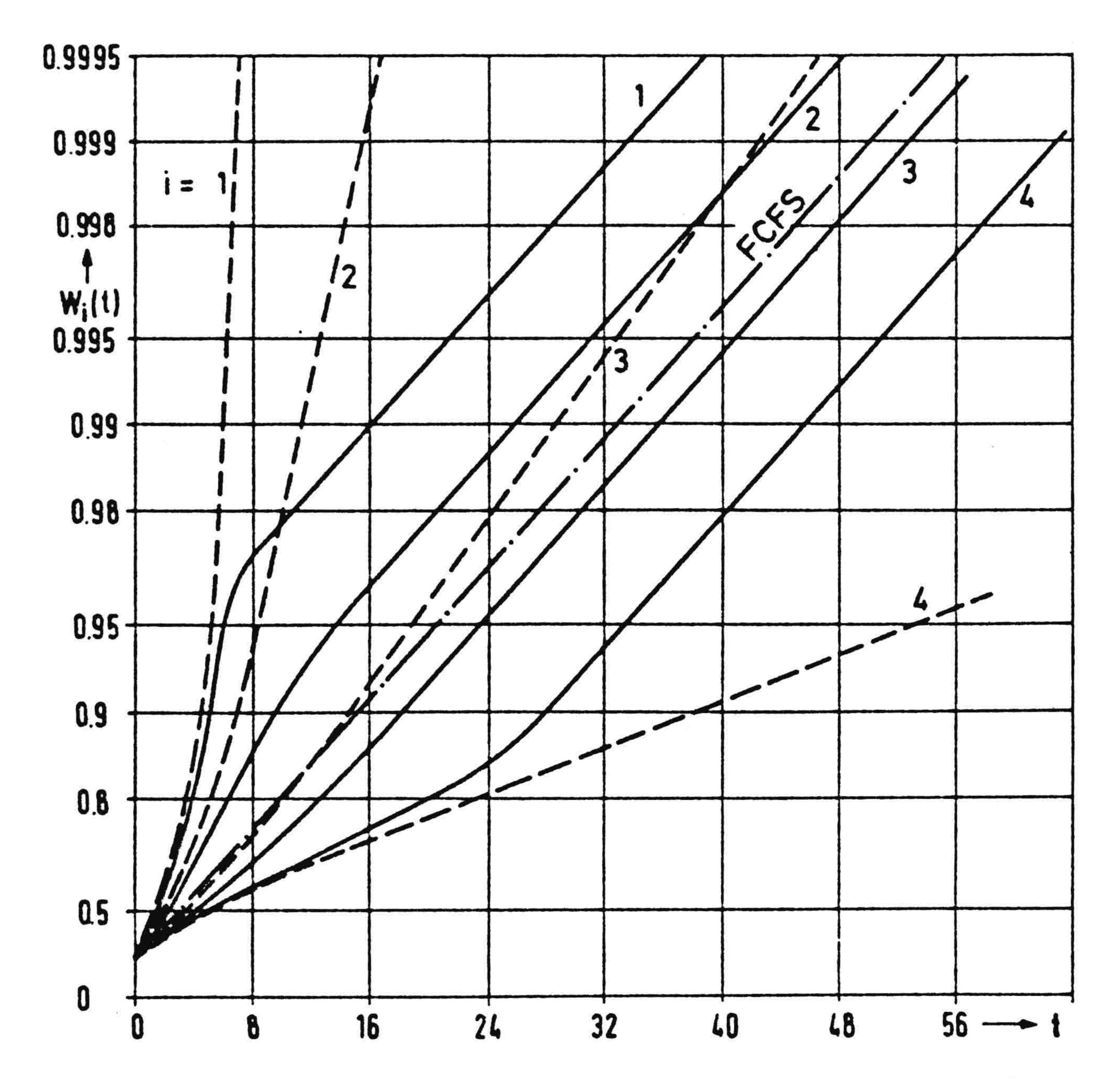
Back in 1962 Jackson [2] considered a model with discrete time and a nonpreemptive RU discipline. The service times of all requests were taken from one geometric distribution, and the requests themselves were taken from a Bernoulli arrival process. The individual urgency  $\omega_i$  is computed for every arriving request according to a given distribution. Model M1, which corresponds to this discrete-time model but is time continuous, would have a Poisson arrival process and negative exponentially distributed service times which all originate from a common d.f. Model M1 differs from the model considered in this paper in that different type requests may have different and arbitrary service-time d.f.'s. Therefore, our model is more general.

In [2] two relations are derived for the discrete time model

$$\lim_{f \to 1} [Q_i(f) - Q_j(f)] = \omega_i - \omega_j \quad [cf., (8)]$$
 (12)

and

$$\lim_{f \to 1} [Q_i(f) - Q(f)] = \omega_i - \overline{\omega}'. \tag{13}$$



g. 6. Waiting-time d.f.'s in the FCFS, RU (solid lines), and HOL (dashed lines) disciplines. Parameters of Example 1,  $\rho = 0.75$ .

 $\overline{\omega}'$  is defined by an expression which for model M1 may be approximated by the mean weighted urgencies  $\overline{\omega}$ , cf. list of abbreviations. Q(f) is the f-quantile of the waiting-time d.f. of a single-queue FCFS model.

### B. Approximation of $\overline{\omega}'$ in (13)

Our simulation results support the suspicion that (12) and (13) seem to be generally applicable, regardless of the service-time d.f.'s assumed. We have found that the difference between the waiting-time d.f. in the FCFS discipline and the curves, constructed from (12) using  $\overline{\omega}$  instead of  $\overline{\omega}$ , always remains small. Only differences up to 9 percent have been observed [4]. However, these differences are considerably larger than those observed by Jackson in his model. This is due to the fact that unlike Jackson's model the mean service times  $\beta_i$  in our examples are not equal, but different.

Our observation is that the difference  $(\overline{\omega} - \overline{\omega}')$  is, depending on the example, greater or smaller than zero. The reason or this can be plausibly explained as the dependancy of the waiting-time d.f.'s in the RU discipline on the urgency number assignment being optimal in Example 1 and worst in Example 2 (see Section II-D).

Note that the load to the RU model is the same in both Examples 1 and 2. Hence  $P(w_{FCFS} \le t)$  is the same for both examples which is not the case for the curves  $W_i(t)$  and W(t). Only the assignment of the same set of urgency numbers is different, which is precisely the reason why  $\overline{\omega}'$  differs from  $\overline{\omega}$ . Our observation is that  $\epsilon = \overline{\omega} - \overline{\omega}'$  for all the examples studied was always very small (a few percent of  $\overline{\omega}$ ). Therefore we decide to accept (13) with  $\overline{\omega}'$  substituted by  $\overline{\omega}$  as an approximation which leads to

$$\lim_{f \to 1} [Q_i(f) - Q(f)] = \omega_i - \overline{\omega}. \tag{14}$$

# C. Approximations of the Waiting-Time d.f.'s in the FCFS and RU Discipline

From Figs. 3, 4, and 5 it can be seen for our two examples that the waiting-time d.f.'s, when plotted semilogarithmically, run as parallel straight lines from about  $f = w_i(t) = 0.98$  on.

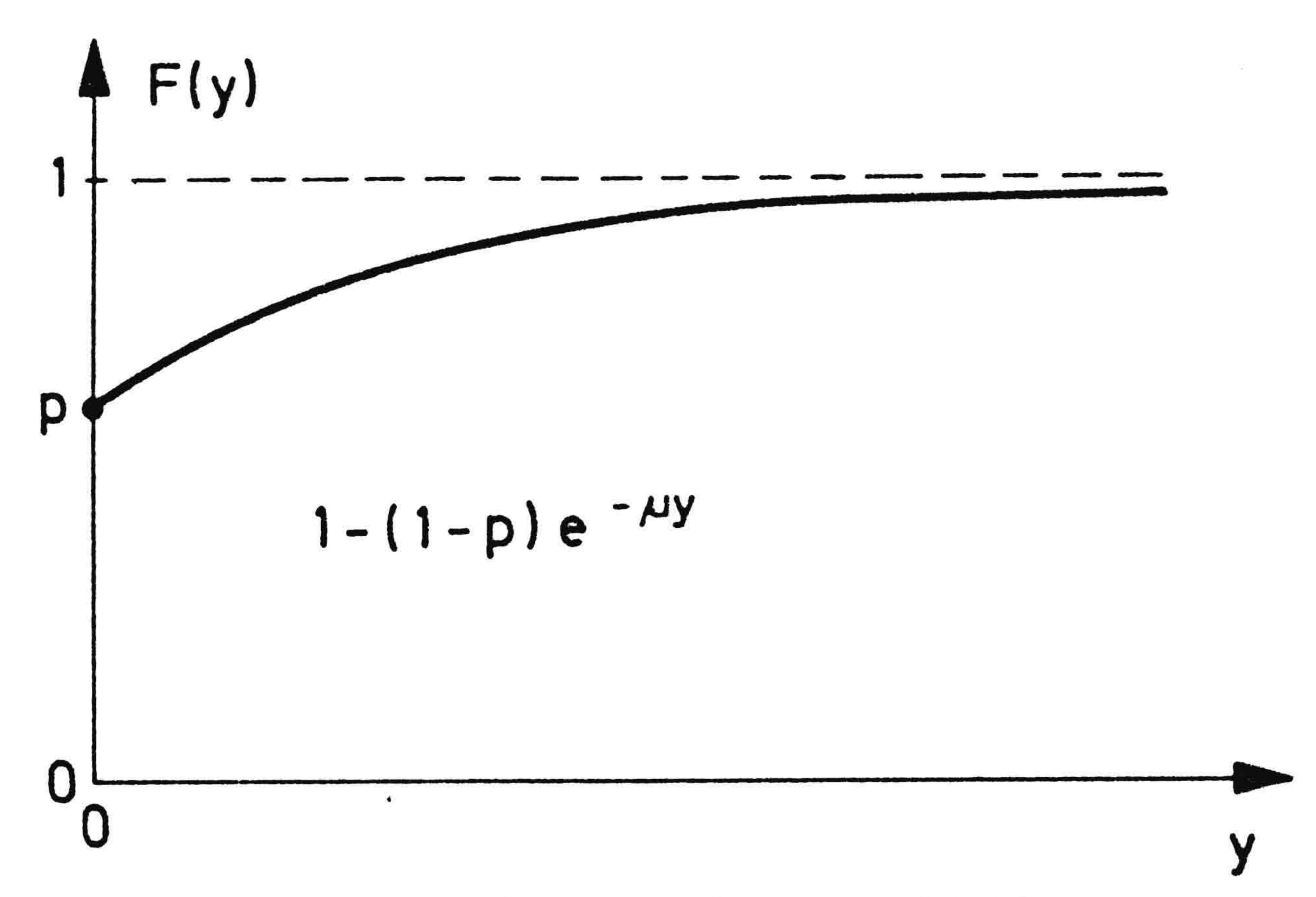


Fig. 7. A degenerated exponential d.f.

We are now going to approximate these functions through proper formulas.

1) The waiting-time d.f. in the FCFS discipline: The waiting-time d.f. in the FCFS discipline may be computed by inversion of their Laplace transform which can be a very difficult task. For the examples introduced in this paper it can be shown that the d.f.'s found by simulation, especially in the tail, can be approximated very well by a degenerated exponential d.f., cf. Fig. 7,

$$F(y) = 1 - (1 - p) e^{-\mu y}. \tag{15}$$

The rth moment of a degenerated d.f. is

$$\gamma^{(r)} = r!(1 - p)/\mu^r$$
.

By setting equal the first two moments of the waiting-time d.f. in the FCFS discipline and the degenerated d.f., namely,  $W_{\text{FCFS}} = \gamma^{(1)}$  and  $W_{\text{FCFS}}^{(2)} = \gamma^{(2)}$ , it is possible to compute the two unknown parameters p and  $\mu$  of (15). Thereby a degenerated d.f. is defined having the same first two moments as the waiting-time d.f. in the FCFS discipline. Using (5) and

$$W_{\text{FCFS}}^{(2)} = 2W_{\text{FCFS}}^2 + \lambda \beta^{(3)}/[3(1-\rho)] \tag{16}$$

1**İ** 

$$C^2 = W_{\text{FCFS}}^{(2)} / W_{\text{FCFS}}^2 - 1 \tag{17}$$

is the squared coefficient of variance we find

$$p = \frac{C^2 - 1}{C^2 + 1}$$
 and  $\mu = \frac{2}{W_{\text{FCFS}}(1 + C^2)}$ . (18)

Equations (15) and (18) define an approximation of  $P(w_{FCFS} \le t)$ , whose first two moments are correct and which belongs to the family of the degenerated exponential d.f.'s

$$P(w_{\text{FCFS}} \le t) \approx P(t) = 1 - (1 - p) e^{-\mu t}.$$
 (19)

In semilogarithmic scaling this function appears as a straight line.

We are now going to study the quality of (19) as an approximation for  $P(w_{FCFS} \le t)$ . Fig. 8 shows simulated (+) and approximated waiting-time d.f.'s in the FCFS discipline using the parameters of Examples 1 and 2. (We have mentioned already that both examples should produce the same waiting-time d.f. in the FCFS discipline.) All our investigations, cf. [4], showed that, especially in the tails of the waiting-time d.f.'s, simulation and approximation fit together very well.

2) The tail of the waiting-time d.f. in the RU discipline: In our first approximation, cf. (14), we assumed the tails of the waiting-time d.f.'s in the FCFS and RU discipline to run parallel. In our second approximation we introduced (19) to ap-

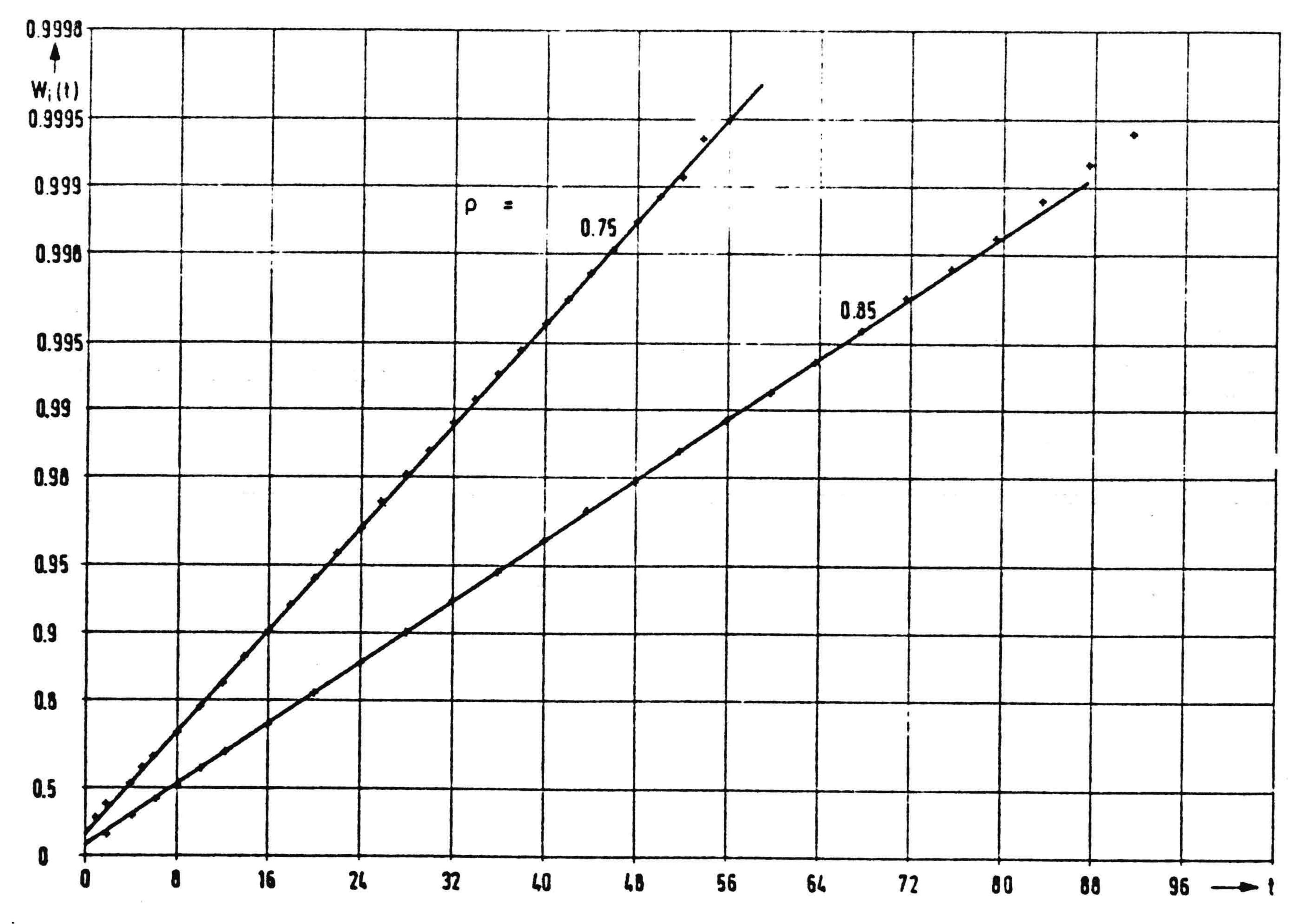


Fig. 8. Waiting-time d.f.'s in the FCFS discipline, cf. Examples 1 and 2. Approximations by means of (19) (lines) and simulation results (+) are shown. Parameter is the total offered traffic.

proximate the tail of the waiting-time d.f. in the FCFS discipline. Combining these two results, we find as the third approximation in this paper a close formula to describe the tails of the waiting-time d.f.'s in the RU discipline.

From (14) it follows immediately that for every pair of f-quantiles  $(f \to 1)$   $Q_i(f) = w_i$  and  $Q_{FCFS}(f) = w_{FCFS}$  of the waiting-time d.f.'s in the RU and FCFS disciplines, respectively, one may write

$$w_i - w_{FCFS} = \omega_i - \overline{\omega}$$

cf. (9). Recalling (19)

$$P(w_{\text{FCFS}} \leq t) \approx 1 - (1 - p) e^{-\mu t}$$

and inserting

$$w_{FCFS} = w_i - \omega_i + \overline{\omega}$$

re have

$$P([w_i - \omega_i + \overline{\omega}] \le t) \approx 1 - (1 - p) e^{-\mu t}$$

$$P(w_i \le t) \approx 1 - (1 - p) e^{-\mu (t - \omega_i + \overline{\omega})}. \tag{20}$$

This function defines the tails of the waiting-time d.f.'s in the RU-NPRE discipline.

A comparison of (19) and (20) reveals that they differ by a factor. We aim at an expression for the waiting-time d.f.'s in the RU discipline which has the form of a degenerated exponential d.f., cf. (15). Instead of (20) we write

$$P(w_i \le t) \approx 1 - (1 - p_i) e^{-\mu t}$$
 (21)

and find, by setting equal the right-hand sides of (20) and (21),

$$p_i = 1 - (1 - p) e^{-\mu(\overline{\omega} - \omega_i)}$$
.

Note that p and  $p_i$  are the probabilities which arise for t = 0 in (15) and (21), cf. Fig. 7. From simulation experiments we know that (21) is a good approximation, if either the total offered traffic is large,  $\rho > 0.6$ , or the urgency numbers have a small variance, or both. Otherwise, the d.f.'s in the RU discipline appear to be similar to those of the HOL discipline which is always the case for small waiting times t. The reason

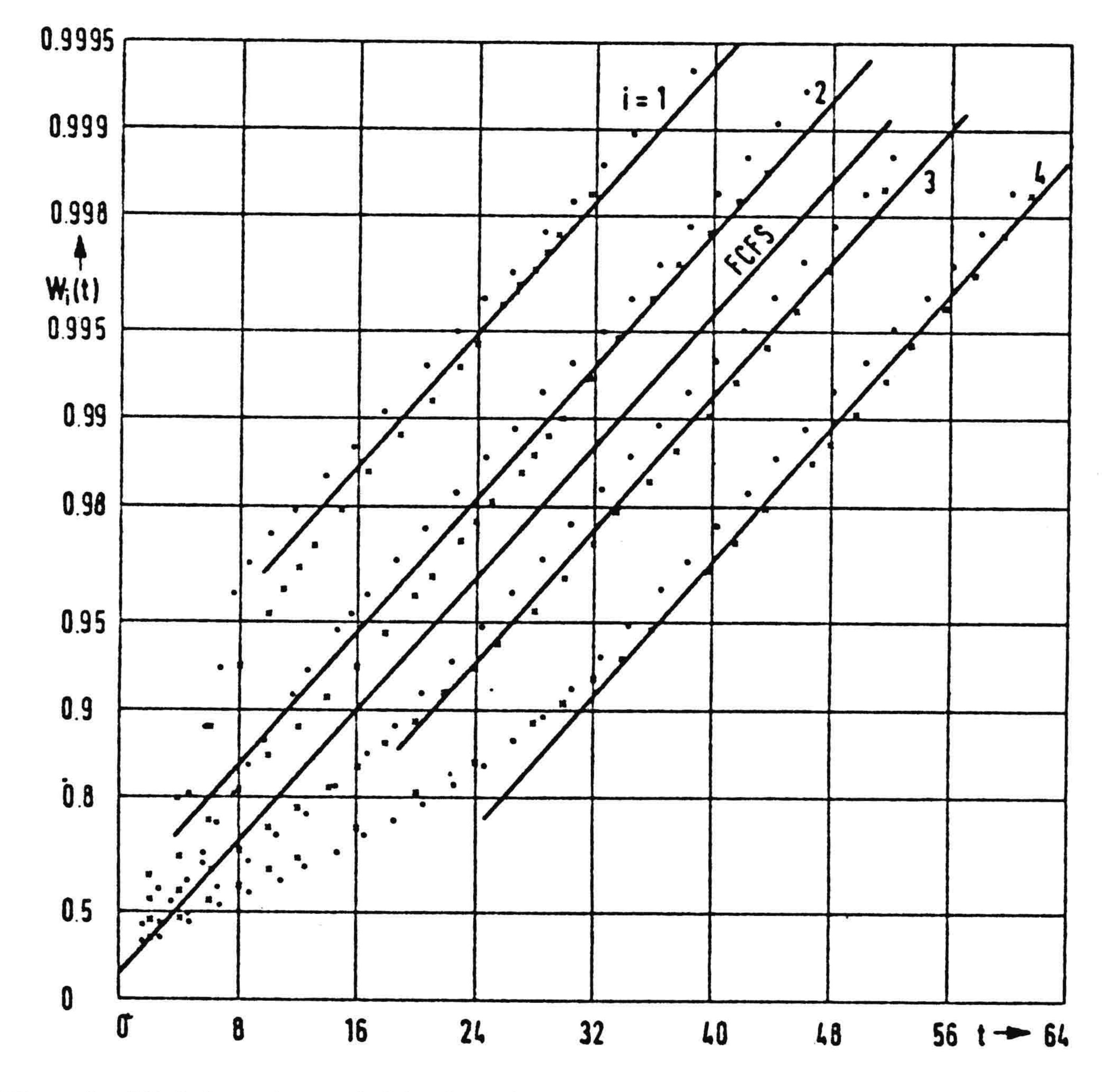


Fig. 9. Waiting-time d.f.'s in the FCFS and RU disciplines, cf. Examples 1 and 2, approximated by means of (21), (lines) and simulation results  $(\cdot, X)$ . Additionally, the FCFS approximation is shown. The total offered traffic is  $\rho = 0.75$ .

for this is that without the conditions mentioned, there is only a very small probability of missing deadlines for any type of request.

Fig. 9 shows both simulation results (·, X) and approximate waiting-time d.f.'s computed from (21), which appear as straight lines for the examples defined in Section II. It can be seen that for large waiting times, simulation and approximation agree very well. Deviations are typically below the 10 percent range.

discipline appear to be similar to those of the HOL discipline In Fig. 9 we once again can study the results of a different which is always the case for small waiting times t. The reason assignment of a set of urgency numbers  $\omega_i$  to a set of mean

service times  $\beta_i$  as defined by Examples 1 and 2. The approximate waiting-time d.f.'s are derived from the d.f. of the FCFS discipline and therefore are the same for both Examples 1 and 2. As a consequence of the optimal urgency number assignment as defined in Section II-D, which is realized in Example 1, the real d.f.'s (·) compare favorably to the computed ones (lines). On the other hand the d.f.'s (X) of Example 2, in which the  $\omega_i$  are especially poorly assigned to the  $\beta_i$ , are less favorable than the computed ones. Note that the approximation by (21) does not take into account the urgency-number assignment. This should be the reason for the deviation between simulated and computed d.f.'s observed.

3) The waiting-time d.f.'s in the RU discipline for small waiting-times: Equation (21) does not apply to small waiting times. In [1] we used a simple approximation of the waiting-time d.f. in the RU discipline to compute the corresponding mean waiting time  $W_i$  of type-i requests

$$P(w_i > t) = \rho e^{-\rho t/W_i}. \tag{22}$$

This simple approximation only depends on the total offered traffic  $\rho$  and  $W_i$ . Fig. 10 shows this approximation (lines) compared to simulation results (X) for Example 1. The most unsatisfactory approximation appears to be for type-1 requests which was not needed in [1]. It can be seen that (22) in general cannot be called a good approximation.

Nevertheless, the error produced by this approximation has a negligible effect on the computation of  $W_i$ ; simulation results agree well with the computed values of  $W_i$ , cf. [1].

# D. Mean Waiting Time of Requests Which Have Missed Their Deadlines

From (21) the approximate probability of a type-i request missing its deadline (at  $t = \omega_i$ ), presuming a sufficiently large t, can be computed

$$P(w_i > \omega_i) \approx (1 - p) e^{-\mu \overline{\omega}}$$
.

The waiting-time d.f. of requests missing their deadlines is approximately given by

$$P(w_i \leq t | t > \omega_i) = [P(w_i \leq t) - P(w_i \leq \omega_i)] / [1 - P(w_i \leq \omega_i)]$$

$$P(w_i \leq t | t > \omega_i) = 1 - e^{\mu \omega_i} e^{-\mu t}.$$

Apparently a degenerated exponential d.f. arises whose mean

$$E(w_i|t>\omega_i)=E_i=e^{\mu\omega_i/\mu}$$

only depends on  $\mu$  and  $\omega_i$ . Note that  $\mu$  can be computed from (18). The approximate mean waiting times  $E_i$  of requests having missed their deadlines differ approximately in the factor  $e^{\omega_i}$ .

# IV. SOME TYPICAL PROPERTIES OF THE WAITING-TIME D.F. IN THE PREEMPTIVE RU DISCIPLINE

In the preemptive RU (RU-PRE) discipline a type-i request may be interrupted (preempted) by a type-j request if  $w_i < \omega_i$  and  $w_j = \omega_j$ . The priority  $q_i(t)$  of a request remains unchanged once its servicing has begun. But if a preemption occurs, the priority of the interrupted request again follows (4) with  $w_i(t)$  being the total time such a request remains in the system without being serviced. The urgency  $\omega_i$  represents the time  $w_i$  a request is willing to wait from its arrival  $T_i$  to its deadline  $t_i$ . For both the nonpreemptive and preemptive versions of the RU discipline a request's deadline may be defined to be reached if  $w_i(t)$  equals  $\omega_i$ . If the waiting time  $w_i$  of a request being serviced is less than its voluntary waiting time  $\omega_i$ , while another type-j request has already waited its voluntary time  $\omega_j$ , then the type-j request preempts the type-i

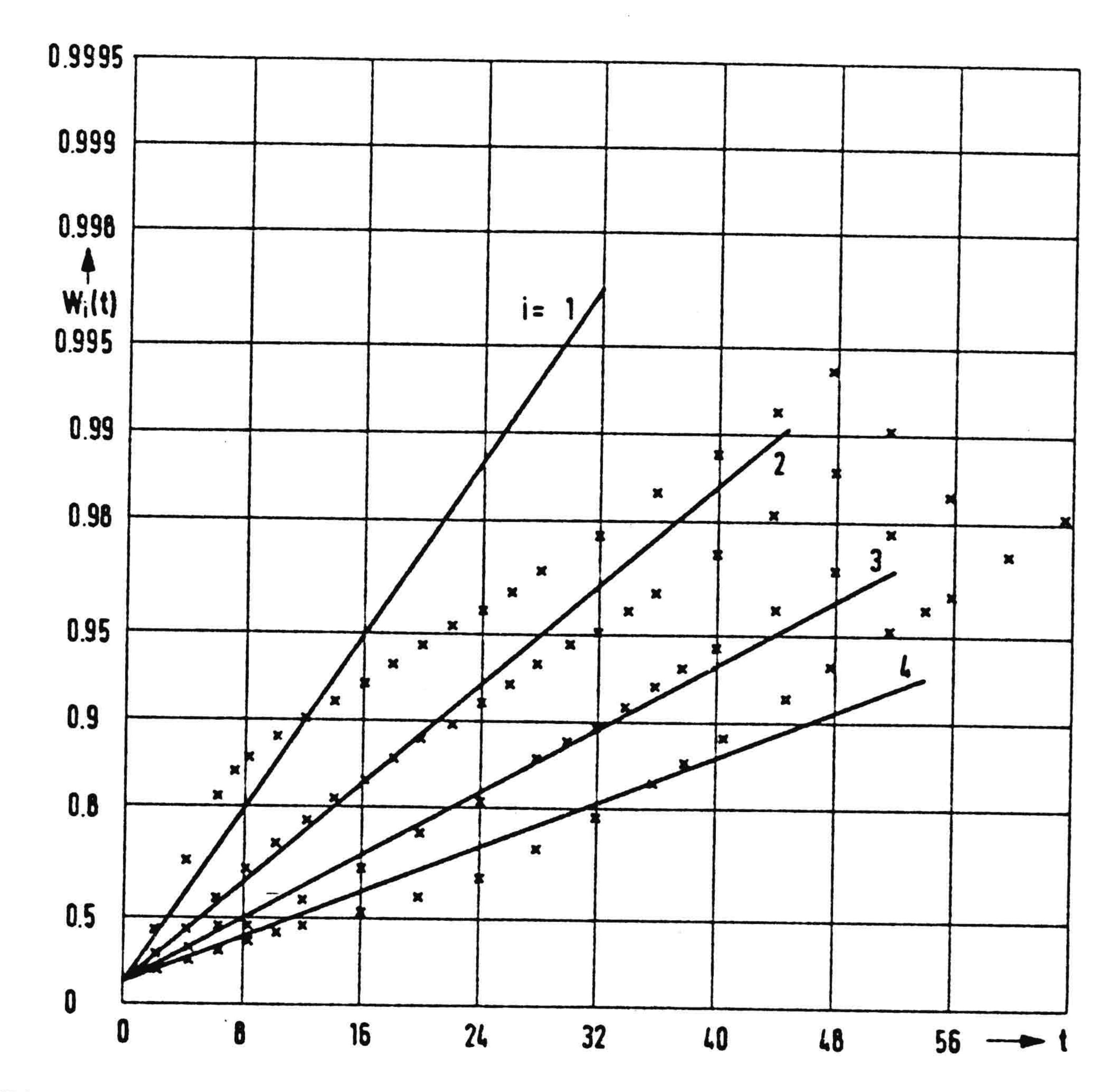


Fig. 10. Waiting-time d.f.'s in the RU discipline, cf. Example 1, approximated by means of (22), (lines) and simulation results ( $\times$ ). The total offered traffic is  $\rho = 0.85$ .

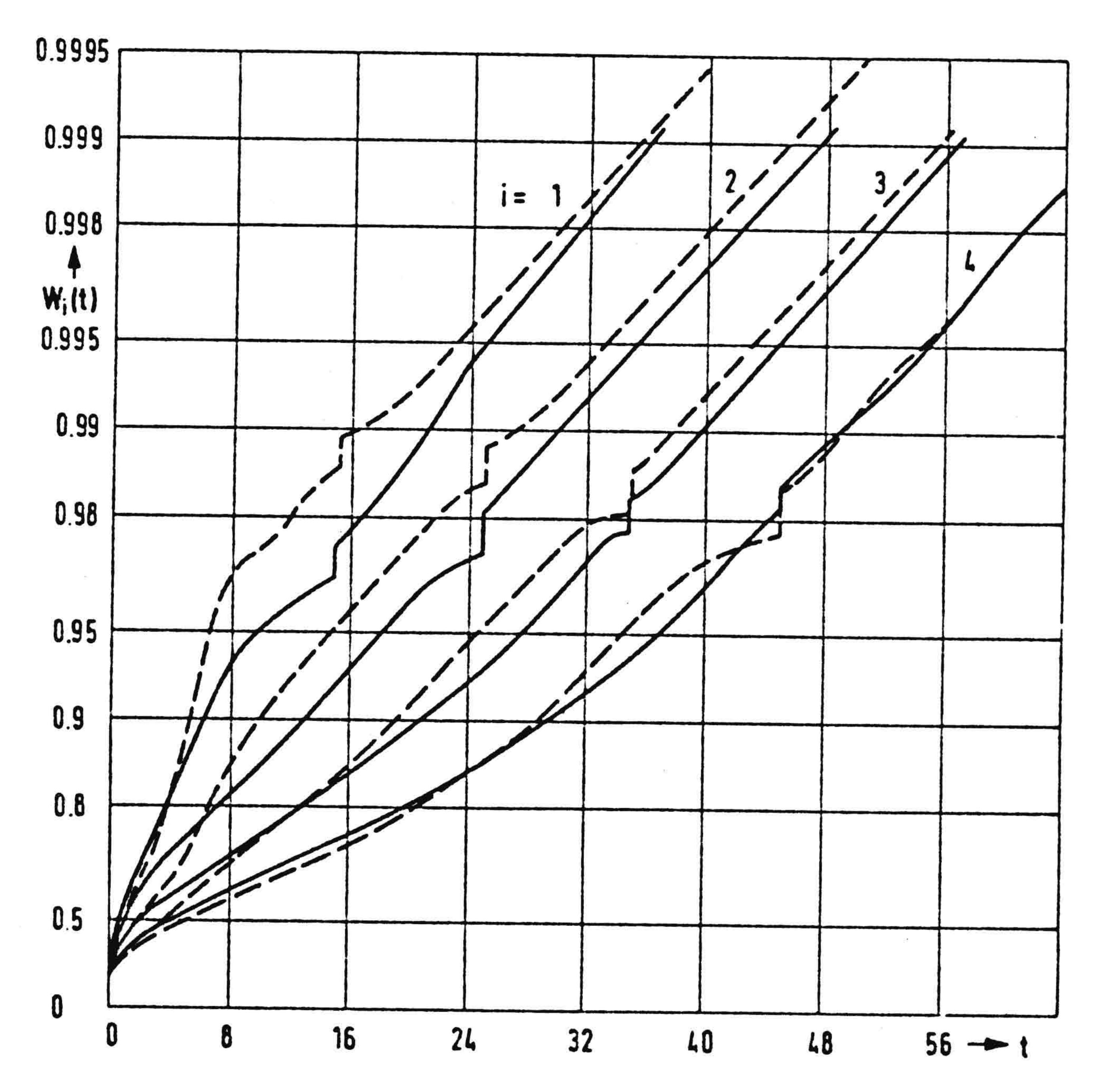


Fig. 11. Simulated waiting-time d.f.'s in the RU-PRE discipline for Examples 2 (solid lines) and 1 (dashed lines). The total offered traffic is  $\rho = 0.75$ . The jumps appear at  $t = \omega_i$ .

request. The case i = j is impossible. Requests, whose servicing just barely began on time or which missed their deadlines, may not be interrupted. Summarizing, it can be said that in the RU-PRE discipline only requests whose servicing began ahead of schedule may be interrupted and only by request whose deadline has been reached.

versions of the RU discipline a request's deadline may be defined to be reached if  $w_i(t)$  equals  $\omega_i$ . If the waiting time  $w_i$  of a request being serviced is less than its voluntary waiting time  $\omega_i$ , while another type-j request has already waited its voluntary time  $\omega_j$ , then the type-j request preempts the type-i

both disciplines arise. Especially worth noting is the jump at  $t = \omega_i$  which results from interrupting prematurely started servicing for such requests which would otherwise miss their deadlines. The size of the jump exactly defines the probability for interrupting other types of requests. It can be observed that the waiting-time d.f.'s in the tail run parallel, as was the case for the RU-NPRE discipline (cf. [6]). For all examples considered the slope of these tails in the RU-PRE and RU-NPRE disciplines was not found to be the same (the same offered traffic  $\rho$  is presumed).

It appears that interruptions are useful for increasing the probabilities of meeting deadlines if the coefficients of variance of the service-time d.f.'s are large, i.e.,  $C_i > 1$ . For small  $C_i$ 's our simulation supports the conjecture that interruptions of requests in favor of others, having a larger expected service time, is disadvantageous for meeting deadlines. The deadlines of some types of service requests are then met better with the RU-NPRE discipline than with the RU-PRE discipline.

#### REFERENCES

- [1] B. Walke, "Improved bounds and an approximation for a dynamic priority queue," in *Proc. 3rd Int. Symp. Modeling and Perform. Eval. Comput. Syst.*, Bonn, Germany, Oct. 1977. Amsterdam, The Netherlands: North-Holland, 1978, pp. 321-346.
- [2] J. R. Jackson, "Waiting time distributions for queues with dynamic priorities," Naval Res. Logist. Quart., vol. 9, pp. 31-46, Mar. 1962.
- [3] R. W. Conway, W. L. Maxwell, and L. W. Miller, Theory of Scheduling. Reading, MA: Addison-Wesley, 1967.
- [4] B. Walke and W. Rosenbohm, "Waiting-time distributions for deadline oriented serving," in *Proc. 4th Int. Symp. Modeling and Perform. Eval. Comput. Syst.*, Vienna, Austria, Feb. 6-8, 1979. New York: North-Holland, 1979.
- [5] E. B. Veklerov and I. M. Dukhovnyi, "Queueing systems with time-limited categorical priority," *Probl. Inform. Transm.*, vol. 12, pp. 73-79, Jan.-Mar. 1976.
- [6] H. M. Goldberg, "Jackson's conjecture on earliest due date scheduling," Dep. Finance and Management Sci., Univ. Alberta, Edmonton, Alberta, Canada, Sept. 1978.