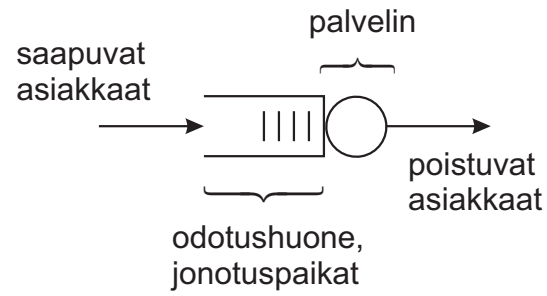


Waiting systems

Now we turn our focus on waiting systems. These are the genuine queues where there is a waiting room and the customers may have to wait for the service.

The basic elements of a (single server) queue are as shown in the figure.



Double time axis (in a single server system)

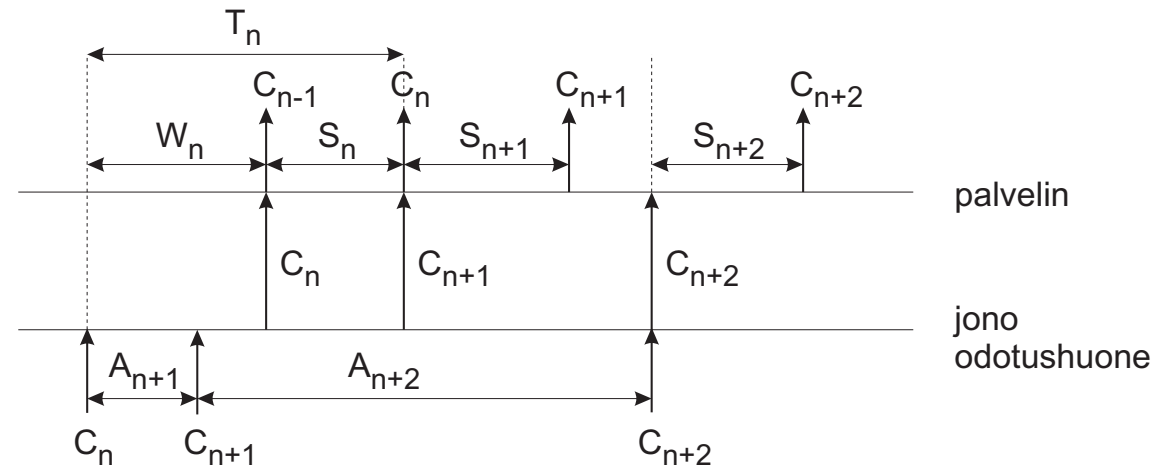
- C_n customer n
- S_n service time of customer n (time it takes to discharge the work)
- X_n service requirement of customer n (the work required)
- W_n waiting time of customer n
- T_n $W_n + S_n$ the total time spent in the system by customer n
time in system, sojourn time
- A_n (or t_n) the interarrival time between customers $n - 1$ and n
- C the service rate or capacity of the server (also denoted by c)

The service time depends on the service requirement (work) and the service rate: $S_n = X_n/C$.

In telecommunication applications the service may mean transmission of a packet on the line.

Then the work may be measured e.g. in units of kbit and the service rate is measured in kbit/s.

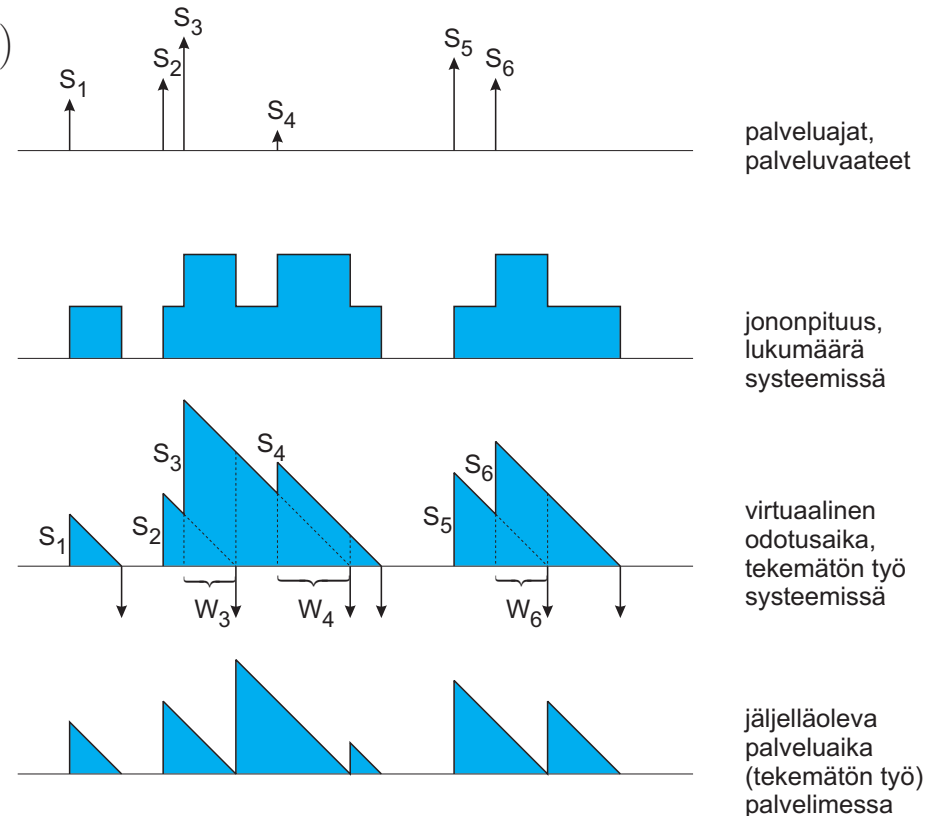
By inspection, one sees that for FIFO



$$\boxed{W_{n+1} = (W_n + S_n - A_{n+1})^+} \quad \text{where } (x)^+ = \max(x, 0)$$

Queue length, unfinished work and virtual waiting time

- N_t (or Q_t or L_t) number of customers in the system (“number in system”, “queue length”)
- S_n service time of customer n (time to discharge the work)
- X_t unfinished work in the queue at time t
- V_t virtual waiting time at time t
- W_n the real waiting time of customer n
- C the service rate or capacity of the server (also denoted by c)



- Virtual waiting time V_t means the time which a customer would have to wait for service if the customer happened to arrive at time t (in a FIFO queue).

V_t is the time it takes to discharge the unfinished work in the queue, X_t , i.e., $V_t = X_t/C$.

- In the case of Poisson arrivals the distribution of W_n is by the PASTA property the same as the stationary distribution of V_t .

The $M/M/1$ queue

Number of customers in an $M/M/1$ queue



By the method of a cut, one gets the balance condition

$$\lambda\pi_{n-1} = \mu\pi_n \quad \text{or} \quad \pi_n = \rho\pi_{n-1} \quad \text{where } \rho = \lambda/\mu \text{ (traffic intensity, load),}$$

wrom which we get recursively

$$\pi_n = \rho^n \pi_0 \quad (\text{in order for the queue to be stable, we have to require } \rho < 1)$$

The probability of an empty queue π_0 is obtained from the normalization condition $\pi_0 + \pi_1 + \pi_2 + \dots = 1$

$$\pi_0 = 1 / \sum_{n=0}^{\infty} \rho^n = 1 - \rho \quad \begin{array}{l} \text{(the probability that the server (ant the queue) is empty} \\ = 1 - \rho \Rightarrow \\ \text{probability that the server is busy} = \rho) \end{array}$$

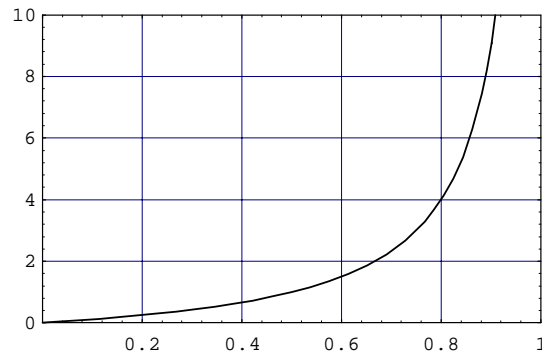
The queue length distribution of an $M/M/1$ queue, $\pi_n = P\{N = n\}$,

$$\boxed{\pi_n = (1 - \rho) \rho^n} \quad n = 0, 1, \dots \quad \text{Geom}_0(\rho) \text{ distribution (starts from 0)}$$

The average number of customers in the system

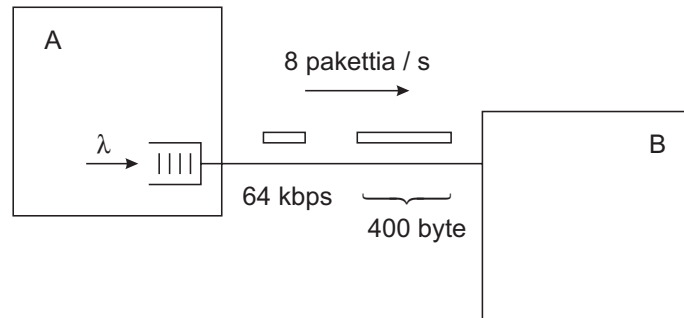
$$\begin{aligned} E[N] &= \sum_{i=0}^{\infty} i \pi_i = (1 - \rho) \sum_{i=0}^{\infty} i \rho^i = (1 - \rho) \rho \frac{d}{d\rho} \sum_{i=0}^{\infty} \rho^i \\ &= (1 - \rho) \rho \frac{d}{d\rho} \left(\frac{1}{1 - \rho} \right) = \frac{\rho}{1 - \rho} \quad \text{the mean of the } \text{Geom}_0(\rho) \text{ distribution (starts} \\ &\quad \text{from 0)} \end{aligned}$$

$$E[N] = \frac{\rho}{1 - \rho} = \underbrace{\rho}_{\substack{\text{customers in} \\ \text{the server}}} + \underbrace{\frac{\rho^2}{1 - \rho}}_{\substack{\text{waiting} \\ \text{customers}}}$$



The tail probability: the probability that there are at least n customers in the system,

$$P\{N \geq n\} = \sum_{i=n}^{\infty} \pi_i = (1 - \rho) \sum_{i=n}^{\infty} \rho^i = \rho^n$$

Example.

- Router A sends 8 packets per second, on the average, to router B.
- The mean size of a packet is 400 byte (exponentially distributed).
- The line speed is 64 kbit/s.

How many packets are there on the average in router A waiting for transmission or being transmitted and what is the probability that the number is 10 or more?

The utilization of the line (server) is

$$\rho = 8 \text{ s}^{-1} \times 400 \times 8 \text{ bit} / 64 \times 10^3 \text{ bit s}^{-1} = 0.4.$$

This can be also calculated in the form λ/μ , where

$$\lambda = 8 \text{ packets/s}, \quad \mu = 64 \text{ kbit/s} / (400 \times 8 \text{ bit/packet}) = 20 \text{ packets/s} \Rightarrow \lambda/\mu = 8/20 = 0.4$$

Thus $E[N] = 0.4 / (1 - 0.4) = 0.67$.

The probability that the number of packets is 10 or more is $0.4^{10} = 10^{-4}$.

Sojourn and waiting times in the $M/M/1$ queue

Little's result:

The average sojourn time (time in system) $E[T] = E[N]/\lambda$

The average waiting time $E[W] = (E[N] - \rho)/\lambda$

$$E[T] = \frac{1}{1 - \rho} \cdot \frac{1}{\mu} = \frac{1}{\mu - \lambda}$$

$$E[W] = \frac{\rho}{1 - \rho} \cdot \frac{1}{\mu} = \frac{\rho}{\mu - \lambda}$$

Independence of the scheduling discipline

For the $M/M/1$ -FIFO queue we have derived the queue length distribution $\pi_n = (1 - \rho)\rho^n$.

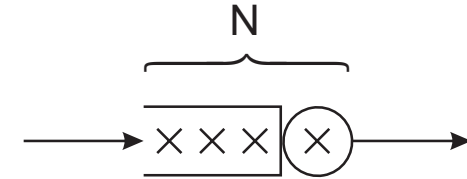
- This distribution is independent of the scheduling discipline (FIFO, LIFO, PS),
 - all these scheduling disciplines lead to exactly the same balance equations (proof is left as an exercise)
- Thus also the mean time in system, $E[T] = 1/(\mu - \lambda)$, is independent of the discipline (by Little's result the mean time in system equals the mean queue length divided by λ .)
- In the contrast, e.g. the distributions of W and T do depend on the discipline.

Note. The queue length distribution is not insensitive to the service time distribution in an $M/M/1$ -FIFO queue. However, in LIFO and PS queues the insensitivity holds.

The distribution of the sojourn time

Assume that an arriving customer finds N customers in the system (including the customer in the server, if any).

By the memoryless property of the exponential distribution also the remaining service time of the customer in service (if any) is distributed as $\sim \text{Exp}(\mu)$.



The time T spent by a customer in the system consists of the time it takes to serve the customers ahead in the queue and the customer's own service time

$$T = \underbrace{S_1 + S_2 + \dots + S_N}_{\text{customers ahead}} + \underbrace{S_{N+1}}_{\text{OWN}} \quad \text{sum of } (N + 1) \text{ rvs with } \text{Exp}(\mu) \text{ distribution}$$

$$\begin{cases} S_i \sim \text{Exp}(\mu) & \text{independent} \\ N \sim \text{Geom}_0(\rho) & \text{equilibrium distribution of the queue length (starts from 0), PASTA!} \end{cases}$$

$$\begin{aligned} f_T(t) &= \sum_{n=0}^{\infty} f_{T|N}(t, n) P\{N = n\} = \sum_{n=0}^{\infty} \overbrace{\mu \frac{(\mu t)^n}{n!} e^{-\mu t}}^{\text{Erlang}(n+1, \mu)} (1 - \rho) \rho^n \\ &= \mu(1 - \rho) e^{-\mu t} \sum_{n=0}^{\infty} \frac{(\mu \rho t)^n}{n!} = \mu(1 - \rho) e^{-\mu(1-\rho)t} \end{aligned}$$

$$\boxed{f_T(t) = (\mu - \lambda) e^{-(\mu - \lambda)t}} \quad \text{exponential distribution } \text{Exp}(\mu - \lambda)$$

The distribution of the sojourn time (continued)

The same result can be derived also by using the result for the Laplace transform of a random sum.

$$\left\{ \begin{array}{l} \mathcal{G}_{N+1}(z) = \frac{(1-\rho)z}{1-\rho z} \quad N+1 \sim \text{Geom}(1-\rho), \text{ starts from } 1 \\ f_S^*(s) = \frac{\mu}{\mu+s} \end{array} \right.$$

$$\begin{aligned} f_T^*(s) &= \mathcal{G}_{N+1}(f_S^*(s)) = \frac{(1-\rho)\frac{\mu}{\mu+s}}{1-\rho\frac{\mu}{\mu+s}} = \frac{\mu-\lambda}{(\mu+s)-\lambda} = \frac{(\mu-\lambda)}{(\mu-\lambda)+s} \\ &\Rightarrow \sim \text{Exp}(\mu-\lambda) \end{aligned}$$

Distribution of the waiting time

The waiting time W consists of the service times of the customers in the system upon the arrival

$$W = S_1 + \cdots + S_N, \quad \text{where } S_i \sim \text{Exp}(\mu) \text{ and } N \sim \text{Geom}_0(1 - \rho) \text{ (starts from 0)}$$

If $N = 0$ there are no terms in the sum and $W = 0$.

The tail distribution of W is derived by conditioning

$$\begin{aligned} \text{P}\{W > t\} &= \underbrace{\text{P}\{W > t \mid N = 0\}}_0 \text{P}\{N = 0\} + \text{P}\{W > t \mid N > 0\} \underbrace{\text{P}\{N > 0\}}_\rho \\ &= \rho \cdot \text{P}\{W > t \mid N > 0\} \end{aligned}$$

By the memoryless property of the geometric distribution N conditioned on $N > 0$ is distributed as $\text{Geom}(\rho)$ (starts from 1).

Thus the sum $S_1 + \cdots + S_N$ conditioned on $N > 0$ is distributed precisely as $S_1 + \cdots + S_{N+1}$ before and obeys the distribution $\text{Exp}(\mu - \lambda)$.

$$\text{P}\{W > t\} = \rho e^{-(\mu - \lambda)t}$$

The waiting time is 0 with a finite probability $PW = 0 = 1 - \text{P}\{W > 0\} = 1 - \rho$. This, of course, is equal to the empty queue probability $\text{P}\{N = 0\}$.

Finite queue: the $M/M/1/K$ system

Let there be K system places (waiting room + server)

The equilibrium equations across the cuts are the same as before

$$\pi_n = \rho^n \pi_0, \quad n = 0, 1, \dots, K$$

The only difference is in the normalization

$$\sum_{n=0}^K \pi_n = 1 \quad \Rightarrow \quad \pi_0 = (1 + \rho + \dots + \rho^K)^{-1} = \frac{1 - \rho}{1 - \rho^{K+1}}$$

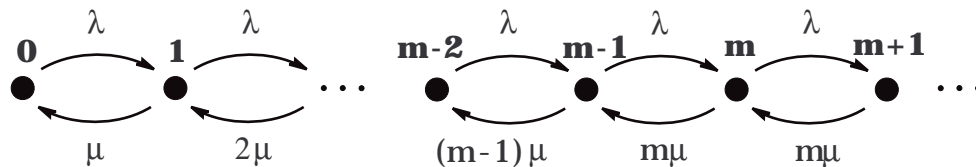
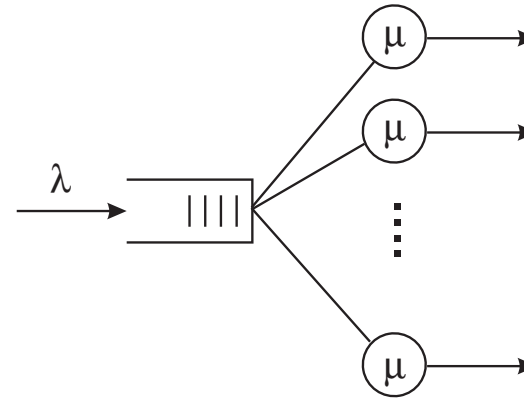
$$\pi_n = \frac{\rho^n}{1 + \rho + \dots + \rho^K} = \frac{1 - \rho}{1 - \rho^{K+1}} \rho^n$$

- The probability π_K of state K is the probability that an arriving customer finds the system full (“the buffer overflows”).
- When $K = 1$, we have a single server loss system ,

$$\pi_n = \frac{\rho^n}{1 + \rho}, \quad n = 0, 1.$$

The $M/M/m$ queue (Erlang's waiting system)

- m parallel servers
- Poisson arrivals
- Exponential service time distribution



- The state transition diagram is, up to state m the same as in the loss system.
- Beyond that state, it is identical with the diagram of an $M/M/1$ queue where the capacity of the server is $m\mu$.

The balance equations can again be written by using the method of a cut:

$$\begin{cases} \lambda\pi_{n-1} = n\mu\pi_n, & n \leq m \\ \lambda\pi_{n-1} = m\mu\pi_n, & n > m \end{cases}$$

The solution up to a constant factor π_0 is

$$\begin{cases} \pi_n = \pi_0 \frac{(m\rho)^n}{n!}, & n \leq m \\ \pi_n = \pi_0 \frac{m^m \rho^n}{m!}, & n > m \end{cases} \quad \begin{array}{l} a = \lambda/\mu \quad \text{traffic intensity} \\ \rho = \lambda/m\mu = a/m \quad \text{traffic intensity per server.} \end{array}$$

The probability π_0 of state 0 is determined by the normalization condition $\sum_n \pi_n = 1$,

$$\pi_0 = \left(\underbrace{\sum_{n=0}^{m-1} \frac{(m\rho)^n}{n!}}_u + \underbrace{\frac{(m\rho)^m}{m!(1-\rho)}}_v \right)^{-1}$$

The probability P_q that upon an arrival all servers are busy and the customer has to wait is

$$P_q = C(m, a) = \sum_{n=m}^{\infty} \pi_n = \sum_{n=m}^{\infty} \frac{\pi_0 m^m \rho^n}{m!} = \frac{\pi_0 (m\rho)^m}{m!(1-\rho)} = \frac{v}{u+v}$$

Erlang's C formula
 $a = m\rho; \rho = a/m$

The mean number of waiting customers \bar{N}_q

$$\bar{N}_q = \sum_{n=0}^{\infty} n \pi_{m+n} = \sum_{n=0}^{\infty} n \pi_0 \frac{m^m \rho^{m+n}}{m!} = P_q \sum_{n=0}^{\infty} n (1 - \rho) \rho^n$$

The sum is of the same form as the mean queue length in an $M/M/1$ queue. Thus

$$\boxed{\bar{N}_q = P_q \frac{\rho}{1 - \rho}} \quad \bar{N} = m\rho + \bar{N}_q \quad \Rightarrow \quad \boxed{\bar{N} = m\rho + P_q \frac{\rho}{1 - \rho}}$$

By Little's result we obtain the mean waiting and sojourn times:

$$\boxed{\begin{cases} \bar{W} = \frac{\bar{N}_q}{\lambda} = P_q \cdot \frac{1}{m\mu - \lambda} \\ \bar{T} = \frac{\bar{N}}{\lambda} = \frac{1}{\mu} + \bar{W} = \frac{1}{\mu} + P_q \cdot \frac{1}{m\mu - \lambda} \end{cases}}$$

The distribution of the waiting time

$$P\{W > t\} = \underbrace{P\{W > t | N < m\}}_0 P\{N < m\} + P\{W > t | N \geq m\} \underbrace{P\{N \geq m\}}_{P_q}$$

$$P\{W > t\} = P_q e^{-(m\mu - \lambda)t}$$

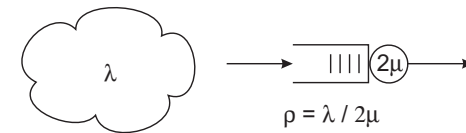
When $N \geq m$ the system behaves as an $M/M/1$ queue with capacity $m\mu$.

Example 1

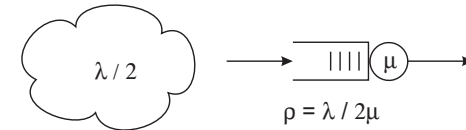
A printer is attached to the LAN of the department. The printing jobs are assumed to arrive with a Poissonian intensity λ and the actual printing times are assumed to obey the distribution $\text{Exp}(\mu)$.

The capacity of the printer has become insufficient with regard to the increased load. In order to improve the printing service, there are three alternatives:

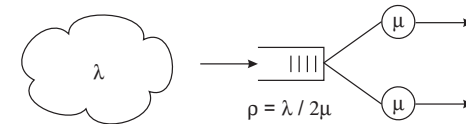
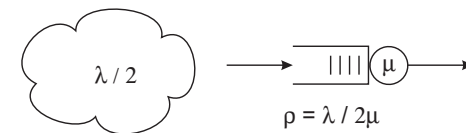
1. Replace the old printer by a new one twice as fast, i.e. with service rate 2μ .



2. Add another similar printer (service rate μ) and divide the users in two groups of equal size directing the works in each group to their own printer. The arrival rate of jobs to each printer is $\lambda/2$.



3. The same as alternative 2, but now there is a common printer queue where all jobs are taken and the job at the head of the queue is sent to whichever printer becomes free first.



Example 1 (continued)

Lut us compare the performance of the alternatives at different loads. As measure of performance we use the mean sojourn time of a job \bar{T} (time in system, from the arrival of the printing job to the full completion of the job).

1. In this case we have an $M/M/1$ queue with parameters λ and 2μ .

$$\rho = \frac{\lambda}{2\mu} \qquad \bar{T}_1 = \frac{1}{2\mu - \lambda} = \frac{1}{1 - \rho} \cdot \frac{1}{2\mu}$$

2. Now we have two separate $M/M/1$ queues with parameters $\lambda/2$ and μ .

$$\rho = \frac{\lambda/2}{\mu} = \frac{\lambda}{2\mu} \qquad \bar{T}_2 = \frac{1}{\mu - \lambda/2} = \frac{1}{1 - \rho} \cdot \frac{1}{\mu}$$

The load per server is the same as before. Now just everything happens two times slower (both arrivals and the service).

3. In the case of a common printing queue, an appropriate model is the $M/M/2$ queue with parameters λ and μ .

$$\rho = \frac{\lambda}{2\mu} \qquad \bar{T}_3 = \frac{1}{\mu} + P_q \frac{1}{2\mu - \lambda} \approx \begin{cases} \frac{1}{\mu} & \rho \ll 1 \\ \frac{1}{1 - \rho} \cdot \frac{1}{2\mu} & \rho \approx 1 \end{cases}$$

Example 1: Summary of the comparison

Take case 1 as the reference: calculate the sojourn times in cases 2 and 3 in relation to that in case 1.

	T_2/T_1	T_3/T_1
$\rho \ll 1$	2	2
$\rho \approx 1$	2	1

- Alternative 1, i.e. one fast printer is the best one.
- In alternative 2, the sojourn time is twice as long as in case 1.
- In case 3, the second printer does not help at all at low loads: each job is taken directly into the service (without waiting) but the actual printing takes twice the time as with the fast printer.
- At heavy loads, the mean sojourn time of case 3 is the same as in case 1 (in both cases it consists mainly of the waiting). Two slow printers fed by a common queue discharge the work in the queue as efficiently as one fast printer.
- This is not the case for the alternative 2. When the queues are separate, it is possible that one printer stays idle while there are jobs waiting in the queue for the other printer. This deteriorates the overall performance in such a way that also at high loads alternative 2 is on the average two times slower than alternative 1.

Example 2

- A telephone switch is modelled as an $M/M/m$ system (when all lines are busy, the callers are let to wait by signalling them the ring tone)
- How many lines (m) are needed that the probability that a caller has to wait longer than time t_{\max} is less than 1 % ?

$$P_q e^{-(m\mu - \lambda)t_{\max}} < 0.01 \quad \Rightarrow \quad m > \frac{\log(100P_q) + \lambda t_{\max}}{\mu t_{\max}}$$

P_q is a function of m (monotonically decreasing); thus the inequality is still an implicit one.

It can be solved by trying sequentially values $m = 1, 2, 3, \dots$ until the inequality is satisfied.

By letting the callers to wait for a free line for a while before blocking them, the number of blocked calls can be reduced or, conversely, the load of the system ρ can be increased in comparison with a loss system with the same blocking probability.