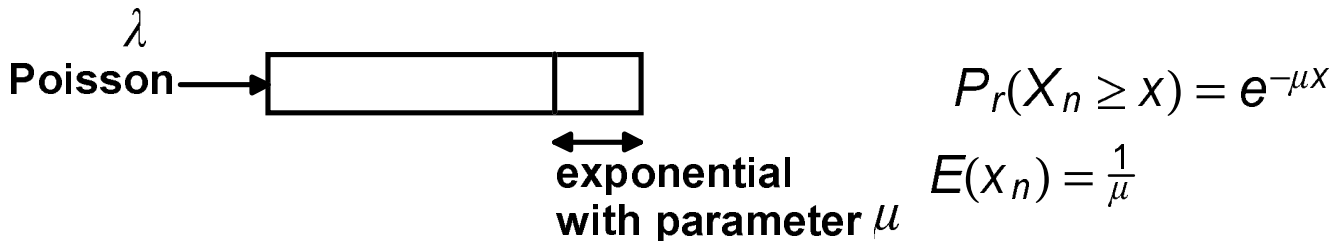


## The M/M/1 Queueing System

**M (Poisson arrival process) / M (exponential service times) / 1 (1 server)**



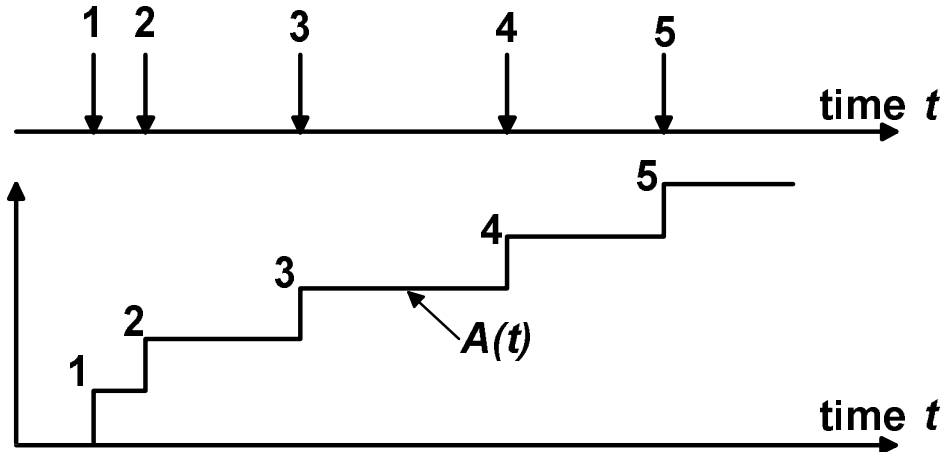
**Notation:** . / . / . / .

Arrival process	Service times	# of servers	Max. # of customers in the system
M: Poisson	M: Exponential		
D: Deterministic	D: Deterministic		
G: General	G: General		

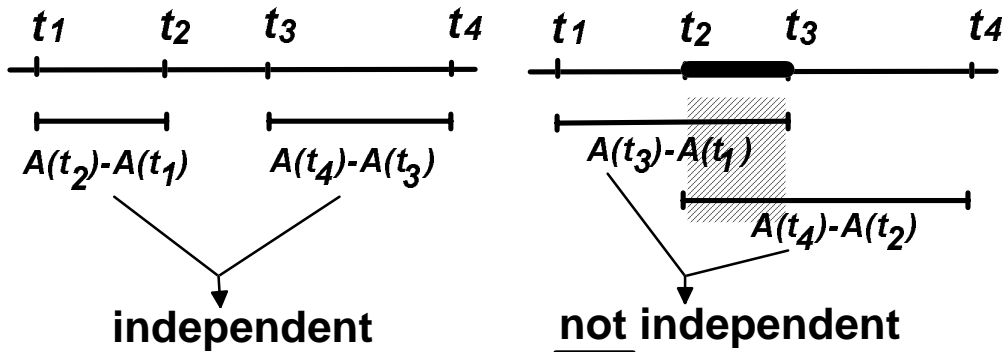
- **We also assume that service times are mutually independent and also independent of all interarrival times.**
- **The server is always serving a customer if any customer is in the system. Assume FCFS service to be specific.**

# Poisson Process of Rate $\lambda$

A Poisson process  $A(t)$  is a counting process.  
 For each  $t \geq 0$ ,  $A(t)$  is a random variable denoting the number of arrivals from 0 to  $t$ .



Number of arrivals in disjoint time intervals are independent.



Number of arrivals in any interval of length  $\tau$  is Poisson with parameter  $\lambda \cdot \tau$

$$P\{A(t+\tau) - A(t) = n\} = e^{-\lambda\tau} \frac{(\lambda\tau)^n}{n!}, n = 0, 1, \dots$$

$$E\{A(t+\tau) - A(t)\} = \lambda \cdot \tau \quad \lambda : \text{arrival rate}$$

## Properties of Poisson process

Let  $t_n =$  time of  $n$  th arrival

$\tau_n = t_{n+1} - t_n =$  interarrival time

$$\bullet P(\tau_n \geq s) = P\{A(t_n + s) - A(t_n) = 0\} = e^{-\lambda s}$$

( Interarrival times are exponentially distributed with parameter  $\lambda$  , mean  $\frac{1}{\lambda}$ , variance  $\frac{1}{\lambda^2}$  )

$$\bullet P(\tau_n \geq r + t \mid \tau_n \geq t) = \frac{P(\tau_n \geq r+t, \tau_n \geq t)}{P(\tau_n \geq t)} = \frac{e^{-\lambda(r+t)}}{e^{-\lambda t}} = e^{-\lambda r} = P(\tau_n \geq r)$$

(memoryless)

• For any  $t$ , and any (small)  $\delta$ :

$$P\{A(t+\delta) - A(t) = 0\} = 1 - \lambda\delta + o(\delta)$$

$$P\{A(t+\delta) - A(t) = 1\} = \lambda\delta + o(\delta)$$

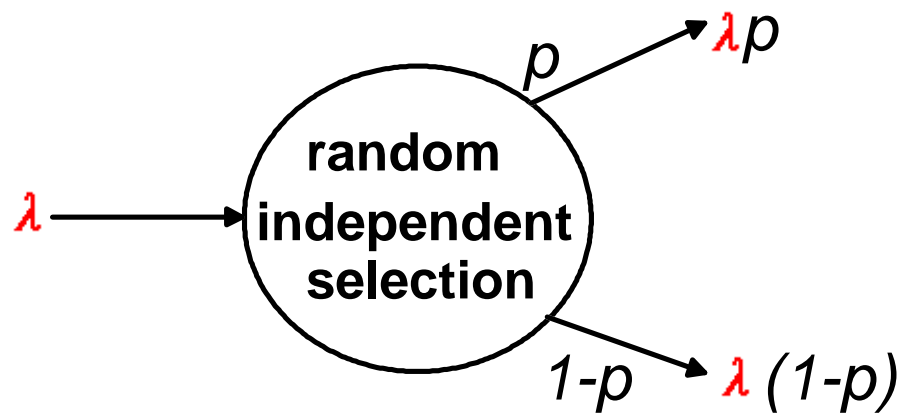
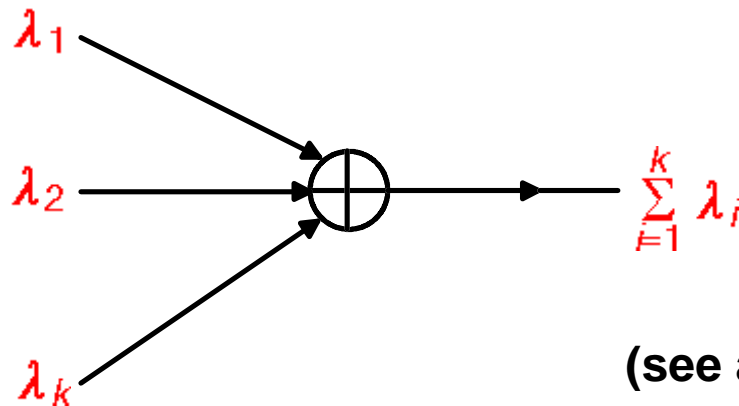
$$P\{A(t+\delta) - A(t) \geq 2\} = o(\delta)$$

where  $\lim_{\delta \rightarrow 0} \frac{o(\delta)}{\delta} = 0$

These follow from definition:

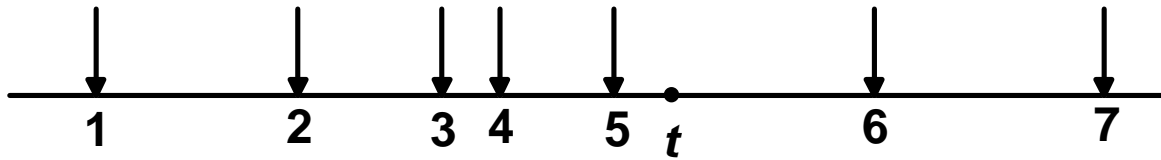
$$P\{A(t+\delta) - A(t) = n\} = \frac{e^{-\lambda\delta} (\lambda\delta)^n}{n!}$$

- If  $A_1(t), A_2(t), \dots, A_k(t)$  are independent Poisson processes of rate  $\lambda_1, \lambda_2, \dots, \lambda_k$ , then  $A_1(t) + A_2(t) + \dots + A_k(t)$  is a Poisson process of rate  $\lambda_1 + \lambda_2 + \dots + \lambda_k$

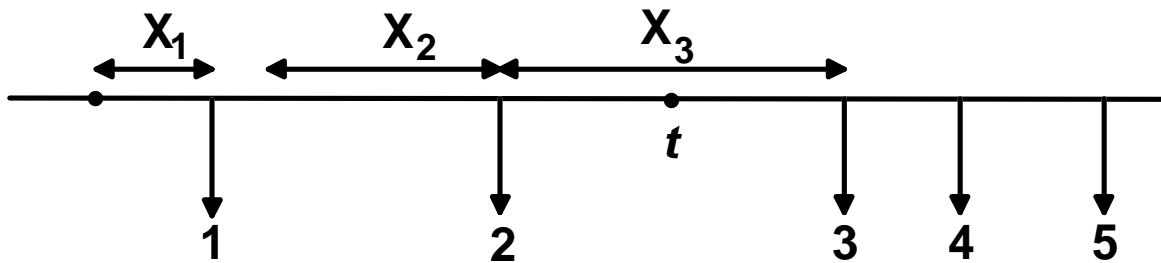


If each arrival of a Poisson process is independently sent to system 1 with prob.  $p$  and system 2 with prob.  $1-p$ , the arrivals to each system are Poisson and independent.(see also Ex.3.11a)

### Arrivals



### Departures



Starting at a particular time  $t$ , the subsequent arrivals do not depend on what has happened in the past, and the subsequent departures depend only on the number  $N(t)$  of customers in the system at time  $t$ .

In particular, since service time exponential, it makes no difference, how long the current customer has been in service; the remaining time until departure is still exponential.

Future # of customers in the system depends on past numbers only through the present number  $N(t)$ .

**We focus at times  $0, \delta, 2\delta, 3\delta, \dots, k\delta, \dots$  ( $\delta$  small).**

$N_k \stackrel{\text{def}}{=} N(k\delta) = \#$  of customers in the system at time  $k\delta$ .

$P_{ij} = P\{N_{k+1} = j \mid N_k = i\}$  (transition probabilities)

$P_{00} = 1 - \lambda\delta + o(\delta)$  (no arrival)

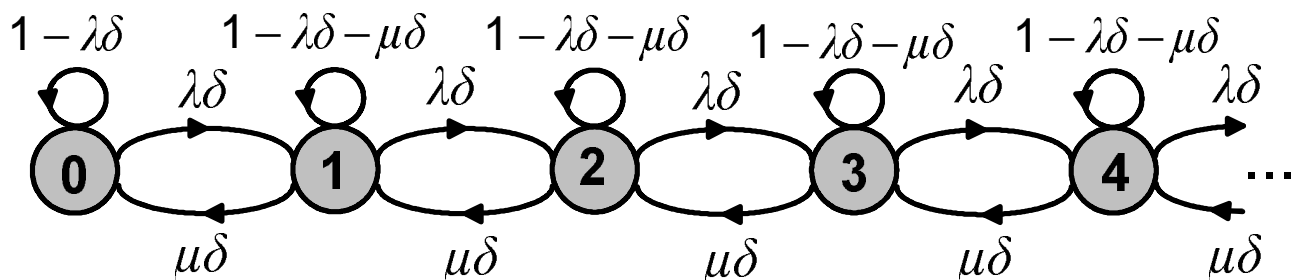
$P_{ii} = 1 - \lambda\delta - \mu\delta + o(\delta), i \geq 1$  (no arrival/departure)

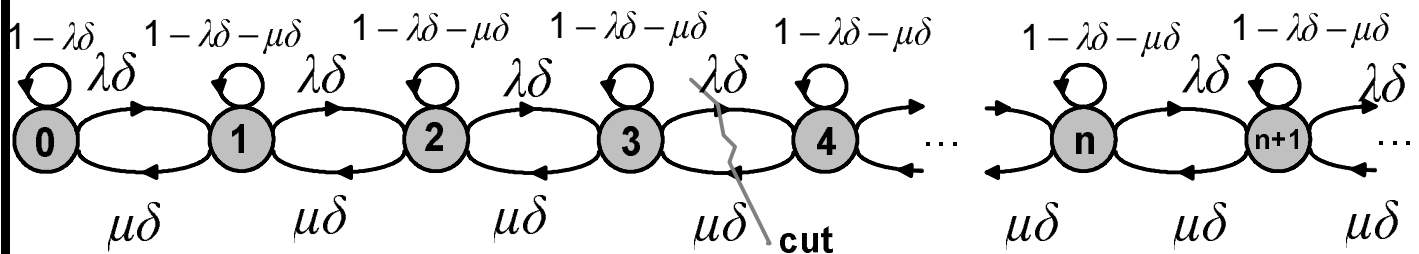
$P_{i,i+1} = \lambda\delta + o(\delta), i \geq 0$  (one arrival)

$P_{i,i-1} = \mu\delta + o(\delta), i \geq 1$  (one departure)

**[Note: for any state  $n \geq 1$ , the server is busy and probability of departure is  $P_r(X \leq \delta) = 1 - e^{-\mu\delta} = \mu\delta + o(\delta)$ ]**

**$P_{ij} = o(\delta), j \neq i, i+1, i-1$  (i.e. the probability of multiple arrivals/departures is negligible.)**





Let  $p_n$  be the “steady-state” probability that the system is in state  $n$ . [ i.e.  $p_n = \lim_{t \rightarrow \infty} P(N(t) = n)$  ]

**Note:** over an arbitrarily long period of time, the number of transitions from  $n$  to  $n+1$  is the same as from  $n+1$  to  $n$  (plus or minus one).

Thus for any  $n$ :

$$p_{n-1} \lambda \delta = p_n \mu \delta \Rightarrow p_n = \left(\frac{\lambda}{\mu}\right) p_{n-1} = \left(\frac{\lambda}{\mu}\right)^2 p_{n-2} = \dots = \left(\frac{\lambda}{\mu}\right)^n p_0$$

Define

$$\rho = \frac{\lambda}{\mu} \text{ (“utilization factor”)}$$

$$\Rightarrow p_n = \rho^n p_0, \quad n = 1, 2, \dots$$

**To find  $p_0$ :**

$$\sum_{n=0}^{\infty} p_n = 1 \Rightarrow \sum_{n=0}^{\infty} \rho^n p_0 = 1 \Rightarrow p_0 \cdot \frac{1}{1-\rho} = 1 \Rightarrow p_0 = 1 - \rho$$

$$\Rightarrow p_n = (1 - \rho) \cdot \rho^n, \text{ for } n \geq 0 \quad \rho = \frac{\lambda}{\mu} < 1$$

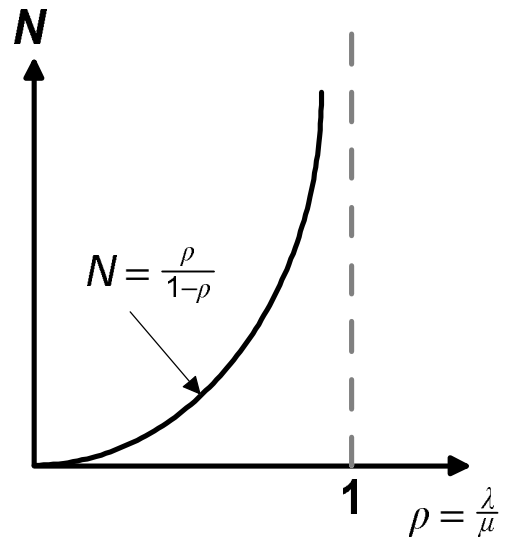
$\rho$  = **Probability that the system has at least one customer (=  $1 - p_0$ )**  
= **Probability server is busy**

**The expected number  $N$  of customers in the system is**

$$N = \sum_{n=0}^{\infty} n p_n = \sum_{n=0}^{\infty} (1 - \rho) n \cdot \rho^n = \frac{\rho}{1 - \rho}$$



Number in the system blows up as  $\rho \rightarrow 1$ ,  $N \rightarrow \infty$ ; i.e. as arrival rate  $\lambda$  approaches service rate  $\mu$ .



From Little's theorem, average customer delay  $T$  is

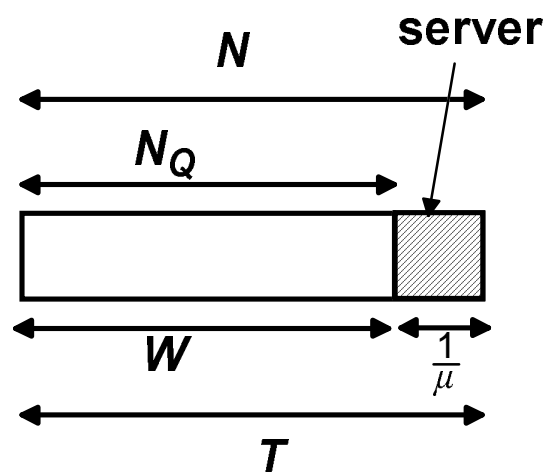
$$T = \frac{N}{\lambda} = \frac{\rho}{\lambda(1-\rho)} = \frac{1}{\mu-\lambda} \quad , \lambda < \mu$$

Average time in queue  $W$  is

$$W = \frac{N}{\lambda} - \frac{1}{\mu} = \frac{1}{\mu-\lambda} - \frac{1}{\mu} = \frac{\rho}{\mu-\lambda}$$

Average number of customers in queue  $N_Q$  is

$$N_Q = \lambda W = \frac{\rho^2}{1-\rho}$$



### Example 1 (Scaling of an $M/M/1$ queue)

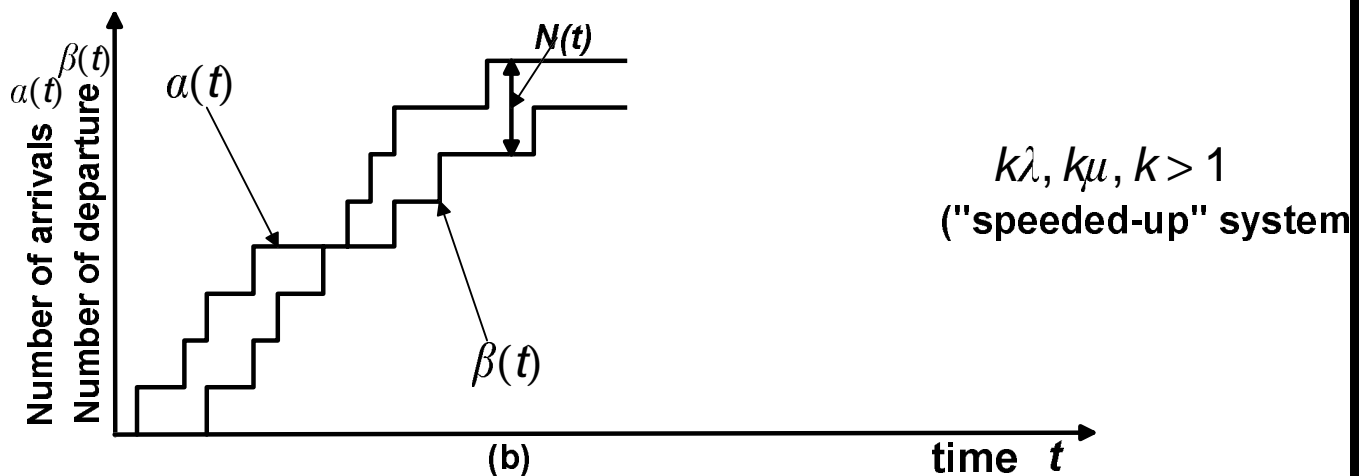
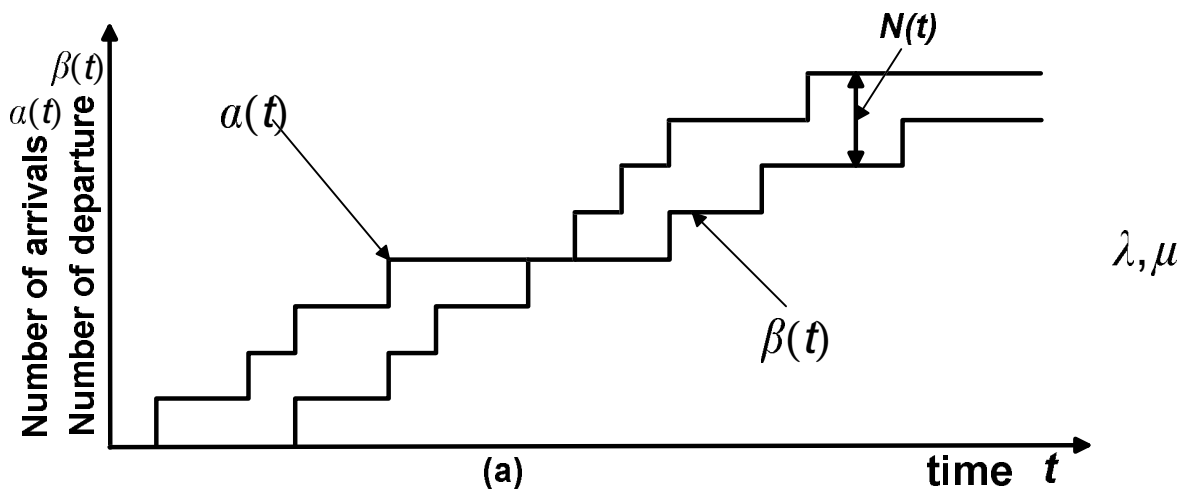
$$N = \frac{\rho}{1-\rho}, \rho_n = (1-\rho)\rho^n \text{ for } n \geq 0, N_Q = \frac{\rho^2}{1-\rho}$$

where  $\rho = \frac{\lambda}{\mu}$

If one scales the arrival rate  $\lambda$  and the service rate  $\mu$ , by a constant factor  $k$ ,  $N$ ,  $N_Q$ , and  $\rho_n$  are unchanged

$$T = \frac{1}{\mu-\lambda}, \quad W = \frac{\rho}{\mu-\lambda}$$

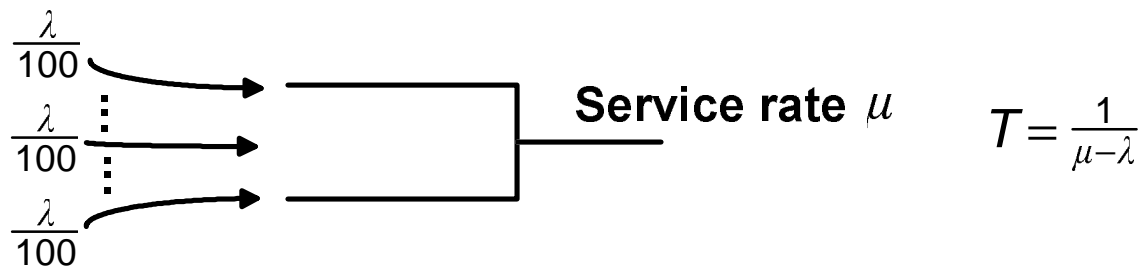
System delay  $T$  and queueing delay  $W$  vary as  $\frac{1}{k}$



## Example 2 (statistical multiplexing verses FDM)

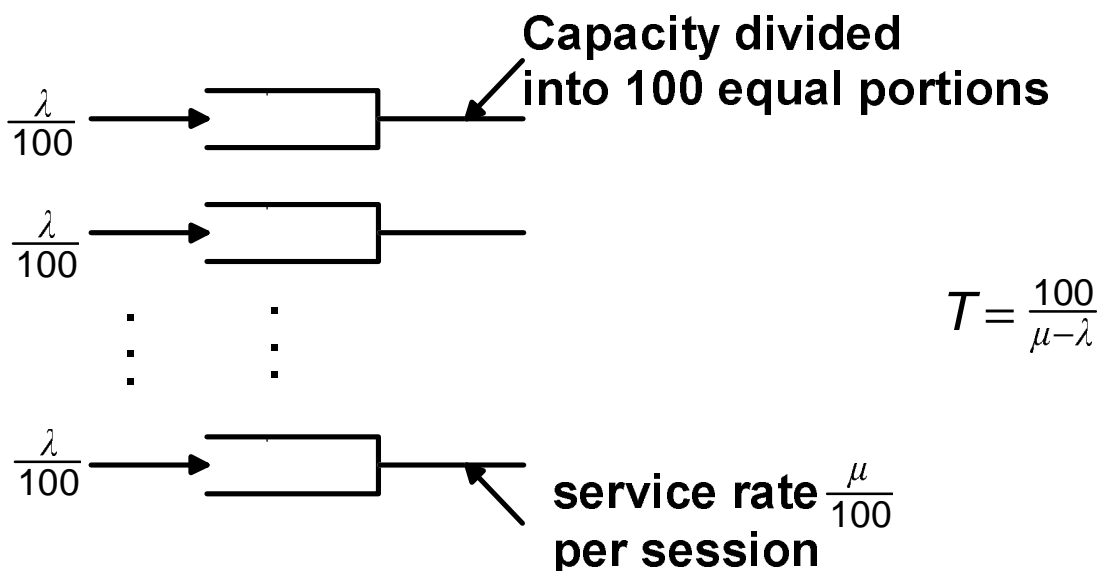
Consider 100 sessions with Poisson arrivals of combined rate  $\lambda$  and exponentially distributed packet lengths sharing a link with service rate  $\mu$  packets/sec.

### Statistical multiplexing



### Frequency Division Multiplexing (FDM)

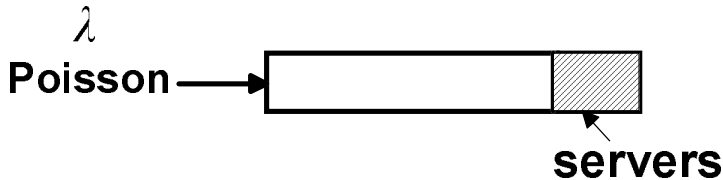
If FDM is used, each session has rate  $\frac{\lambda}{100}$  and “sees” service rate  $\frac{\mu}{100}$



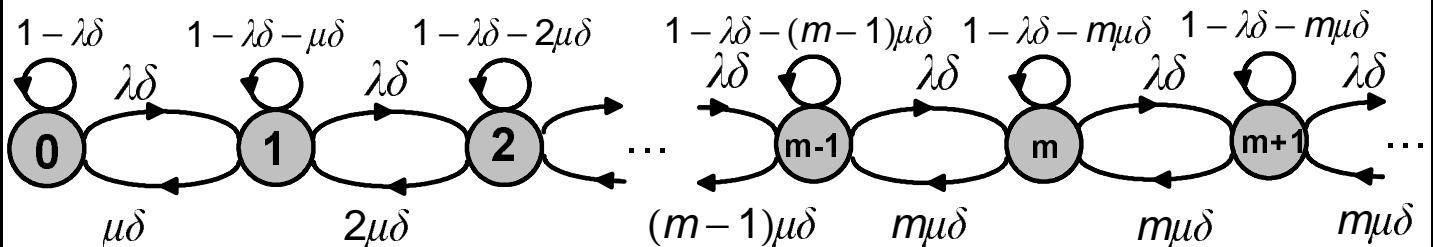
## M/M/m queue

Poisson arrivals of rate  $\lambda$

$m$  servers, each exponentially distributed with rate  $\mu$



- Given  $n$  customers in the system,  $n \leq m$ , a new arrival will occur in an increment  $\delta$  with probability  $\lambda\delta$ . A departure will occur with probability  $n\mu\delta$ .
- For  $n > m$  a departure occurs with probability  $m\mu\delta$ .



$$\lambda p_{n-1} = n\mu p_n, \quad n \leq m$$

$$\lambda p_{n-1} = m\mu p_n, \quad n > m$$

If  $n \leq m$ :

$$p_n = \left(\frac{\lambda}{n\mu}\right)p_{n-1} = \frac{\lambda^2}{n(n-1)\mu^2}p_{n-2} = \dots = \frac{\lambda^n}{n!\mu^n}p_0$$

If  $n > m$ :

$$p_n = \left(\frac{\lambda}{m\mu}\right)p_{n-1} = \left(\frac{\lambda}{m\mu}\right)^2 p_{n-2} = \dots = \left(\frac{\lambda}{m\mu}\right)^{n-m} p_m = \frac{\lambda^n p_0}{m^{n-m} m! \mu^n}$$

Let  $\rho = \frac{\lambda}{m\mu} < 1 \Rightarrow p_n = \begin{cases} p_0 \frac{(m\rho)^n}{n!}, & n \leq m \\ p_0 \frac{m^m \rho^n}{m!}, & n > m \end{cases}$

$$\sum_{n=0}^{\infty} p_n = 1 \Rightarrow \dots p_0 = \left[ \sum_{n=0}^{m-1} \frac{(m\rho)^n}{n!} + \frac{(m\rho)^m}{m!(1-\rho)} \right]^{-1}$$

The probability an arriving customer will find all servers busy (and will have to wait) is

$$P(\text{all servers busy}) = P_Q = \sum_{n=m}^{\infty} p_n = \dots = \frac{p_0 (m\rho)^m}{m!(1-\rho)}$$

**Erlang C formula** used in telephony

**The expected number of customers in queue is**

$$N_Q = \sum_{n=0}^{\infty} n p_{m+n} = \sum_{n=0}^{\infty} n p_0 \cdot \frac{m^m \rho^{m+n}}{m!} = \frac{p_0 (m\rho)^m}{m!} \underbrace{\sum_{n=0}^{\infty} n \rho^n}_{\frac{\rho}{(1-\rho)^2}}$$

$$= P_Q \cdot \frac{\rho}{1-\rho} \Rightarrow \frac{N_Q}{P_Q} = \frac{\rho}{1-\rho} \quad (\rho = \frac{\lambda}{m\mu})$$

**Waiting time :**

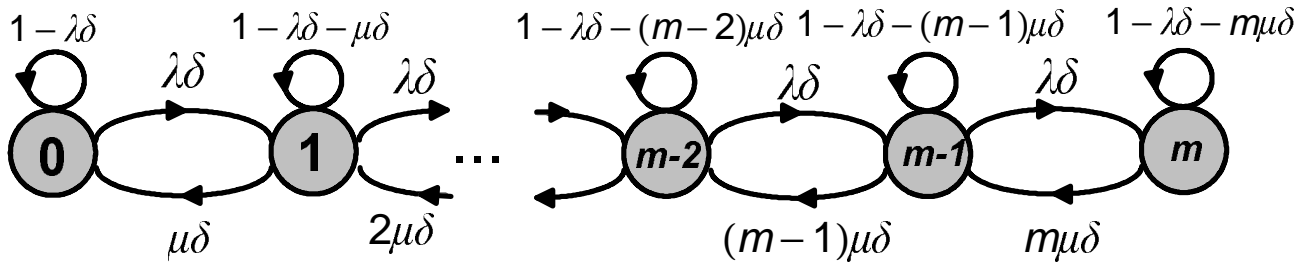
$$W = \frac{N_Q}{\lambda} = \frac{\rho P_Q}{\lambda(1-\rho)}$$

**Total time in system:**

$$T = W + \frac{1}{\mu}$$

## M/M/m/m Queue

Customers arriving when all m servers busy are thrown away, never to return.



$$\begin{aligned} \lambda p_{n-1} &= n\mu p_n, \quad n = 1, 2, \dots, m \\ p_n &= \frac{p_0}{n!} \cdot \left(\frac{\lambda}{\mu}\right)^n, \quad n = 1, 2, \dots, m \\ \sum_{n=0}^m p_n &= 1 \rightarrow p_0 = \left[ \sum_{n=0}^m \left(\frac{\lambda}{\mu}\right)^n \frac{1}{n!} \right]^{-1} \end{aligned}$$

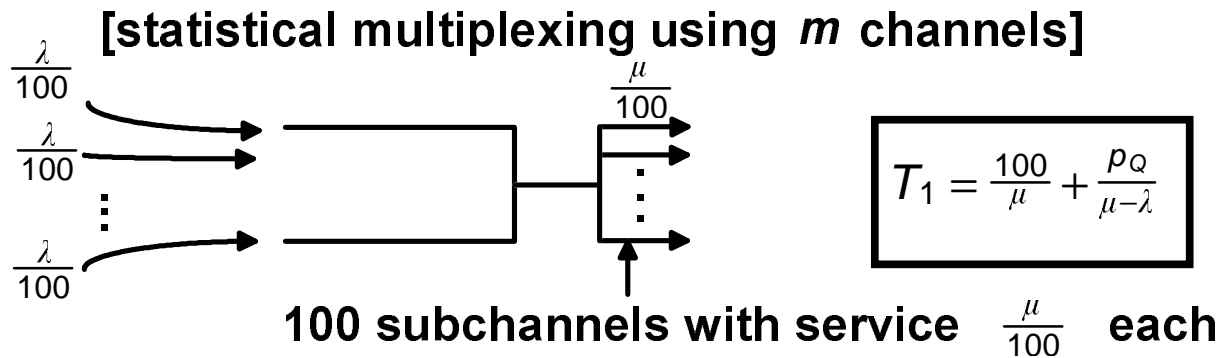
Probability that a customer finds all servers busy is

$$p_m = \frac{\frac{(\frac{\lambda}{\mu})^m}{m!}}{\sum_{n=0}^m \frac{(\frac{\lambda}{\mu})^n}{n!}}$$

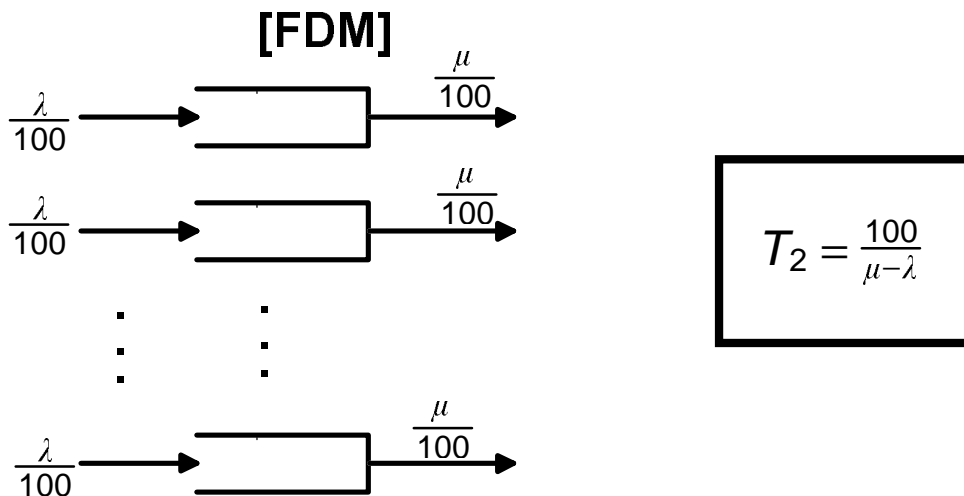
**Erlang B formula** ( $N, T$  less than for  $M/M/m$  but not all customers get served)

If one uses these models for session arrivals in a voice network, one sees that customer behavior is somewhat between  $M/M/m$  and  $M/M/m/m$  behavior—some customers go away if they can't get through and some keep trying.

**Example: Assume  $m=100$  sessions sharing a link. Assume 100 frequency bands, but packets are assigned to any available band. This is an  $M/M/100$ .**

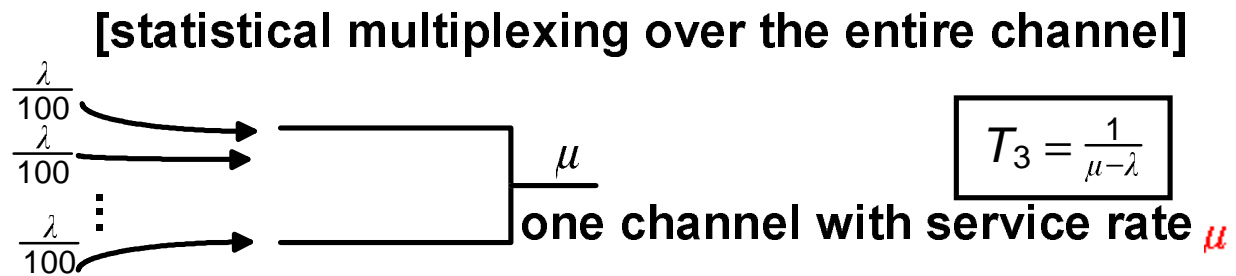


**Under light load is almost the same as FDM**





**Under heavy load delay is almost the same as statistical multiplexing**



$\lambda$  **small**:  $T_1 \approx T_2 \approx 100T_3$ ,  $\lambda$  **large** ( $\lambda \approx \mu$ ):  $T_1 \approx T_3 = \frac{T_2}{100}$

## M/M/∞ queue.

As in M/M/m, but with  $m = \infty$  servers.

Expressions can be found by taking the limit  $m \rightarrow \infty$  in the M/M/m expressions:

$$\rho_n = \frac{1}{n!} \left(\frac{\lambda}{\mu}\right)^n \rho_0$$

$$\rho_0 = \left[1 + \sum_{n=1}^{\infty} \left(\frac{\lambda}{\mu}\right)^n \frac{1}{n!}\right]^{-1} = e^{-\lambda/\mu}$$

$$\rho_n = \frac{e^{-\lambda/\mu}}{n!} \left(\frac{\lambda}{\mu}\right)^n$$

Average # in the system  $N = \frac{\lambda}{\mu} (= \sum_{n=0}^{\infty} n \rho_n)$

$$T = \frac{N}{\lambda} = \frac{1}{\mu}$$