

# Networking

Prof. Dr.-Ing. Holger Hermanns

Dependable Systems & Software  
Saarland University

Summer 04

*Intros*  
*Tue, 10:30*

Session D: Queueing Basics

# Queueing Basics

Study chapter 3 of  
[Bertsekas/Galagher]

## Our goal today:

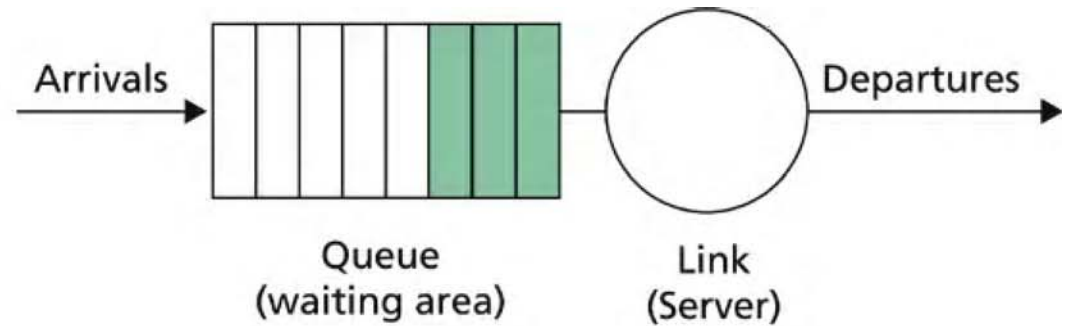
- ❑ develop an understanding of the fundamentals of capacity and delay behaviour of the Internet core

- ❑ approach:
  - straight into the heart of the matters

## Overview:

- ❑ Basic Queueing model
- ❑ Little's law
- ❑ Markov everywhere:  $M/M/...$
- ❑ From single queues to networks of queues
- ❑ Beyond Markov properties

# Context



## □ Why study queues?

- Framework for analyzing network queueing delay.

## □ Typical measures of interest:

- The average number of customers in the system.
- The average delay per customer.

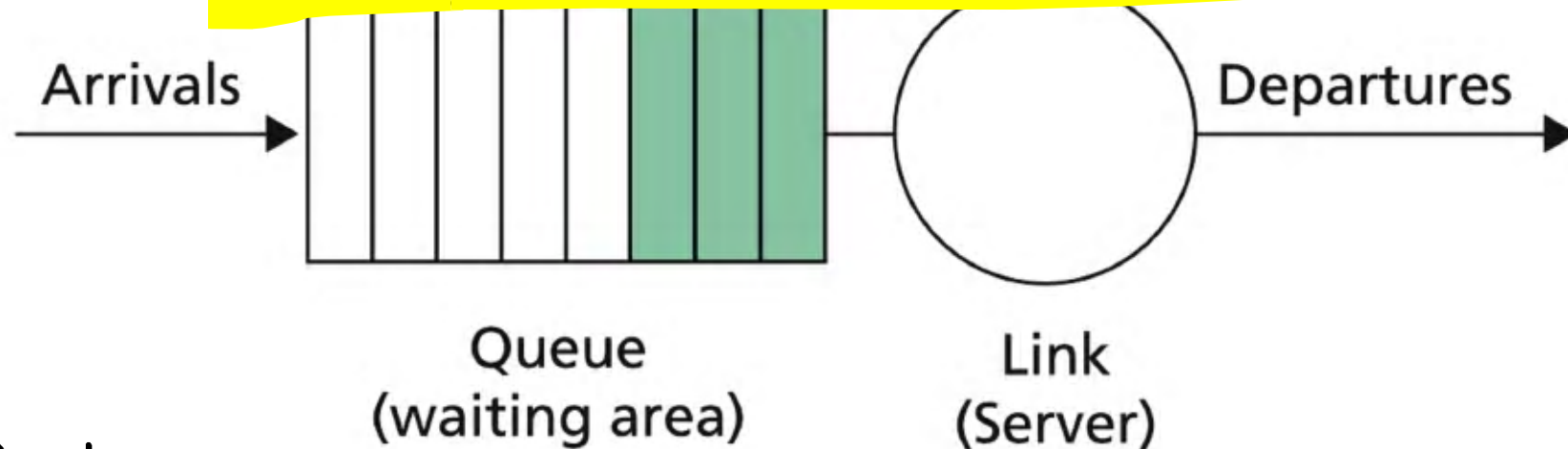
## □ What do we need as input?

- The customer arrival rate  $\lambda$ .
- The customer service rate  $\mu$ .

Think of:

- Server: Router
- Customer: Packet
- $\lambda$  arriving packets per second
- $L$  bits per packet
- $C$  bits per second (link bandwidth)
- $\mu = L/C$

# Queues in the Internet



## □ Router

- e.g. buffering of incoming packets (waiting to be routed)
- e.g. buffering of outgoing packets (waiting to get on the link)

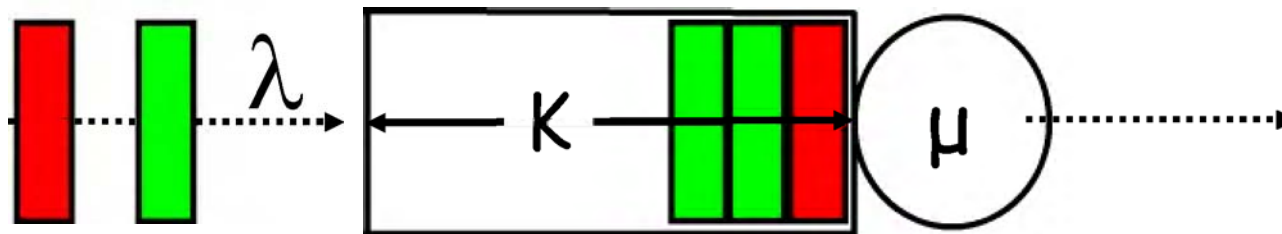
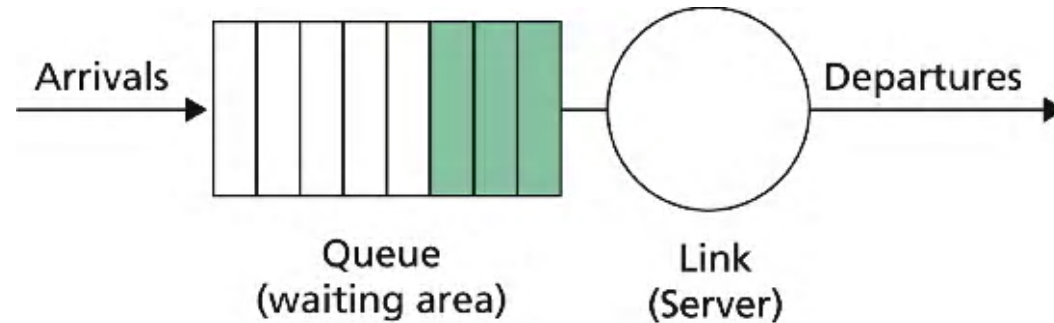
## □ Client

- e.g. buffering of streaming media

## □ Server

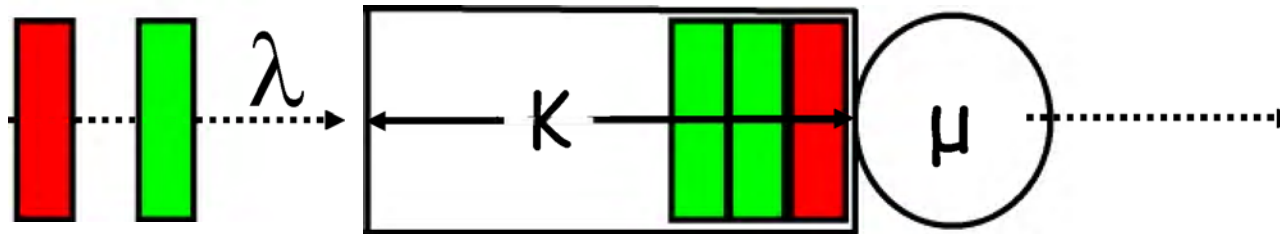
- e.g. buffering requests to be processed

# The generic model of a queue



- Buffer of size  $K$   
(# of customers in system)
- Customers arrive at rate  $\lambda$
- Customers are served at rate  $\mu$
  
- $\lambda$  and  $\mu$  are average rates  
but we don't bother much about this for now

# Queues: General Observations



## □ Increase in $\lambda$ :

- more customers in queue (on average),
- longer delays to get through queue;

## □ Decrease in $\mu$ :

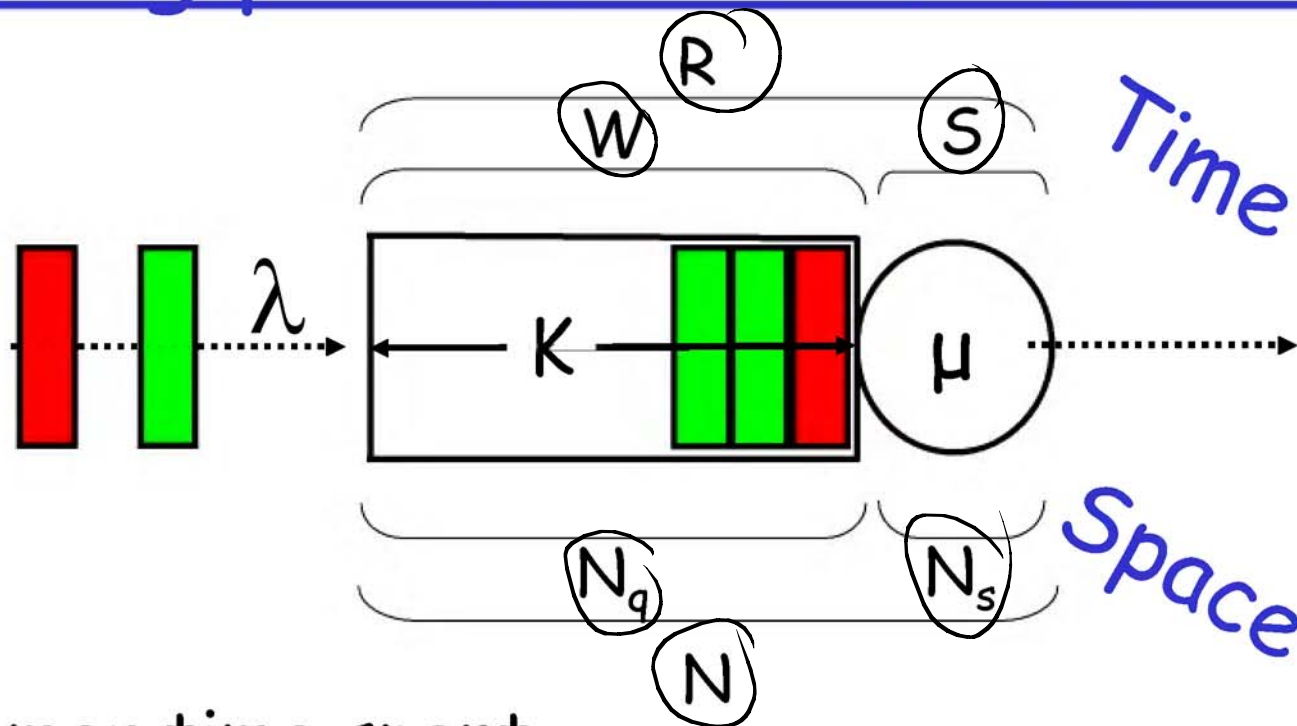
- longer delays to get processed,
- leads to more customers in queue;

## □ Decrease in $K$ :

- customer drops more likely,
- less delay for the "average" customer accepted into the queue.



# Queueing parameters of interest



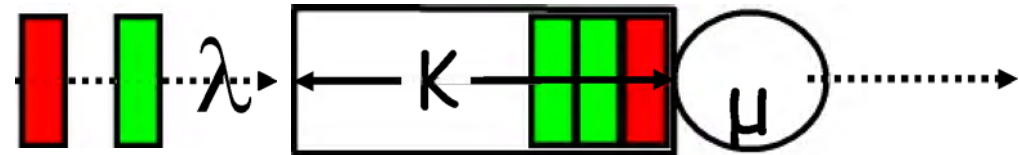
## □ Customer time spent

- in queue:  $W$  (waiting time)
- in service:  $S$
- in the complete system:  $R$  (response time)

## □ Number of customers

- in queue:  $N_q$
- in service:  $N_s$  ('utilisation')
- in the complete system:  $N$

# Little's Law



- Let  $c_i$  be the  $i^{\text{th}}$  customer arriving to the queue
- Let  $N_i$  be the number customers already in the queue when  $c_i$  arrives
- Let  $R_i$  be time spent by  $c_i$  in the system
  - both in the queue and while being served
- If  $K = \infty$  (unlimited queue size) then

*expectation*  
*(or 'average' if you like)*

$$E[N] = \lim_{i \rightarrow \infty} E[N_i] = \lambda \lim_{i \rightarrow \infty} E[R_i] = E[R]$$



# Little's Law in various flavours

On the long run ( $t \rightarrow \infty$ ),

## □ System

- $E[N]$  number of customers in system, **on average**.
- $E[R]$  response time, **on average**.

$$E[N] = \lambda E[R]$$

3/sec

## □ Queue

- $E[N_q]$  queue length, **on average**.
- $E[W]$  waiting time, **on average**.

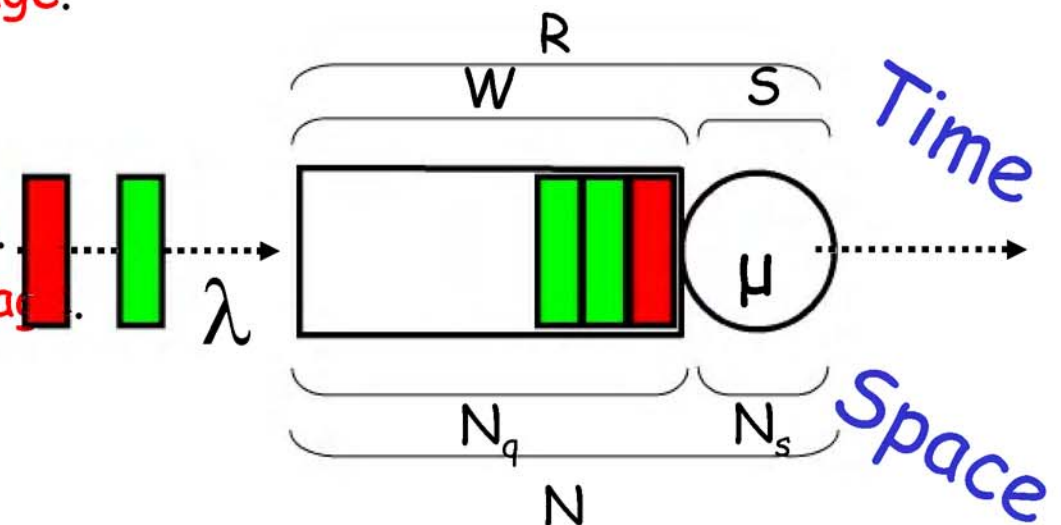
4.5 =  $E[N_q] = \lambda E[W]$

1.5 sec

## □ Server



- $E[N_s]$  utilisation **on average**.
- $E[S]$  servicing time **on average**.

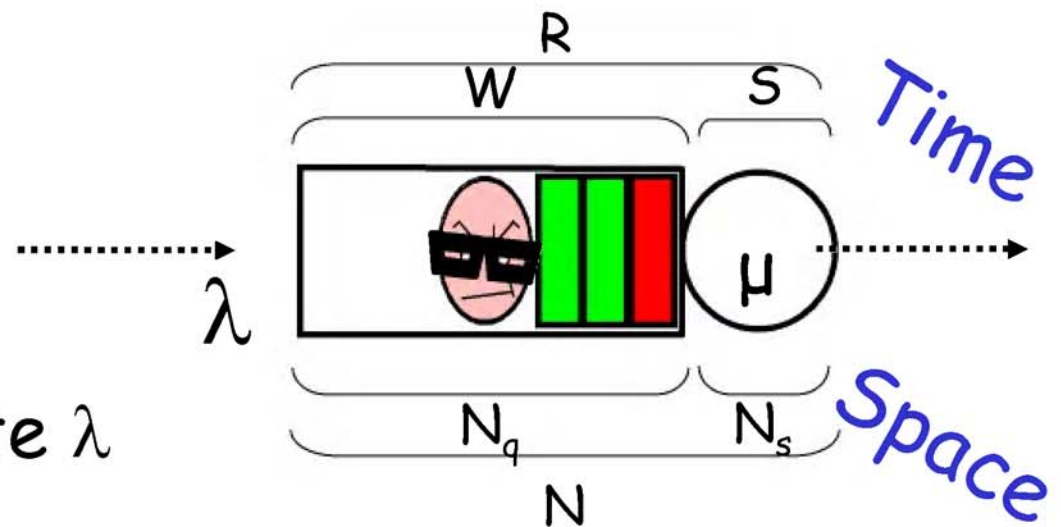
$$E[N_s] = \lambda E[S]$$



# Little's Law


How to understand the law:

- An average customer  sees an average queue length  $E[N_q]$ .
- It takes him  $E[W]$  on average to travel through the queue.
- When  leaves the queue, on average the queue still has length  $E[N_q]$ , being refilled with rate  $\lambda$  during  $E[W]$ .




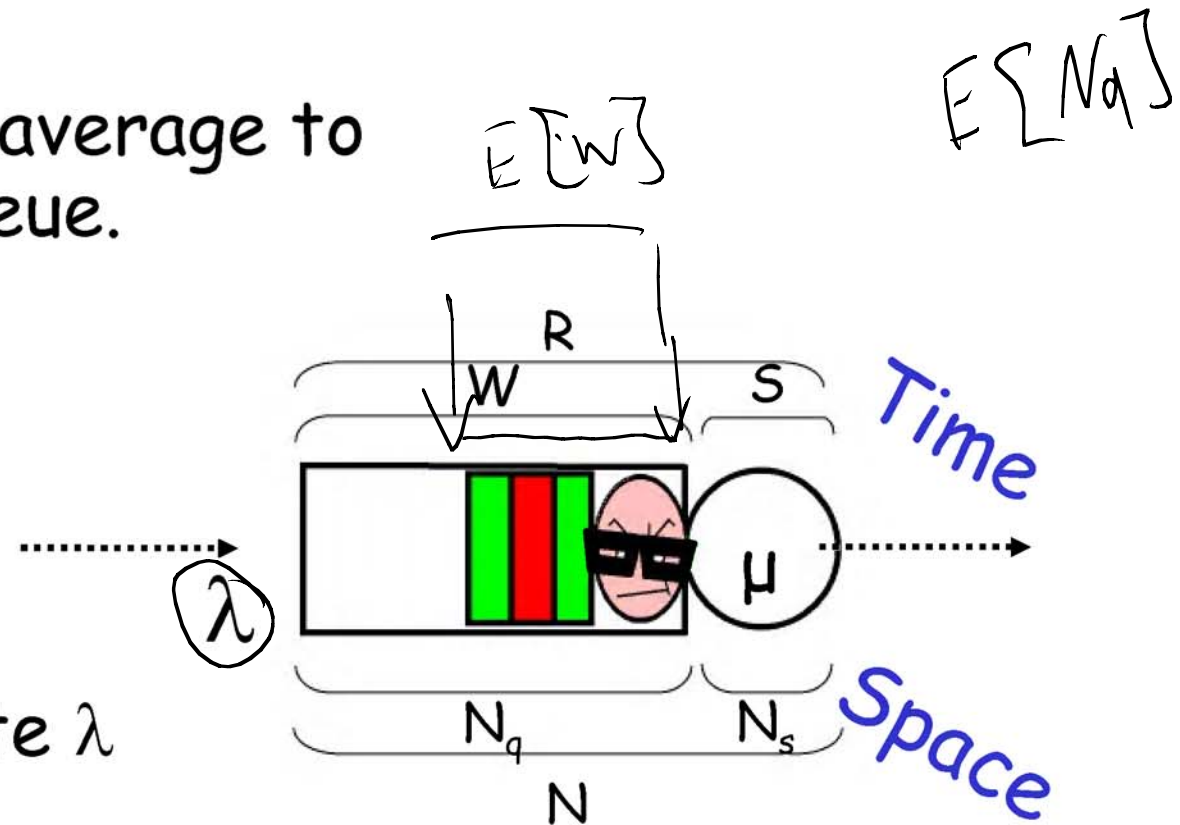
# Little's Law

How to understand the law:

□ An average customer  sees an average queue length  $E[N_q]$ .

□ It takes him  $E[W]$  on average to travel through the queue.

□ When  leaves the queue, on average the queue still has length  $E[N_q]$ , being refilled with rate  $\lambda$  during  $E[W]$ .



# Little's law

- Related average number of some entity to
  - time spent in this entity
  - average number of entities in this entity
  
- No assumptions about
  
- Only based on expected values, only provides long-run expectations
  
- Independent of
  - number of servers (could have more servers per queue),
  - used scheduling discipline (could serve in reverse order),
  - queue length etc etc.

*Very very general, holds for (almost) all networks of queues etc etc.*



# An application of Little's law

□ Consider a window flow control system for packet transmission, with window size  $W$ .

□ There cannot be more than  $W$  packets in the session at any time, i.e.  $W \geq E[N]$ .

□ According to Little

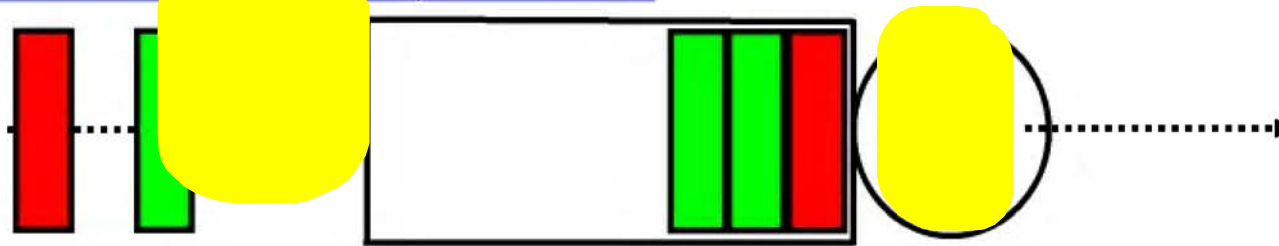
$$W \geq \lambda E[R]$$

where  $R$  is the response time of the system.

□ Now, if congestion builds up and  $R$  increases,  $\lambda$  must eventually decrease.

Q: How?

# The M/M/1 Queue



- So far we have been quite relaxed about the arrival time distributions and the service time distributions.
- Now let's study a specific case more closely.

the so-called **M/M/1** queue.

This type of queue has been a classical model for telecommunication network dimensioning, and continues to serve as a model for Internet traffic -- though being somewhat less appropriate. (More on this later).



# The M/M/1 Queue



- The M/M/1 queue is obtained by fixing
  - M: Markov (i.e. memoryless) arrival rate  $\lambda$
  - M: Markov (i.e. memoryless) service rate  $\mu$
  - 1: a single service station.

*both  $\lambda$  and  $\mu$   
are parameters of  
exponential distributions.*

- For given  $\lambda$  and  $\mu$  we know  $E[N]=\lambda E[R]$  etc.
- How do we determine  $E[N]$  or  $E[R]$ ?

# On the appropriateness of the exponential distribution

□ Statistical independent behaviours naturally lead to the exponential distribution.

□ Examples

○

○ cars passing bridges  
(except rush hour etc)

○ length of files on an arbitrary computers  
(prior to mp3)

...

usually referred to as 'Poisson' arrivals  
for reasons that you are invited to find out yourself

# On the appropriateness of the exponential distribution

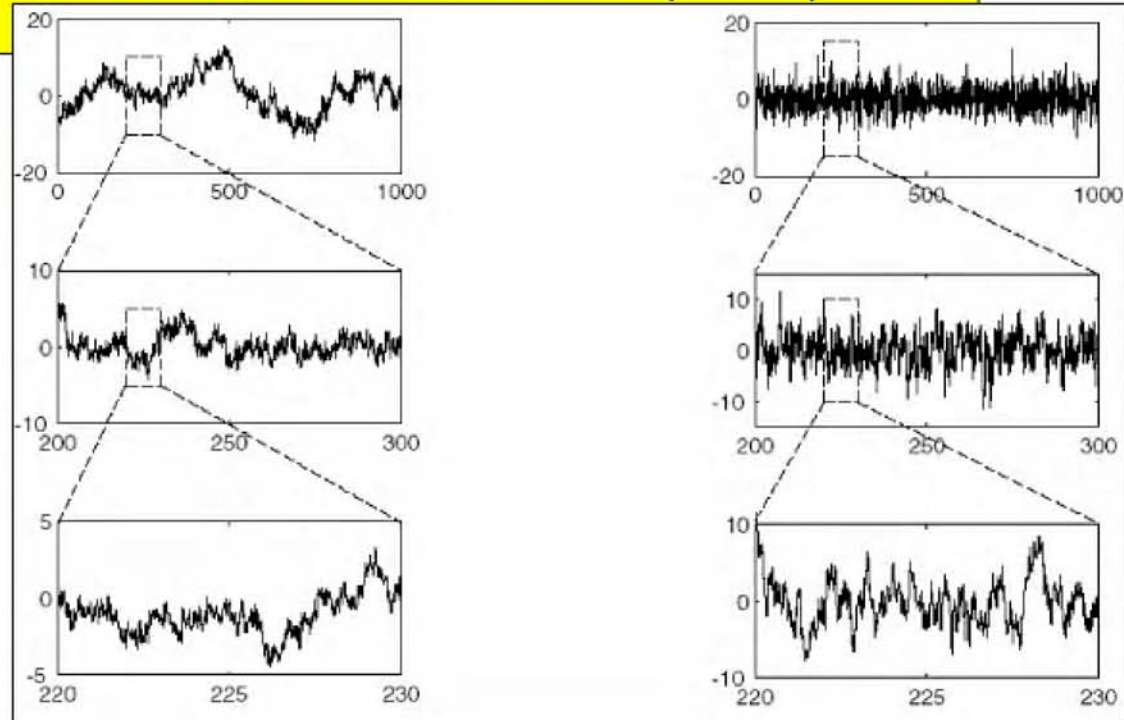
□ Lucent Press Release, June 6, 2001

<http://www.lucent.com/press/0601/010606.bla.html>

Cao, Cleveland, Lin, and Sun, "Internet traffic tends toward Poisson and Independent as the Load Increases" (2002)

'Packet interarrival times (of aggregate traffic) become exponentially distributed and independent as the link load increases'

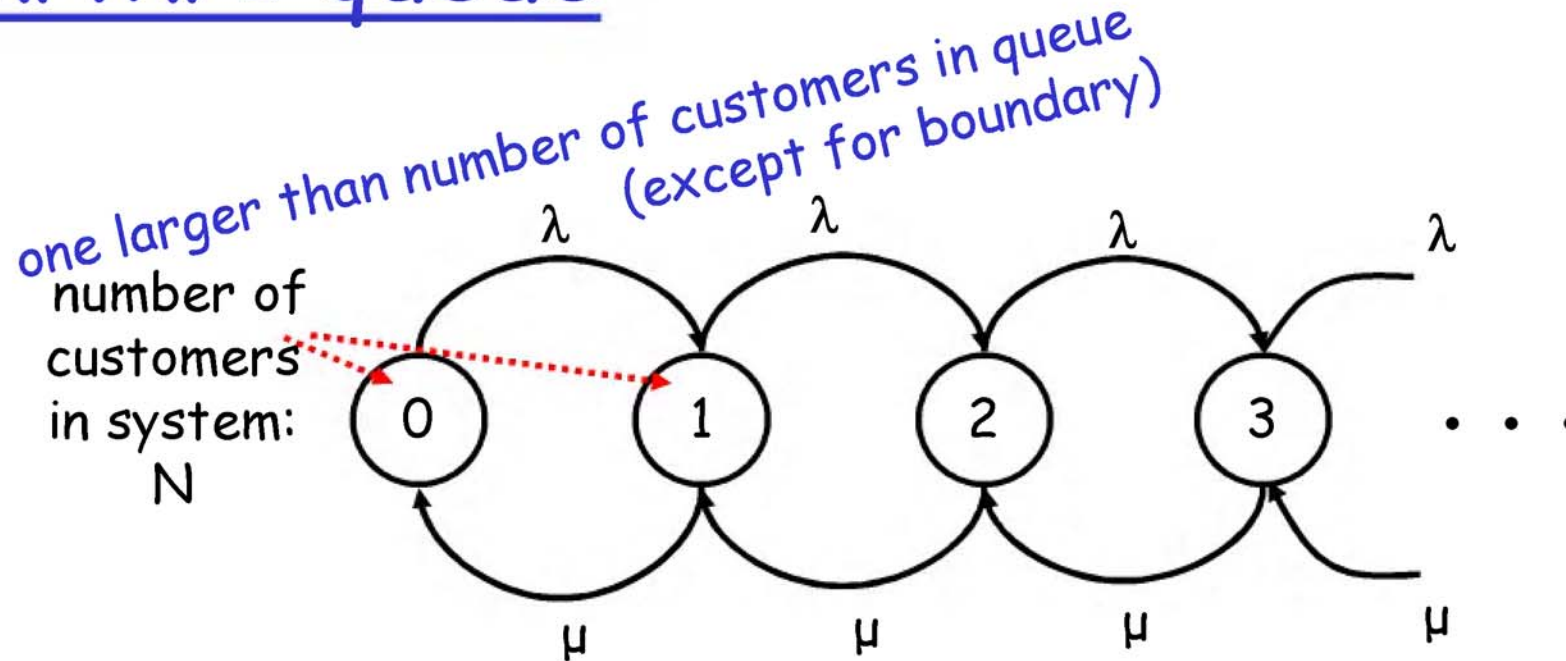
Not so sure about this!



heavy-tailed (self-similar)

exponential

# A state diagram for the M/M/1 queue



□ Is a 'Continuous Time Markov Chain'

□  $E[N] \equiv \sum_{n=0}^{\infty} n P(N=n)$

$$0 = -\lambda \pi_0 + m \pi_1 \Rightarrow \lambda \pi_0 = m \pi_1$$

$$(\lambda + m) \pi_1 = \lambda \pi_0 + m \pi_2$$

$$\pi_0 = \frac{m}{\lambda} \pi_1$$

$$(\lambda + m) \pi_{n+1} = \lambda \pi_n + m \pi_{n+2}$$

$$\pi_2 = \frac{\lambda}{m} \pi_0$$

$$(\lambda + m) \frac{\lambda}{m} \pi_0 = \lambda \pi_0 + m \pi_2$$

$$\left( \frac{(\lambda + m)\lambda}{m} - \frac{\lambda m}{m} \right) \pi_0 = m \pi_2$$

$$\frac{\lambda^2}{m} \pi_0 = m \pi_2$$

$$\pi_2 = \left( \frac{\lambda}{m} \right)^2 \pi_0$$

$$\pi_{n+2} = \left( \frac{\lambda}{m} \right) \pi_n =$$

$$\left( \frac{\lambda}{m} \right)^{n+1} \pi_0$$



# M/M/1 queue analysis

Let's assume an infinite buffer capacity  
(we are not the only one to do so).

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

$$\sum_{n=0}^{\infty} \pi_n = 1$$

We are dealing with the following matrix

$$Q = \begin{pmatrix} -\lambda & \lambda & 0 & 0 & \dots \\ \mu & -\lambda - \mu & \lambda & 0 & \dots \\ 0 & \mu & -\lambda - \mu & \lambda & \dots \\ \vdots & \vdots & \dots & \dots & \dots \end{pmatrix}$$

for which we need to solve

$$\vec{\pi} Q = 0$$

to obtain the steady-state state probabilities  $P(N=n) = \pi_n$ .

If  $\lambda < \mu$  we get:  $\pi_n = \rho^n (1 - \rho)$



# M/M/1 queue analysis

*instable if  $\lambda > \mu$*

As long as  $\lambda < \mu$ , queue has the following (long-run) limit properties

- $P(N=n) = \rho^n(1-\rho)$ 
  - (indicates fraction of time spent with  $n$  customers in queue)
  - Utilization =  $1 - P(N=0) = \rho$

## □ Notation

- $\rho = \lambda/\mu$ 
  - ratio of arriving/departing traffic ('traffic intensity')
- $N = \#$  customers in system
- $R =$  customer time in system

- $E[N] \equiv \sum_{n=0}^{\infty} n P(N=n) = \rho/(1-\rho)$

- $E[R] = E[N] / \lambda$  (Little)  
 $= \rho/(\lambda(1-\rho)) = 1/(\mu - \lambda)$

# M/M/1 queue analysis

20

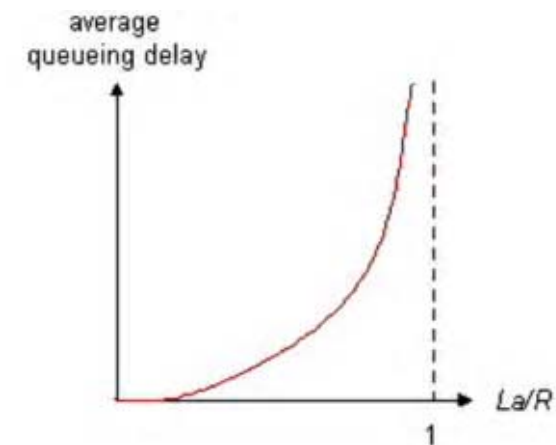
1: Introduction

## Queueing delay (revisited)



- $R$ =link bandwidth (bps)
- $L$ =packet length (bits)
- $a$ =average packet arrival rate

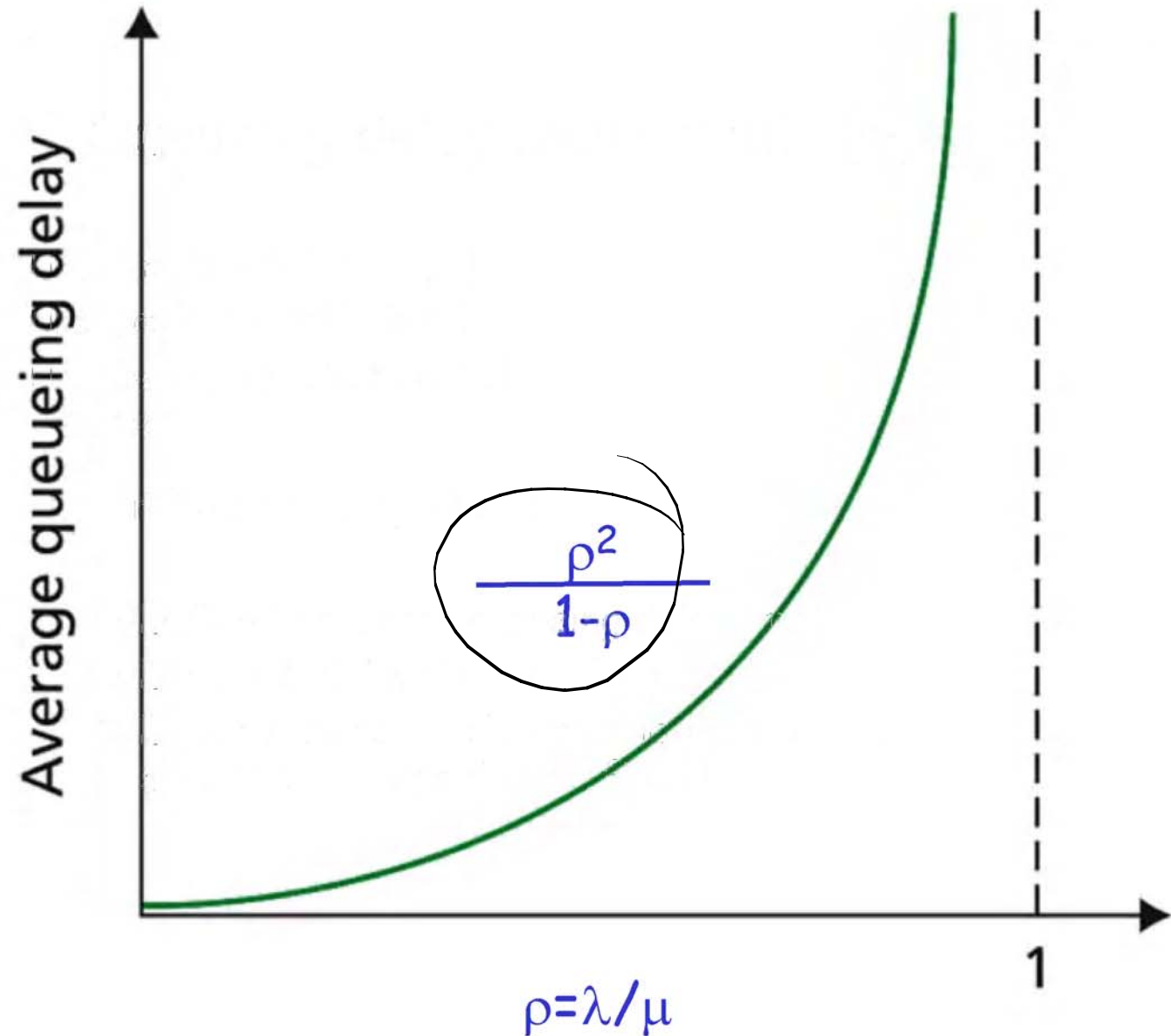
$$\text{traffic intensity} = La/R$$



- $La/R \sim 0$ : average queueing delay small
- $La/R \rightarrow 1$ : delays become large
- $La/R > 1$ : more "work" arriving than can be serviced, average delay infinite!

# M/M/1 queue analysis

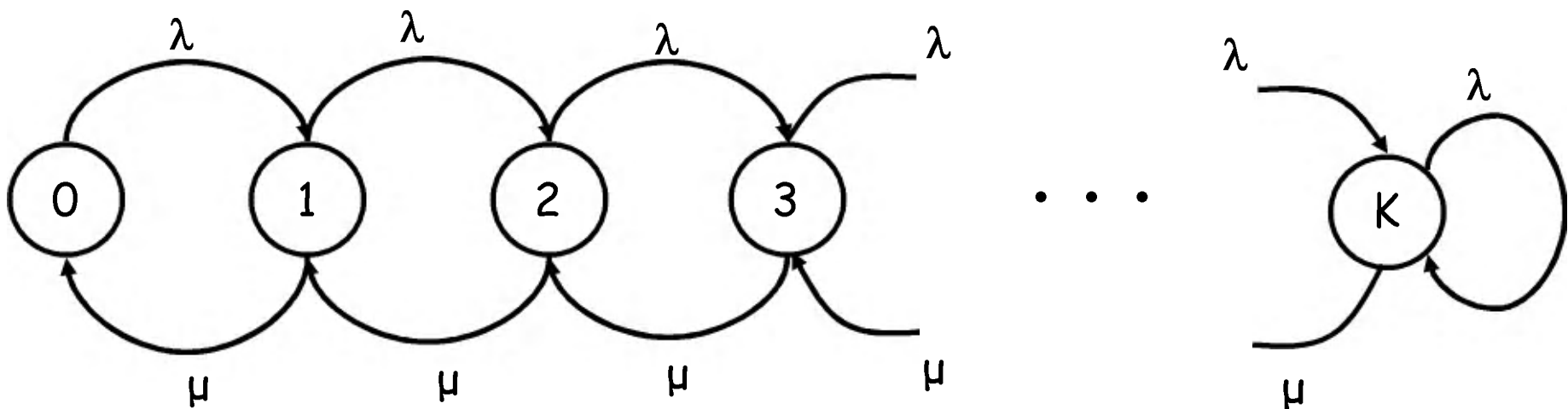
Q: Why this?



# Bounded buffer M/M/1 queue

Note:  $\rho \geq 1$  permitted

- Also can be modeled as a CTMC
  - requires  $K+1$  states for a model (queue + server) that holds  $K$  packets
  - stay in state  $K$  upon a customer arrival



# Bounded buffer queue properties

$$\square P(N=n) = \begin{cases} \rho^n(1-\rho) / (1 - \rho^{K+1}), & \rho \neq 1 \\ 1 / (K+1), & \rho = 1 \end{cases}$$

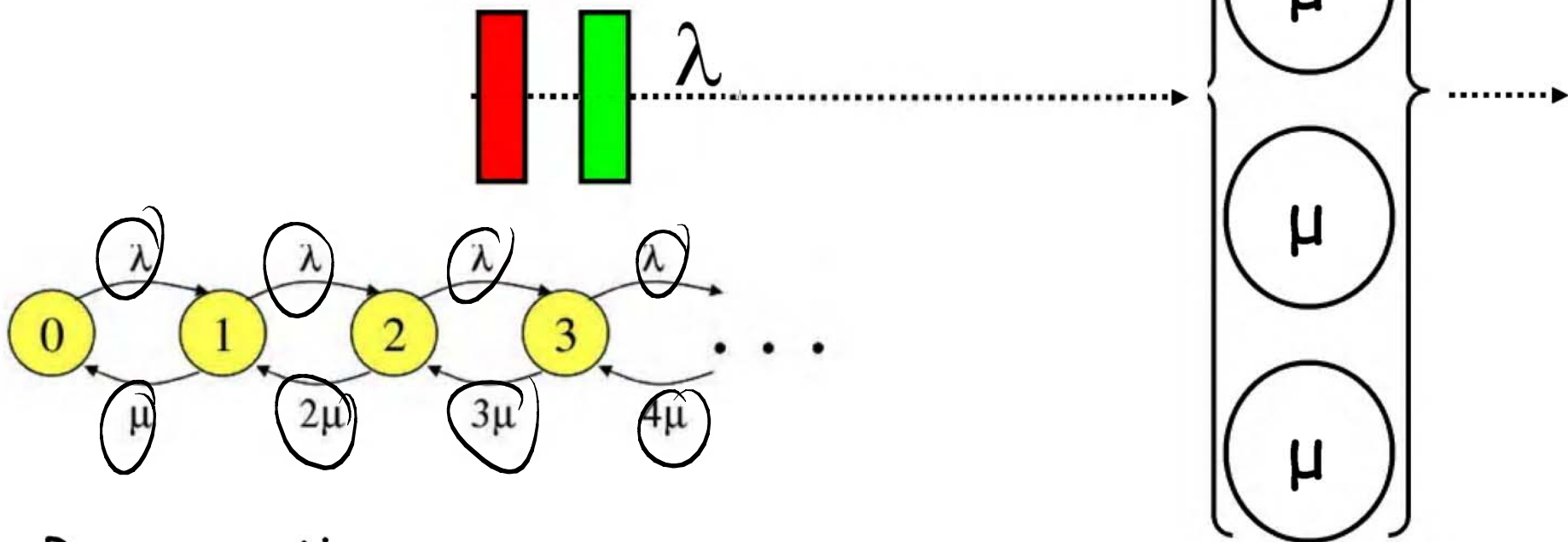
$$\square E[N] = \begin{cases} \rho / ((1-\rho)(1 - \rho^{K+1})), & \rho \neq 1 \\ 1 / (K+1), & \rho = 1 \end{cases}$$

divide unbounded buffer values by  $(1 - \rho^{K+1})$

$$\square \text{Utilisation} = 1 - P(N=0) = \rho(1 - \rho^K) / (1 - \rho^{K+1})$$

# What if many servers serve the same queue ?

- The extreme case:  $M/M/\infty$
- No waiting, always someone around to serve you!



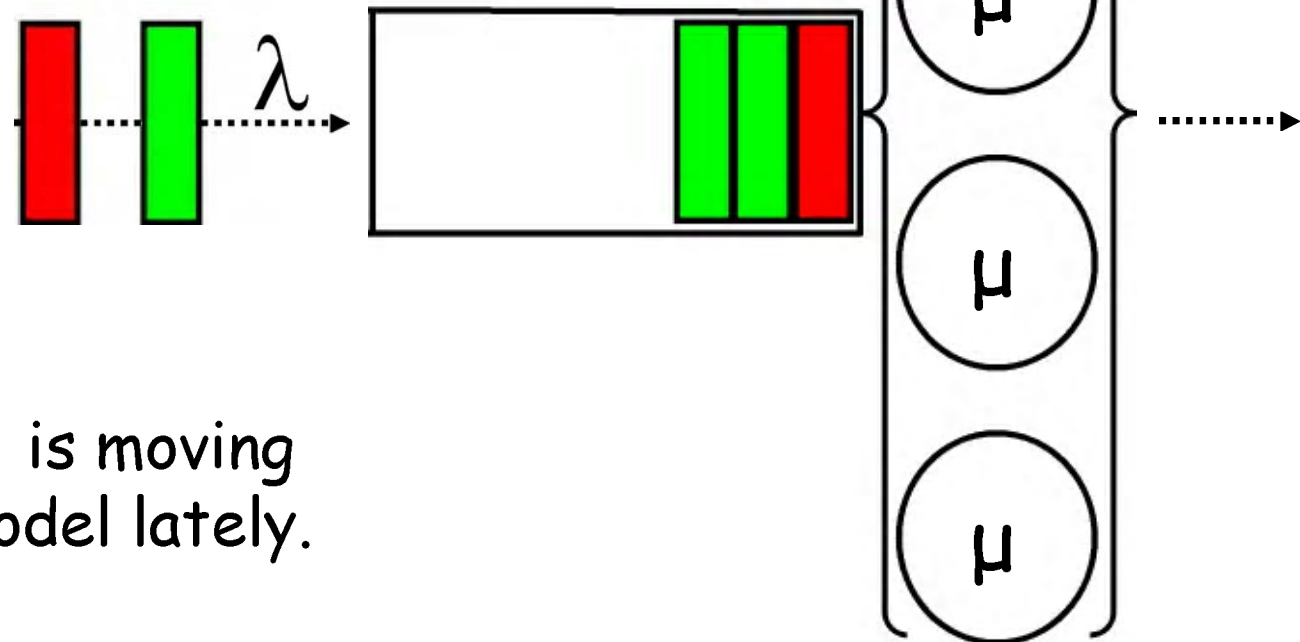
- Response time:  
 $R=1/\mu$



# What if many servers serve the same queue ?

- The not-so extreme case:

- $m$  servers at your service, serving a single queue

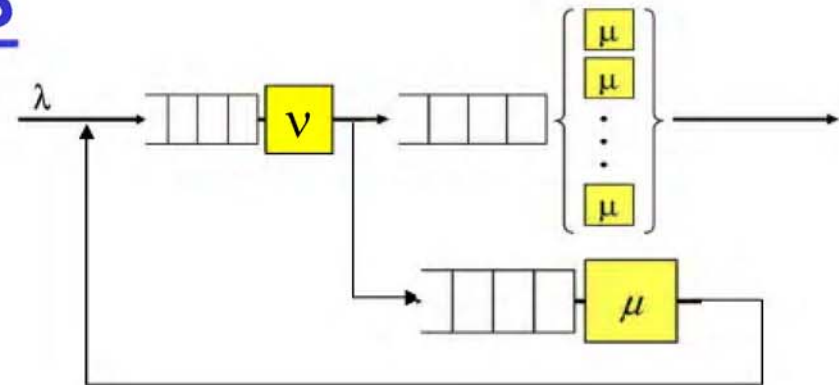


- *Deutsche Bahn* is moving towards this model lately.

Q: Why?

# Networks of queues

- Queueing networks consist of
  - queues
  - routes of customers
- If all component queues are  $M/M/...$ ,  
*you now know* how to analyse them!



# Networks of queues

- Queueing networks consist of
  - queues
  - routes of customers
- If all component queues are  $M/M/...$ , you now know how to analyse them!
- A simple example:
  - 'Tandem' queue:
    - 1st:  $M/M/1$  with queue size 3
    - 2nd:  $M/M/5$  with queue size 3

