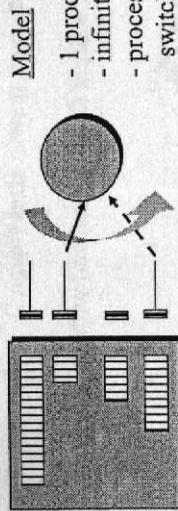


## Processor Sharing



Model

- 1 processor
- infinite buffer capacity
- processor can switch between requests, switchover times assumed negligible

153

*J. Koenigst., Operational Research/Queueing systems*

## Processor Sharing (2)

Switching between requests can take various forms:

- preemption with *abandonment* – destroys work!
- ...with *rollback* – creates work!
- ...with *resumption* – work-conserving

Problems:

- feasibility: service may not be "divisible"
- cost: switchover times may be significant, storage of requests' attained service

*J. Koenigst., Operational Research/Queueing systems*

154

**Processor Sharing (3)**



Why ask for trouble?

- stronger delay differentiation – favoring small-size requests
- ... possibly without information on request sizes (magic!)

**Basic evaluation criterion:** mean waiting delay normalized with respect to requested service time.

For a request with service time  $x$  it is  $\frac{w_m^{(x)}}{x}$ .

In the case of time sharing,  $w_m^{(x)} \geq \tau_m$  (waiting can't be shorter than residual service found in progress on arrival), so

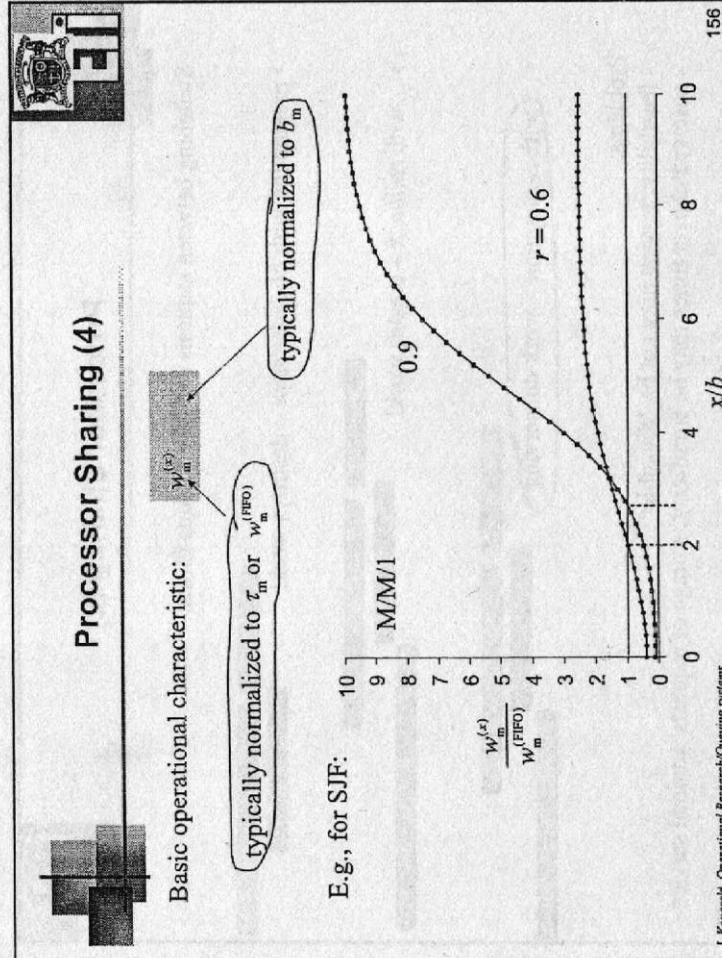
$$\lim_{x \rightarrow 0} \frac{w_m^{(x)}}{x} = \infty$$

and this can't be overcome by whatever sophisticated QD we devise.

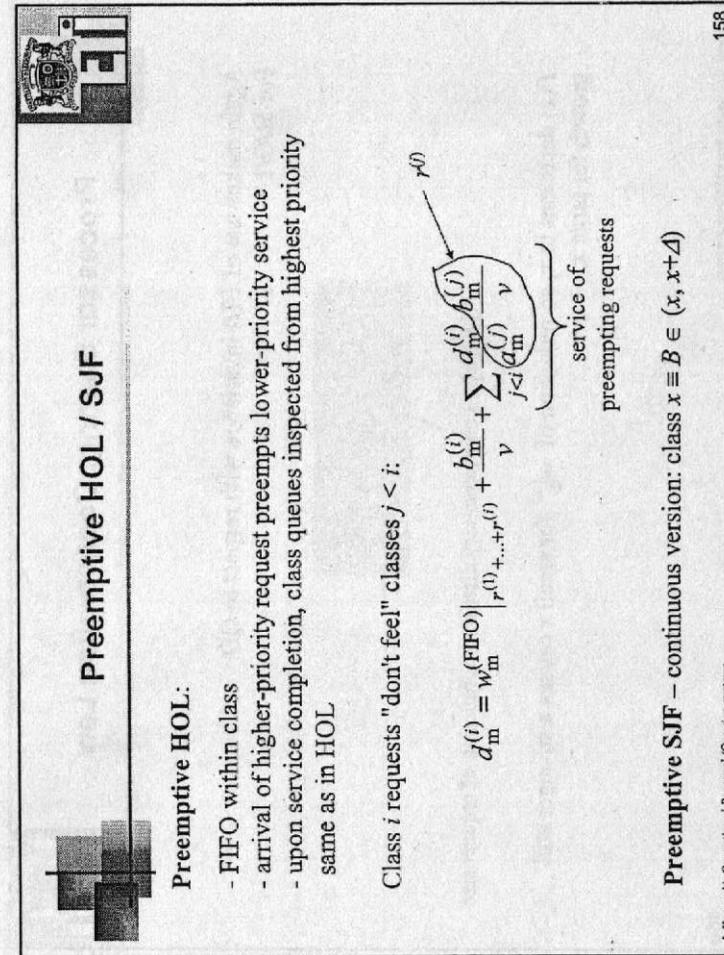
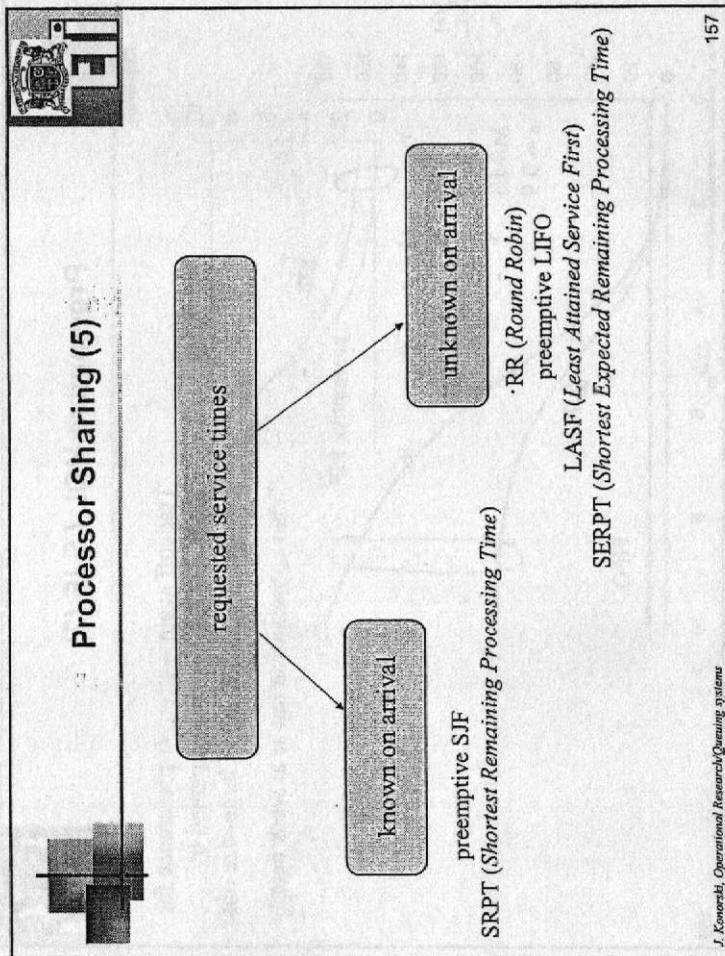
*J. Kornaraki, Operational Research/Queueing systems*

155

*J. Kornaraki, Operational Research/Queueing systems*



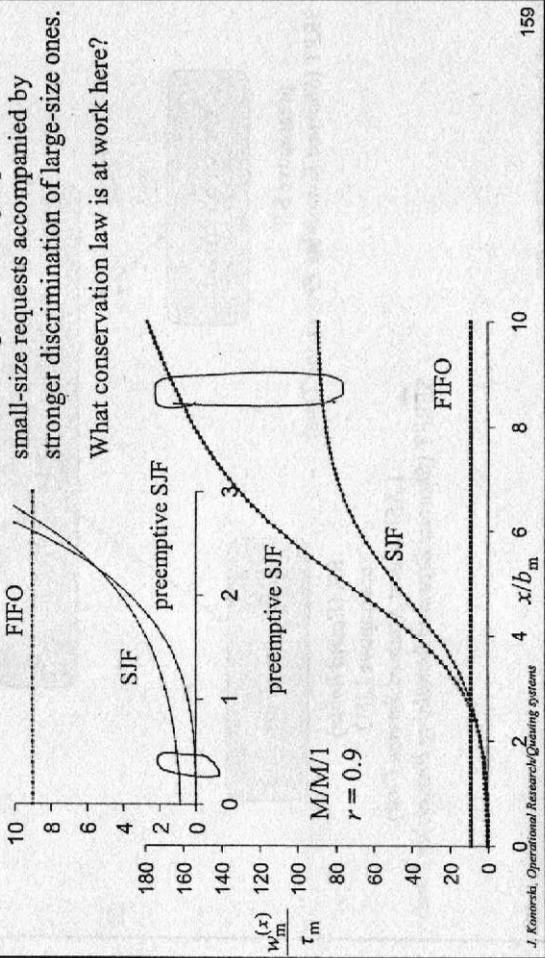
*J. Kornaraki, Operational Research/Queueing systems*



## Preemptive HOL / SJF (2)

Tradeoff again – stronger preference for small-size requests accompanied by stronger discrimination of large-size ones.

What conservation law is at work here?



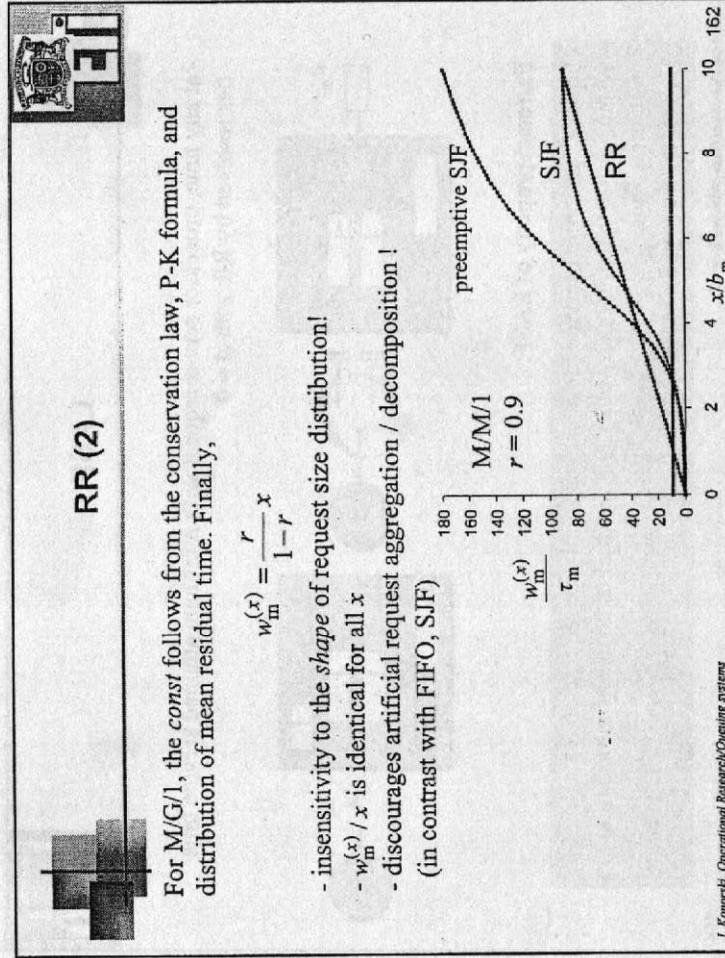
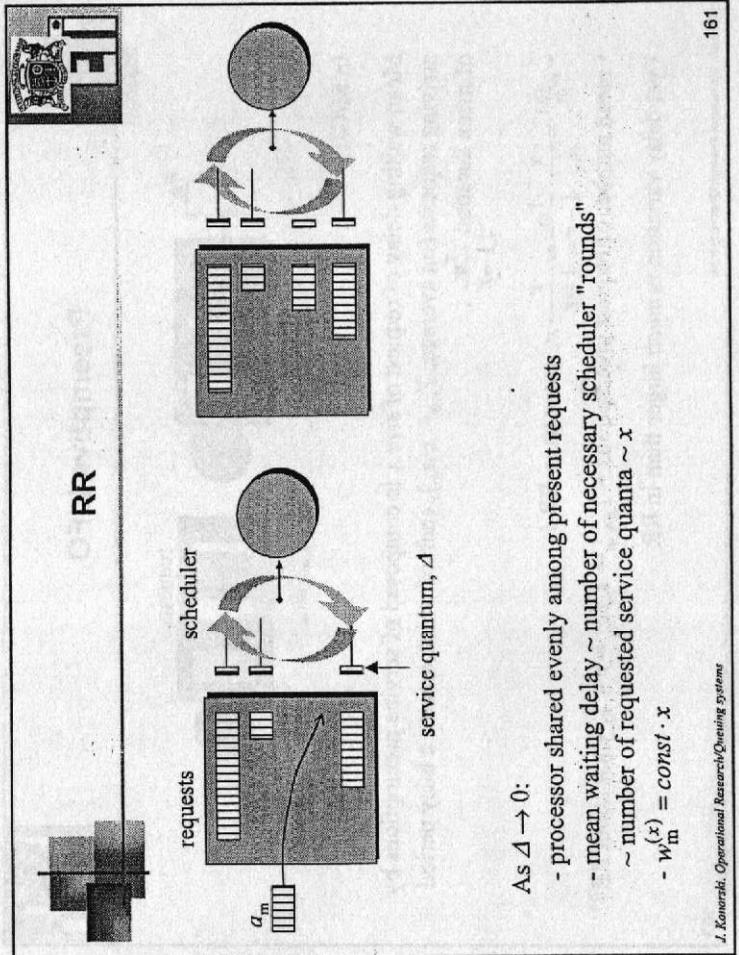
## Processor Sharing: Conservation Law

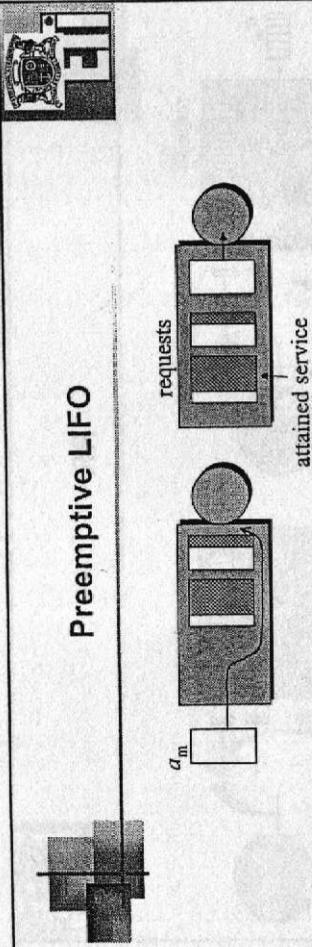
Again makes use of  $U(t)$  invariance with respect to QD.  
For M/G/1,

$$\int_0^{\infty} \frac{w_m^{(x)}}{\tau_m} P(x) dx = w_m^{(FIFO)}$$

$\Pr[B \geq x]$  – complementary distribution function of request size

$P(x)$  decreases in  $x$ , so reduction of  $w_m^{(x)}$  for small  $x$  causes a stronger still growth for large  $x$ .





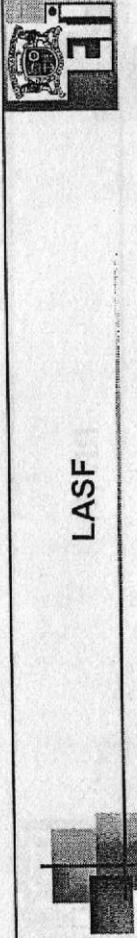
In M/G/1:

Mean waiting delay of request of size  $x$  is composed of service preemptions by arriving requests (on average,  $Y_{x,m} = x/a_m$ ), each of which creates a busy period of mean duration  $\frac{\tau_m}{1-r}$ .

- $w_m^{(x)} = \frac{x}{a_m} \cdot \frac{\tau_m}{1-r} = \frac{r}{1-r} \cdot x$ , same as in RR !
- mean number of request preemptions =  $Y_{B\&M} = \tau_m/a_m = r < 1$ , what about RR ?!
- yet delay variation is *much* larger than in RR

*J. Konarak, Operational Research/Queuing systems*

163



- at any time, processor serves request with minimum attained service time
- ties resolved by RR with  $\Delta = 0$

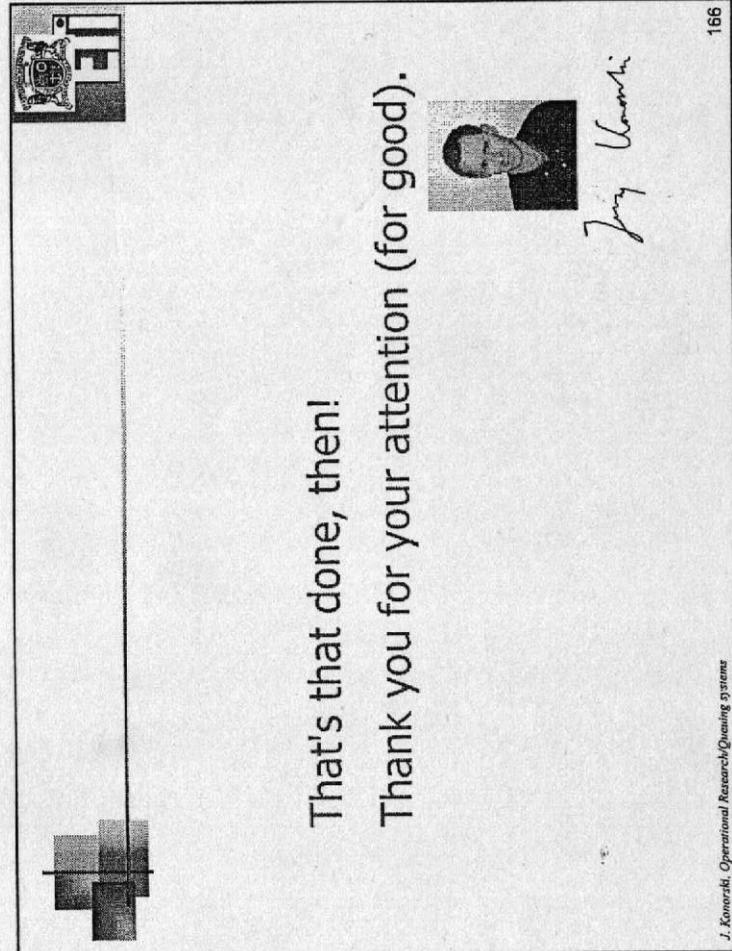
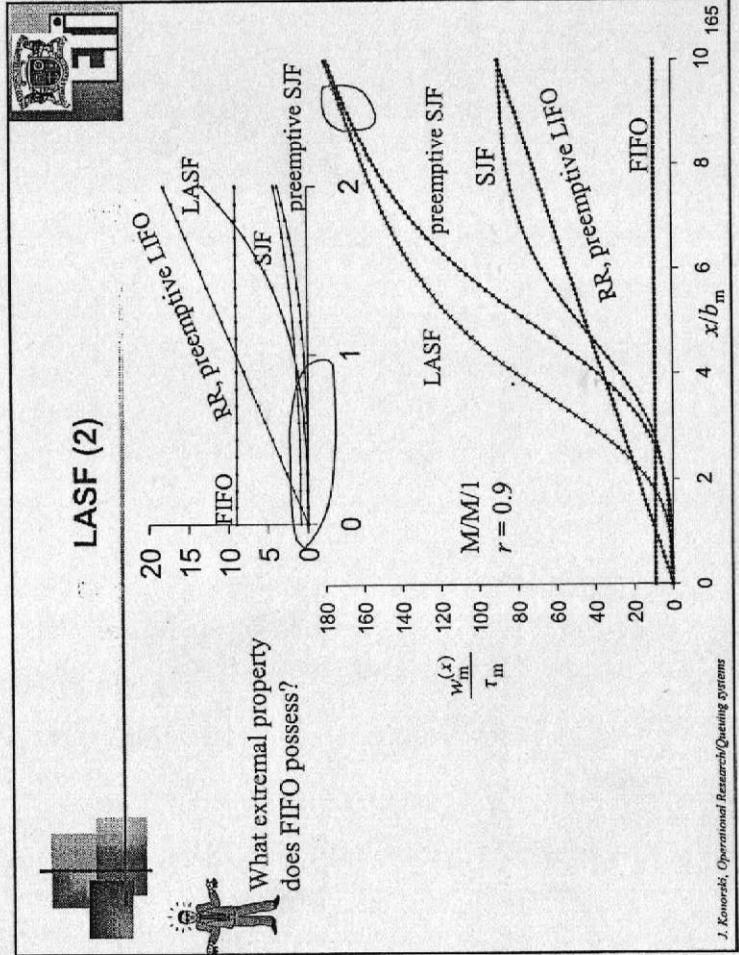


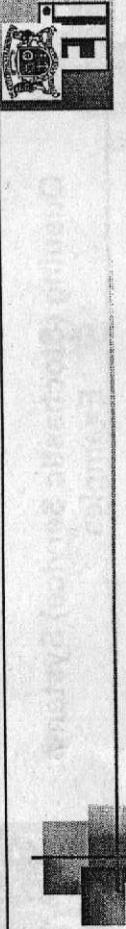
Extremal property of LASF:

Equalizing attained service times causes small-size requests to leave the system earlier than under any other QD that makes no use of information on request sizes.  
Hence, favors small-size requests in the strongest way possible.

*J. Konarak, Operational Research/Queuing systems*

164





# Operational Research

## Queuing Systems 1: Description and Operation

Jerzy Konorski

Room 139 (old bldg)  
office hours: tba  
[jekon@eti.pg.gda.pl](mailto:jekon@eti.pg.gda.pl)

*J. Konorski, Operational Research/Queuing Systems*

1



## Recommended Reading

- L. Kleinrock: *Queuing systems*, vol. I, II, Wiley 1975-1976
- D. Gross, C.M. Harris: *Fundamentals of Queuing Theory*, Wiley 1998
- Joti Lal Jain, W. Boehm, Sri Gopal Mohanty: *A Course on Queuing Models*, Chapman & Hall 2006
- G. Bolch, S. Greiner, H. de Meer, K. S. Trivedi: *Queueing Networks and Markov Chains. Modeling and Performance Evaluation with Computer Science Applications*, 2nd Ed., Wiley-Interscience 2006

2

*J. Konorski, Operational Research/Queuing Systems*

**Queuing (Stochastic Service) Systems:**

**Examples**

- computer system (mainframe / call center / database / Web server):  
interruptions / system tasks / queries / transactions wait to be processed when  
operators / processors / data storage released
- communication device (network card / telephone exchange / link multiplexer):  
data frames / subscriber calls wait for free capacity
- transport infrastructure (toll gate / gas station / harbor quay / runway):  
vehicles await a free "service slot"
- service access point (ATM / supermarket checkout / public office):  
customers / shoppers wait to be served / attended to by clerk / till lady / server

*J. Kamburki, Operational Research/Queuing systems*

3

**Queuing Systems(2)**

Queuing Systems(2)

link

traffic flows

aggregate traffic

router

E1

ATM

POST OFFICE

NO HELMETS

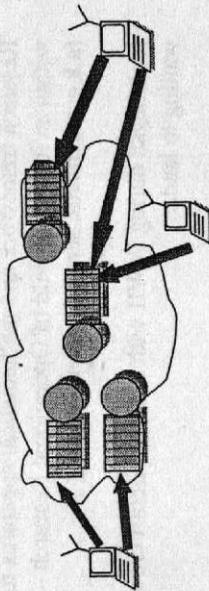
Queues everywhere!  
Must be accounted for in systems design.

*J. Kamburki, Operational Research/Queuing systems*

4

## Queuing Systems(3)

- system has resources – limited, reusable    *2009 by Konoracki*
- perceives events in the form of request arrivals = some entity demands access to the resources
- in response, system assigns resources to request enabling their consumption for a prespecified service time
- resources capable of serving requests = **processors**
- arriving request may find all processors busy serving other requests; then is stored in a **buffer** = waiting area for queuing requests until processor becomes free and service can commence



*J. Konoracki, Operational Research/Queuing systems*

5

## Queuing Theory: Mission

- Population of requests / request sources usually very large.
- Renders pointless optimization of specific request arrival scenarios e.g., scheduling for earliest termination or minimum processor usage.
- Only meaningful is analysis and design of service systems whose input is an **arrival stream** = unpredictable on-the-fly arrivals of successive requests.
- To this end we study trajectories of various queue characteristics over time = **queuing (service) processes**.

*J. Konoracki, Operational Research/Queuing systems*

6

## Queuing Theory: Mission (2)

### *WwII bombing*

With a large request population, instantaneous demand often exceeds instantaneous service supply – this is how queues form.

System designers are supposed to keep resulting damage under control e.g.,

- customer dissatisfaction due to delays / rejections, balking (refusing to join a long queue),  
*flight*
- buffer and queue management burden,  
*flight*

with a view of the economics of processor usage.

Research framework and mathematical apparatus for that were developed within an important field of Operational Research – **queuing theory** (a.k.a. **stochastic service systems theory**).

It all began during WWII with bomber aircraft crowding over the airfield, waiting to land...

J. Konorat, Operational Research/Queuing systems

7

## Simplest Model



The simplest model of a queuing system consists of:

- a processor,
  - a buffer, and
  - an arrival stream. *Stream* *balking*
- 

Characteristics of service process depend  
on those of arrival stream and the way  
buffer & processor system operates.  
How?  
This is what queuing theory is about

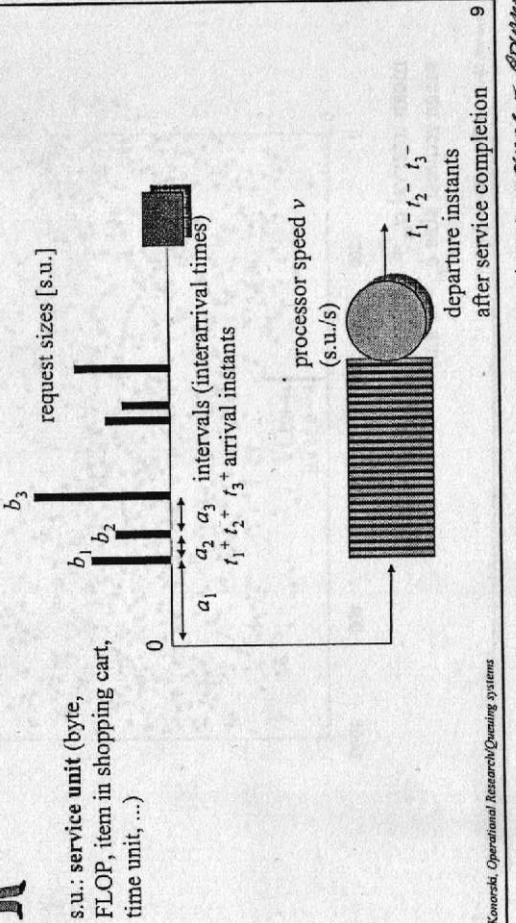
J. Konorat, Operational Research/Queuing systems

8

Simplest Model (2)

Which characteristics of the air system operation are relevant?

Which characteristics of the arrival stream and which rules of queuing system operation are relevant?



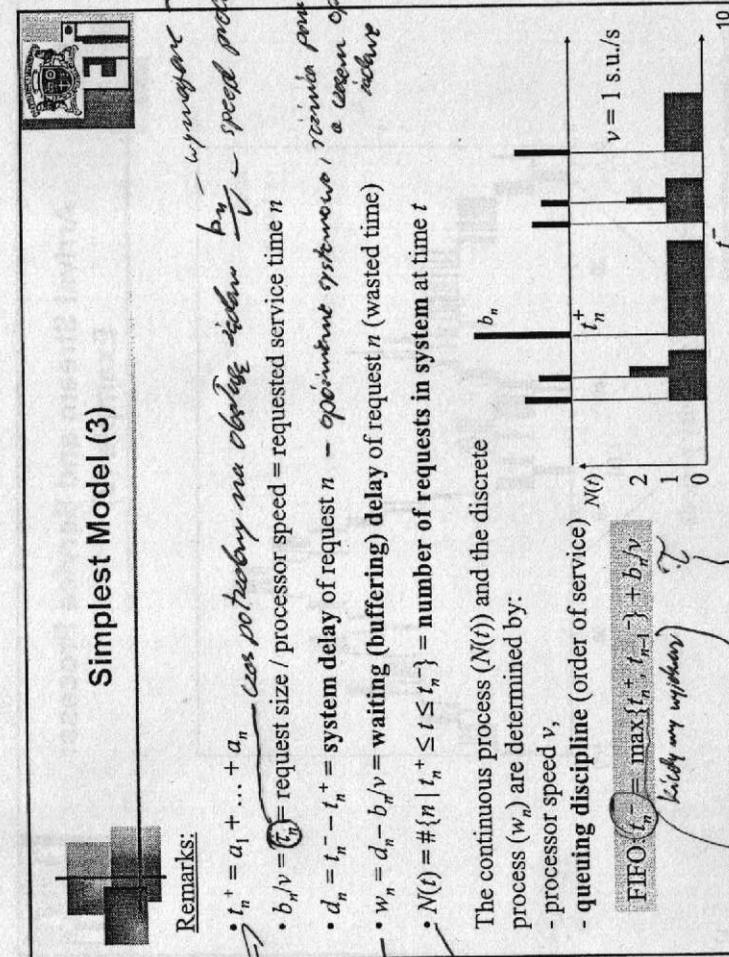
*J. Konarski, Operational Research/Queueing systems*

### Simplest Model (3)

**Remarks:**

- $t_n^+ = a_1 + \dots + a_n$   $\rightarrow$  total size of requests
  - $b_n/v = t_n^+$   $\rightarrow$  processor speed = requested service time  $n$
  - $d_n = t_n^- - t_n^+$   $\rightarrow$  system delay of request  $n$  = request time + a certain overhead  $\rightarrow$  total system time
  - $w_n = d_n - b_n/v$   $\rightarrow$  waiting (buffering) delay of request  $n$  (wasted time)
  - $N(t) = \#\{n \mid t^- \leq t \leq t^+\}$   $\rightarrow$  number of requests in system at time  $t$

Lieska et al.  
w systems  
fe w ~~but~~  
+ he co se too  
l' <sup>l' / l' / l' / l'</sup>

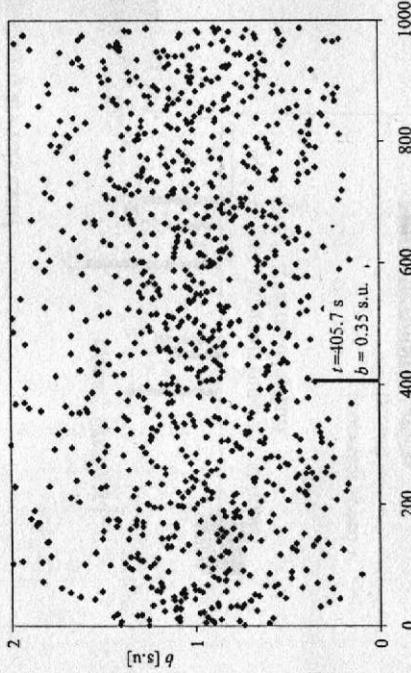


—

Queuing systems

los w/ sick  
populations  
but more  
likely to be sick

## Arrival Stream and Service Process: Example

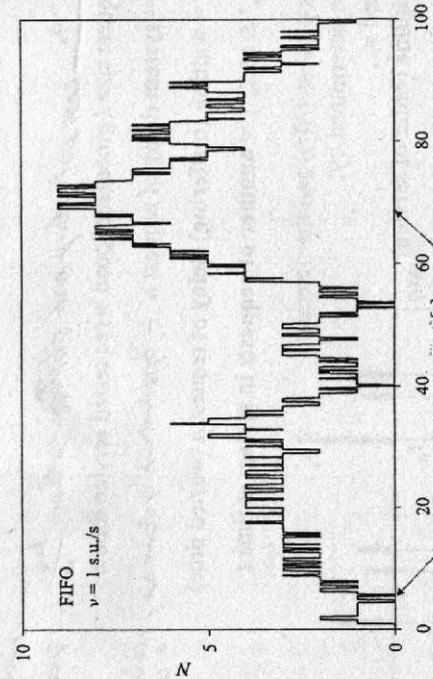


mean interval  $a_m = 1$  s  
mean request size  $b_m = 1$  s.u.

J. Konszki, Operational Research/Queueing systems

11

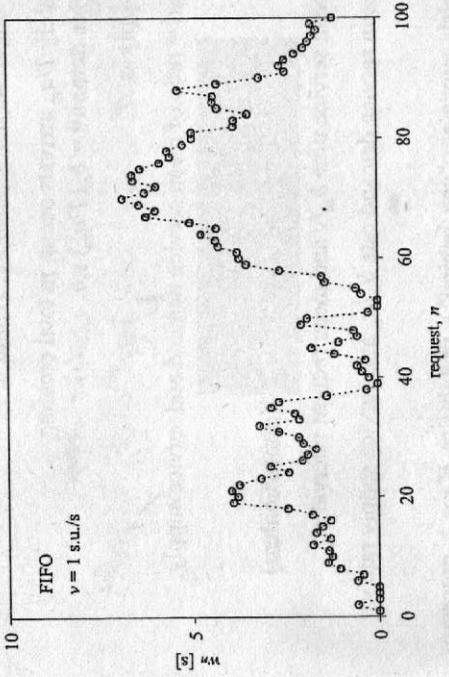
## Arrival Stream and Service Process: Example (2)



J. Konszki, Operational Research/Queueing systems

12

### Arrival Stream and Service Process: Example (3)



J. Kurose, Operational Research/Queueing systems

13

### Properties of "Good" Service Process

*what was absent  
(zero) = process  
no idle*

- from requests' viewpoint:  
small waiting delays, rare buffer overflows *rare as system - processor performance better*
- from system operator's viewpoint:  
high processor utilization (rare idle periods) *do not worry about - manage more tasks*

These are contradictory! Rare idle periods imply:

- occurrences of queuing
- long queues becoming prevalent *least off the load*  
*(rare) / messages*
- systematic queue growth (instability) / avalanche of buffer overflows *processor unable to work in real time!*

Relationship between arrival stream characteristics and processor speed  
determines an important parameter, offered load

*processor  
why task  
systems?*

14

J. Kurose, Operational Research/Queueing systems

## Offered Load

Consider a long observation period  $T$ . Within it,

- approximately  $T/a_m$  arrivals occur, in total creating mean service demand =  $b_m(T/a_m)$  s.u.
- service supply is  $vT$

Offered load = ratio of mean service demand and service supply:

$$r = \frac{b_m(T/a_m)}{vT} = \frac{b_m/v}{a_m} = \frac{b_m/a_m}{v} \quad (\text{dimensionless})$$

= ratio of mean service time  $b_m/v$  and mean request interval  $a_m$ ,

= ratio of mean service demand rate  $b_m/a_m$  and service supply rate  $v$

If  $r > 1$  persists, processor "gets behind" – instability. If  $r < 1$ , processor "keeps pace" – system stable. What happens under  $r = 1$ ?

J. Kouvatsi, Operational Research/Queueing systems

$\frac{b_m/v}{a_m} = \frac{\text{load}}{\text{pace}}$  breaks or  
none!

$r < 1 - \text{spikes left}$

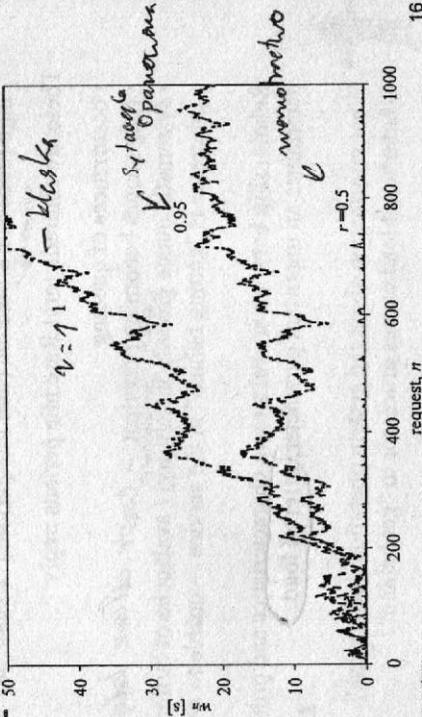
wandering site

## Offered Load (2)

same arrival stream ( $a_m = 1$  s,  $b_m = 1$  s.u.)  
decreasing  $v$



instability at  $r = 1$ !



J. Kouvatsi, Operational Research/Queueing systems

## Impact of Input Speed

So far immediate input assumed of requests from source to system (arbitrary  $a_n$ ).  
In reality, request transfer from source occurs at finite speed.



Can be modeled as a "virtual" input queuing system with processor speed  $v' < \infty$ , and arrival stream with  $b_n$  and arbitrarily small  $a'_n$ ; offered load  $r' = (b_m/a'_m)/v'$ .

Arrival stream at the real system has  $a_n \geq b_n/v'$ ; offered load  $r = (b_m/a_m)/v$ .

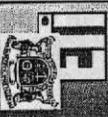
Clearly,  $a_m = a'_m$ , so  $r' = (v/v)r$ .

$$\left. \begin{array}{l} r' = 0 \\ r' = 0.5r \\ r' = r \end{array} \right\} \text{corresponds to } \left\{ \begin{array}{l} v' = \infty \text{ (infinite input speed)} \\ v' = 2v \\ v' = v \text{ (no queues)} \end{array} \right.$$

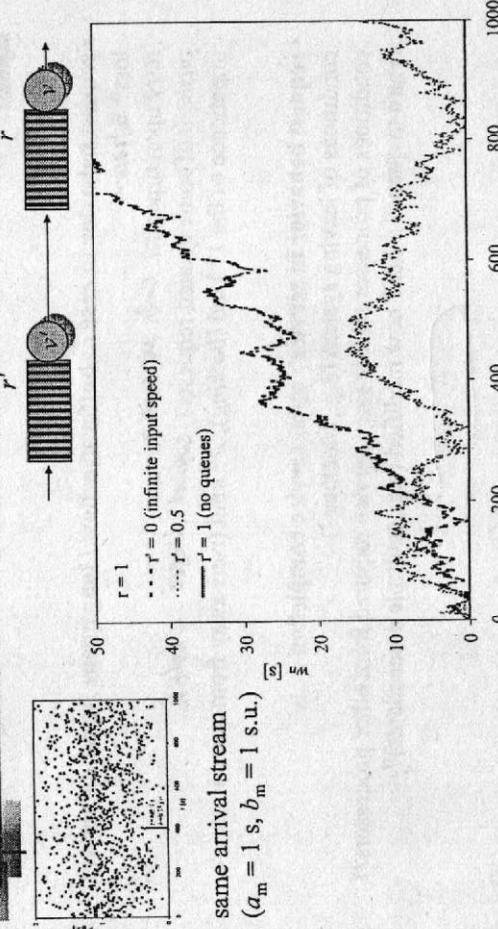
J. Konarski, Operational Research/Queueing systems

17

$T_m =$



## Impact of Input Speed (2)



J. Konarski, Operational Research/Queueing Systems

At  $r = 1$  with finite input speed, system can still remain stable.

At  $r = 1$   
; slowjana maksi' pnyahale yet  
to next tanda, ie positive stable

## Towards Richer Models

How is queuing process affected by other characteristics of the arrival stream, request behavior within system, queuing discipline, service rules?

- **arrival stream**
  - how exactly are  $(a_n)$  and  $(b_n)$  generated?
  - time variability? dependence on queuing process? bulk arrivals?
  - request sizes  $b_n$  – known/unknown on arrival?
- **buffer**
  - finite capacity 
  - limited accessibility? when is a request rejected – *drop-tail*, ...?

J. Konarski, Operational Research & Queuing systems

19

coincidence  
propagation  
behavior?

## Towards Richer Models (2)

- **request behavior in case of buffer overflow / long queue found on arrival**
  - loss? 
  - retry upon timeout? 
  - pushout of some queued requests? 
  - impatience of the 1<sup>st</sup> kind (*balking*), 2<sup>nd</sup> kind (runs away from queue)

coincidence  
propagation  
behavior?  
by blocking process  
by parallel wa  
indirecte process  
process w/ blanking?

- **request behavior in service / upon service completion**
  - conditions of leaving system (e.g., blocking?)
  - conditions of processor release (e.g., service required from other processors?)
  - return to queue? when? how modified (e.g., multiple descendants)?

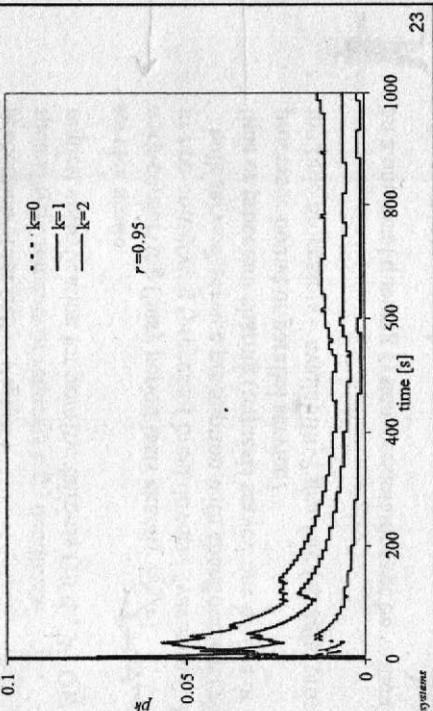
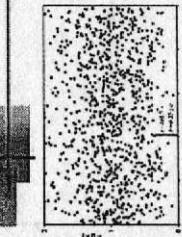
J. Konarski, Operational Research & Queuing systems

20



## Steady State (2)

same arrival stream ( $a_m = 1 \text{ s}$ ,  $b_m = 1 \text{ s.u.}$ )

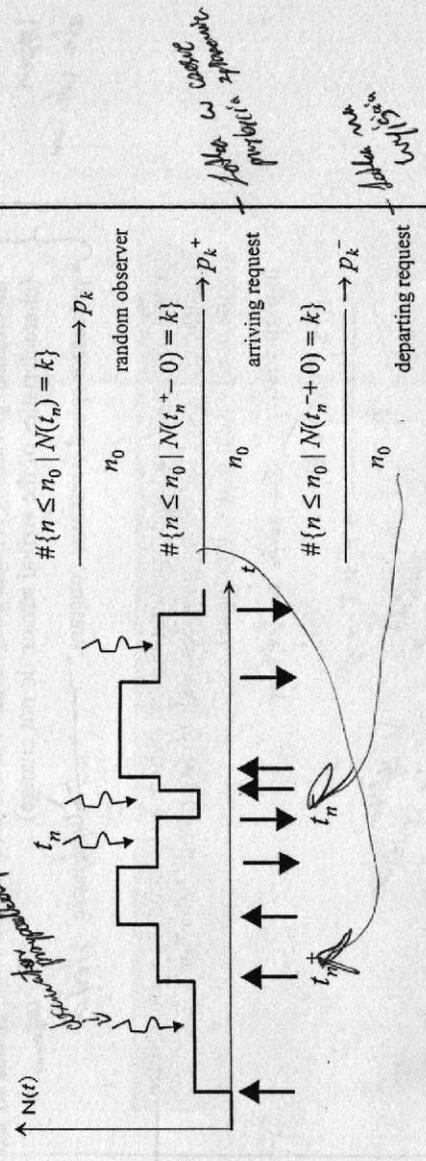


J. Komornik, Operational Research/Queueing systems

## Steady State (3)

Observation of  $N(t)$  – viewpoint matters!

As  $n_0 \rightarrow \infty$ :

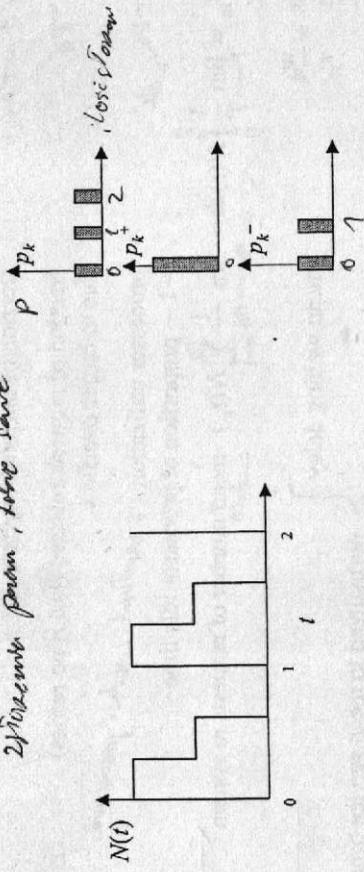


J. Komornik, Operational Research/Queueing systems

24

## Steady State (4)

"Practitioners", beware!



J. Konarak, Operational Research/Queueing systems

25

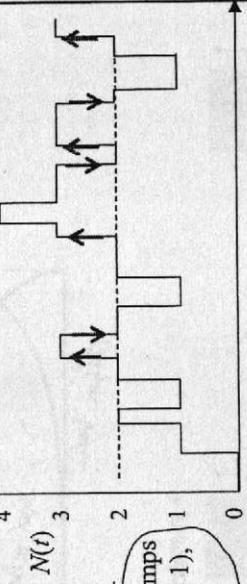
against steady state



## above proportionality

**Unit jump argument:** if  $N(t)$  only changes by  $\pm 1$  then  $p_k^+ \equiv p_k^-$

(% arrivals finding  $k$  in system = % departures leaving  $k$ )



Since  $q_1 < 1$  result  $1$

w downward process

Remark: Only accepted requests count as arrivals. If  $Q < \infty$  (some requests can be rejected due to buffer overflow) then

$$\frac{p_k}{p_{k-1}} = p_k^-, k = 0, \dots, Q-1$$

$1 - p_Q^+$  - probability to be busy path

With more by above property to line sit into to project  
 $\lambda \rho < \infty$   $Q < \infty$  (otherwise loss buffer entry) to or slow down

$$\frac{\alpha + 1}{\beta + T} \approx \frac{\alpha}{\beta}, \alpha \rightarrow \infty$$

26

13

## Steady State (6)

Relevant evaluation criteria:

$\rightarrow$   $p_0$  fraction of requests lost due to lack of service

$p_1^+ + \dots + p_{Q-1}^+$  fraction of buffered requests

$L = p_Q^+$  fraction of requests rejected (lost if no retries)

due to buffer overflow

processor utilization  $\rightarrow$   $\zeta_{idle}$  time spent in processor

$= 1 - \text{proportion of processor idle time}$

$= 1 - \lim_{n_0 \rightarrow \infty} \frac{1}{n_0} \sum_{n=1}^{n_0} N(t_n)$  mean number of requests in system

$$N_m = \lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T N(t) dt = \lim_{n_0 \rightarrow \infty} \frac{1}{n_0} \sum_{n=1}^{n_0} N(t_n)$$

mean waiting delay } normalized to mean service time  
mean system delay }

J. Kumerits, Operational Research/Queueing systems

$\rightarrow$   $\frac{20 \text{ min}}{10 \text{ min}} = 2$   $\rightarrow$  evaluate the performance monitor

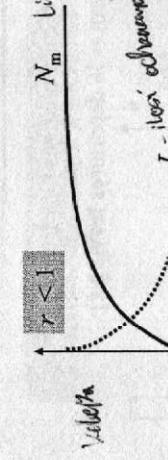
10 wait longer

10 work longer

27

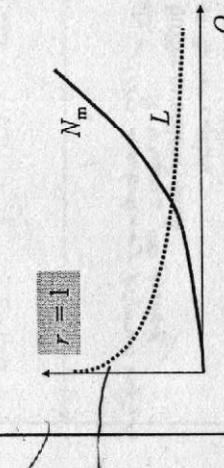


## Steady State (7)

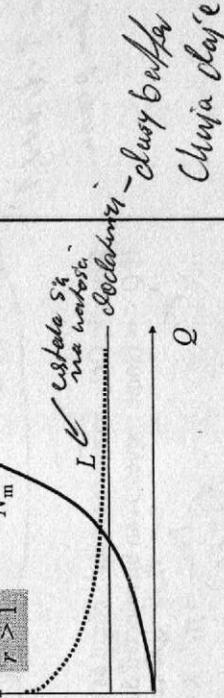


Überlagernd  
System stabil

+ systematic upward drift in time



Nicht überlagernd  
System instabil



Nicht überlagernd  
System instabil

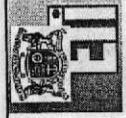
J. Kumerits, Operational Research/Queueing systems

28

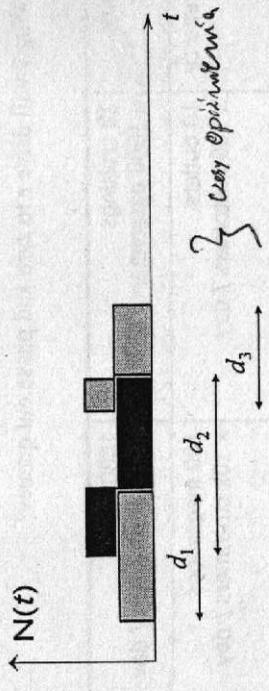
Churia daye

Churia daye

## Oliver's Theorem e optimização



### Little's Theorem



$$\text{as } n, T \rightarrow \infty, \int_0^T N(t) dt = d_1 + \dots + d_n$$

3 way optimisation

lost requests  
produced older  
systems

$(T/a_m)(1-L)d_m = \text{util. prod. system}$

J. Keonan, Operational Research/Optimizing Systems

J. Keonan, Operational Research/Optimizing Systems

Sistemas  
operacionais

29

### Little's Theorem (2)

$$N_m = \frac{1-L}{a_m} d_m$$

process

mean - avg program in  
system

mean system throughput = (mean system delay)  $\times$  (mean system delay)

mean population = (mean circulation)  $\times$  (mean lifetime)

Valid for any part of system:

- processor:  $1-p_0 = \frac{1-L}{a_m} r_m = (1-L)r$  flow conservation equation
  - buffer:  $N_m - (1-p_0) = \frac{1-L}{a_m} w_m$  system instant
- i - steady state:  
low priority o priority

J. Keonan, Operational Research/Optimizing Systems

$$OPV = 1-p_0 = \frac{1-L}{a_m} r_m = (1-L)r$$

but :  $\frac{1-L}{a_m} W_m$

 **Is Queuing Theory Losing Momentum?**

Growing  $v$  will drive  $r$  to zero and phase out queues?



transatlantic transport	10 cruisings x 1000 passengers /20 days	1000 flights x 250 passengers / day	x500
retail commerce	10 outlets x 100 customers / day	100 megashops x 10k customers / day	x1000
Internet access link	10 kb/s	100 Mb/s	x10k
processing power cost	\$15m/GFLOPS (1984)	\$0.5/GFLOPS (2007)	x30m
	<a href="http://en.wikipedia.org/wiki/GFLOPS">http://en.wikipedia.org/wiki/GFLOPS</a>		

Yet the answer is No.

*J. Konarski, Operational Research/Queuing Systems*

31

 **Is Queuing Theory Losing Momentum? (2)**



First, airliners, supermarkets, Internet links, and mainframe computers seem more crowded than ever. Same for online banking, toll-free numbers, hub airports etc.

Service demand rate  $b_m/a_m$  grows in step with service supply rate  $v$ ,  
and so  $r$  isn't dropping any!

*J. Konarski, Operational Research/Queuing Systems*

32

## Is Queuing Theory Losing Momentum? (3)

Second, how do queues form anyway?

Not only because of  $\nu < \infty$ , but above all because of variability of  $a_n$  (arythmic arrivals) and  $b_n$  (capricious demand) exhibited by request source!

Arythmic arrivals cause *instantaneous* offered load to vary between 0 and  $\infty$ . To get rid of queues, even occasional, one needs  $\nu = \infty$ .

Under  $b_m/a_m < \infty$  this gives  $r = 0$  i.e., zero processor utilization !! Highly uneconomical, no matter what progress technology and management make.

What is economical? Keep  $r < 1$  i.e.,  $\nu > b_m/a_m$ , but *not much*. Meaning, allow queues at times.

J. Kumoriski, Operational Research/Queuing Systems

33

## Comments on Queuing Theory

Queuing theory is a mathematical analysis tool. When designing a queuing system, perhaps one could do better with a prototype / simulator?

- credible estimates of troublesome characteristics  
rare events – queue length crosses threshold,  
long busy period – do we have instability here?
- qualitative (rather than scenario-specific) influence of parameter settings  
upon relevant characteristics of queuing process  
– saves a lot of unnecessary experimenting
- ! - universal (qualitatively, often also quantitatively) impact of results for simple models – carry over to much more realistic ones  
*mainly*

Contrary to what might seem, mathematical analysis is very costly.  
Only pays off if provides answers that would be hard to get otherwise.

J. Kumoriski, Operational Research/Queuing Systems

34

## Comments on Queuing Theory (2)



Agner K. Erlang (1878-1929)

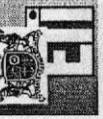
Danish mathematician and engineer,  
was the first to appreciate that modern telephony cannot do without probability

today's teletraffic unit = 1 *erlang*

*J. Konorski, Operational Research/Queuing systems*

35

## Comments on Queuing Theory (3)



How is queuing theory related to the theories of:

### job scheduling

finding an optimum schedule for a fixed job set  
vs. unpredictable on-the-fly arrivals of requests

### concurrent processes

deterministic analysis of specific event scenarios  
vs. massive population of random events, where only statistical  
characteristics are worth studying

### stochastic processes

similar calculus  
queuing process = nonlinear, infinite-memory transformation of arrival stream

*J. Konorski, Operational Research/Queuing systems*

36



# Operational Research

## Queuing Systems 2: Stochastic Models and Characteristics

Jerzy Konorski

Room 139 (old bldg)

jerkon@eti.pg.gda.pl

*J. Konorski, Operational Research/Queuing systems*

37



## Random Variables and Stochastic Arrival Streams

For an arrival stream observed over a finite time,  $a_n$ ,  $b_n$  and any other useful characteristics can be calculated.

Yet for prediction of characteristics over infinite observation periods, or computer imitation of the arrival stream, one needs a model of generation of  $(a_n)$  and  $(b_n)$ .

With rather impractical exceptions, deterministic models are of no interest:

- impractical – arrival instant and size of the next request rarely known in advance,
- carry no information (what is known in advance doesn't ever come as a surprise),
- pose no design challenge.

From now on we focus on stochastic (models of) arrival streams i.e., consider relevant quantities to be **random variables**:

- described by a **probability distribution** over a set of values (**realizations**),
  - this probability distribution exists, is either known or can be derived somehow (**the Bayesian approach**).

*J. Konorski, Operational Research/Queuing systems*

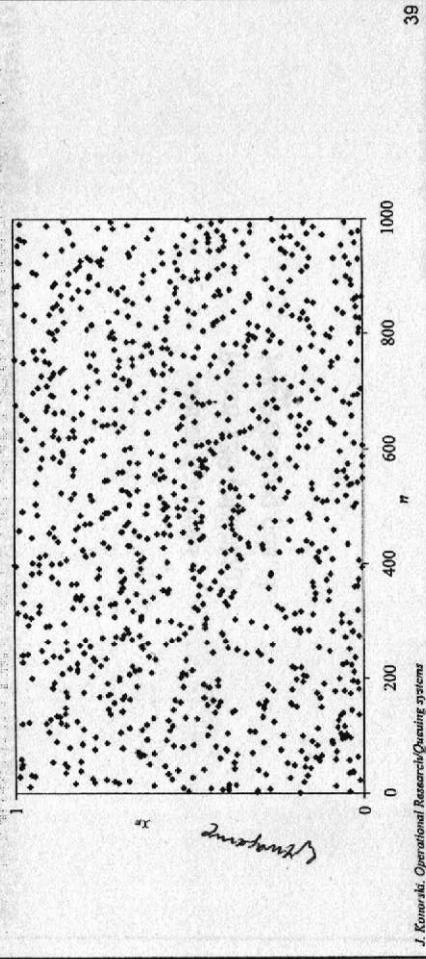
38

## Random Variables and Stochastic Arrival Streams (2)



An example of a random variable is value returned by the function `rrandom` (if we are deliberately oblivious to the algorithm of pseudorandom number generation). Its probability distribution is uniform on  $[0,1]$ .

The first 1000 observed realizations:



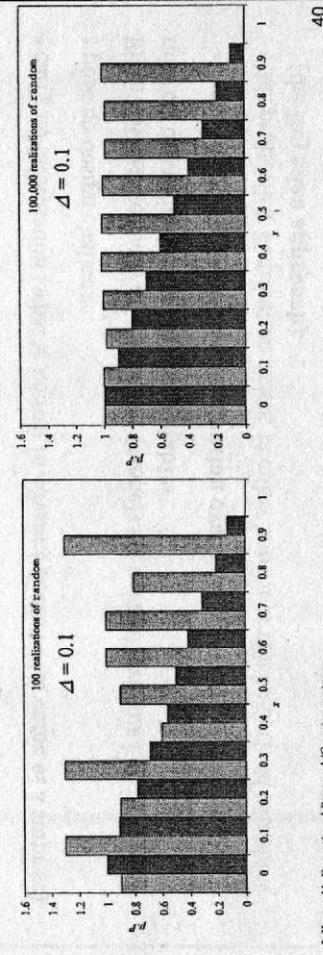
## Empirical Distributions



Having realizations  $x_1, \dots, x_N$  one constructs a histogram:

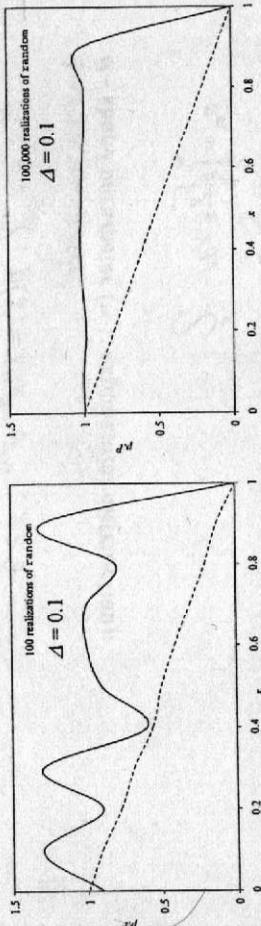
- divide the range of possible realizations into bins of length  $\Delta$ ,
- count realizations falling into  $i$ th bin:  $k_i = \#\{n \mid i\Delta \leq x_n \leq (i+1)\Delta\}$ ,
- at  $i\Delta$ , draw a bar of width  $\Delta$  and height  $p_i = (k_i/N)/\Delta$ .

Cumulative histogram constructed analogously, with bars of height  $C_i = \sum_{j \geq i} p_j$ .



## Empirical Distributions (2)

...or for readability, use smooth lines instead of bars:



$$\text{Mean value: } x_m = \frac{x_1 + \dots + x_N}{N}$$

$$\text{Standard deviation (dispersion around mean): } \sigma_x = \sqrt{\frac{1}{N-1} \sum_{n=1}^N (x_n - x_m)^2}$$

*Alternative Skewness*

J. Konečný, Operational Research/Queueing systems

41

## Theoretical Distributions

Probability density function and complementary distribution function

- histograms one would obtain taking  $N \rightarrow \infty$  and  $\Delta \rightarrow 0$ .

$$P(x) = \lim_{\Delta \rightarrow 0} \frac{\Pr[x \leq X < x + \Delta]}{\Delta}$$

$$\text{Probability density function: } P(x) = \Pr[x' \leq X < x''] = \int_x^{x''} p(x) dx.$$

$$\text{Complementary distribution function: } P(x) = \Pr[X \geq x] = \int_x^{\infty} p(y) dy$$

$$x_m = \int x p(x) dx, \quad \sigma_x = \sqrt{\int_{-\infty}^{\infty} (x - x_m)^2 p(x) dx}$$

J. Konečný, Operational Research/Queueing systems

42

2

## Theoretical Distributions (2)

Modeling for engineering applications often uses Weibull distribution:

$$P(x) = e^{-\lambda x^\theta}, \quad p(x) = \lambda \theta x^{\theta-1} e^{-\lambda x^\theta}, \quad x \geq 0$$

$\lambda$  – scale parameter,  
 $\theta$  – shape parameter (= 1: exponential distribution).

$$x_m = \int_0^{\infty} \sqrt[\theta]{\frac{y}{\lambda}} e^{-y} dy$$

$$\sigma_x = \sqrt{\int_0^{\infty} \frac{y^2}{\lambda^2} e^{-y} dy - x_m^2}$$

J. Konečný, Operational Research/Queuing Systems

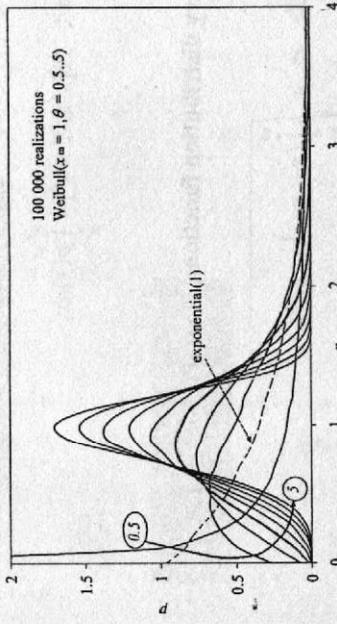
43

## Computer Generation of Arrival Streams

Given random generator that returns values  $(z_n)$ ,  
 how to generate pseudorandom numbers  $(x_n)$  with arbitrary  $P(x)$ ?

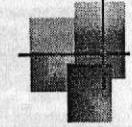
$$x_n \text{ solves } P(x) = z_n \quad \text{e.g., for Weibull distribution: } x_n = \theta \sqrt{-\frac{\ln z_n}{\lambda}}$$

(method of inverted distribution function; many others exist).



J. Konečný, Operational Research/Queuing Systems

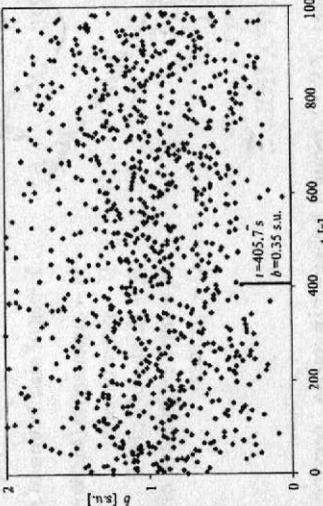
44



## Computer Generation of Arrival Stream (2)

Introduce random variables:

- $A$  – interarrival interval (realizations:  $a_n$ )
- $B$  – request size ( $b_n$ )



Arrival stream we met before had been generated as  $A, B \sim \text{Weibull}(1, 2.5)$

J. Kowalski, Operational Research/Queuing systems

not good & not efficient, too busy & no structure

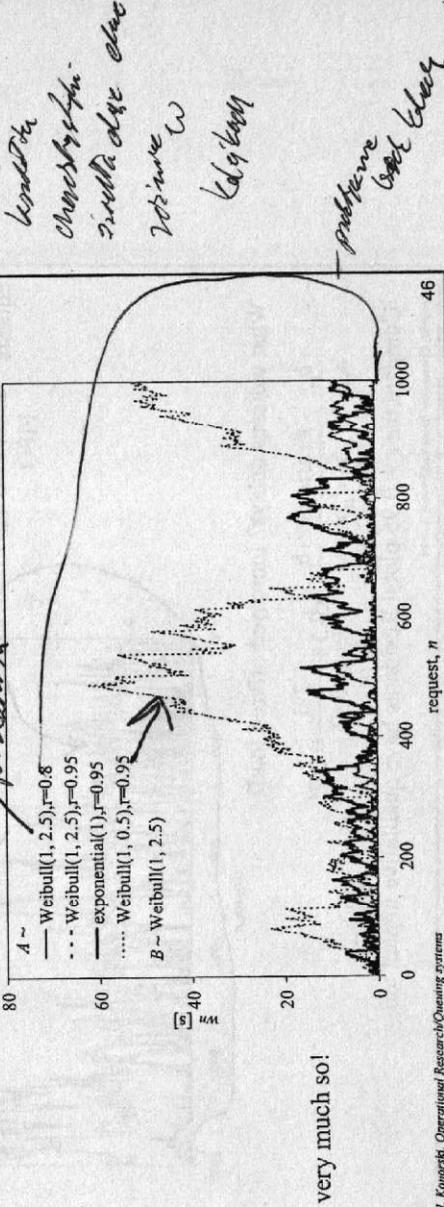
process takes time  
depends on all  
kinds of variables  
medium load  
intensity  $\rightarrow$

## Impact of Distribution of $A$ and $B$

We have seen the impact of mean values  $a_m$  and  $b_m$  upon the queuing process  
(through offered load).

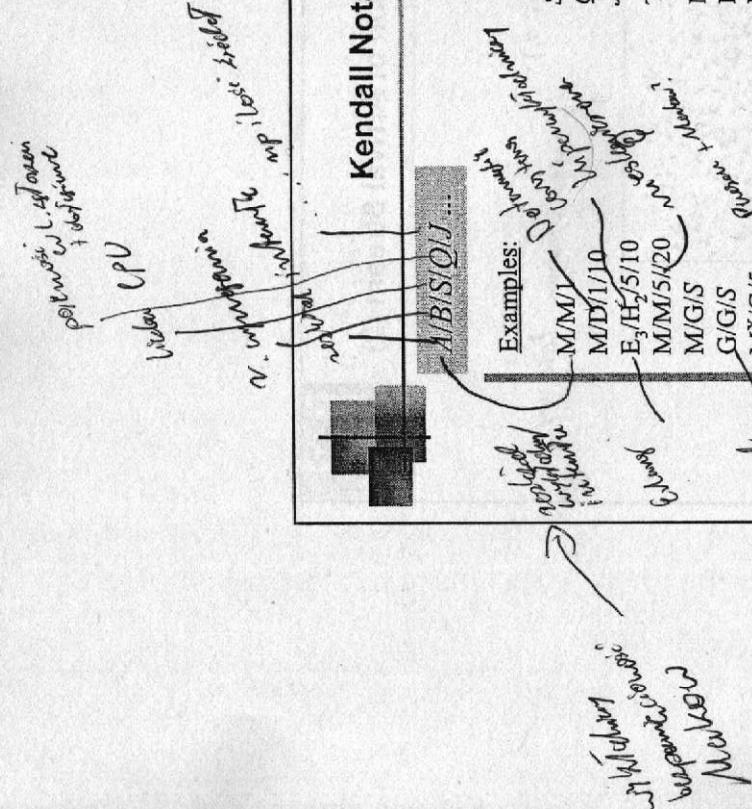
Is the shape of distributions of  $A$  and  $B$  equally relevant?

P. Krok



very much so!

J. Kowalski, Operational Research/Queuing systems



Kendall Notation

Examples:

M/M/1  
M/D/1/10  
E<sub>3</sub>/H<sub>2</sub>/5/10  
M/M/5/20  
M/G/S  
G/G/S  
M<sup>X</sup>/G/5  
MMPP/G/1

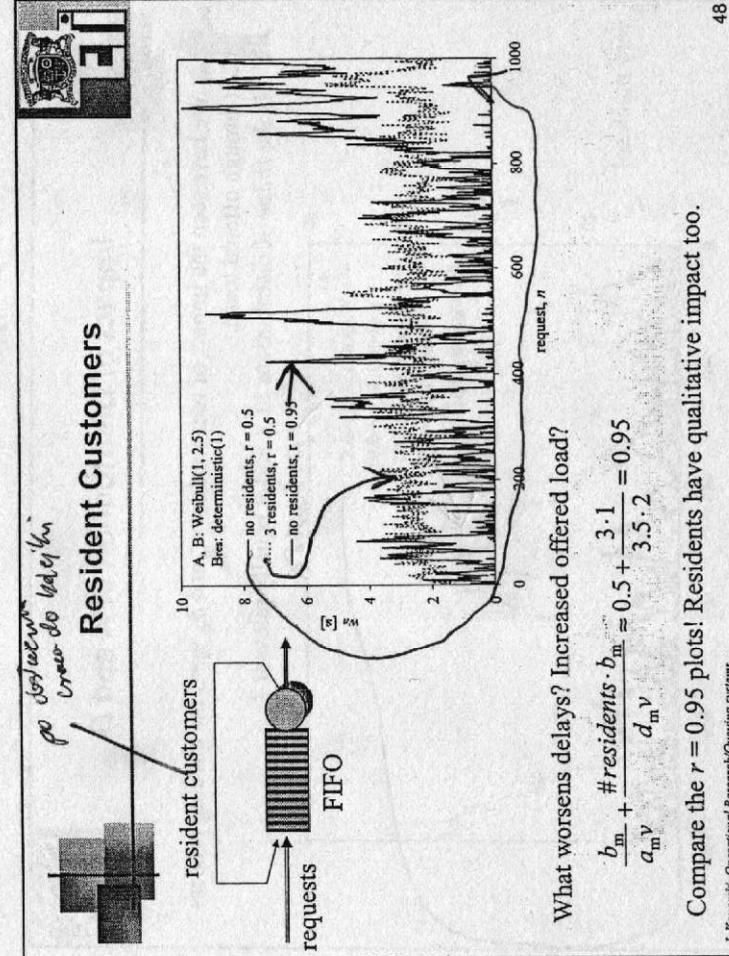
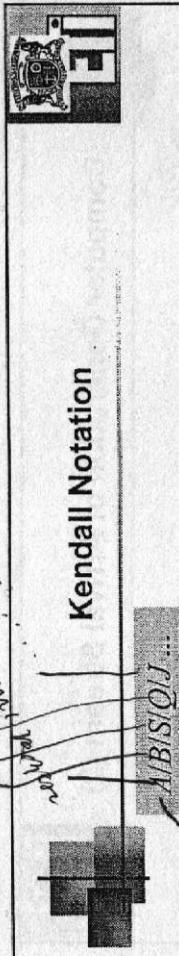
Examples:

- $S$  - # processors
- $Q$  - buffer capacity
- $J$  - request source population size
- Types of distribution of  $A, B$ :
- M - exponential (Markovian)
- D - deterministic
- $E_k$  - Erlang of order  $k$
- $H_k$  - hyperexponential of order  $k$
- G - general

CENTRAL INSTITUTE OF ENGINEERING / COMMENCEMENT RULES?

二

!



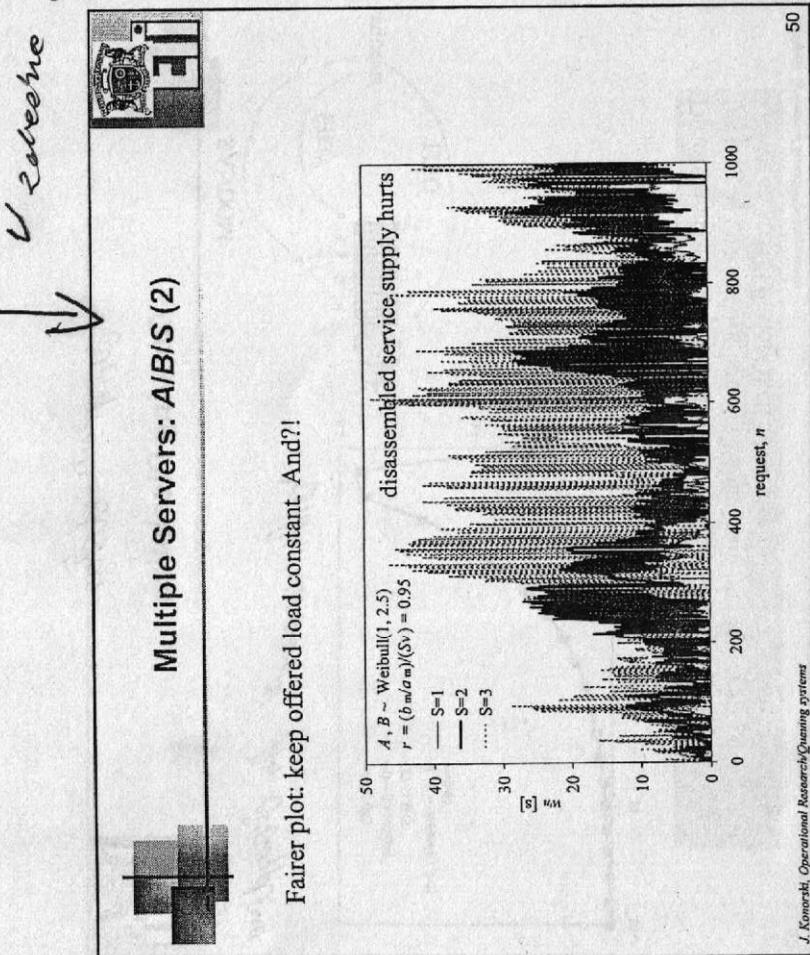
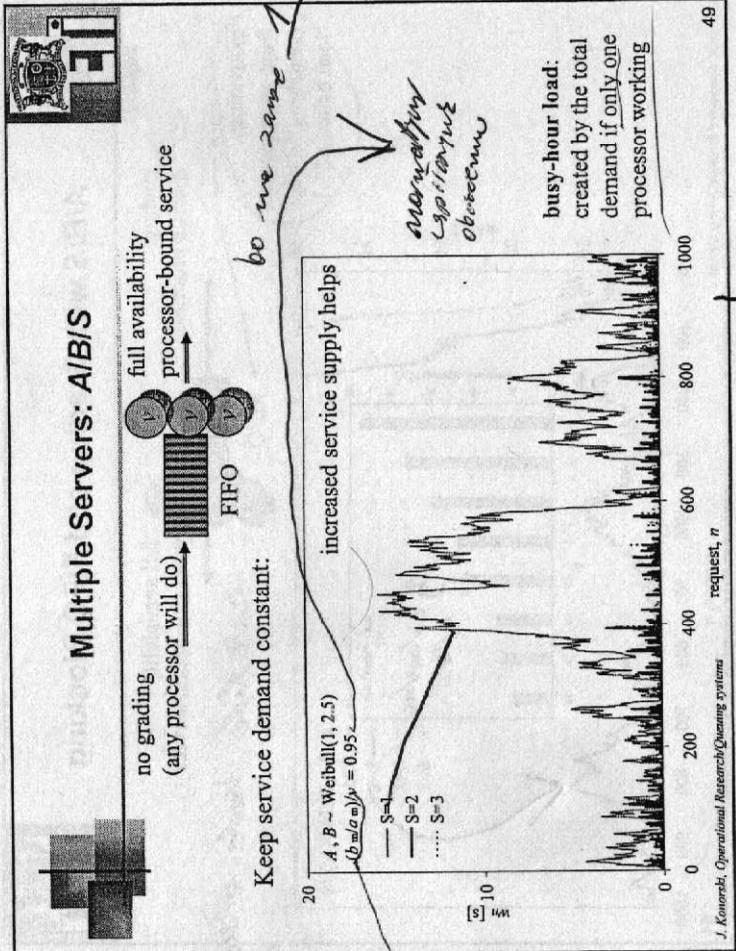
What worsens delays? Increased offered load?

$$\frac{b_m}{a_v} + \frac{\# \text{residents}}{d_v} \cdot b_m = 0.5 + \frac{3.1}{3.5 \cdot 2} = 0.95$$

Compare the  $r = 0.95$  plots! Residents have qualitative impact too.

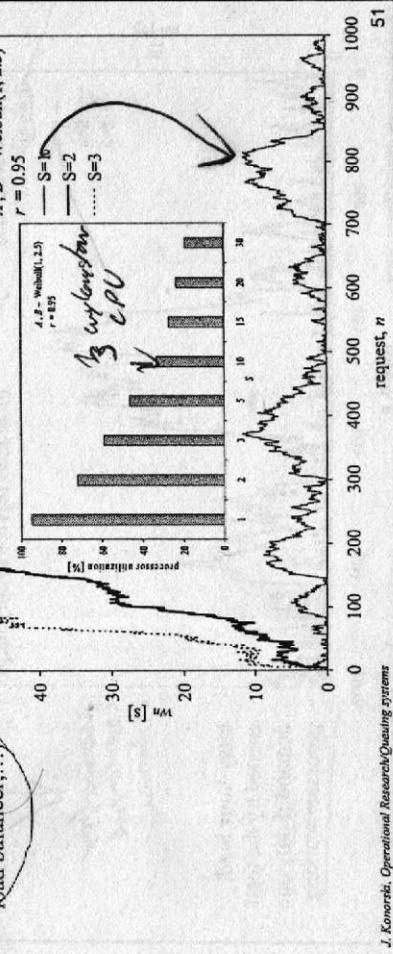
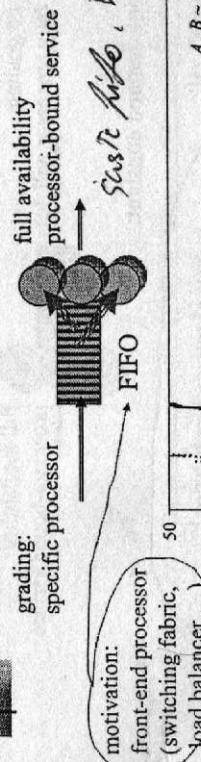
J. Komorski, Operational Research/Queuing Systems

48



A/B/S with Grading and FIFO Blocking

meester  
paroleer  
pale Venetiaan  
lambertus duynen  
deutse

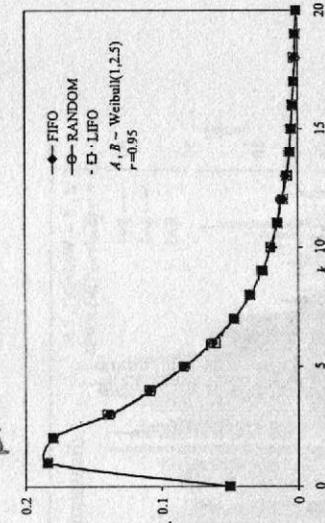
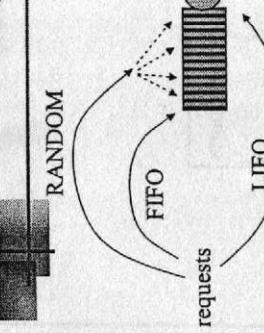


J. Konc

one  
by  
size  
process;  
where no  
new laboratory

Queuing Discipline

Queuing Discipline



W. H. Dyer  
1911

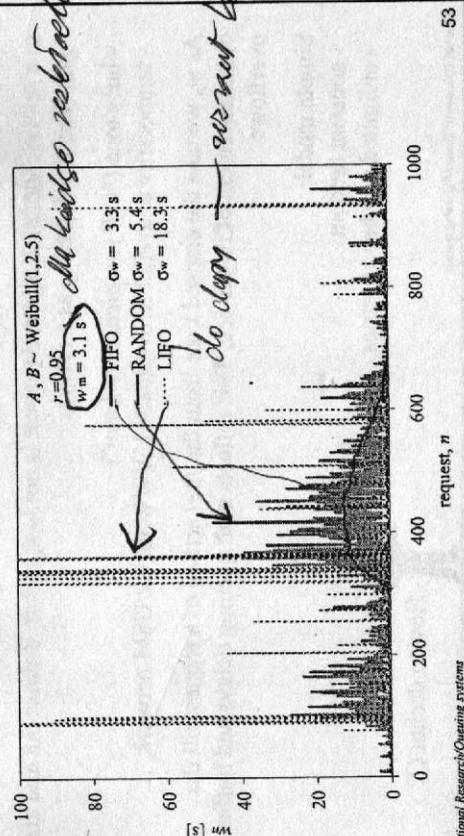
Distribution of  $N(t)$  is QD invariant if (as long as QD is work-conserving) known on arrival or not taken advantage.

J. Konorski, Operational Research/Queueing Systems

Dopóki jak Lektorówka i Absynt na doszły we wykładzie "Zajęcia z wojen"

## Queuing Discipline (2)

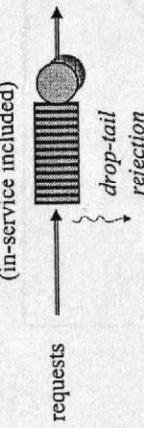
Which QD cause largest  $w_m$ ? largest  $\sigma_w$ ?  
With known real usage pattern  
to also refine



J. Konarak, Operational Research/Queuing systems

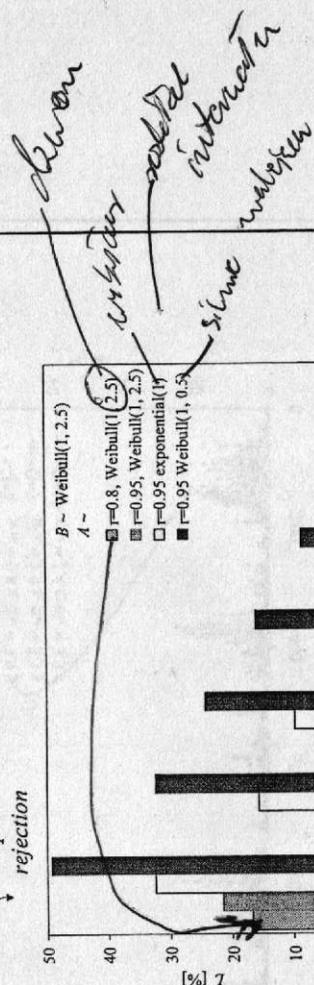
## Finite Buffer Capacity: A/B/1/Q

maximum  $Q$  requests  
(in-service included)



Motivation:  
- memory cost!  
- delays, management cost!

Low waiting delays ( $w_n$ ) traded for increased loss  $L$  due to buffer overflow.



J. Konarak, Operational Research/Queuing systems

*Polish version is going well after optimization*

## A/B/1/Q with Retries

In a realistic model, a rejected request is not lost; rather, it times out and arrives again (this is referred to as a **retry**).

- busy tone ("will you please try later")
- temporarily unavailable/overloaded WWW server/GSM network, ...

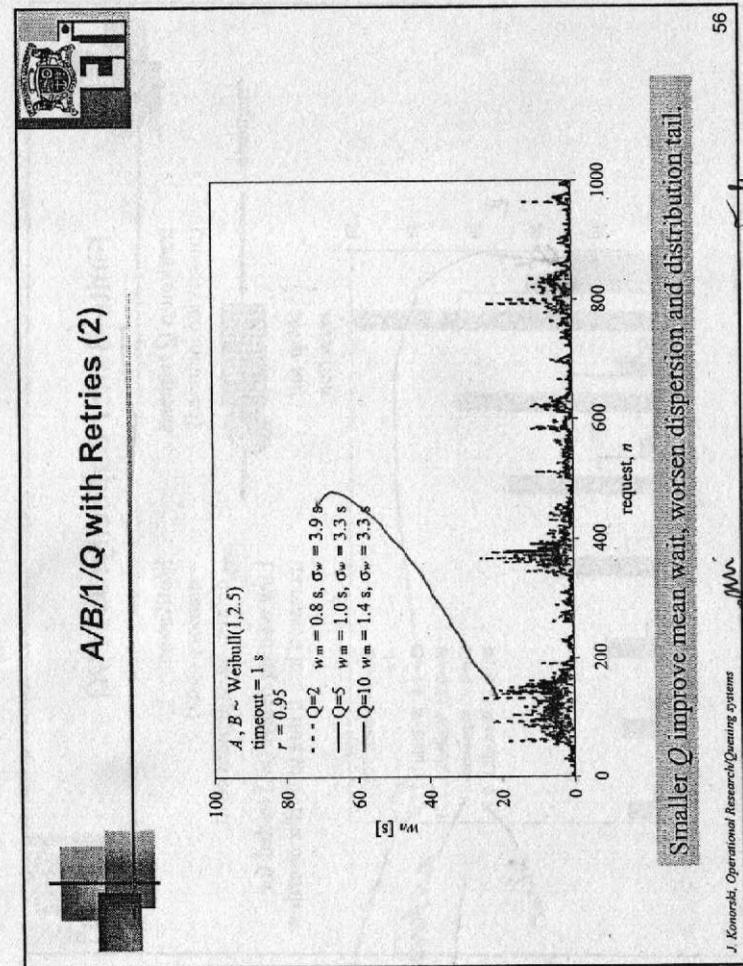
As  $w_n$  we take the elapsed time from the *first* arrival of a request till the commencement of its service. This reflects both queuing delays and buffer overflows.

Simple model:

- constant timeout
- unlimited number of retries.

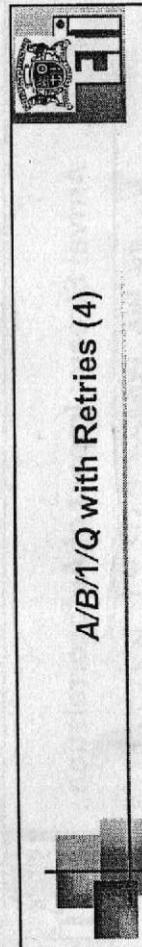
How to design  $Q$ ? 55

J. Komorski, Operational Research/Queueing systems





core backlog  
increasingly



**Common Models of Arrival Stream**

• Bernoulli  
• Erlang  
• gamma  
• Weibull  
• renewal streams, often named after distribution of  $A$

more complex e.g., time-varying distributions of  $A$  and  $B$ , *Markov Modulated Poisson Process*, *Batch Markovian Arrival Process*, fractal Brownian process, ...  
– model nonstationarity, dependence on the queuing process, bulk arrivals, internal correlation, ...

*J. Komorník, Operational Research/Queuing systems*

59

**Arrival Stream: Impact of Autocorrelation**

• Under what condition  $A$  is at least  $\sim \langle a_m \rangle$ : Are distributions of  $A$  and  $B$  enough to determine the queuing process (given fixed  $v$ ), or do we need information on the internal correlation in  $(a_n)$ ?

$$\text{corr}_a(l) = \frac{1}{\sigma_a^2} \sum_{m=1}^M (a_m - \bar{a})(a_{m+l} - \bar{a}), M \rightarrow \infty, l = 0, 1, 2, \dots$$

(autocorrelation function = how correlated are intervals  $l$  requests apart, correlation normally vanishes for larger  $l$ )

Renewal arrival stream is uncorrelated (white noise-like):

$$\text{corr}_a(l) = \begin{cases} 1, & \text{if } l = 0 \\ 0, & \text{if } l \neq 0 \end{cases}$$

60

**Arrival Stream: Impact of Autocorrelation (2)**

Generate  $(a_n)$  and  $(b_n)$  according to Weibull(1, 2.5) distribution using the method of inverted distribution function. Input the obtained renewal arrival stream to a queuing system with  $r = 0.95$ .

Observe the queuing process  $(w_n)$ .

Next, shuffle i.e., apply random permutation to the  $(a_n)$ , use the same  $(b_n)$  and again observe  $(w_n)$ .

Is there any impact from shuffling?

**Experiment 1***reliable solution observed**idrone low**reduces low**Shuffle**shuffles**shuffles*

61

**Arrival Stream: Impact of Autocorrelation (3)**

Shuffling makes almost no difference.

Legend:

- $A \sim$  — Weibull(1, 2.5)
- shuffled
- $B \sim$  Weibull(1, 2.5)
- shuffled

$r=0.95$

*problem less problem*  
*we see*  
*working*  
*reducing*

62

## Arrival Stream: Impact of Autocorrelation (4)

### Experiment 2

Take  $A \sim \text{Weibull}(1, 2.5)$ , and generate  $(a_n)$ :

- as iid intervals from successive random numbers – renewal stream,
- by repeating each successive interval  $R$  times, where  $R$  is drawn from  $\text{Weibull}(1, 0.5)$  distribution.

In both variants, distribution of  $A$  is the same.

However, the second variant yields  $(a_n)$  with long-range autocorrelation – a **self-similar** arrival stream.

Using the same  $(b_n)$  as before, input the obtained arrival stream to a queuing system with  $r = 0.95$ .

*Arrival stream with autocorrelation*

J. Konarzki, Operational Research/Queuing systems

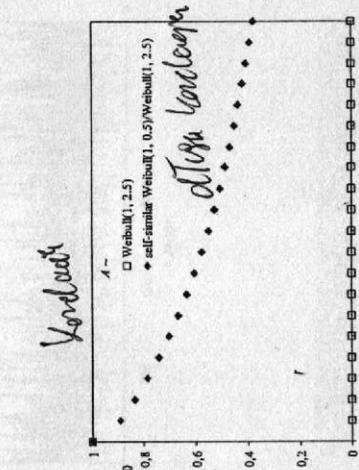
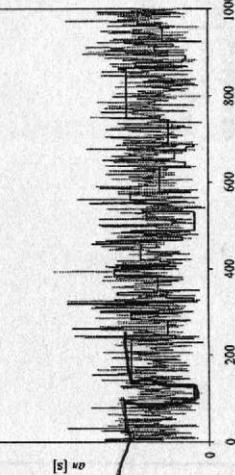
*Arrival stream without autocorrelation*

## Arrival Stream: Impact of Autocorrelation (5)

$A \sim$

— Weibull(1, 2.5)

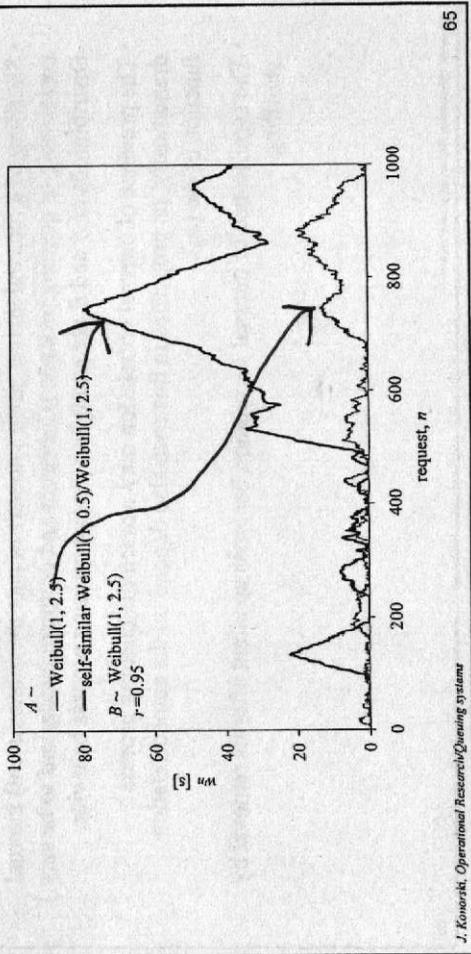
→ self-similar Weibull(1, 0.5) \* Weibull(1, 2.5)



J. Konarzki, Operational Research/Queuing systems

## Arrival Stream: Impact of Autocorrelation (6)

Comparison of the queuing process for renewal and self-similar arrivals is quite spectacular... (note again: distributions of  $A$  and  $B$  are same in both variants!)



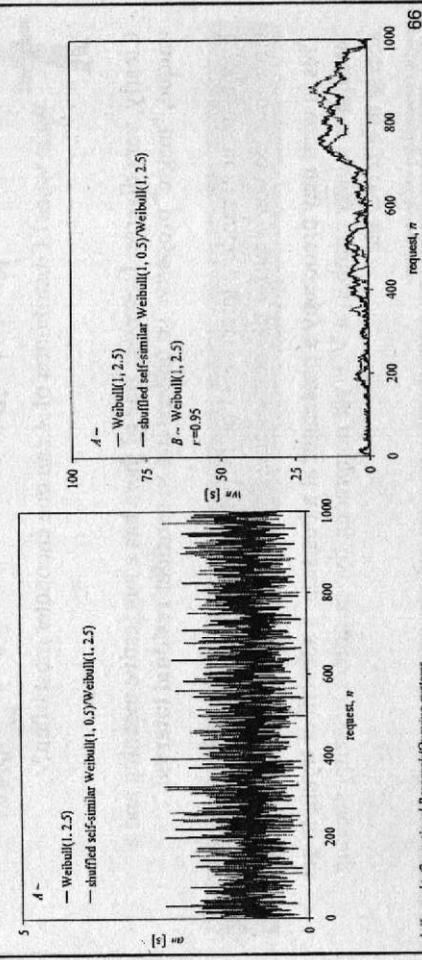
Konaraki PST  
BARD2 O 272-  
Tasmane  
Zimberto  
Lovelett

## Arrival Stream: Impact of Autocorrelation (7)

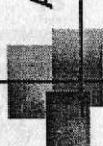
### Experiment 3

Shuffle the obtained self-similar ( $a_n$ ), which of course removes long-range autocorrelation, but preserves distributions of  $A$  and  $B$ .

Use the same ( $b_n$ ) and again compare with the renewal stream variant...



Modulation  
with no server  
latency



## Arrival Stream: Impact of Autocorrelation (8)



### Conclusions:

- Shuffling of a renewal arrival stream doesn't impact the (nonexistent) internal correlation, or queuing process. (Construct and compare histograms to be sure.)
- Distributions of  $A$  and  $B$  are enough to predict queuing process behavior.
- The presence of internal correlation may worsen the queuing process dramatically, its properties in this case also depend on the autocorrelation function of the  $(a_n)$ .
- The significance of internal correlation becomes apparent after its removal by shuffling.

J. Konsztat, Operational Research/Queuing systems

67

11

U 2 problems A  
with no latency.  
How long do we have?  
We have no random



## Renewal Arrival Stream: Residual Interval



Except special fields of research (heartbeat anomalies, overflows of the Nile, Web traffic analysis etc.), renewal streams model real-world arrival streams adequately. Henceforth we focus upon them.

When they will make overflow?

What types of distribution of  $A$  can one encounter most often?

Clearly, very diverse. However, one of them has a suggestive meaning and a unique, "magic" property. To understand it, consider **residual interval**.

Events occur at random intervals. You arrive at a random instant.

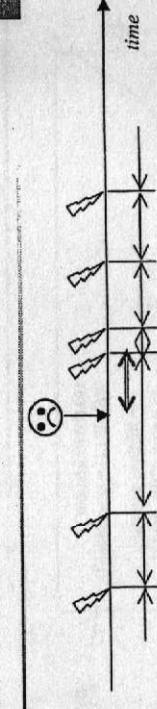
How long do you wait for the next event?

This is what may preoccupy a passenger at a bus stop, a subscriber trying to get through to a busy number, a VIP yet nonpreemptive customer urgently seeking access to a server etc.

J. Konsztat, Operational Research/Queuing systems

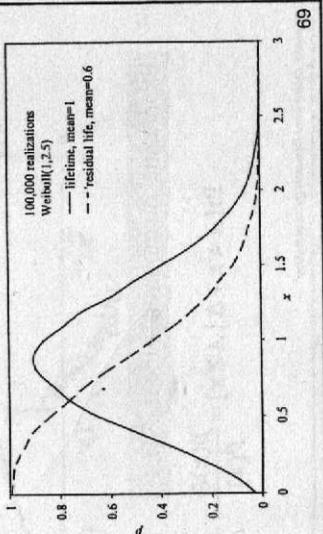
68

## Renewal Arrival Stream: Residual Interval (2)



Buses run every hour, on average. "Statistical passenger" waits for half an hour?

Given the distribution of inter-bus interval (lifetime), what distribution does the residual interval have? Shifted towards smaller realizations?

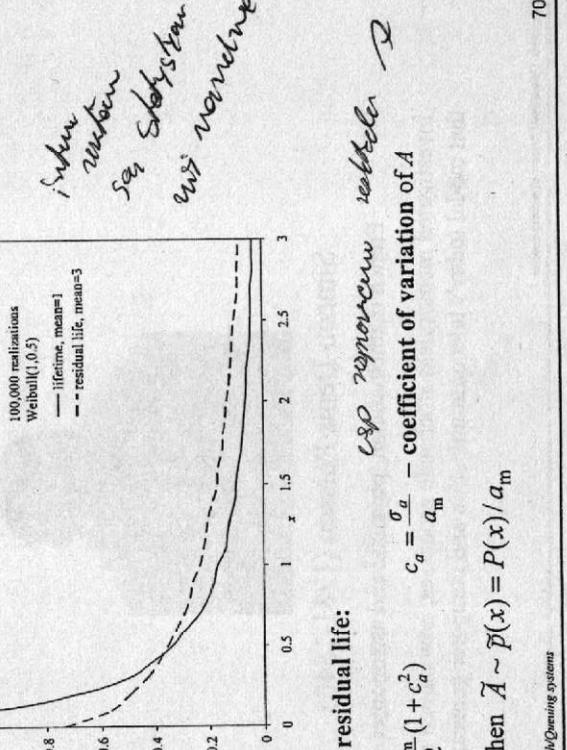


J. Konarak, Operational Research/Queuing systems

## Renewal Arrival Stream: Residual Interval (3)



Infinite waiting for system until next arrival



J. Konarak, Operational Research/Queuing systems

Paradox of residual life:

$$\tilde{a}_m = \frac{a_m}{2} (1 + c_a^2) \quad c_a = \frac{\sigma_a}{a_m} - \text{coefficient of variation of } A$$

If  $A \sim P(x)$  then  $\tilde{A} \sim \tilde{P}(x) = P(x)/a_m$

70

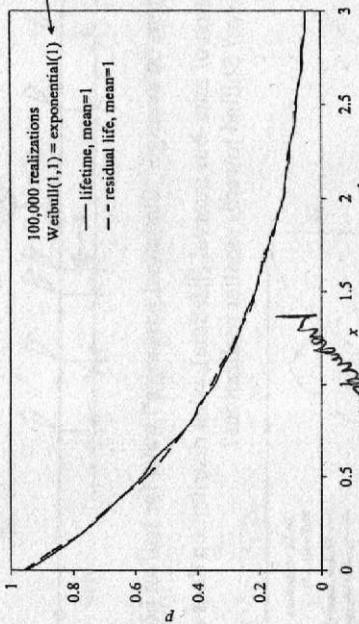
Wieland | Wieland (1) | Wieland  
Wieland | Wieland -

### Renewal Arrival Stream: Residual Interval (4)

100,000 realizations  
Weibull(1,1) = exponential(1)  
— lifetime, mean=1  
— residual life, mean=1

$P(x) = e^{-x/a_m}$

Exponential distribution is memoryless



J. Komorski, Operational Research/Queueing systems

71

Stream arrival  
not measured in  
topping interval  
now



### Poisson Arrival Stream



### Simeon-Denis Poisson (1781-1840)

French mathematician, physicist and astronomer investigated memoryless stochastic processes, now named after him, that model today's telecommunication and computer generated traffic

J. Komorski, Operational Research/Queueing systems

72

Stream arrival  
not measured in  
topping interval  
now

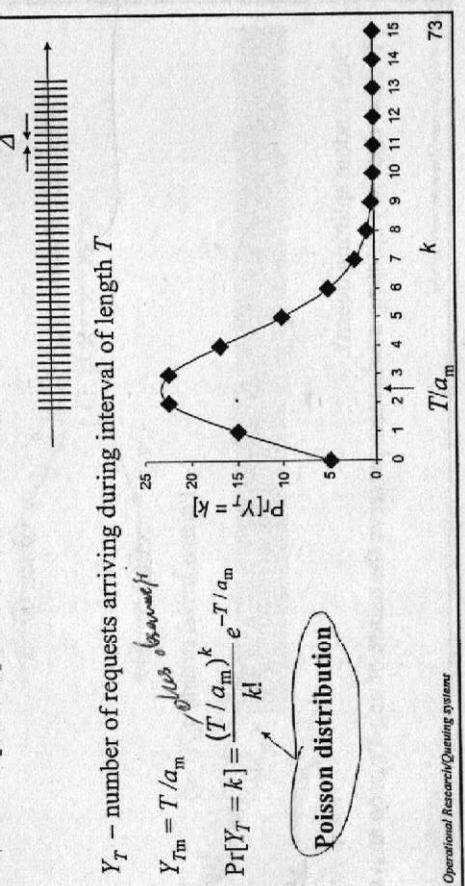
the probability of there is delay older ( older ) than  
while passing

36

## Poisson Arrival Stream (2)

At any instant of time, new request arrival occurs with constant probability.  
In Kendall notation:  $M/\dots$  systems.

$$\Pr[\text{new request in } (t, t + \Delta)] = \Delta/a_m + o(\Delta)$$



## Poisson Arrival Stream (3)

For stationary stochastic models, time average-type characteristics of queuing processes are determined using probability theory.

Steady state is then referred to as **statistical equilibrium**.

Consists in the time averages of interest stabilizing over time e.g., system state probabilities, loss probability, waiting time distributions etc.

*bedienungsdauer*

*p. B. wickeln  
p. verarbeiten (as 2  
schaffern mit  
heute & loswerden*

**PASTA (Poisson Arrivals See Time Averages):** for Poisson arrivals,  $P_k^+ \equiv p_k$

(requests arriving according to a Poisson stream "see" the same queue length  $L$  as does a random observer).

$$\Pr[N(t) = k] \frac{\Delta}{a_m} = \Pr[N(t) = k]$$

Hence, in  $M/G/S/Q$ : loss probability due to buffer overflow equals  $L = p_Q^+ \equiv p_Q$

*J. Kemerati, Operational Research/Queuing systems*

*zur Zeit kommt  
im Kreis laufen  
wiederholen sich  
die gleichen  
Vorgänge*

*zur Zeit kommt  
im Kreis laufen  
wiederholen sich  
die gleichen  
Vorgänge*

*zur Zeit kommt  
im Kreis laufen  
wiederholen sich  
die gleichen  
Vorgänge*

*zur Zeit kommt  
im Kreis laufen  
wiederholen sich  
die gleichen  
Vorgänge*

**Poisson Arrival Stream (4)**

Random splitting:

Random splitting preserves arrival stream's Poissonian nature!

Neither does any splitting mechanism preserve the nature of non-Poisson arrivals.

J. Konarski, Operational Research/Queueing systems

75

**Aggregation of Renewal Arrival Streams**

Independent components

- renewal arrival streams  $a_m^{(j)}$  (input to backbone link, mainframe / Web server...)

reduced at backbone

$$\frac{T}{a_m} = \sum_{j=1}^J \frac{T}{a_m^{(j)}} \quad (\text{in particular, for identical components, } a_m = \frac{a_m^{(j)}}{J})$$

What interval distribution does the system "see"? At least a renewal stream?  
Not necessarily. In general, analytic calculation difficult if not impossible.

Handwritten notes include 'reduced at backbone' and 'residual intervals:  $\bar{A} = \min\{\bar{A}^{(1)}, \dots, \bar{A}^{(J)}\}$ '

J. Konarski, Operational Research/Queueing systems

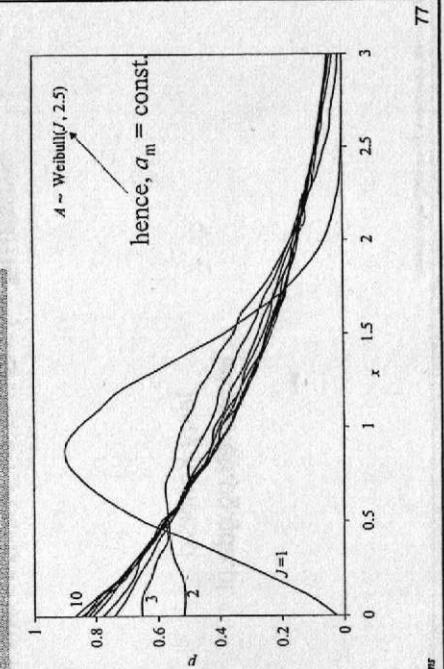
76

## Aggregation of Renewal Arrival Streams (2)

With  $J \rightarrow \infty$ , but  $a_m > 0$  (a practical model) and independent component streams:

$\sum_{j=1}^J A^{(j)}$  *looks like we've been nice*

Aggregated arrival stream is Poisson (Palm theorem).



J. Konsoli, Operational Research/Queueing Systems

Stochastic  
note observe  
process each  
interval  
David John  
Ugure  
 $\sum$  Poisson  
in next  
line will  
have a/  
with  $J$  have  $a_j$

## Aggregation of Renewal Arrival Streams (3)

Proof: omitted :)

$$\tilde{A} = \min\{\tilde{A}^{(1)}, \dots, \tilde{A}^{(J)}\}, \text{ so } \Pr[\tilde{A} \geq x] = \prod_{j=1}^J \Pr[\tilde{A}^{(j)} \geq x]$$

$J \rightarrow \infty$ , but  $a_m > 0$  (a practical model). That is,  $a_m^{(j)} \rightarrow \infty$ .

For any finite  $x$ ,  $x/a_m^{(j)} \rightarrow 0$  and we can neglect  $Y_x^{(j)} > 1$ .

$$\Pr[\tilde{A}^{(j)} \geq x] = 1 - \Pr[Y_x^{(j)} > 0] \approx 1 - \Pr[Y_x^{(j)} = 1] \approx 1 - Y_{x,m}^{(j)} \approx 1 - \frac{x}{a_m^{(j)}}$$

Finally,

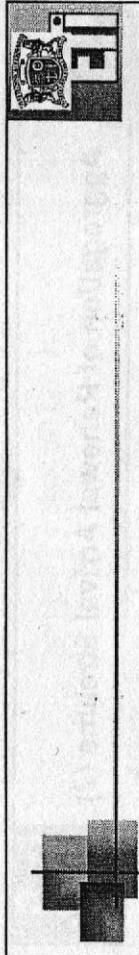
$$\Pr[\tilde{A} \geq x] = \prod_{j=1}^J \left( 1 - \frac{x}{a_m^{(j)}} \right) \approx \exp \left( - \sum_{j=1}^J \frac{x}{a_m^{(j)}} \right) = e^{-x/a_m}$$

since for  $x, y, z, \dots \ll 1$ , the following holds:  $(1-x)(1-y)(1-z)\dots \approx e^{-(x+y+z+\dots)}$

Since residual interval in the aggregated stream is exponentially distributed, so is interval itself!

J. Konsoli, Operational Research/Queueing systems

78



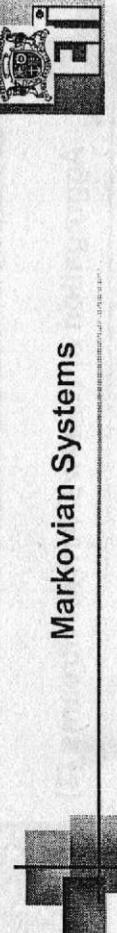
# Operational Research

## Queuing Systems 3: Markovian Models

Jerzy Komorski  
jekon@eti.pg.gda.pl

J. Komorski, Operational Research/Queuing systems

79



### Markovian Systems

Recall that queuing theory deals with systems and processes that can be observed, measured, and simulated.

Mathematical analysis may be useful too, but only if leads to *simple* and *insightful* results.

Take an  $A/B/\dots$  system. Is it easy to predict its characteristic theoretically?

Yes, if necessary simplifications are made:

- not too drastic  
keep models close to reality!  
(or else face charges of shaping the lock to fit the key!)
- yet bold enough  
keep problems tractable!  
get universal insight!

Example: Markovian queuing systems.

J. Komorski, Operational Research/Queuing systems

80

## Markovian Systems (2)

Exhibit the apparently unusual, but most useful **Markov property**.

(which, however, they share with a huge number of real-world dynamical systems – technical, physical, social, biological, economic, ...)

*What does system share  
 $t + \Delta$  - change workspace?*

$\text{state}(t + \Delta) = f(\text{state}(t), \phi)$

*Install(t) + release leave  $\emptyset$  - standardize Zentrale  
 60 new salary of  
 Wissens*

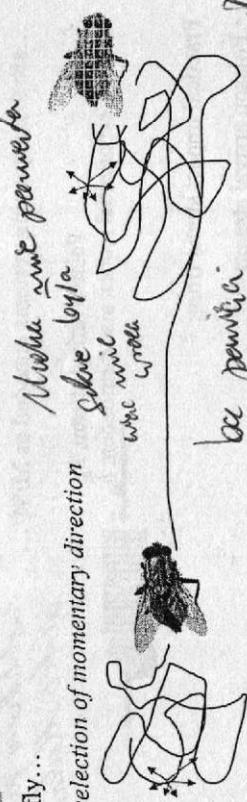
"noise" (random external input at time  $t$ ,  
 in general dependent on the current state,  
 but independent of earlier ones)

J. Konecny: Operational Research/Queueing systems

81

## Markovian Systems (3)

• The fly...



- card shuffling: card order in the deck (*selection of cut point*)
- gambler's capital / population (*current interest / growth rate*)
- $trend(t + \Delta) = (1 - c) \cdot trend(t) + c \cdot \phi(t)$  (*current observation*)
- $market\_share(t + \Delta) = \phi \cdot market\_share(t)[1 - market\_share(t)]$  (*current management performance*)
- Internet topology (*number and points of attachment of new networks*)

J. Konecny: Operational Research/Queueing systems

82



## Markovian Systems (4)



Andrej - Ivanovič Markov

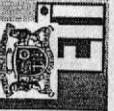
**Andrei A. Markov (1856-1922)**

Russian mathematician

investigated stochastic processes of finite memory, now named after him,  
that model many natural and man-made phenomena

*J. Konorat, Operational Research/Queueing Systems*

83



## Markovian Systems (5)

...are those queuing systems encoded as M/M/...  
Markov processes

Poisson arrival stream,  $a_m$   
exponential request size distribution,  $b_m$



$v$

Practical impact stems from:

- Poisson arrival stream
- Palm theorem
- random splitting
- PASTA
- pessimistic (meaning: fail-safe) performance characteristics – *expensive / slow and unreliable processes may pay off*
- exponential request size distribution: crude approximation of
  - call holding time, Web / P2P file transfer
  - batch processing time
  - ...

*J. Konorat, Operational Research/Queueing Systems*

84

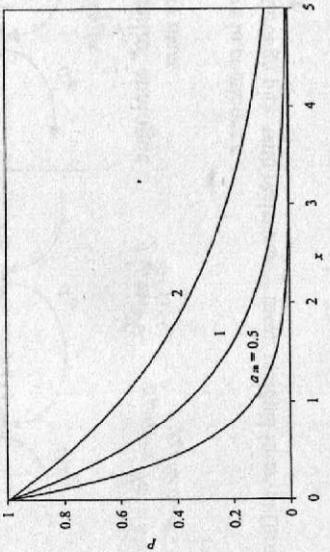
## Markovian Systems (6)

$\nu$  – (constant) processor speed

$A \sim M = \text{exponential}(a_m)$ :  $\Pr(A \geq x) = e^{-x/a_m}$

$B \sim M = \text{exponential}(b_m)$ :  $\Pr(B \geq x) = e^{-x/b_m}$

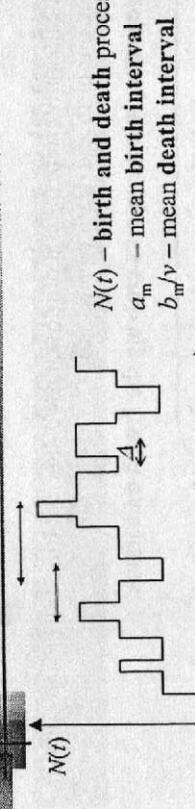
$\Pr[\text{service time} \geq x] = e^{-x/(b_m/\nu)}$



J. Konarak, Operational Research/Queuing Systems

85

## M/M/1 Queuing System



Let  $p_k(t) = \Pr[N(t) = k]$ . We know that  $\lim_{t \rightarrow \infty} p_k(t) = p_k$  (statistical equilibrium). State distribution  $(p_k)$  can be derived from birth and death equations.

- Exponential distribution has no probability mass at 0  $\Rightarrow$  we reside above or just 0  
suppose  $N(t) = k$ , what can happen between  $t$  and  $t + \Delta t$ ?  
practically, only one of the following: nothing / 1 birth / 1 death (if  $k > 0$ ),  
(state transitions between neighbor states only)
- Exponential distribution is memoryless (Markov property)  
 $\Rightarrow$  residual interval (time to occurrence) of next birth / death is statistically the same as the interval between consecutive births / deaths

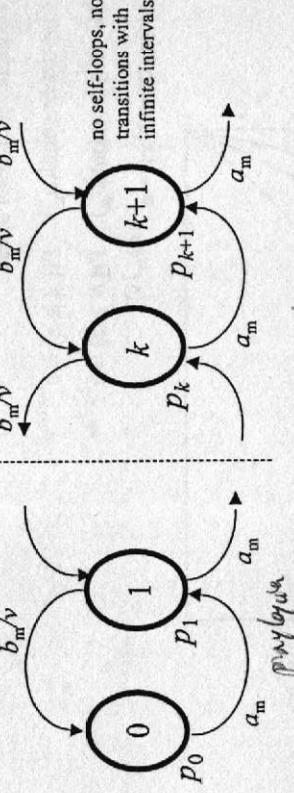
86

Now, why does it always decrease instead  
it was not my initial thought

43

M/M/1 Queuing System (2)

Zelotes 10, 2



- suggestive "hydraulic" analogue
- liquid ~ probability mass
- state ~ container
- $p_k$  ~ liquid pressure in container  $k$
- transition ~ flow through pipe with respect of respective event

*Leishman* positive to either a precipitin = *Leishman*  
*coccidioides* disease and infant *Cystisoma*, from which  
child died.

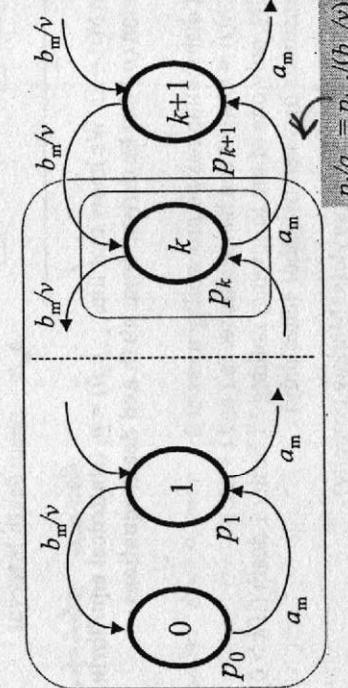
28

M/M/1 Queuing System (3)

Velby north  
landed  
at Weymouth  
Tunis

In statistical equilibrium, in- and outflow must balance out for *any* closed contour.  
(These are our birth and death equations.)

For convenience, select contours crossing the fewest transitions!



$$P_{k+1}^{\text{em}} = P^{k+1}(\mathcal{C}^{\overline{m}, \gamma})$$

J. Konorski, Operational Research/Queueing Systems

reduced storage  $\rightarrow$  posted secondary  
to system

## M/M/1 Queuing System (4)

$$P_r = p_0 \cdot r^k \quad p_0 + p_1 + p_2 + \dots = 1 \Rightarrow p_0 = 1 - r$$

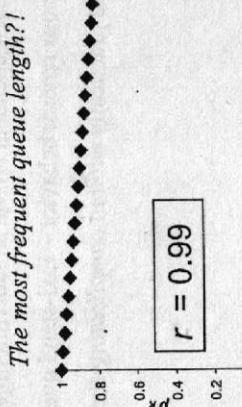
Hence mean queue length and further, by Little's theorem, mean waiting delay:

$$N_m = \frac{r}{1-r}$$

$$d_m = a_m N_m = \tau_m \left( \frac{1}{1-r} \right)$$

$$w_m = d_m - \tau_m$$

very suggestive!

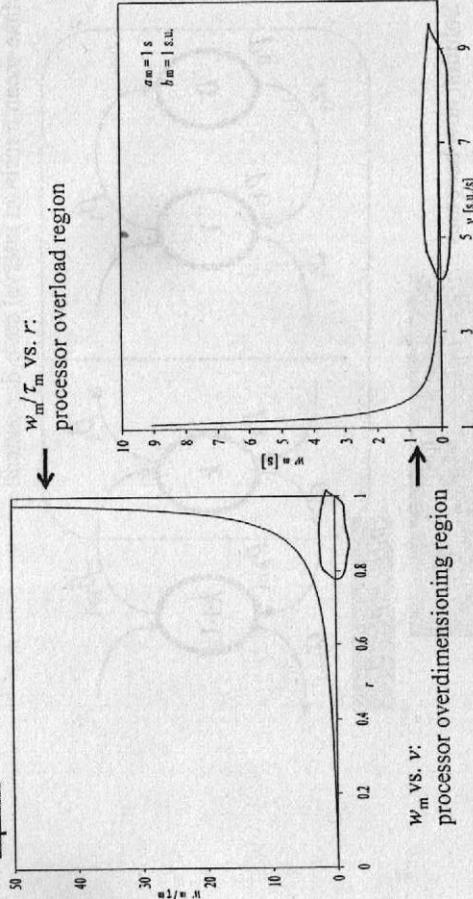


89

J. Kemerer, Operational Research/Queuing systems

## M/M/1 Queuing System (5)

$w_m / \tau_m$  vs.  $r$ :  
processor overload region



Distribution of waiting time does depend on QD.  
For FIFO, will be found with a little richer math... shortly.

J. Kemerer, Operational Research/Queuing systems

90

## M/M/... Systems

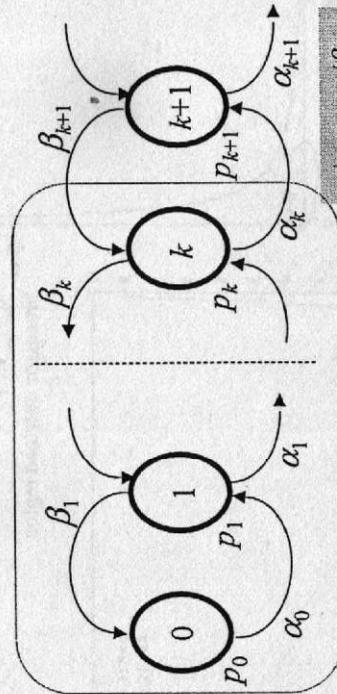
Birth-and-death processes help to analyze far richer and more realistic Markovian models of queuing systems featuring e.g.,

- finite buffer capacity (no-retry, *drop-tail*)
- multiple processors (no grading), perhaps in a queue-dependent number,
- queue-dependent arrival stream (intelligent terminal-type request sources)
- various request behavior – taxi-stand queue / token bucket, impatience, ...  
...and practically *without complicating the math!*.

J. Kononidis, Operational Research/Queuing systems 91

## M/M/... Systems (2)

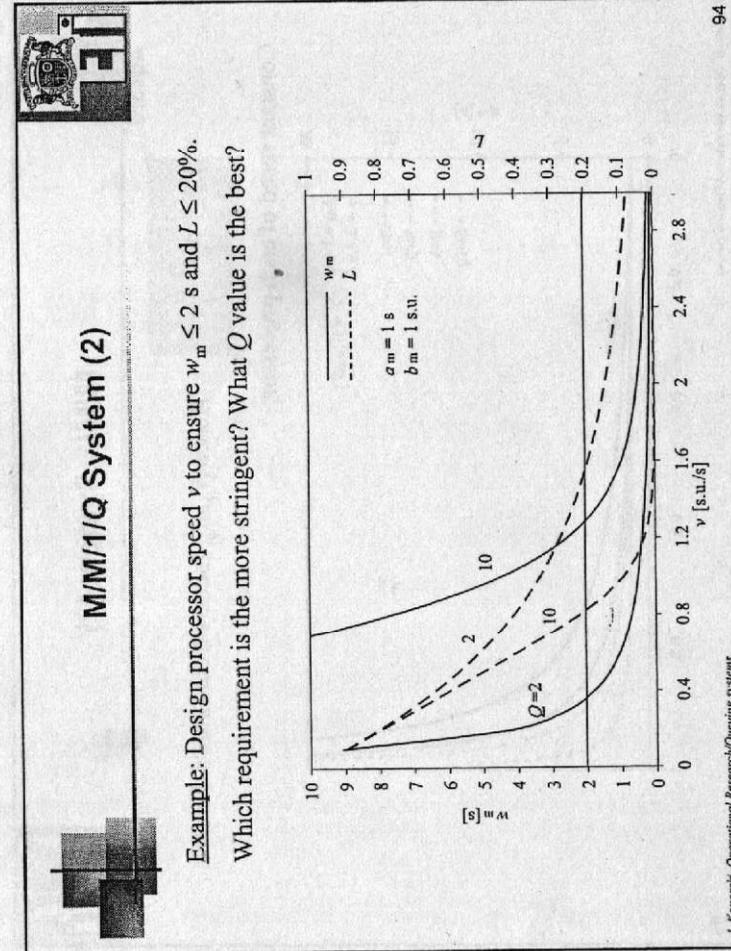
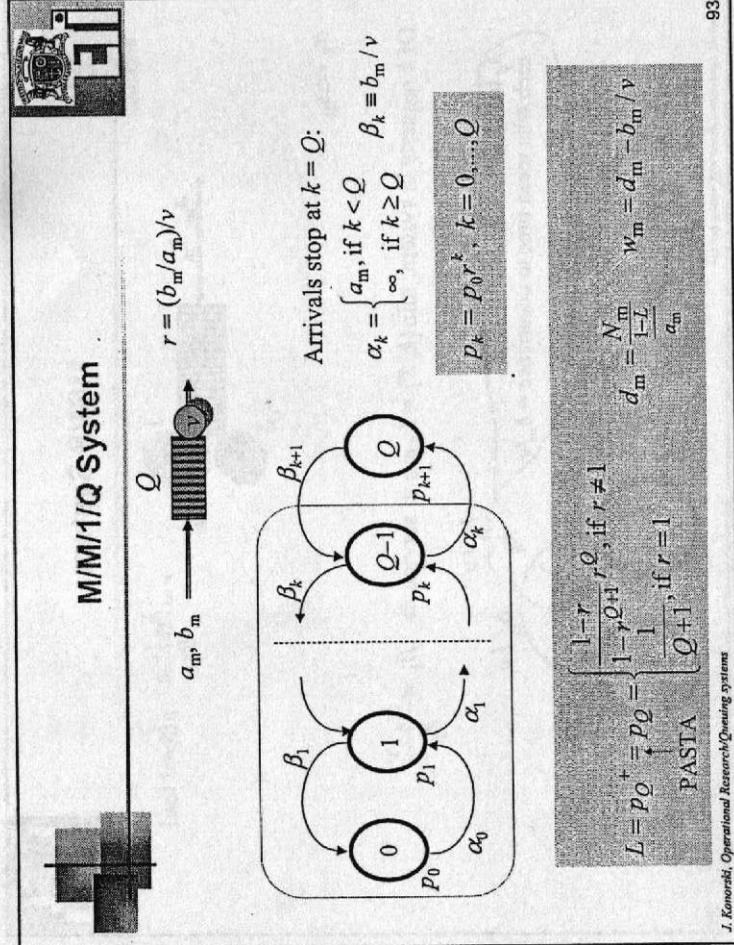
Make mean transition interval state-dependent:

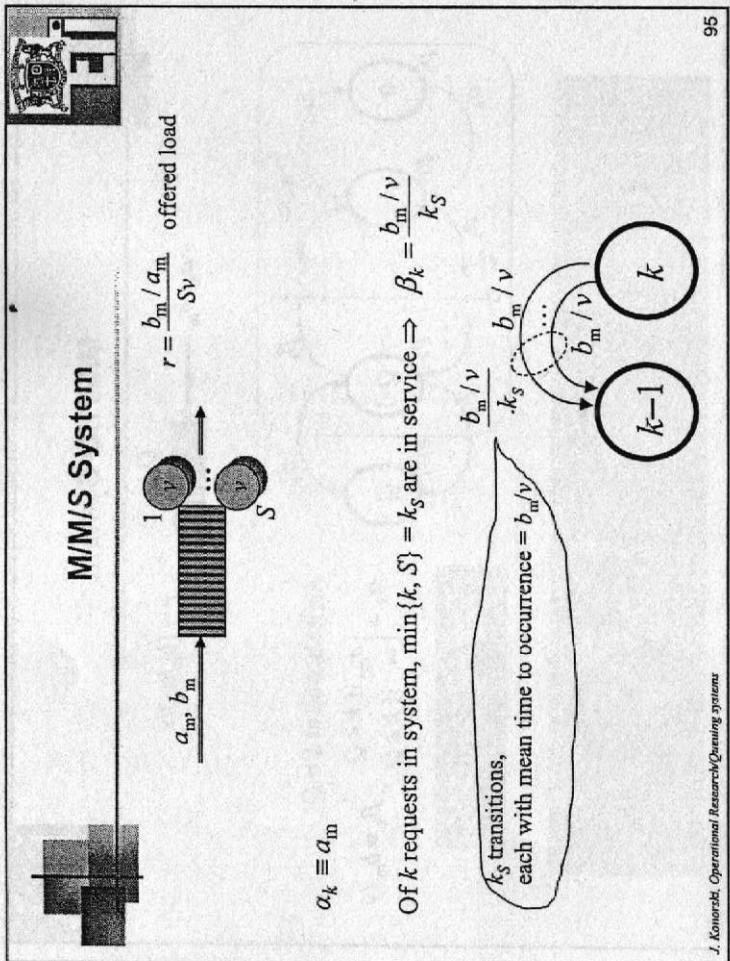


$$p_k \alpha_k = p_{k+1} \beta_{k+1}$$
$$p_k = p_0 \frac{\beta_1 \dots \beta_k}{\alpha_0 \dots \alpha_{k-1}}, \quad k = 0, 1, 2, \dots$$

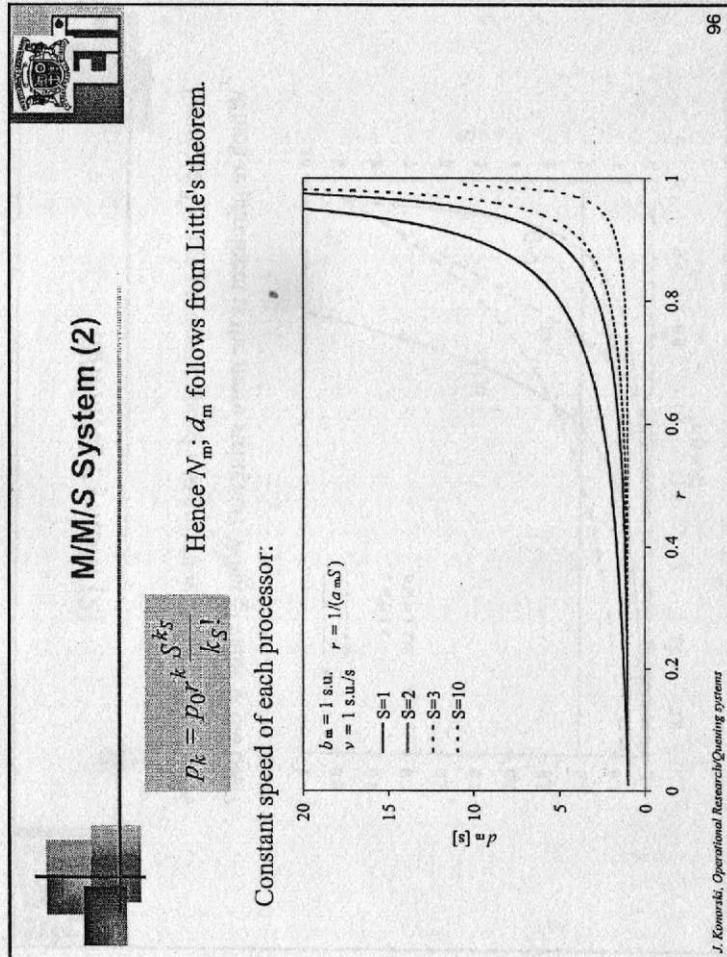
whence all the interesting characteristics:  $L, p_0, N_m, d_m, w_m, \dots$

J. Kononidis, Operational Research/Queuing systems 92

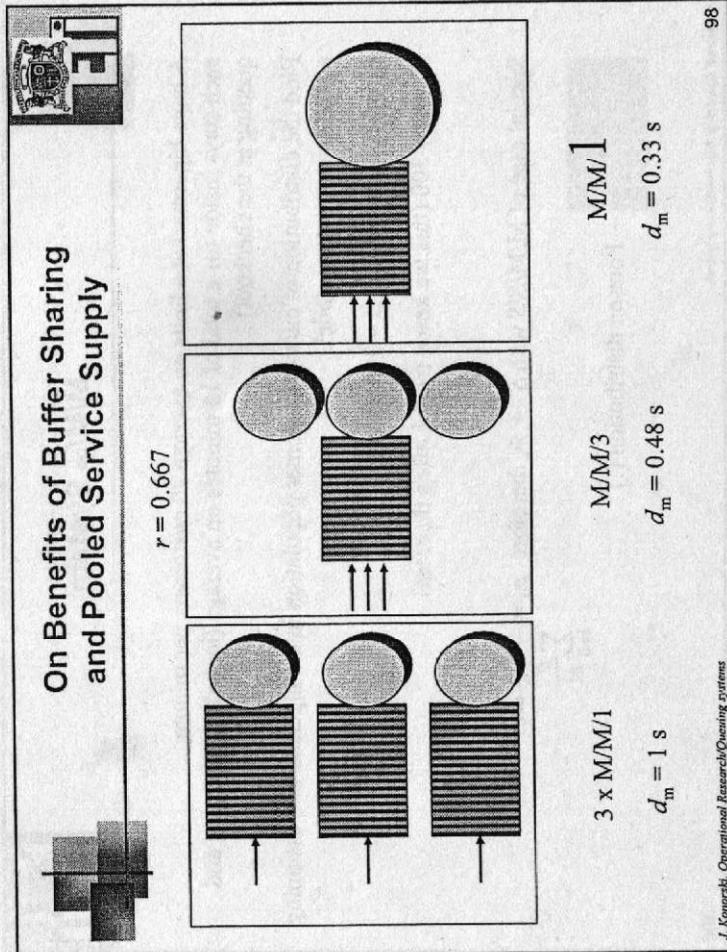
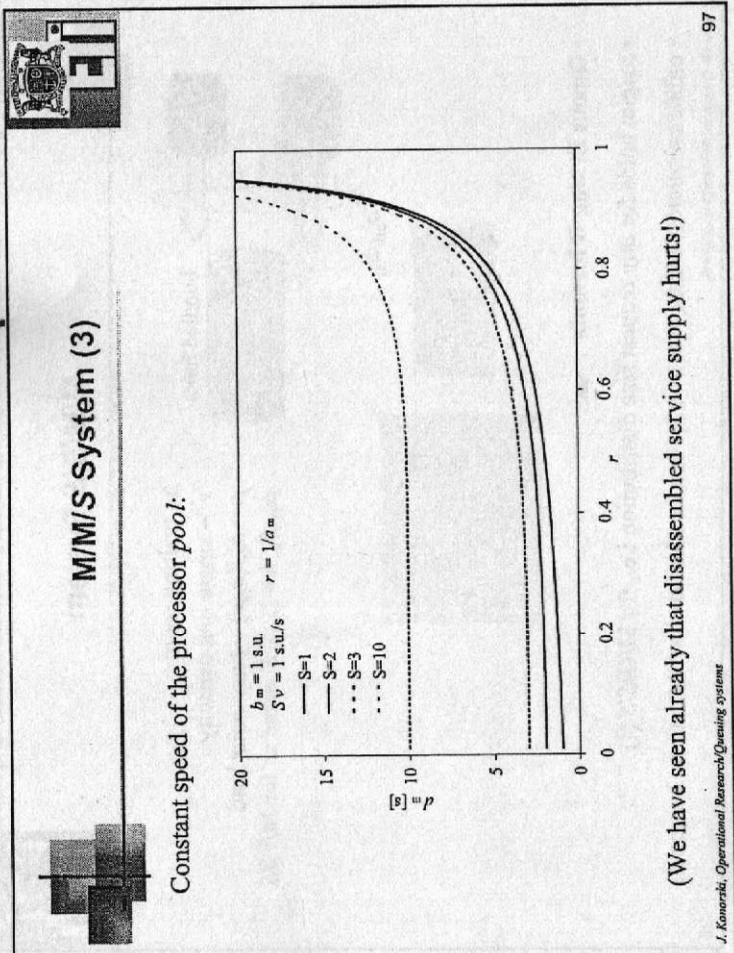


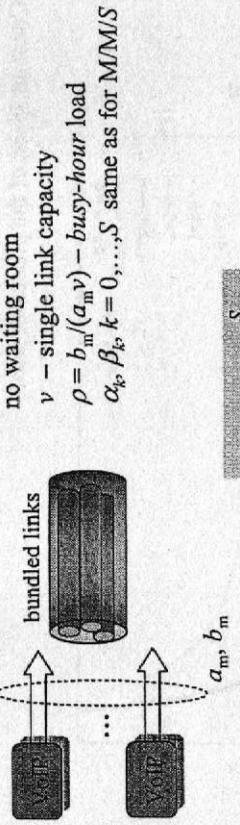
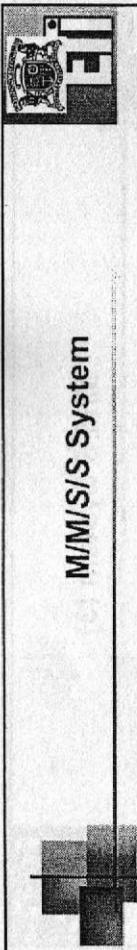


95



96





$$p_k = p_0 \frac{\rho^k}{k!}$$

$$p_S = L = \frac{\rho^S}{S!} \sum_{k=0}^S \frac{\rho^k}{k!}$$

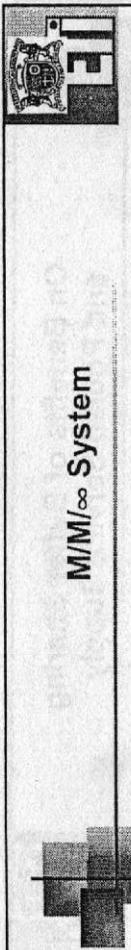
- famous Erlang B formula

- magic: holds for *any* request size distribution i.e., for M/G/S/S (!)

- online calculators exist ([www.voip-calculator.com/calculator/](http://www.voip-calculator.com/calculator/))

J. Kononov, Operational Research/Queueing systems

99



A huge supermarket admits on average 20 customers per minute,  
 each stays inside for a total of 15 minutes on average (including shopping and  
 queuing at the checkout).

Find the distribution of current customer population in the supermarket, assuming  
 a Markovian system model.

$a_m = 3$  s,  $b_m/v = 900$  s,  $\rho = 300$  erlangs

$N_m = \rho = 300$  (this we know from Little's theorem).

Special case of M/M/S/S with  $S \rightarrow \infty$ , therefore  $p_0 = \frac{1}{\sum_{k=0}^{\infty} \frac{\rho^k}{k!}} = e^{-\rho}$

$$p_k = e^{-\rho} \frac{\rho^k}{k!} - \text{Poisson distribution (!)}$$

J. Kononov, Operational Research/Queueing systems

100



















































