

Operational Research

Queuing Systems 1: Description and Operation

1

Jerzy Konorski
Room 139 (old bldg)
office hours: tba
jekon@eti.pg.gda.pl

J. Konorski, Operational Research/Queuing Systems



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

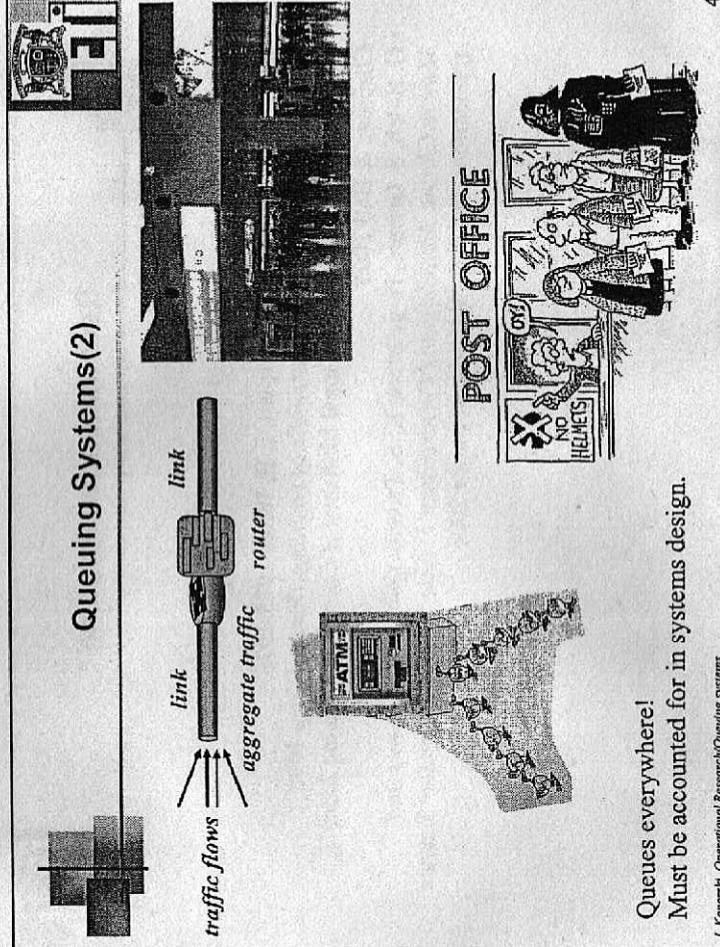
J. Konorski, Operational Research/Queuing Systems



- 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. Kowalski, Operational Research/Queuing systems

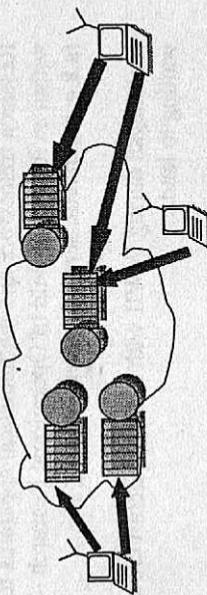
3



4

Queuing Systems(3)

- system has resources – limited, reusable
- 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. Konoriti, Operational Research/Queuing Systems

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. Konoriti, Operational Research/Queuing Systems

6

Queuing Theory: Mission (2)

Wikipedia

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),
if you
- buffer and queue management burden,

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. Konorski, 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* *process*
-
- The diagram illustrates the simplest queuing model. It features a horizontal timeline with an arrow pointing to the right labeled "time". On the left, a vertical bar labeled "arrival stream" has several short vertical tick marks representing individual arrivals. To the right of the arrival stream is a small square labeled "source". An arrow points from the source to a rectangular box labeled "infinite buffer" which contains a hatched pattern. From the right side of the buffer, an arrow points to a circle labeled "processor". Above the processor, an arrow points upwards and to the right, labeled "request generation".

Characteristics of service process depend
on those of arrival stream and the way
buffer & processor system operates.
How?

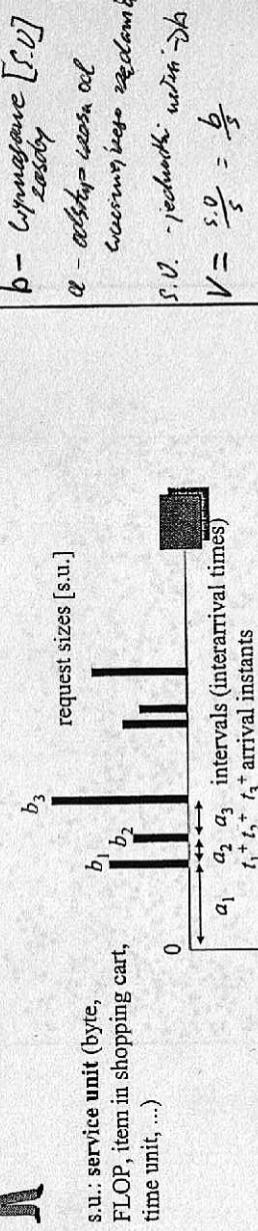
This is what queuing theory is about

J. Konorski, Operational Research/Queuing systems

8

Simplest Model (2)

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



J. Kemerich, Operational Research/Queueing systems

Applique 2 process - spätesten Fristen!
Ces

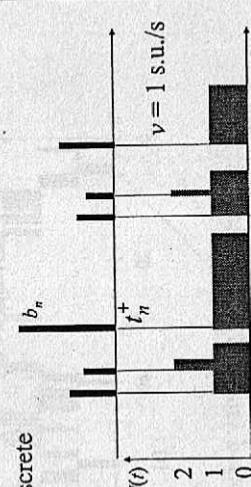
Simplest Model (3)

Remarks:

- $t_n^+ = a_1 + \dots + a_n$ caus potroby na obstaranie byt - spätesten Fristen
- $b_n/v = t_n^+$ caus potroby na obstaranie byt - spätesten Fristen
- $d_n = t_n^- - t_n^+$ = system delay of request n - čas odvodu počas zpracovania požiadavky
- $w_n = d_n - b_n/v$ = waiting (buffering) delay of request n (wasted time)
- $N(t) = \#\{n \mid t_n^+ \leq t \leq t_n^-\}$ = number of requests in system at time t

Lisia potroby
systeme
je v obľúbenom
+ je čo sa teraz
obstaráme

b_n



The continuous process ($N(t)$) and the discrete process (w_n) are determined by:

- processor speed v ,

- queuing discipline (order of service)

$$\text{FIFO: } t_n^+ = \max\{t_{n-1}^+, t_{n-1}^-\} + b_n v$$

(likely my upskaler)

ces výška ces obťažia v CPU

pozemky

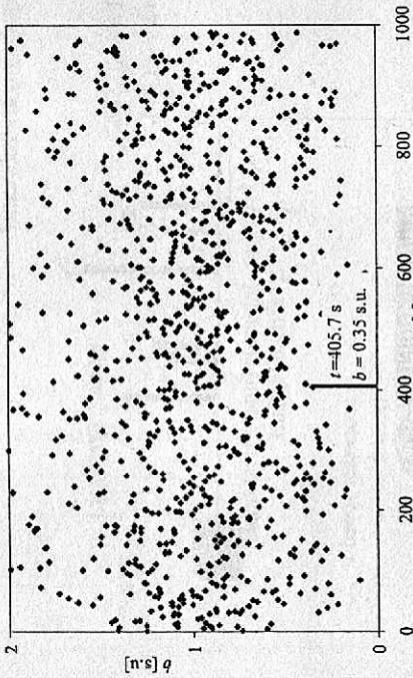
libere

ces požiadavka

libere

ces požiadavka

Arrival Stream and Service Process: Example

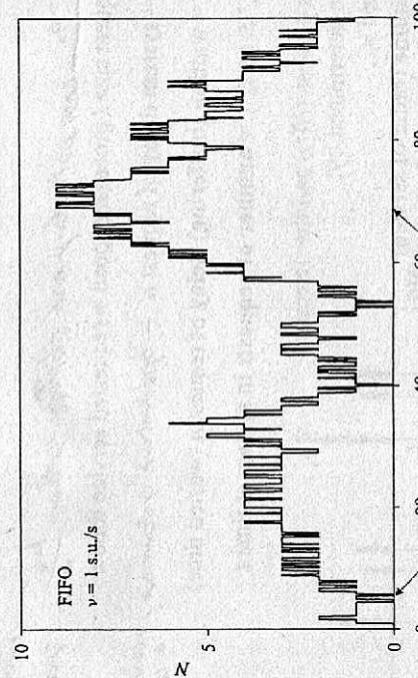


mean interval $a_m = 1$ s
mean request size $b_m = 1$ s.u.

J. Kemerita, Operational Research/Queueing systems

11

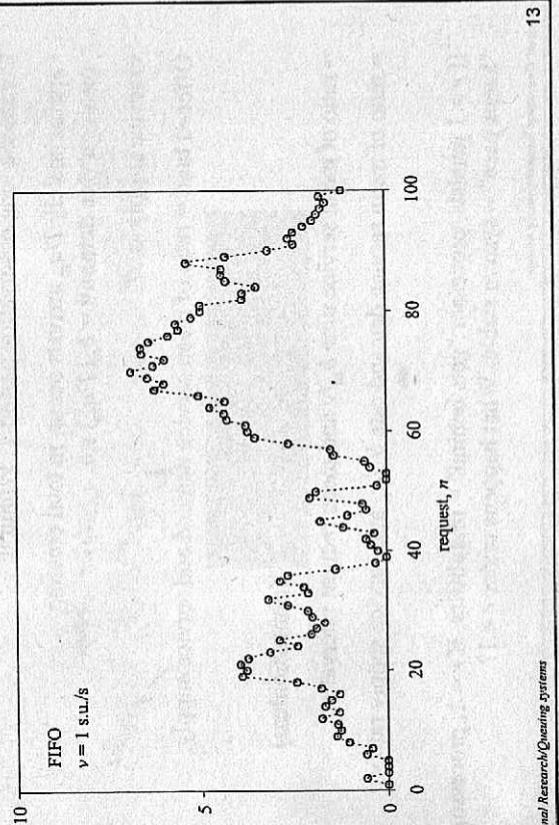
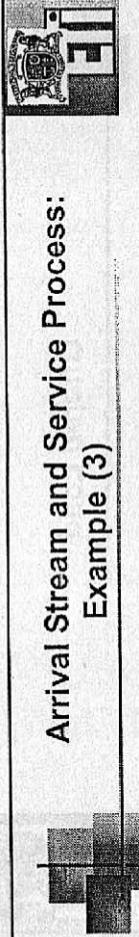
Arrival Stream and Service Process: Example (2)



J. Kemerita, Operational Research/Queueing systems

12

Arrival Stream and Service Process: Example (3)



J. Kemeristi, Operational Research/Queueing systems

Properties of "Good" Service Process

and has sufficient capacity = processor is idle

- from requests' viewpoint:
small waiting delays, rare buffer overflows
- from system operator's viewpoint:
high processor utilization (rare idle periods)

These are contradictory! Rare idle periods imply:

- occurrences of queuing
- long queues becoming prevalent → *queue overflow* / *avalanche of buffer overflows*
- systematic queue growth (instability) / avalanche of buffer overflows
- processor is "getting behind" in the service → unable to work in real time!

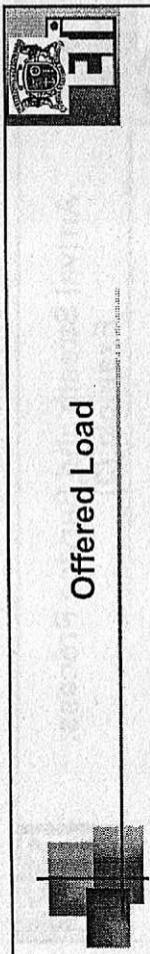
Relationship between arrival stream characteristics and processor speed
determines an important parameter, offered load



- What information about the system operation does it give?
- What offered load causes processor to "get behind"?

J. Kemeristi, Operational Research/Queueing systems

14



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. *second order*
- service supply is vT *first order*

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 ,

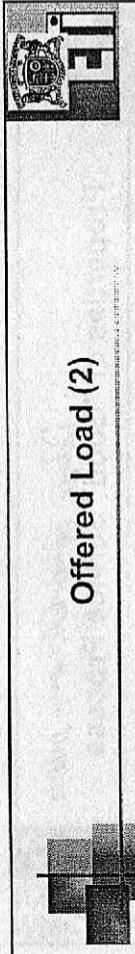
= ratio of mean service demand rate b_m/a_m and service supply rate v

$r > 1$ - *overload*
 $r < 1$ - *underload*

Waking site

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

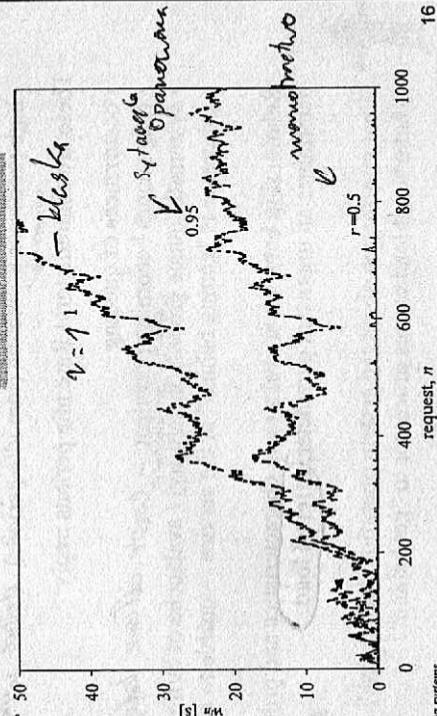
J. Koenraadt, Operational Research/Queueing systems



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



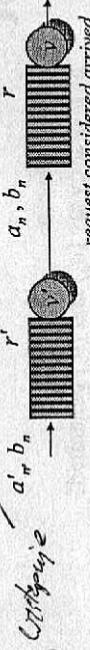
Instability at $r = 1$



J. Koenraadt, 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.



Request is transferred 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 $r = (b_m/a_m)/v$.

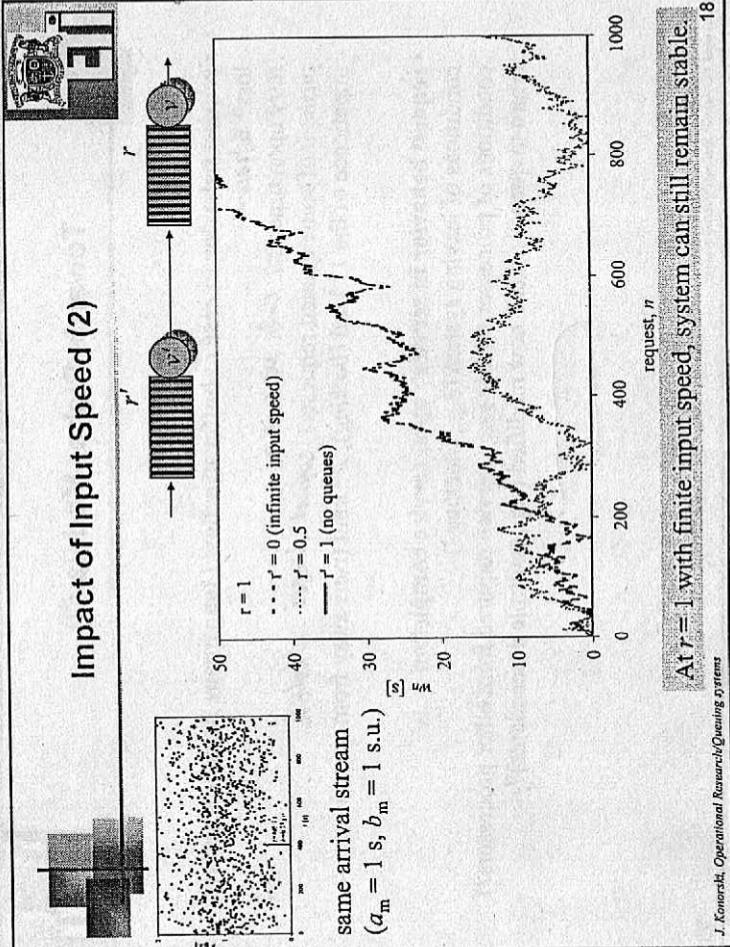
Clearly, $a_m = a'_m$, so $r' = (v/v)r$.

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

J. Konaraki, Operational Research/Queuing systems

17

T = 1



18

At $r = 1$ with finite input speed, system can still remain 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 \mathcal{Q}
limited accessibility? when is a request rejected – drop-tail, ...?

J. Kemeraci, Operational Research/Queuing Systems

19

• *polimorfne
vzorce*

• *polimorfne
vzorce*



Towards Richer Models (2)

- request behavior in case of buffer overflow / long queue found on arrival
loss? *tracery*
retry upon timeout? *with probability*
pushout of some queued requests? *expedit losses* & *balk*
impatience of the 1st kind (*balking*), 2nd kind (runs away from queue)

- 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)?
returning processes
multiple descendants

J. Kemeraci, Operational Research/Queuing Systems

20

Towards Richer Models (3)

- queuing discipline – what decides the order of service

arrival instant (FIFO, LIFO)?
pure chance (RANDOM)? *not me faire faire, nemah liet losen*

request size / current service advancement (SJF, LASF)?

predefined order (RR)? *Leanda*

special requirements of requests e.g., deadlines?

request classification for priority, fairness (HOL, WFQ)?

- service mode

work-conserving (busy period lasts exactly $\sum b_i/v_i$)

or non-conserving (processor breakdowns / "vacations", background arrivals)

bulk service, service preemption with abandonment / rollback?

time or processor sharing (requests served one by one or switched between?)

processor-bound or parallel service?

multiple processors – availability? *grading* (specific processors demanded)?

Can a universal queuing systems simulator ever be written?



J. Kurose, Operational Research/Queuing systems

unigames - IT do GPU nutzlos für den system

nutzlos für den system

Steady State

*lossen
ale lati som
ale lati som*

In a stationary queuing system: *generat losen, ale o gen set mit zonnen*

*{ mechanism governing request generation remains invariable in time zuweile w kum son spott alle
(characteristics of the arrival stream do not change)
processor speed v remains constant } *zyklische* stata*

A well designed stationary queuing system (where $r < 1$) tends to steady state
(characteristics of the queuing process exhibit asymptotic behavior
as the observation period lengthens).

In particular,

lost zonnen = 0

$\{t \leq T \mid N(t) = k\} \rightarrow p_k$ as $T \rightarrow \infty$

*verdutzt v - \rightarrow 2fachen
geno - in system*

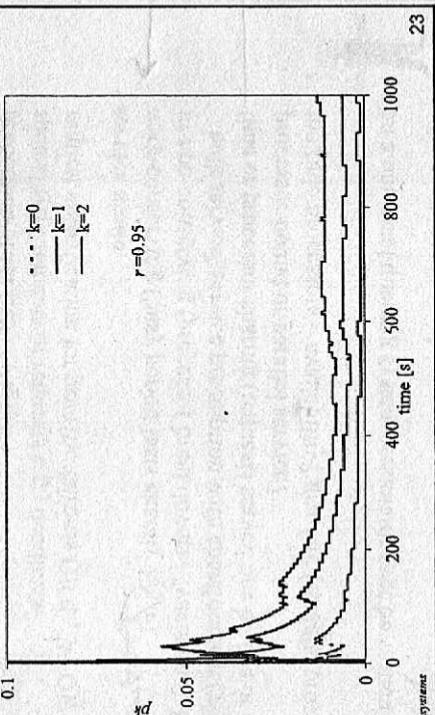
J. Kurose, Operational Research/Queuing systems

Ston woh

*System ohne Zyklenzyklus \rightarrow my stone zebulon, ore 20 another stone obelisk,
siehne leck zebulon w system me bricht, a bishwo sie wohlt raus' hab*

Steady State (2)

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

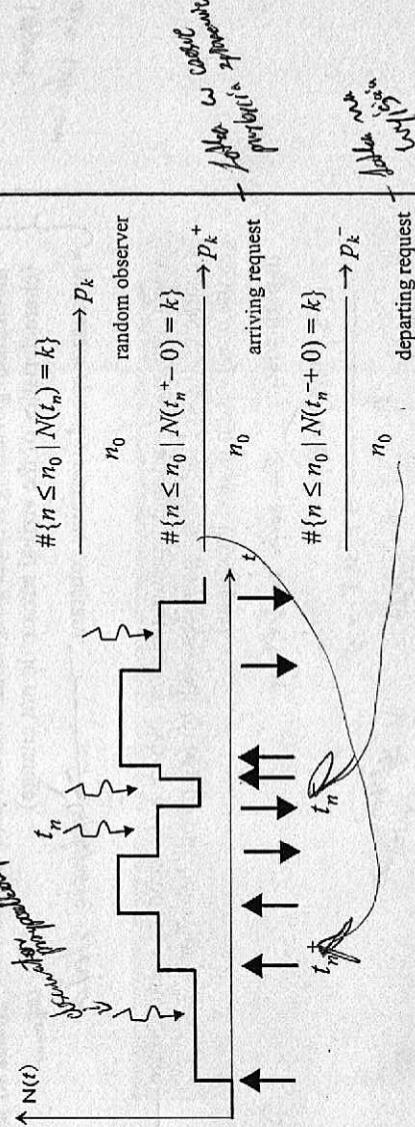


J. Komornik, Operational Research/Queuing systems

Steady State (3)

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

As $n_0 \rightarrow \infty$:



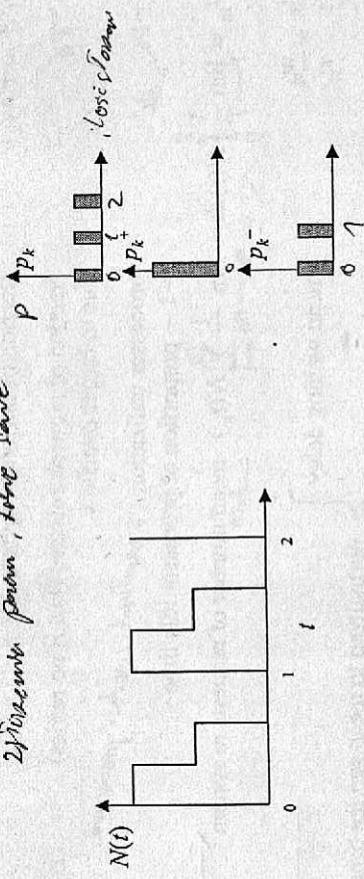
J. Komornik, Operational Research/Queuing systems

24

Steady State (4)

"Practitioners", beware!

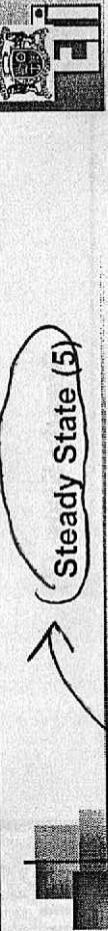
2) *Wear out, have same*



J. Konarak, Operational Research/Queueing systems

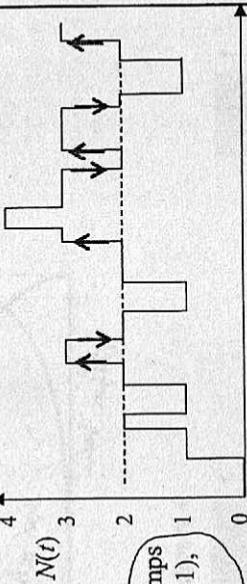
25

alone repository



Unit jump argument: if $N(t)$ only changes by ± 1 then $p_k^+ = p_k^-$

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



Since $q_1 < 1$ rest < 1
in dashed process
 $q_1 = \frac{1}{1-p_0}$
 $(< \infty \text{ } p \neq 1)$

in dashed process

Remark: Only accepted requests count as arrivals. → *the queue is never full*

If $Q < \infty$ (some requests can be rejected due to buffer overflow) then

$$\frac{p_k}{p_0} = p_k^-, k = 0, \dots, Q-1$$

$1 - p_0^+ = \frac{p}{b}$

J. Konarak, Operational Research/Queueing systems

With more than one server process → to link sit into to project
No for $Q < \infty$ (otherwise too much delay) → on the other side

$$\frac{p}{b} \rightarrow \frac{p}{b+T}$$

26

Steady State (6)

Relevant evaluation criteria:

$\lambda \tau_p < 0$ - number of lost requests

$$p_1^+ + \dots + p_{Q-1}^+$$

fraction of buffered requests

$L = p_0^+$ due to buffer overflow

$$1 - p_0^{\text{idle}}$$

processor utilization \rightarrow proportion of processor idle time

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

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

mean waiting delay } normalized to mean service time

mean system delay }

27

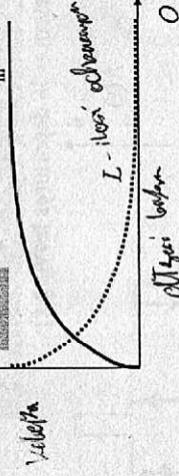
J. Konsztadt, Operational Research/Queueing systems

do drop ... or include the pending notifications
 $\frac{20 \text{ min}}{10 \text{ min}} \dots 0 \text{ or } 10 \text{ min}$

NO

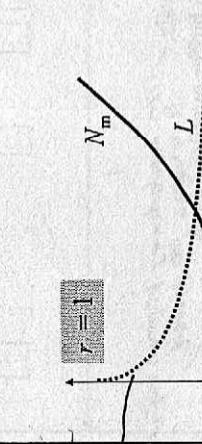
Steady State (7)

$r < 1$ N_m L L_{max}



more than
one request
in queue

+ systematic upward drift in time



less than
one request
in queue

28

J. Konsztadt, Operational Research/Queueing systems

Churn - busy buffer
Churn day

Is Queuing Theory Losing Momentum?

Growing ν 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) http://en.wikipedia.org/wiki/FLOPS	\$0.5/GFLOPS (2007)	x30m

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 ν ,
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 ∞ . 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. Kemerki, 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
hurly
simple models – carry over to much more realistic ones

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

J. Kemerki, 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. Konoracki, 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. Konoracki, Operational Research/Queuing systems

36



Operational Research

Queuing Systems 2: Stochastic Models and Characteristics

Jerzy Konorski

Room 139 (old bldg)

jekon@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_m 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_m) .

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).

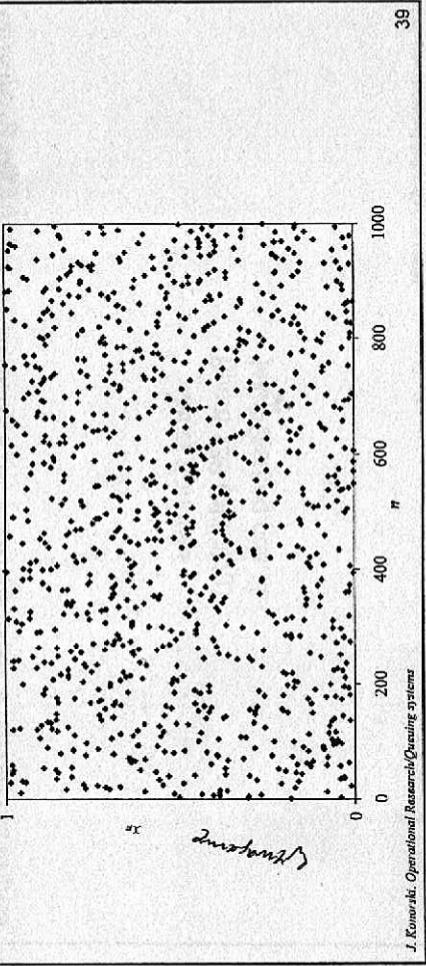
J. Konorski, Operational Research/Queuing systems

38

Random Variables and Stochastic Arrival Streams (2)

An example of a random variable is value returned by the function `r-random` (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:



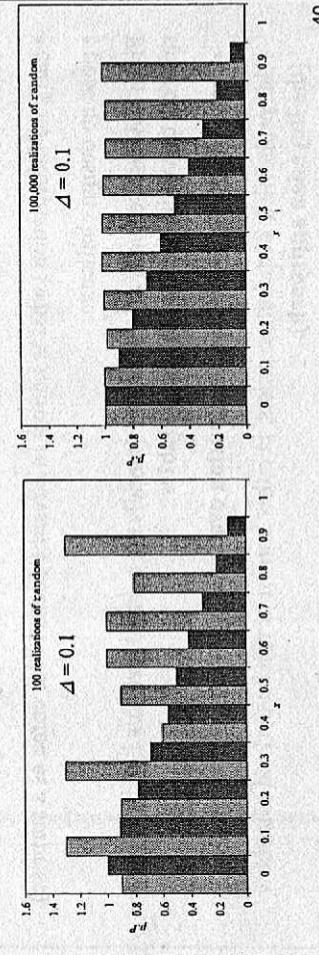
39

Empirical Distributions

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

- divide the range of possible realizations into bins of length Δ ,
- 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 Δ and height $p_i = (k_i/N)/\Delta$.

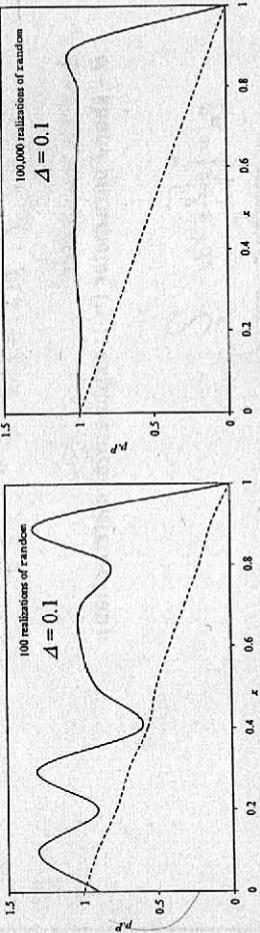
Cumulative histogram constructed analogously, with bars of height $C_i = \sum_{j \geq i} p_j$.



40

Empirical Distributions (2)

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



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

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

Olafurri skellefse

J. Konarik, Operational Research/Queueing systems

41

Theoretical Distributions

Probability density function and complementary distribution function

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

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

$$\text{For any } x' \text{ and } x'', \quad \Pr[x' \leq X < x''] = \int_{x'}^{x''} P(x) dx.$$

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

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

J. Konarik, Operational Research/Queueing systems

42

2

Theoretical Distributions (2)

Modeling for engineering applications often uses Weibull distribution:

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

λ – scale parameter,
 θ – 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/Queueing systems

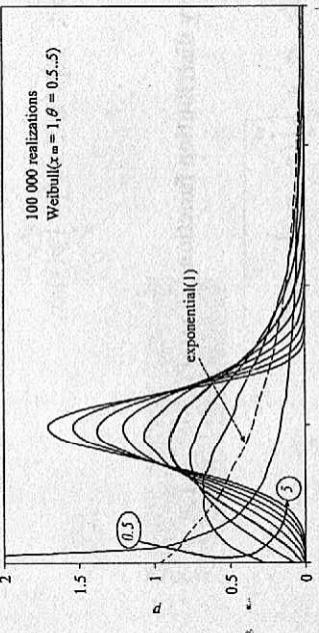
43

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).



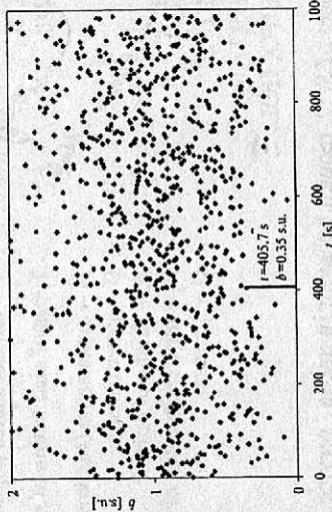
J. Konečný, Operational Research/Queueing systems

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 $A, B \sim \text{Weibull}(1, 2.5)$

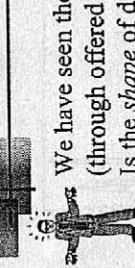
J. Kowarski, Operational Research/Queuing systems

nic send v mo emeemce, bo beezet no strucency

Moes belas
zales k' oel
kinstotka reabilita
reduem
intendue →



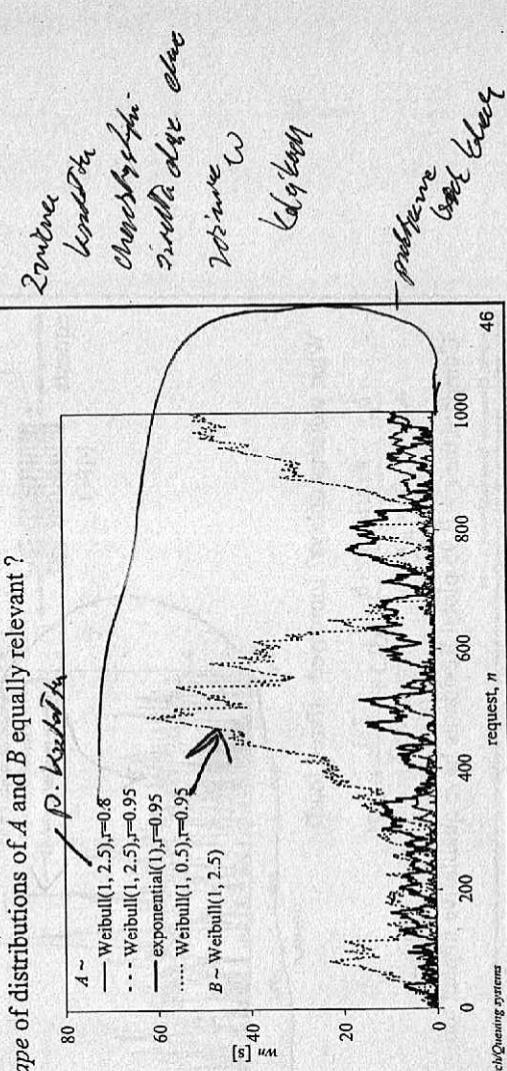
Impact of Distribution of A and B



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

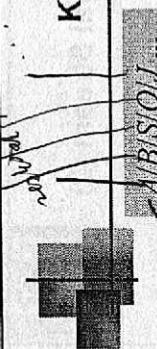
audua siroebud jest
merystresipile



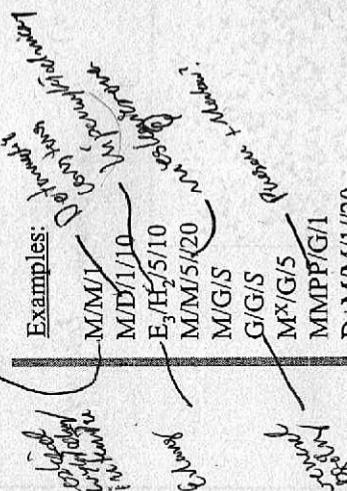
J. Kowarski, Operational Research/Queuing systems

very much so!

لینک
لینک
لینک



Kendall Notation



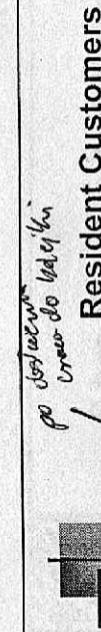
Examples:

- 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 - hiperexponential of order k
- G - general

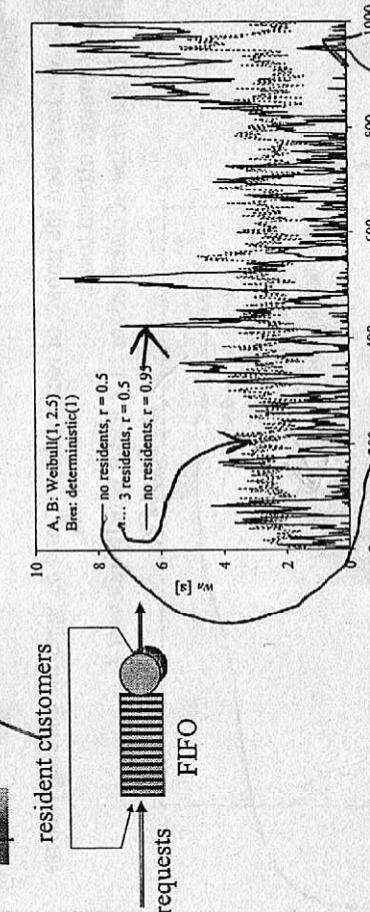
What is the impact of S, Q , request behavior in queue / service, service rules?

J. Kowalski, Operational Research/Queueing Systems

47



Resident Customers



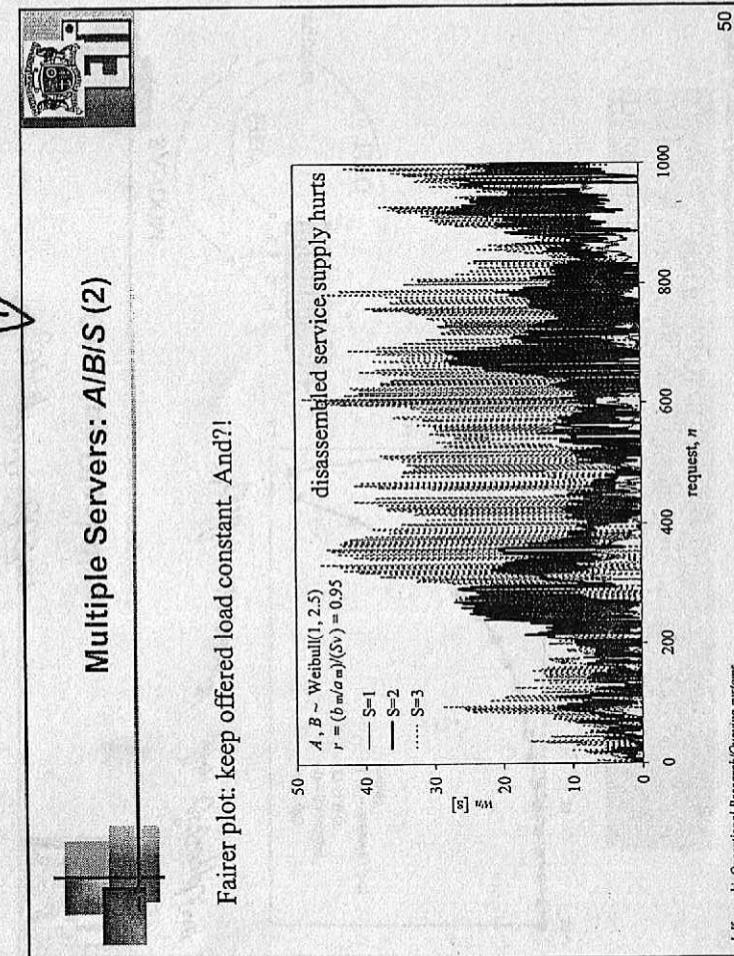
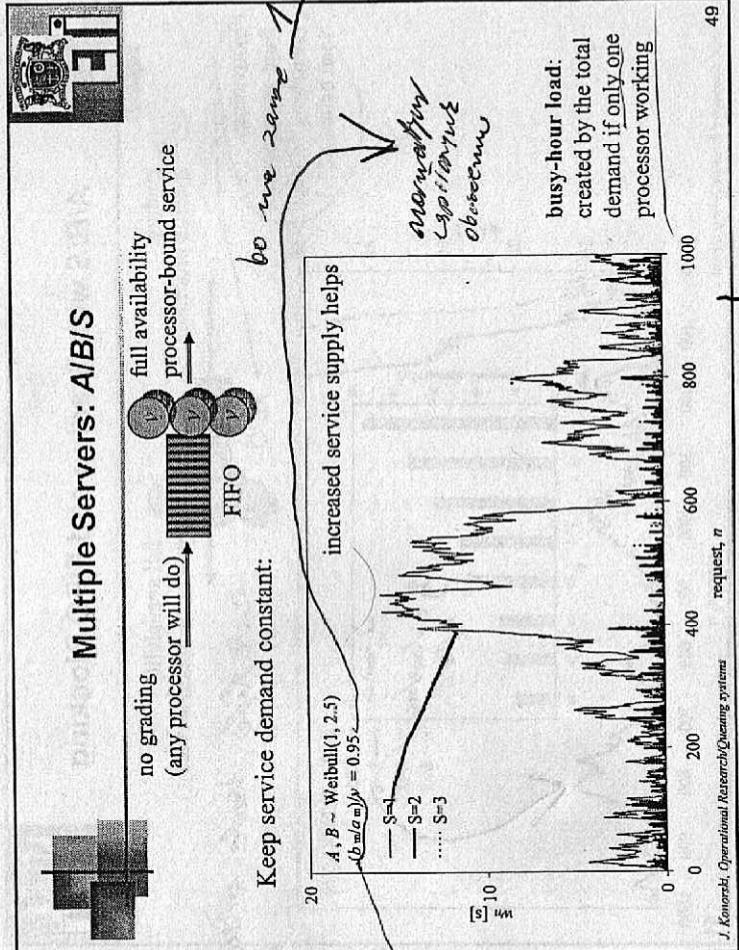
What worsens delays? Increased offered load?

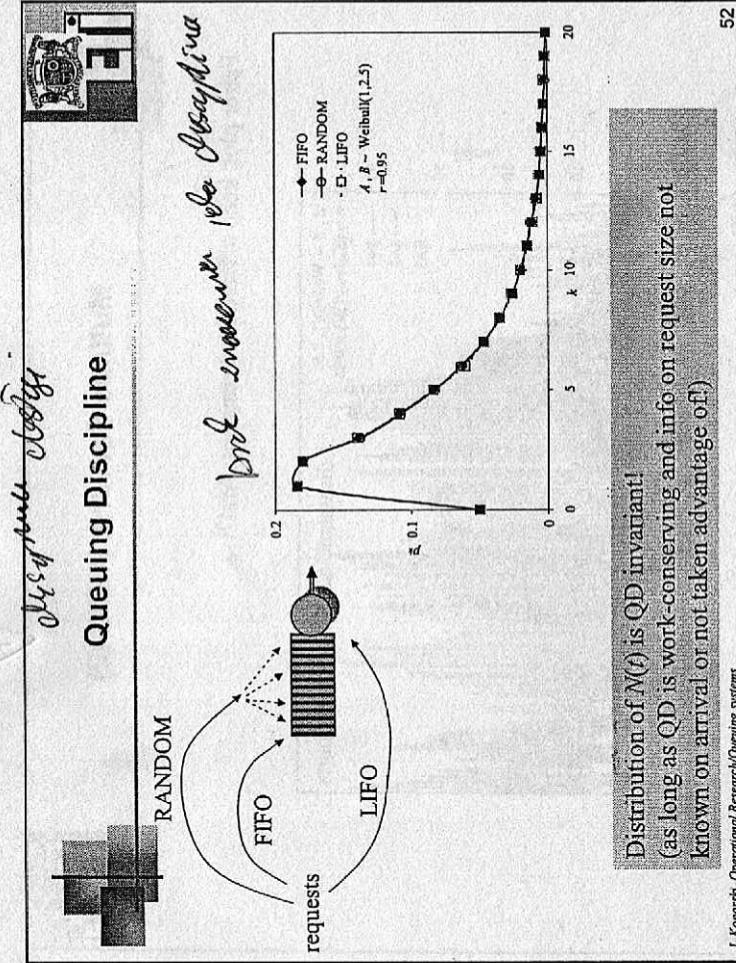
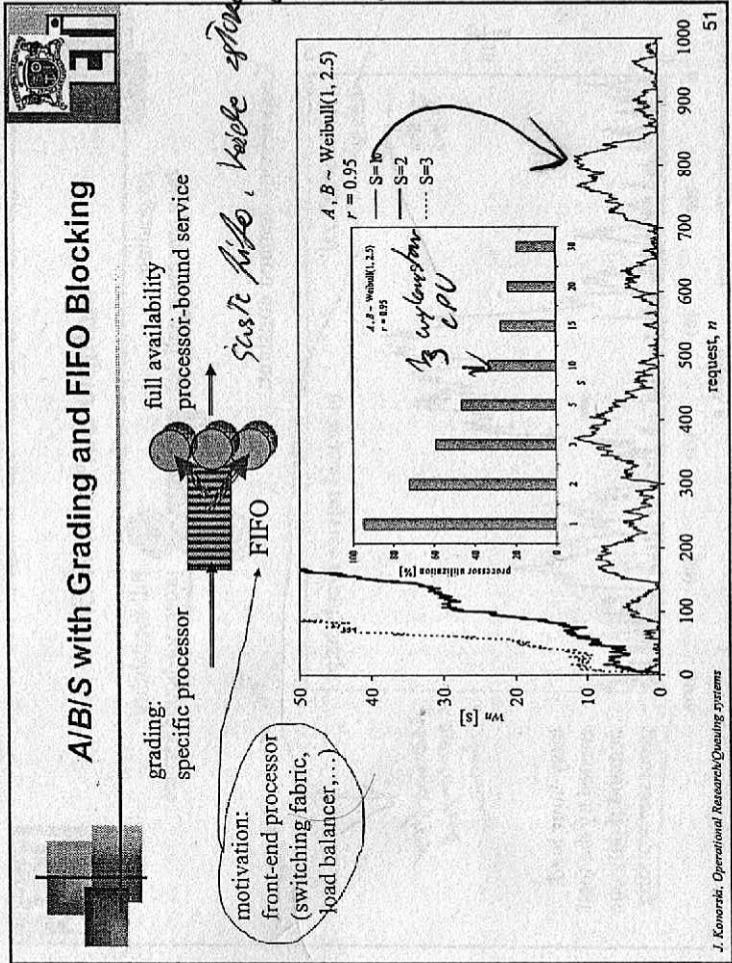
$$\frac{b_m}{a_m v} + \frac{\# \text{residents}}{d_m v} = 0.5 + \frac{3 \cdot 1}{3.5 \cdot 2} = 0.95$$

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

J. Kowalski, Operational Research/Queueing Systems

48





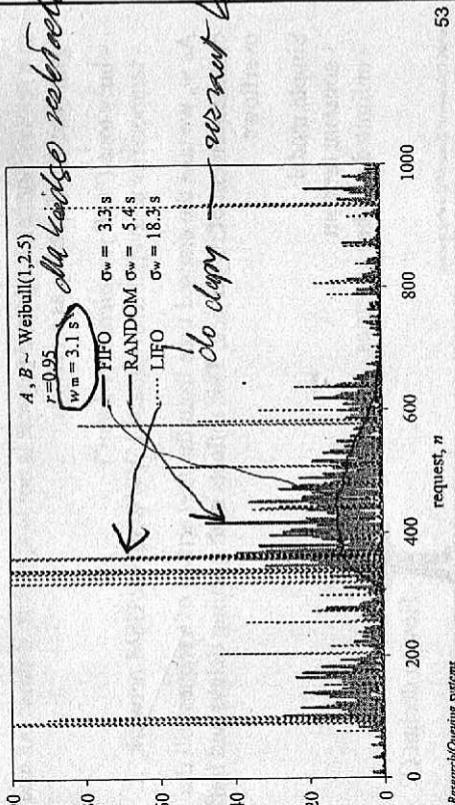
وقت $M(1/\mu)$ از
نیازمندی

Distribution of $M(t)$ is QD invariant!
(as long as QD is work-conserving and info on request size not
known on arrival or not taken advantage of!)

Popkiri چه لذتی داشتند اگر نمی‌توانستند اینجا بخواهند!

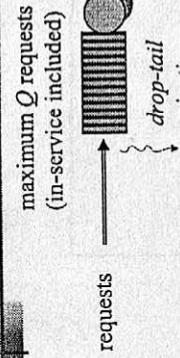
Queuing Discipline (2)

Which QD cause largest w_m ? largest σ_w ?
to allow server to serve
the lowest service rates
for same request thresholding



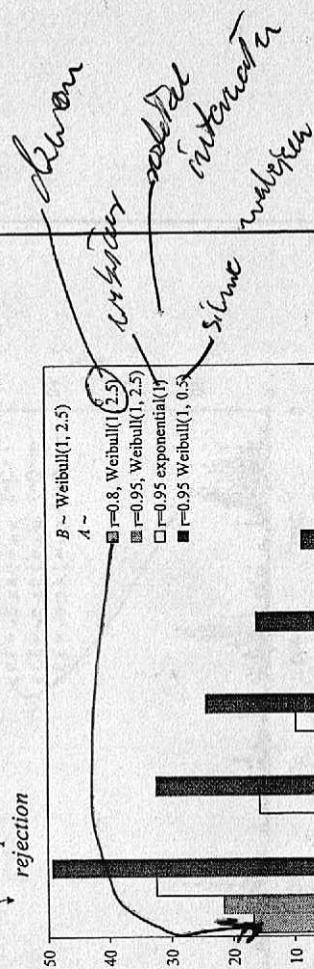
J. Konsakis, Operational Research/Queuing Systems

Finite Buffer Capacity: A/B/1/Q



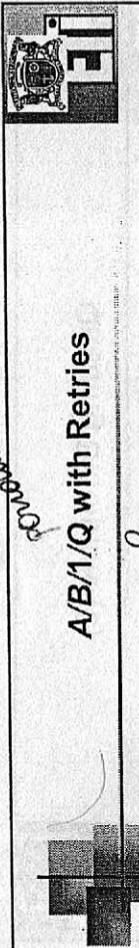
- Motivation:
 - memory cost!
 - delays, management cost!

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



J. Konsakis, Operational Research/Queuing systems

54



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.

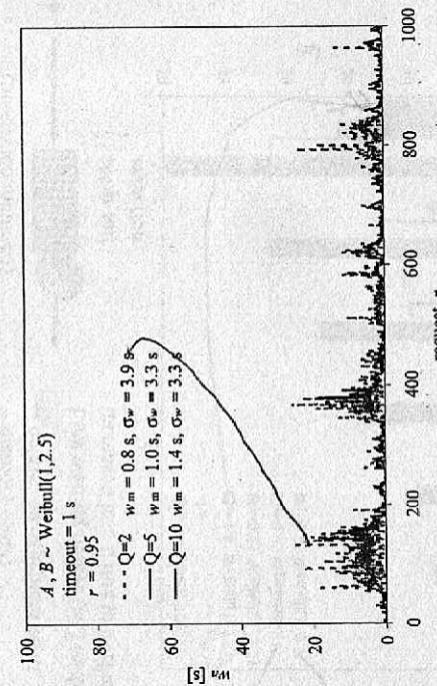
J. Konschik, Operational Research/Queueing systems



55



A/B/1/Q with Retries (2)

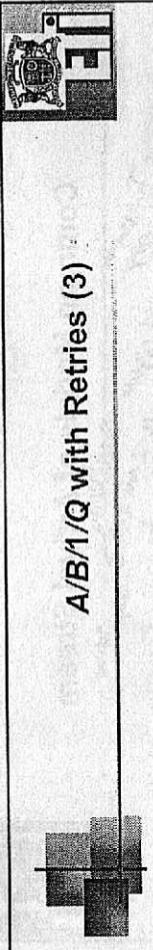


Smaller Q improve mean wait, worsen dispersion and distribution tail.

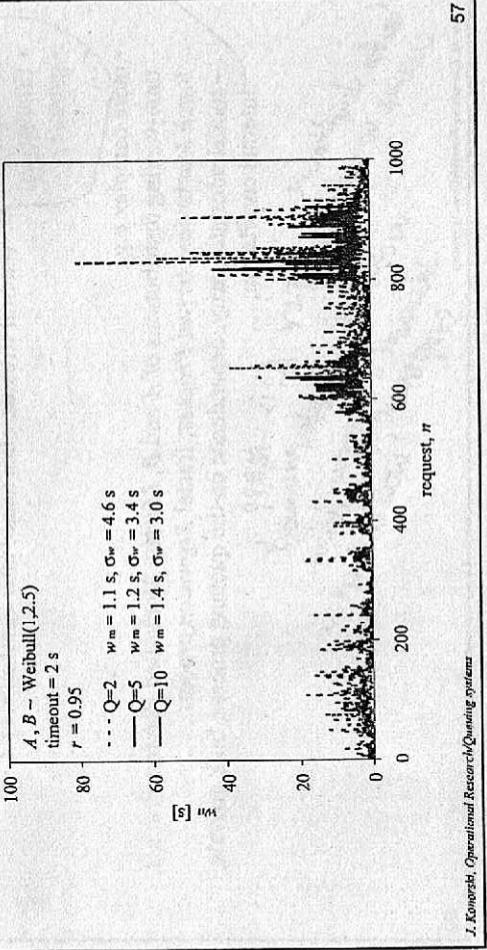
J. Konschik, Operational Research/Queueing systems

56

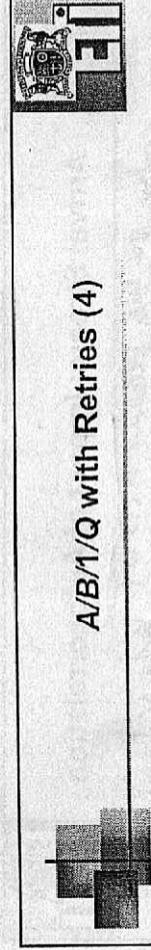
*problem of
request waiting is getting really bad when
Q is too small*



This becomes even more visible with more "patient" requests...

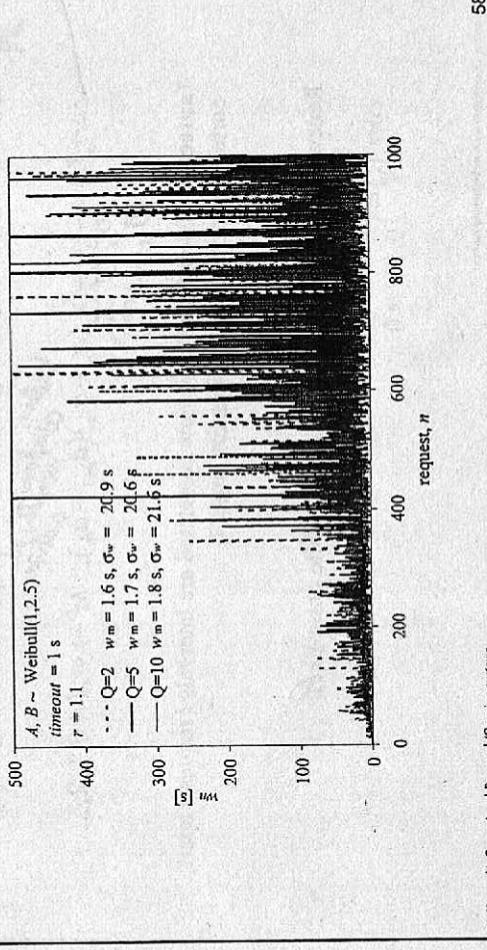


57



... whereas at heavier offered load, distribution tails become dramatic:

core becomes increasingly



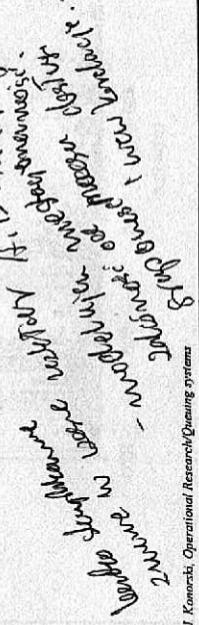
58

Common Models of Arrival Stream



- Bernoulli
- Weibull
- Erlang
- gamma
- ...

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. Kurose, Operational Research/Queuing systems 59

Arrival Stream: Impact of Autocorrelation



U{ Operate per your medium A :S oft land ~ lan
Are distributions of A and B enough to determine the queuing process (given fixed ν), or do we need information on the internal correlation in (a_n) ?

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

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

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

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

J. Kurose, Operational Research/Queuing systems

60



Arrival Stream: Impact of Autocorrelation (2)

Experiment 1

variable arrival distribution

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 difference?

J. Koenraadt, Operational Research/Queuing systems

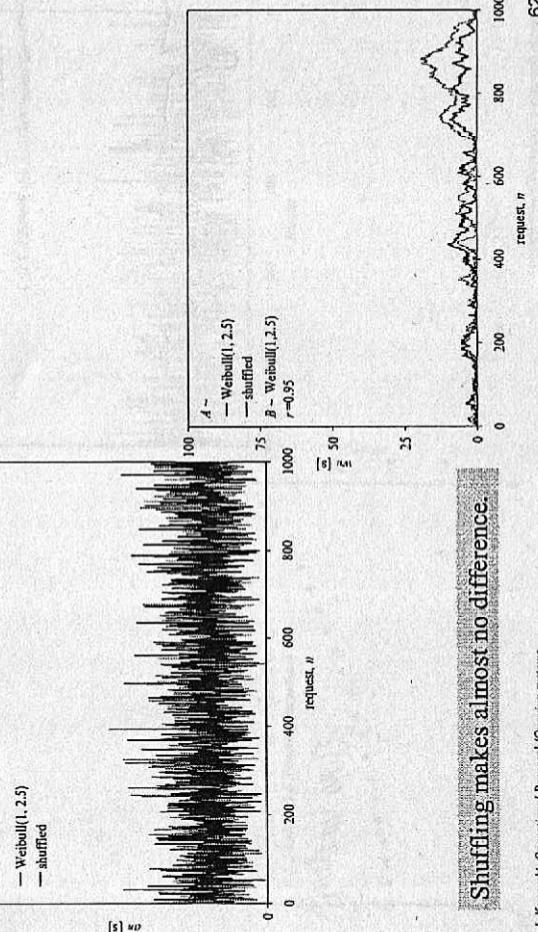
61

done by
reorder
Shuffling
theory

power
log profile
we be
walking
reinvy



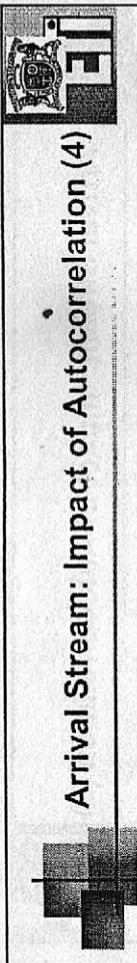
Arrival Stream: Impact of Autocorrelation (3)



J. Koenraadt, Operational Research/Queuing systems

62

Shuffling makes almost no difference.



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.

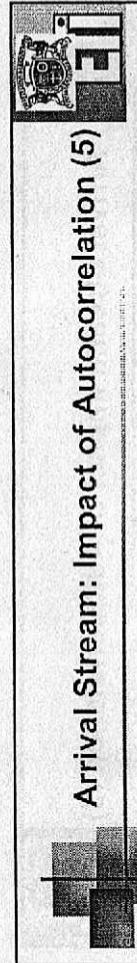
However, the second variant yields (a_n) with long-range autocorrelation – a **self-similar** arrival stream.

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

Yannick & Hervé summ vert à droite

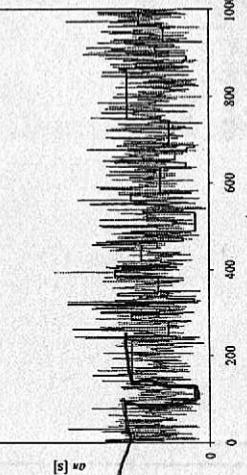
J. Konarski, Operational Research/Queuing Systems

63

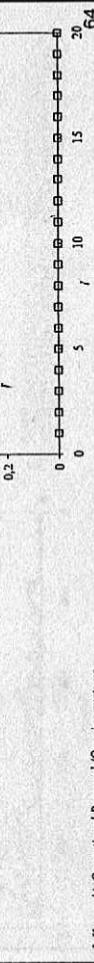


$A \sim$ — Weibull(1, 2.5)

↔ self-similar Weibull(1, 0.5) / weibull(1, 2.5)

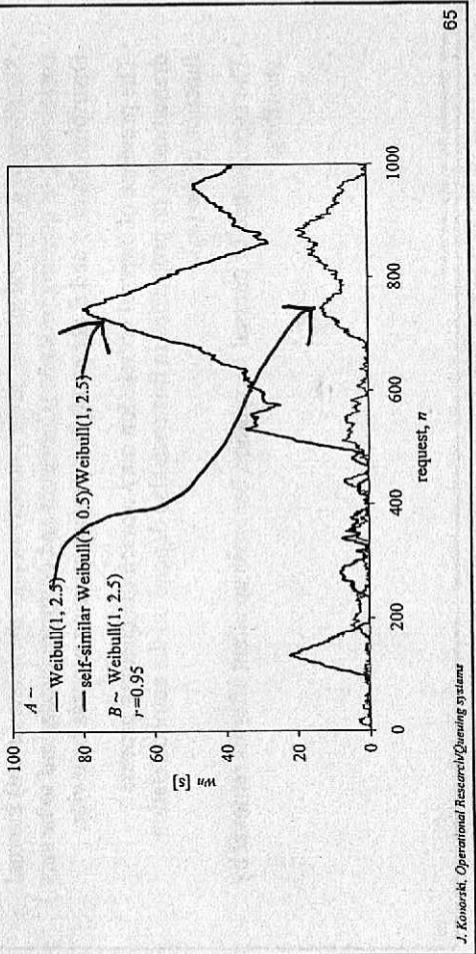


J. KonarSKI, Operational Research/Queuing Systems



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!)



65

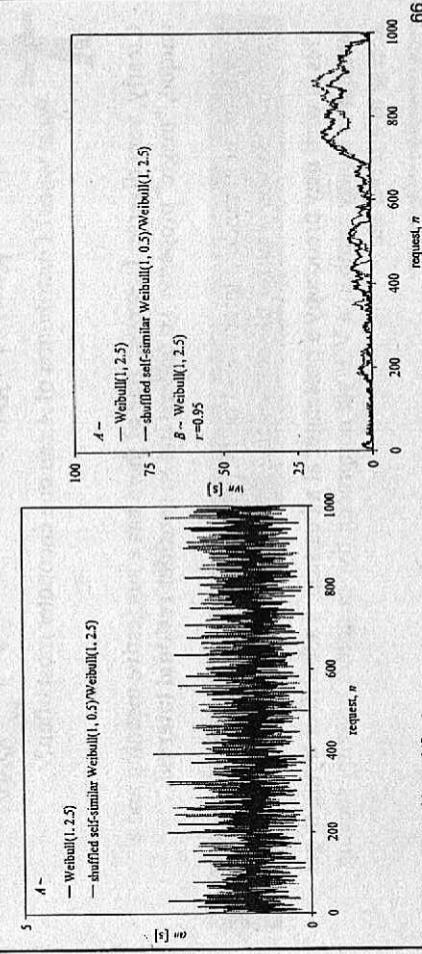
Konarati PST
BARD 202721
Yasamine
Zimmetto
Lovecraft

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 .

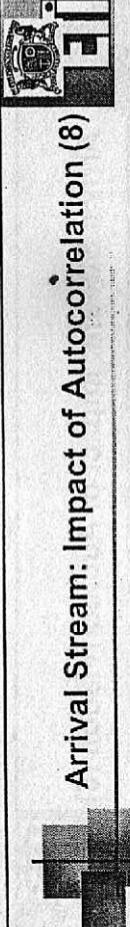
Use the same (b_n) and again compare with the renewal stream variant...



66

Modellieren
mit Web-Services
und Prozessmodellen
bei Verkäufen

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. Konsztadt, Operational Research/Queuing systems

67

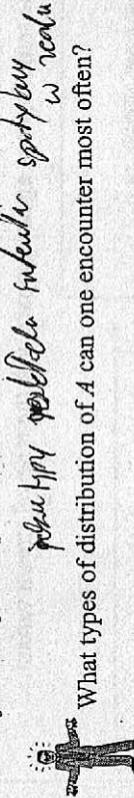
11

U 2 methods A
and B
do not have
the same
internal
correlation



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.



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.

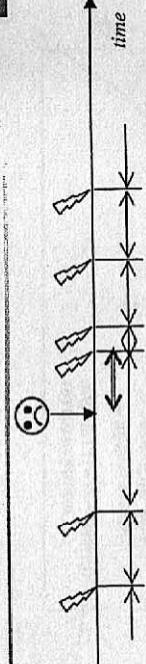
J. Konsztadt, Operational Research/Queuing systems

68

34

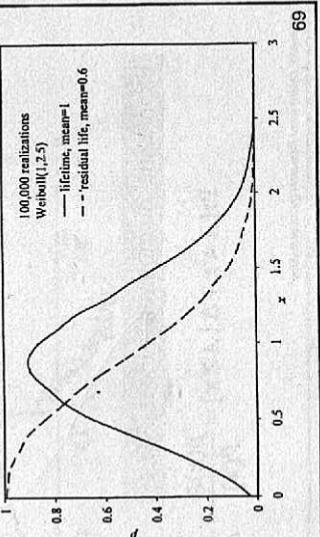


Renewal Arrival Stream: Residual Interval (2)



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

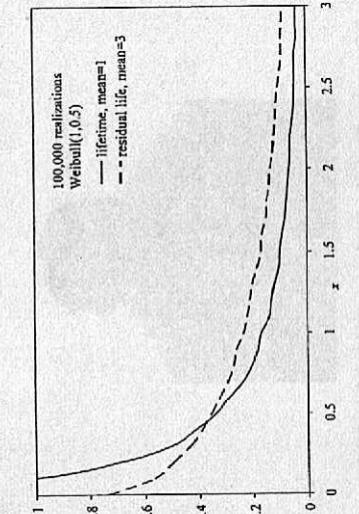
Given the distribution of inter-bus interval (lifetime), what distribution does the residual interval have? Shifted towards smaller realizations?



J. Konarak, Operational Research/Queuing systems



Renewal Arrival Stream: Residual Interval (3)



Paradox of residual life:

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

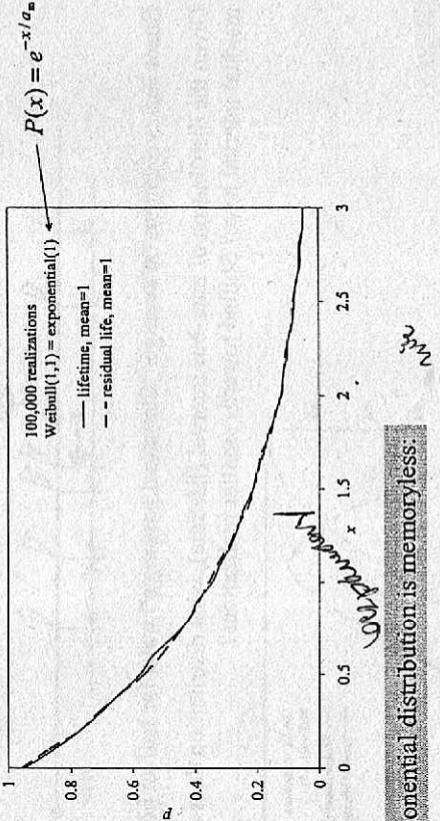
If $A \sim P(x)$ then $\bar{A} \sim \bar{P}(x) = P(x)/a_m$

J. Konarak, Operational Research/Queuing systems

70

Stream of renewal events
with mean =

Renewal Arrival Stream: Residual Interval (4)



J. Konarski, Operational Research/Queueing systems

71

$$\Pr[A < x + \Delta | A \geq x] = \frac{P(x + \Delta)}{P(x)} = \frac{\Delta}{a_m} \quad \text{regardless of } x$$

Stream of renewal events
with mean = 1 second
approximately
many intervals



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. Konarski, Operational Research/Queueing systems

72

Stream of renewal events
with mean = 1 second
approximately
many intervals

ile poisson se d'appelle (souvent) essai !

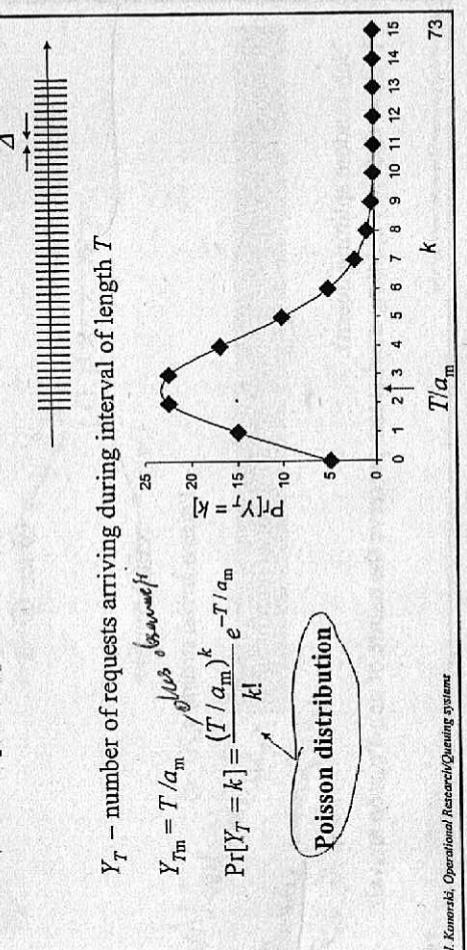
simeon-denis poisson

36

Poisson Arrival Stream (2)

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

bedien darunter

*p 12 wischen
prosekusse (ca 2
stundenrum
nach 2 logogram*

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

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

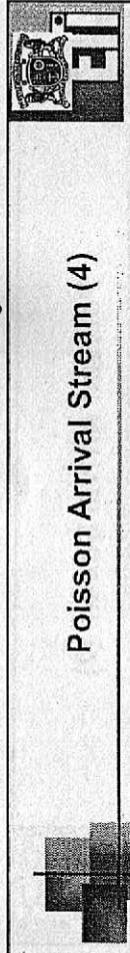
$$\Pr[N(t) = k] \xrightarrow{\text{Poisson}} \frac{\Delta}{a_m} \Pr[N(t) = k]$$

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

J. Kemerer, Operational Research/Queueing systems

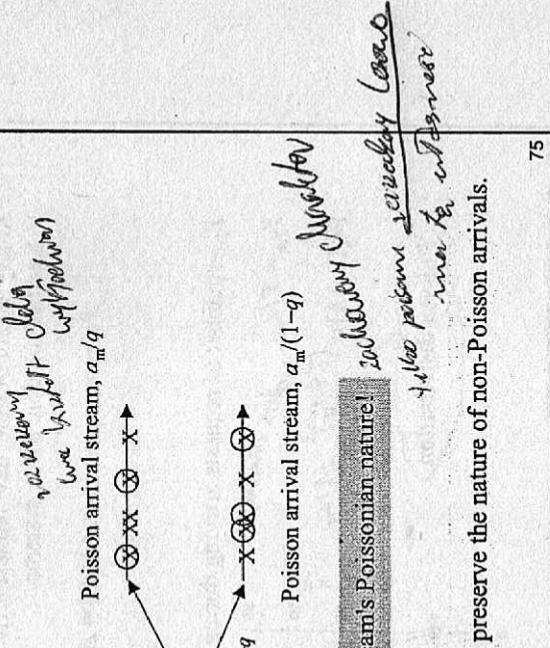
→ look part
i. m. derzeit
→ m. m. derzeit
→ m. m. derzeit

*cup short rest propagandie do
hoch kann es sein w. share
polym*



Poisson Arrival Stream (4)

Random splitting:



J. Konarski, Operational Research/Queueing systems

75



Aggregating streams
which should
not have
overlap

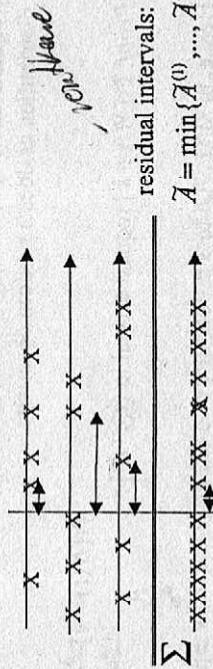
Aggregation of Renewal Arrival Streams

$\frac{T}{a_m} = \sum_{j=1}^J \frac{T}{a_m^{(j)}}$ (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.

1/10 renewal stream



J. Konarski, Operational Research/Queueing systems

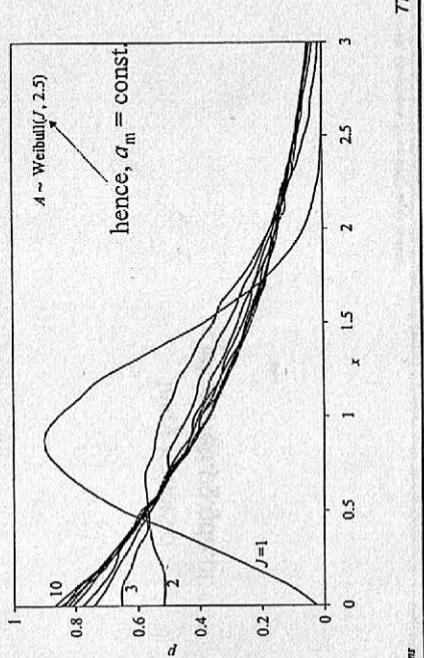
76

Aggregation of Renewal Arrival Streams (2)

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

$\sum_{j=1}^J A^{(j)}$

Aggregated arrival stream is Poisson (Palm theorem).



J. Komorowski, Operational Research/Queueing Systems

Stochastische
metode obserwacji
processu
zincydentów
Dario John
Ugore
 \sum Poisson
im wypadku
dla danego a_j
wysokości 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 $Y_x^{(j)} > 1$.

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

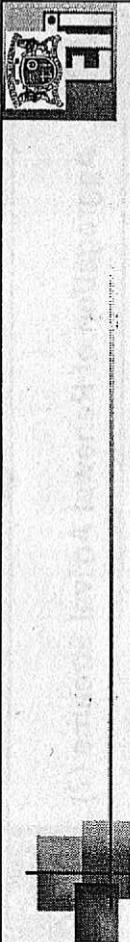
Finally,

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

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

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

J. Komorowski, Operational Research/Queueing Systems



Operational Research

Queuing Systems 3: Markovian Models

Jerzy Konorski
jekon@eti.pg.gda.pl

J. Konorski, 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. Konorski, Operational Research/Queuing systems

80

Markovian Systems (2)

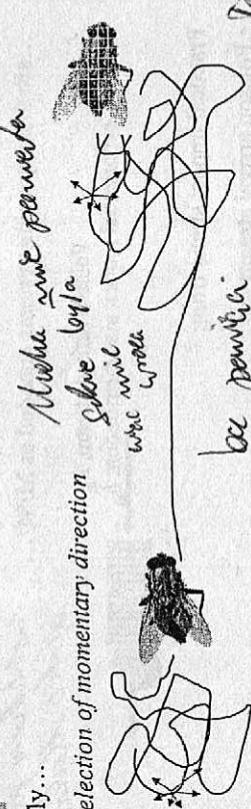
Exhibit the apparently unusual, but most useful **Markov property**.
(which, however, they share with a huge number of real-world dynamical systems – technical, physical, social, biological, economic, ...)

What's the state system given
 $\overline{t + \Delta}$ - change unknown?
 $\text{state}(t) + \text{noise} \xrightarrow{\phi} \text{state}(t + \Delta)$
"noise" (random external input at time t ,
in general dependent on the current state,
but independent of earlier ones)

J. Konečný, Operational Research/Queueing systems

81

Markovian Systems (3)

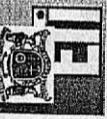
- The fly...


Market share predictor
Share by/a
Share with
etc.
- card shuffling: card order in the deck (*selection of cut point*)
• gambler's capital / population (*current interest / growth rate*)
$$\text{trend}(t + \Delta) = (1 - c) \cdot \text{trend}(t) + c \cdot \phi(t) \quad (\text{current observation})$$

$$\text{market_share}(t + \Delta) = \phi \cdot \text{market_share}(t)[1 - \text{market_share}(t)] \quad (\text{current management performance})$$
- Internet topology (*number and points of attachment of new networks*)

J. Konečný, 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. Koenraadt, Operational Research/Queueing Systems

83



Markovian Systems (5)

...are those queuing systems encoded as M/M/ ∞ ...
Markov chains process

Poisson arrival stream, a_m
exponential request size distribution, b_m



Practical impact stems from:

- Poisson arrival stream
 - Palm theorem
 - random splitting
 - PASTA
 - pessimistic (meaning: fail-safe) performance characteristics — *expressing / else just upholding promises / my property*
- exponential request size distribution: crude approximation of
 - call holding time, Web / P2P file transfer
 - batch processing time
 - ...

J. Koenraadt, Operational Research/Queueing Systems

84

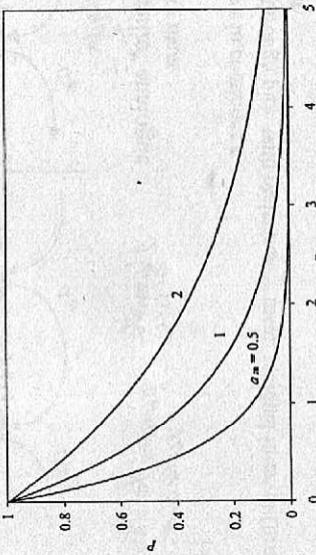
Markovian Systems (6)

ν – (constant) processor speed

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

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

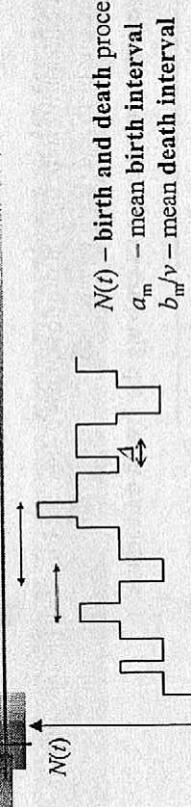
$\Pr[\text{service time} \geq x] = e^{-x/(b_m/\nu)}$



J. Kornai, 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
 \Rightarrow 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

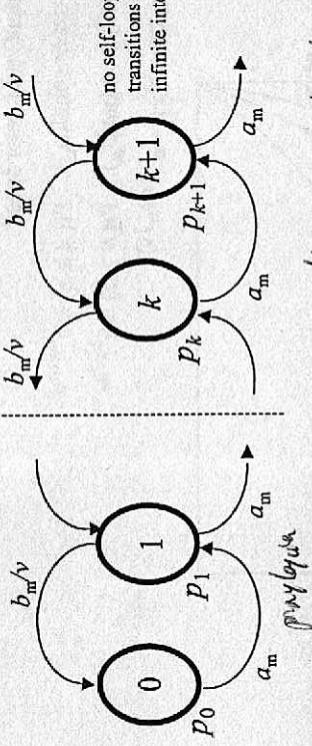
Mivel, ugyen leg elég a gyakorlati interval
 i minden részben megfelelő

43



M/M/1 Queuing System (2)

Hydrodynamic analogy



Suggestive "hydraulic" analogue :

- liquid \sim probability mass
- state \sim container
- p_k \sim liquid pressure in container k
- transition \sim flow through pipe with resistance = mean residual interval (time to occurrence) of respective event

J. Koenig, Operational Research/Queueing systems

87

Hydrodynamic analogy
pressure to flow \propto pressure = fastest release
and instant overflow from tank
closure



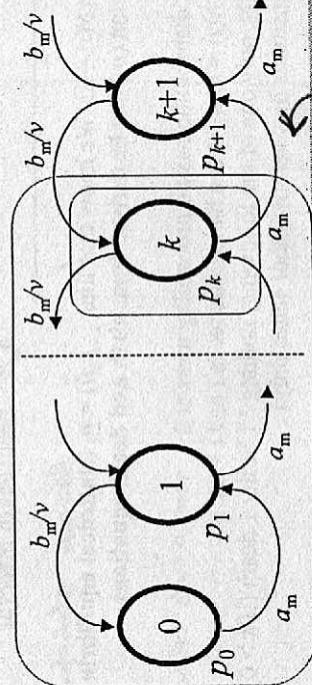
M/M/1 Queuing System (3)

Financial situation looks like water reservoir
get enough. When people must take it away

Water balance
level limited
water inflow
outflow

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!



J. Koenig, Operational Research/Queueing systems

$$\frac{p_0 a_m}{p_{k+1}} = \frac{p_{k+1} (b_m/v)}{a_m}$$

$$\Rightarrow \frac{p_{k+1}}{p_k} = \frac{b_m}{a_m} = r$$

Actual storage \rightarrow post up security ring
 \hookrightarrow system

M/M/1 Queuing System (4)

$$P_k = P_0 \cdot r^k \quad P_0 + P_1 + P_2 + \dots = 1 \Rightarrow P_0 = 1 - r$$

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

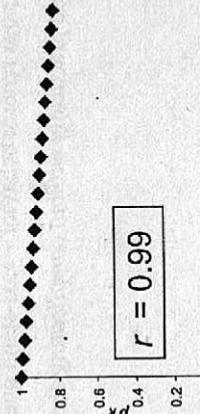
$$N_m = \frac{r}{1 - r}$$

$$d_m = a_m N_m = \tau_m \cdot \frac{1}{1 - r}$$

$$w_m = d_m - \tau_m$$

very suggestive!

The most frequent queue length?!

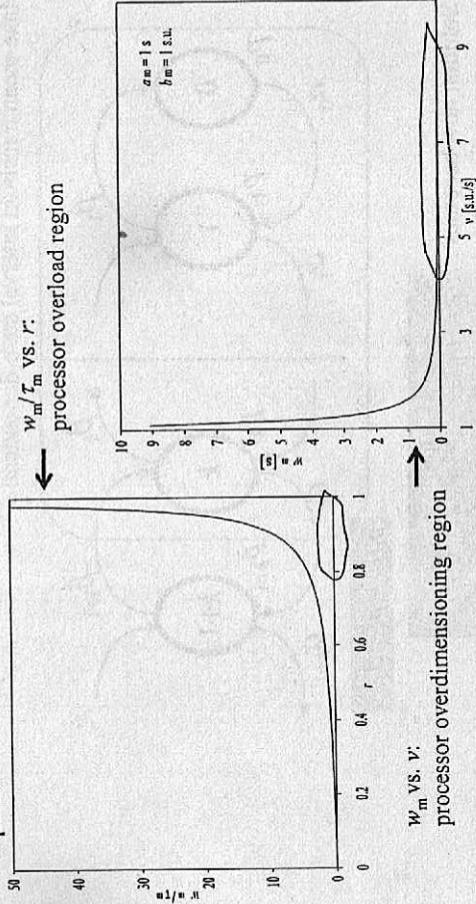


89

J. Koenigs, Operational Research/Queuing systems

M/M/1 Queuing System (5)

w_m / τ_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.

90

J. Koenigs, Operational Research/Queuing systems

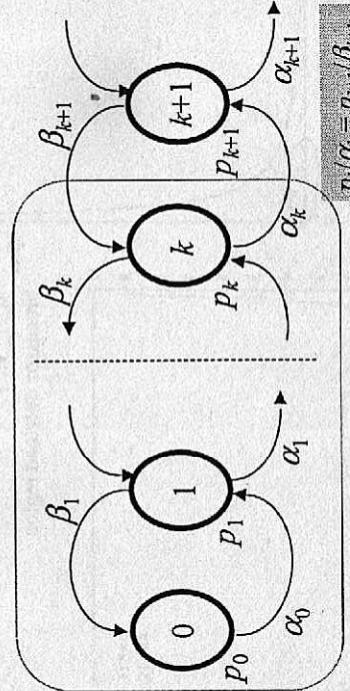
M/M/... Systems

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!*.

M/M/... Systems (2)

Make mean transition interval state-dependent:

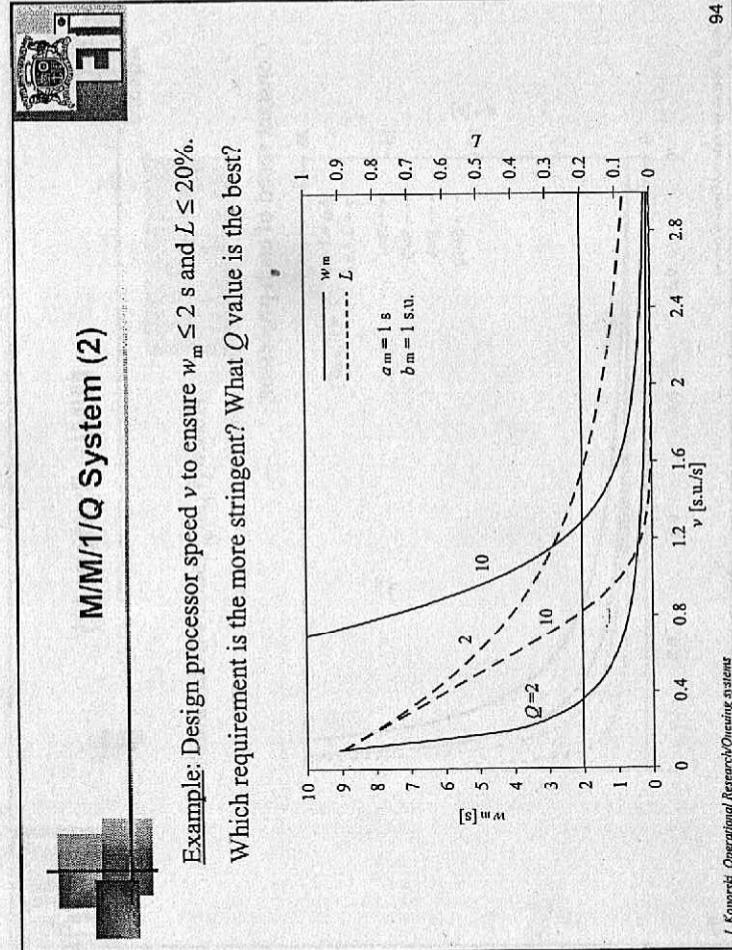
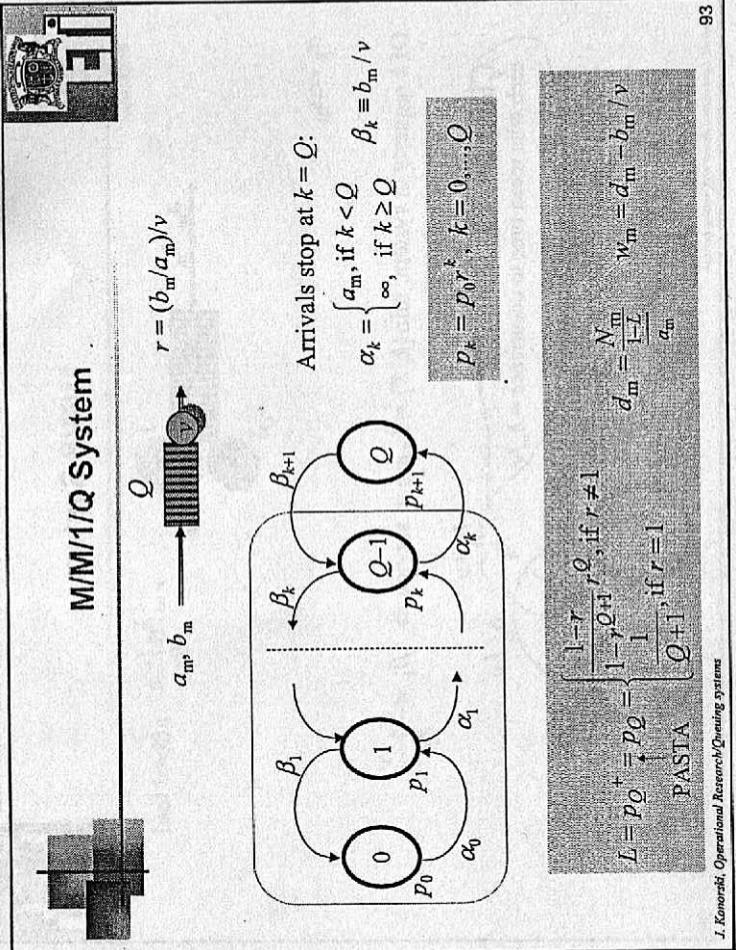


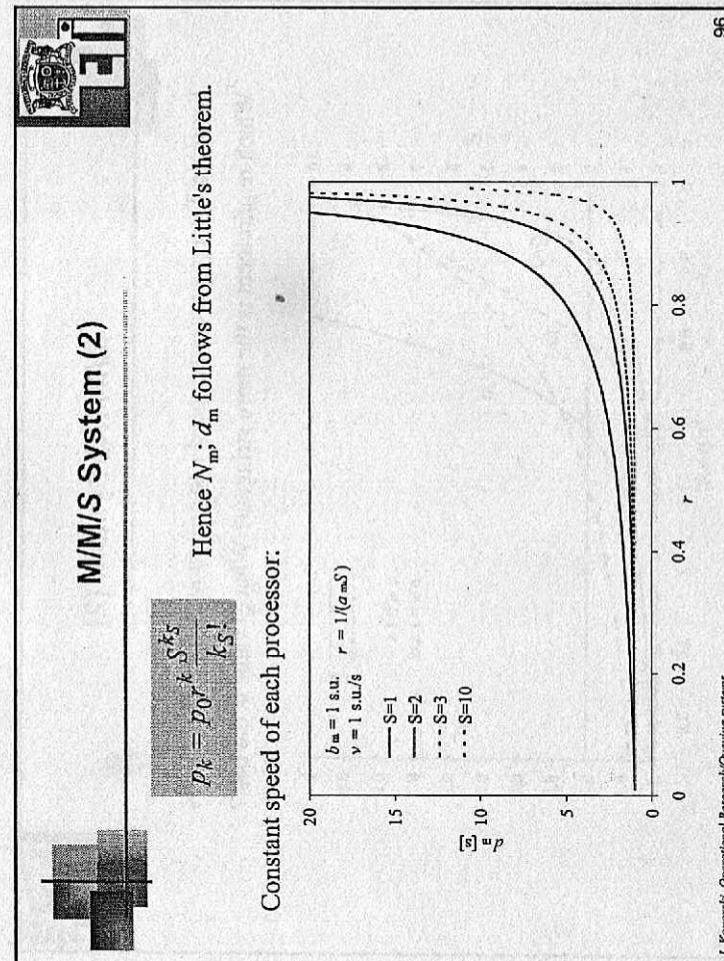
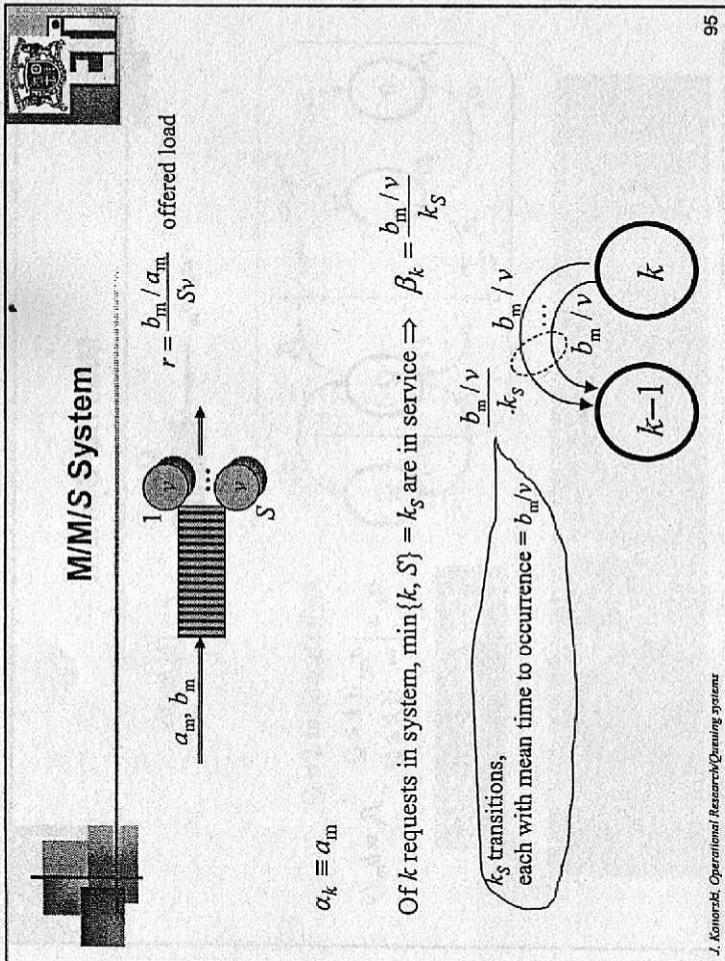
$$p_k \alpha_k = p_{k+1} \beta_{k+1}$$

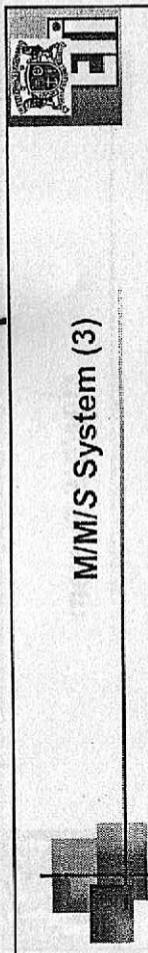
Solution now becomes :

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

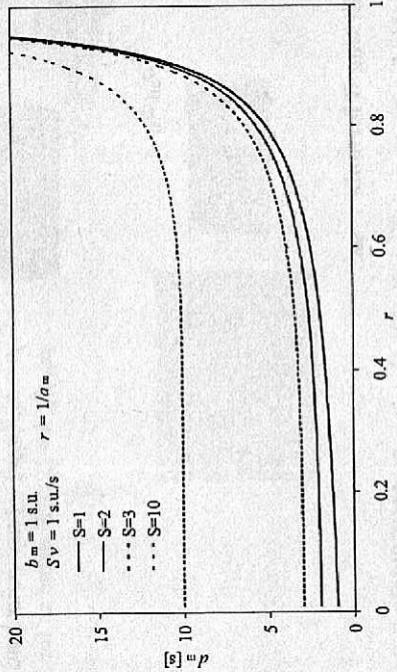
whence all the interesting characteristics: $L, p_0, N_m, d_m, w_m, \dots$





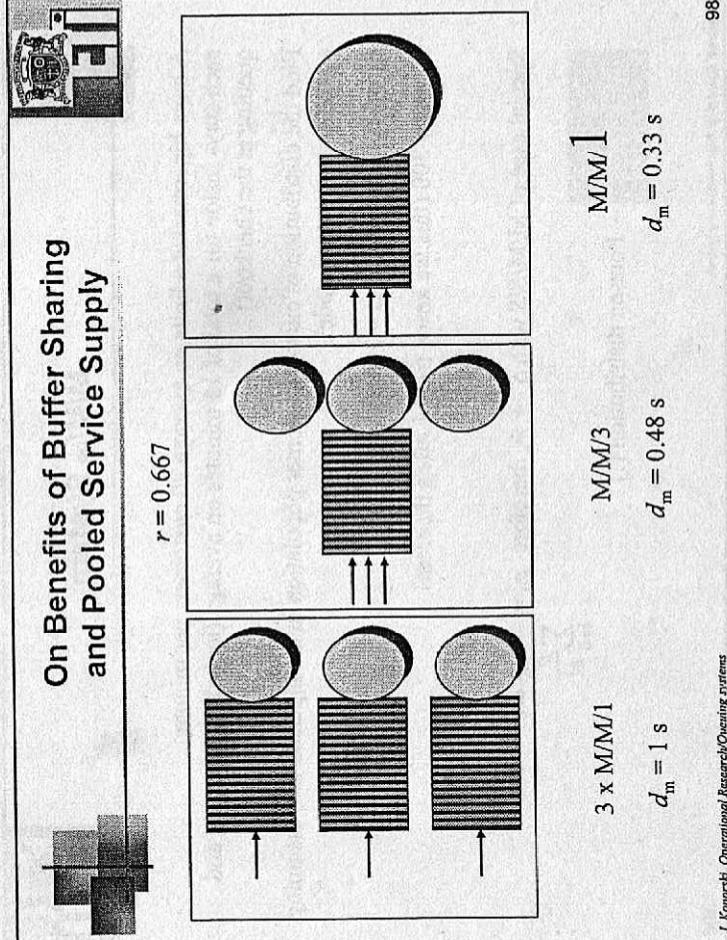


Constant speed of the processor pool:



(We have seen already that disassembled service supply hurts!) 97

J. Koenigs, Operational Research/Queuing Systems



J. Koenigs, Operational Research/Queuing Systems

M/M/S System

no waiting room
 ν – single link capacity
 $\rho = b_m / (\alpha_m \nu)$ – busy-hour load
 $\alpha_k, \beta_k, k = 0, \dots, S$ same as for M/M/S

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

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

- 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/)

J. Koenigs, Operational Research/Queueing systems

99

M/M/ ∞ System

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.

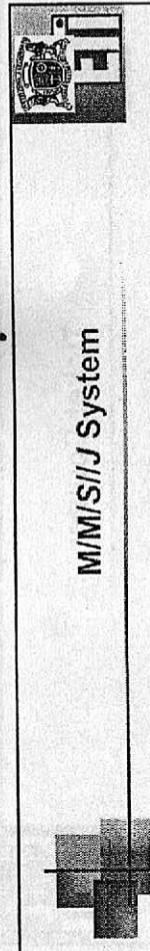
$\alpha_m = 3$ s, $b_m / \nu = 900$ s, $\rho = 300$ erlangs
 $N_m = \rho = 300$ (this we know from Little's theorem).

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

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

J. Koenigs, Operational Research/Queueing systems

100

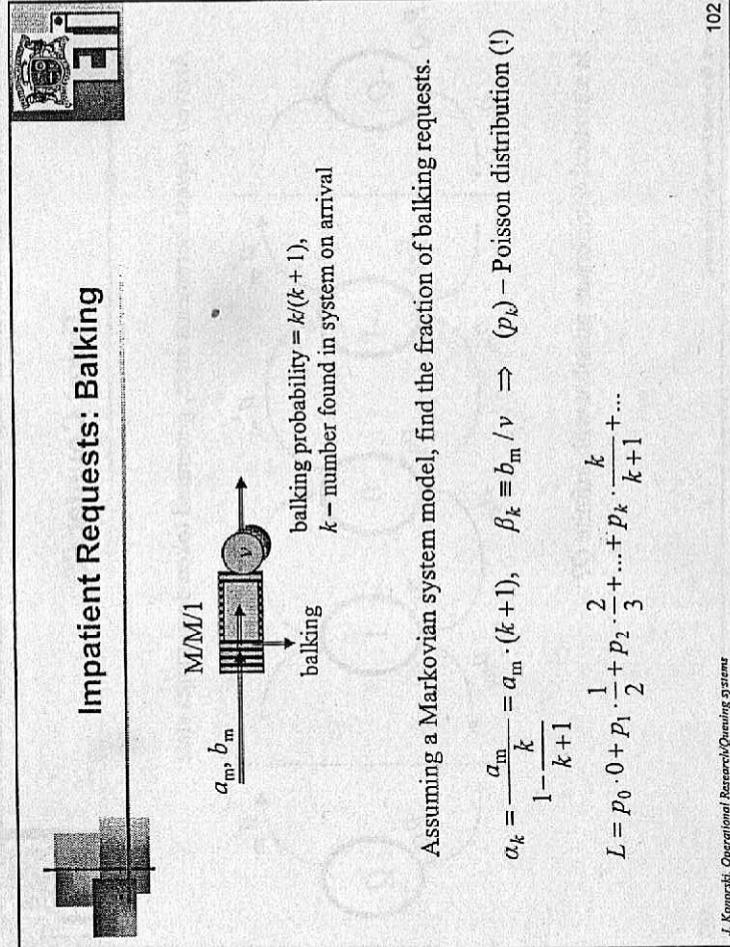


When k requests in system,

- $J - k$ terminals during think time $\Rightarrow \alpha_k = \frac{h_m}{J - k}$
- k_S in service $\Rightarrow \beta_k$ same as for M/M/S
- as $J \rightarrow \infty, h_m \rightarrow \infty, h_m/J \rightarrow a_m$, the system becomes M/M/S

J. Kambouri, Operational Research/Queueing Systems

101



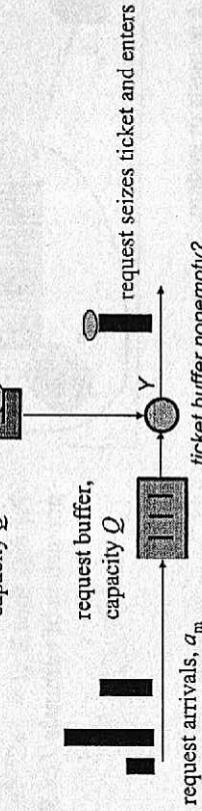
J. Kambouri, Operational Research/Queueing Systems

102



ticket generation, a'_m

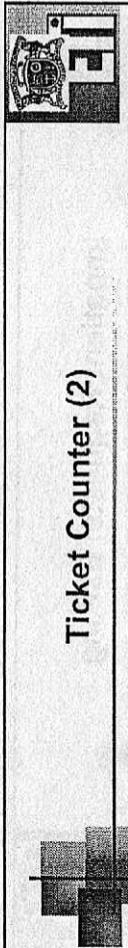
ticket buffer,
capacity Q



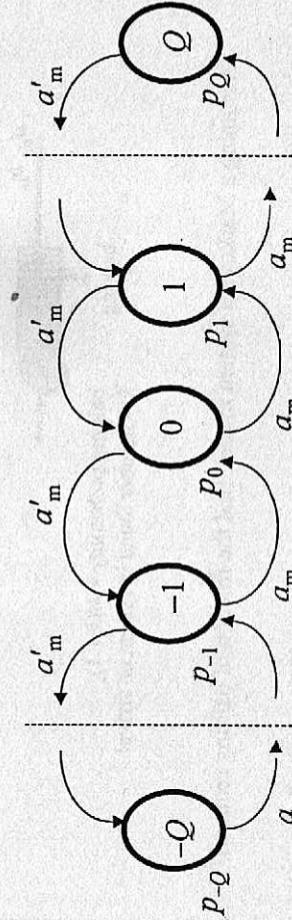
For statistical equilibrium, find the probability of tickets waiting for requests / requests waiting for tickets.

J. Karellović, Operational Research/Queueing systems

103



Arrived request increments state, generated token decrements state:



Is statistical equilibrium possible with infinite Q ?

J. Karellović, Operational Research/Queueing systems

104

A/B/2/Q with Grading and No FIFO Blocking

Flow balance of probability mass:

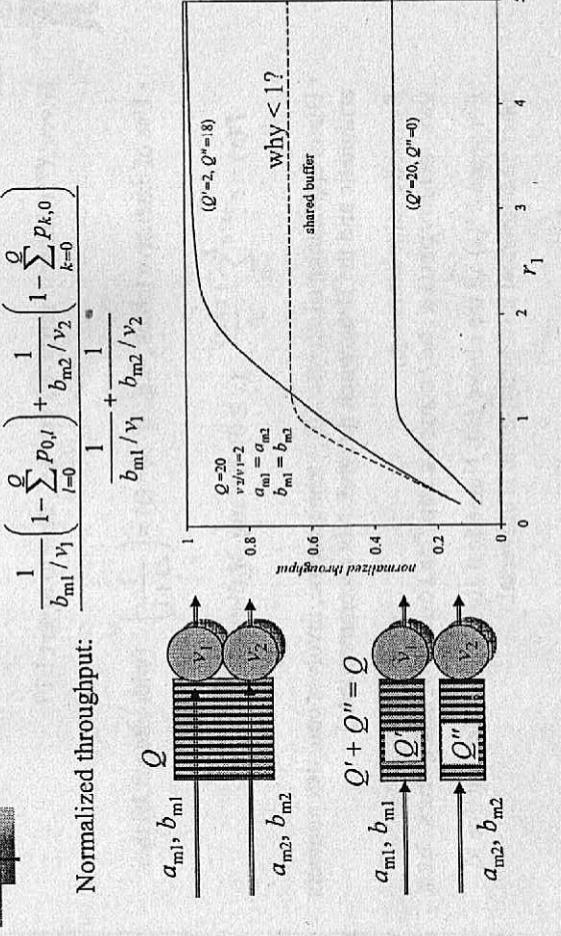
$$P_k/a_{ml} + P_{k,l}/a_{ml} + P_{k,l}/(b_{ml}/v_1) + P_{k,l}/(b_{m2}/v_2) \\ = P_{k-1,l}/a_{ml} + P_{k,l-1}/a_{ml} + P_{k+1,l}/(b_{ml}/v_1) + P_{k+1,l}/(b_{m1}/v_1)$$

$$\text{Solution by guessing and inspection: } P_{k,l} = P_{0,0} \eta^k \eta'_l \quad \text{where } \eta = \frac{b_{ml}}{a_{ml} v_1}, \eta'_2 = \frac{b_{m2}}{a_{m2} v_2}$$

J. Kemerer, Operational Research/Queueing systems

On Benefits of Buffer Partitioning

Normalized throughput:

$$\frac{\frac{1}{b_{ml}/v_1} \left(1 - \sum_{l=0}^Q p_{0,l} \right) + \frac{1}{b_{m2}/v_2} \left(1 - \sum_{k=0}^Q p_{k,0} \right)}{\frac{1}{b_{ml}/v_1} + \frac{1}{b_{m2}/v_2}}$$


J. Kemerer, Operational Research/Queueing systems

Laplace Transform

Finding delay distribution is a harder task and requires more advanced math.
In particular, Laplace transform.

For random variable X with complementary probability function $P(x)$, define:

$$X^*(s) = \int_0^\infty e^{-sx} (-dP(x))$$

E.g., for exponential distribution, $P(X \geq x) = e^{-x/c} \Rightarrow X^*(s) = \frac{1}{cs+1}$

LT is a linear operator.

If random variables X_1 and X_2 are independent and $X = X_1 + X_2$ then

$$X^*(s) = X_1^*(s)X_2^*(s)$$

For sums of 2, 3, ... exponential random variables:

$$X^*(s) = \left(\frac{1}{cs+1}\right)^2, \left(\frac{1}{cs+1}\right)^3, \dots$$

J. Konszak, Operational Research/Queueing systems

Laplace Transform (2)

Given $X^*(s)$, how to retrieve $P(x)$? I.e., how to invert LT?

- Use of extensive tables, e.g., if $X^*(s) = \left(\frac{1}{cs+1}\right)^M$ (with integer M) then

$$P(x) = e^{-x/c} \sum_{i=0}^{K-1} \frac{(x/c)^i}{i!} \quad (x \geq 0) \text{ Erlang-}M \text{ distribution}$$

- Direct application of inverse LT – troublesome, involves complex numbers arithmetic and the Bromwich integral. Not recommended :)
- Symbolic calculators e.g.,
www.educypedia.be/education/calculators/algebra.htm
- For some $X^*(s)$ all the above fail. Numerical algorithms exist, but due to inherent numerical instability, none is universal.

J. Konszak, Operational Research/Queueing systems

108

Laplace Transform (3)

[www.pe.tamu.edu/blasingame/data/P620_reference/P620_Lectures_\(pdf\)/P620_Mdl_Math/P620_Mod1_ML_05_LaplaceTrans.pdf](http://www.pe.tamu.edu/blasingame/data/P620_reference/P620_Lectures_(pdf)/P620_Mdl_Math/P620_Mod1_ML_05_LaplaceTrans.pdf)

- The Gaver formula for numerical Laplace transform inversion is

$$f_{Gaver}(n,t) = \frac{\ln(2)}{t} \frac{(2n)!}{(n-1)!} \sum_{k=0}^n \frac{(-1)^k}{(n-k)k!} f\left[\frac{\ln(2)}{t} (n+k)\right]$$

- The Gaver-Stehfest formula for numerical Laplace transform inversion is

$$f_{Gaver-Stehfest}(n,t) = \frac{\ln(2)}{t} \sum_{i=1}^n V_i f\left[\frac{\ln(2)}{t} i\right]$$

and the Stehfest extrapolation coefficients are given

$$V_i = (-1)_2^{n-i+1} \sum_{k=\lceil \frac{i+1}{2} \rceil}^n \frac{\frac{n-k}{2}}{k_2} \frac{(k-1)!(i-k)!(2k-i)!}{k! (k-1)!(i-k)!(2k-i)!}$$

J. Komar/H. Operational Research/Queuing Systems 109

M/M/1 FIFO: Distribution of System Delay

System delay of request finding k in system on arrival is composed of $k+1$ service times (including one residual, if $k > 0$), each exponentially distributed: $P(x) = e^{-x/(b_m s)}$.

$$\text{Hence, LT of system delay is } \left(\frac{1}{b_m s + 1} \right)^{k+1}.$$

PASTA applies, so $p_k^+ \equiv p_k$

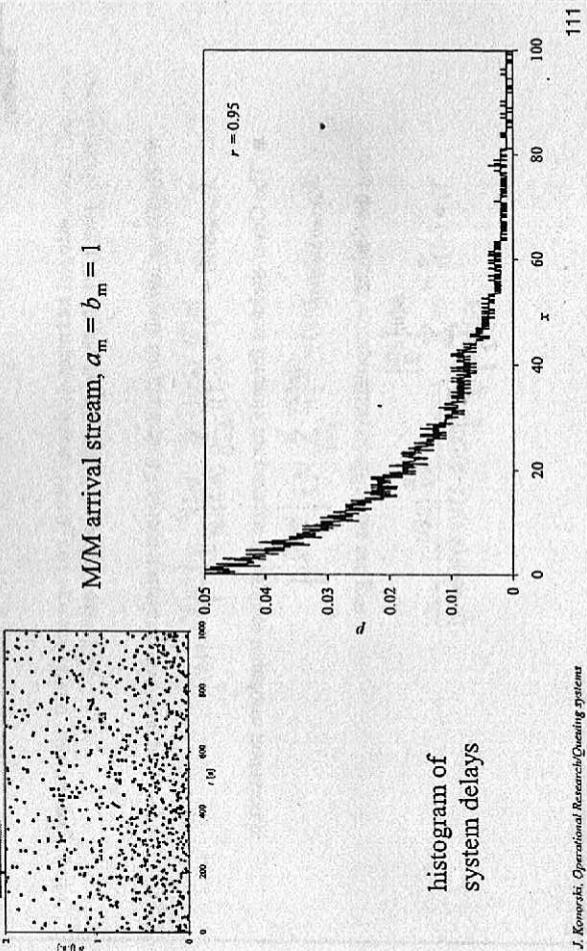
Averaging over k and using LT linearity gives:

$$D^*(s) = \sum_{k=0}^{\infty} p_k \left(\frac{1}{b_m s + 1} \right)^{k+1} = \sum_{k=0}^{\infty} (1-r)r^k \left(\frac{1}{\frac{b_m}{r} s + 1} \right)^{k+1} = \frac{1}{\frac{r}{1-r} s + 1}$$

This corresponds to exponential distribution (whose mean is already known to us)!

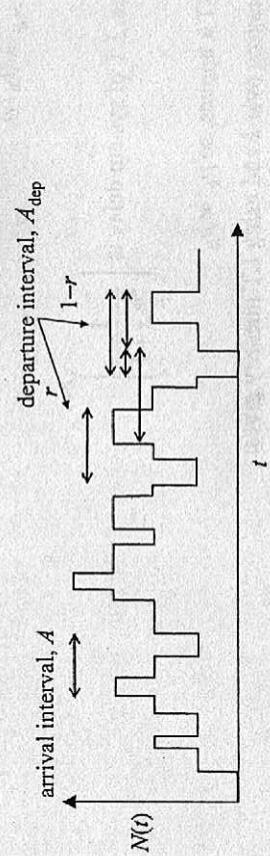
J. Komar/H. Operational Research/Queuing systems 110

M/M/1 FIFO: Distribution of System Delay (2)



M/M/1: Departure Stream

What does departure stream look like in M/M/1?



PASTA and unit jump argument apply, so $p_k \equiv p_k^+ \equiv p_k^-$. Hence, $p_0^- = 1 - r$.

Departure interval A_{dep} with probability r is a single service time, and with probability $1 - r$ consists of a service time and a residual arrival interval.

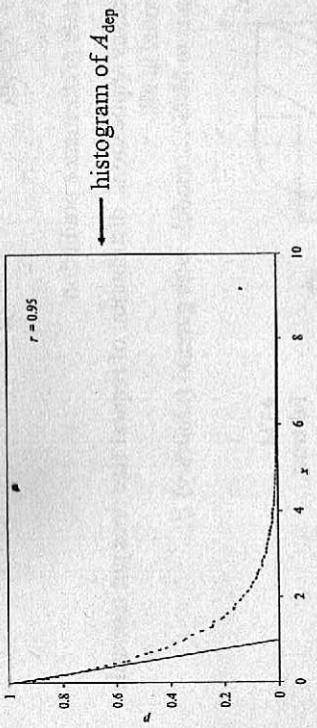
Finding the distribution of A_{dep} is easy once LT has been mastered...

112

J. Konarski, Operational Research & Queueing Systems

M/M/1: Departure Stream (2)

$$A_{wjj}^*(s) = r \frac{1}{b_u s + 1} + (1 - r) \frac{1}{a_{sf} s + 1} \cdot \frac{1}{\frac{b_u}{v} s + 1} = \frac{1}{a_{sf} s + 1} = A^*(s)$$



M/M/1 preserves Poisson stream! (Burke's theorem)

Important in tandem configurations... →

113

J. Konorski, Operational Research/Queuing Systems

Operational Research

Queuing Systems 4: Non-Markovian Models,
Priority Queues, Processor Sharing

Jerzy Konorski
jekon@eti.pg.gda.pl

114

J. Konorski, Operational Research/Queuing Systems

Motivation

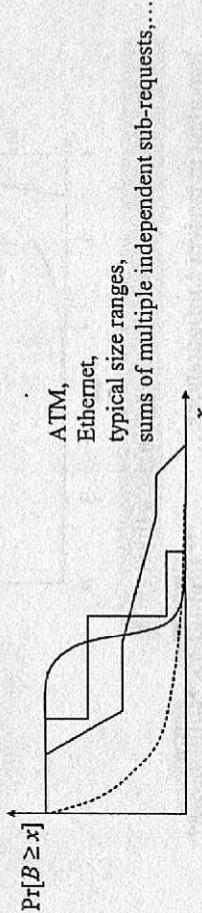
Poisson arrival stream still assumed: $A \sim \Pr[A \geq x] = e^{-x/a_m}$.

Its modeling clout has been pointed out:

- random splitting
- Palm theorem
- PASTA
- fail-safe performance prediction

Yet in many applications, distribution of request size does not resemble exponential at all!

This leads to M/G/... models with general (arbitrary) B .



J. Komoradi, Operational / Research / Queueing systems

115

M/G/1 Analysis

Poisson arrival stream, a_m
general request size distribution, b_m

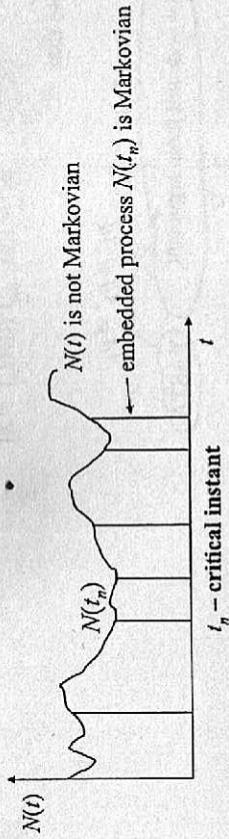
- M/G/1: 1 processor, infinite buffer capacity
- $B \sim P(x)$, arbitrary distribution with mean b_m
- $r = b_m/(a_m v)$, where v – processor speed
- $N(t)$ does *not* possess the Markov property
- two-dimensional process $[N(t), \text{residual service time}(t)]$ does possess it, but is analytically much more complex!

J. Komoradi, Operational / Research / Queueing systems

116

M/G/1 Analysis (2)

However, the situation is far from hopeless: invoke embedded Markov chain.



How should one choose the critical instants (t_n)?

117

J. Karelka, Operational Research/Queueing Systems

M/G/1 Analysis (3)

For M/G/1 the choice is easy:

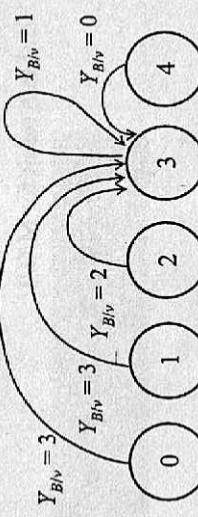
- $t_n \equiv t_n^- + 0$ i.e., just after departure of n th request
- $N(t_n) =$ number of requests left by n th request on departure

Let $Y_{B|v} =$ number of request arrivals during one service time ($B|v$).

Clearly, $Y_{B|v,m} = (b_m/v)/a_m = r$.

(a stochastic equation)

$$N(t_n) = \max[N(t_{n-1}) - 1, 0] + Y_{B|v}$$



118

J. Karelka, Operational Research/Queueing Systems

M/G/1 Analysis (4)

Stochastic equations translate into probability distributions of the involved random variables. Most conveniently at statistical equilibrium:

$$p_k^- = \lim_{n \rightarrow \infty} \Pr[N(t_n) = k]$$

Note that

$$p_k^- \equiv p_k^+ \equiv p_k$$

$\circlearrowleft Q = \infty \text{ & unit jump argument}$ $\circlearrowright \text{PASTA}$

hence at statistical equilibrium, the distribution of $N(t_n)$ is also the distribution of $N(t)$ sampled at random critical instants!

J. Konsztak, Operational Research/Queueing systems

119

M/G/1 Analysis (5)

Method of Laplace Transform (LT) yields (p_k) in a closed form, albeit not as digestible as for M/M/...

1. Find LT of service size distribution: $B^*(s) = \int_0^\infty e^{-sx} (-dP(x))$

2. Desired (p_k) envisaged as a polynomial in a symbolic variable z :

$$p_0 + p_1 z + p_2 z^2 + p_3 z^3 + \dots$$

(Its derivative at $z = 1$ equals N_m – easy to verify.)

3. Main result for M/G/1:

$$p_0 + p_1 z + p_2 z^2 + p_3 z^3 + \dots = (1-r) \frac{1-z}{1-z/B^*(\frac{1-z}{p_0})}$$

J. Konsztak, Operational Research/Queueing systems

120

M/G/1 Analysis (6)

The RHS of the latter relationship is the probability generating function of the distribution (p_k) .

Can be inverted by expansion into a power series around $z = 0$ or using rich tables of the so-called Z-transform.

LTs can be calculated directly or using tables of LTs. Direct LT inversion is not recommended to beginners ;), is easy using tables of LTs given a suitable entry is found, also possible via various approximate numerical algorithms, cf. e.g., [J. Abate and P. P. Valkó: *Multi-precision Laplace transform inversion*, Int. J. Numer. Meth. Engng 2004; 60:979–993].

J. Konarak, Operational Research/Queueing Systems

121

Special Case: M/M/1

To cross-check, find (p_k) for M/M/1, which we have derived before using birth-and-death-equations:

$$P(x) = e^{-x/b_m}, \quad B^*(s) = \frac{1}{b_m s + 1},$$

$$(1-r) \frac{1-z}{1-z/B^*(\frac{1-z}{a_m})} = (1-r) \frac{1}{1-rz} = (1-r) + (1-r)rz + (1-r)r^2z^2 + \dots$$

power series expansion

So $p_k = (1-r)r^k$, as yielded by Markovian analysis.

J. Konarak, Operational Research/Queueing systems

122

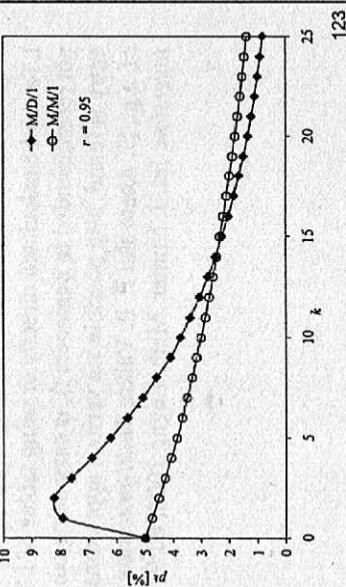


Special Case: M/D/1

For M/D/1, Markovian analysis fails. Method of LT doesn't:

$$P(x) = \begin{cases} 1, & \text{if } x < b_m \\ 0, & \text{if } x \geq b_m \end{cases}, \quad B^*(s) = e^{-b_m s},$$

$$(1-r) \frac{1-z}{1-z/B^*(\frac{1-z}{\sigma_m})} = (1-r) \frac{1-z}{1-ze^{r(1-z)}} =$$



Numerical power series expansion
is not a problem – cf.
www.educypedia.be/education/calculatorsalgebra.htm
or other abundant software.

J. Konorski, Operational Research/Queueing systems



P-K Formula

Differentiating the main result for M/G/1 at $z = 1$, one gets N_m^* and further,
by Little's law:

$$\text{M/M/1} \\ w_m^* = \frac{1 + c_b^2}{2} \cdot \frac{\tau_m}{c_m \cdot 1 - r}$$

Cognitive value:

- magic $1 - r$ again in the denominator, this is no accident!
- comparison with M/M/1 straightforward
- shape of request size distribution only reflected through standard deviation
- variation of request size distribution worsens mean waiting delay
- recalling the residual lifetime formula, one gets an alternative, and very suggestive, formula:

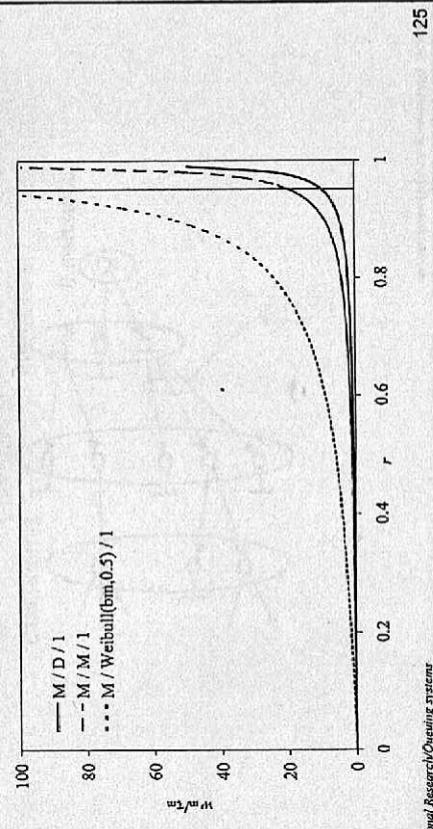
$$w_m^* = \tau_m \cdot \frac{r}{1 - r} \quad \text{e.g., } r = 0.9: 90\% \text{ chance of waiting for 10 residual service times}$$

J. Konorski, Operational Research/Queueing systems

P-K Formula (2)

Example: for $r = 0.95$, compare:

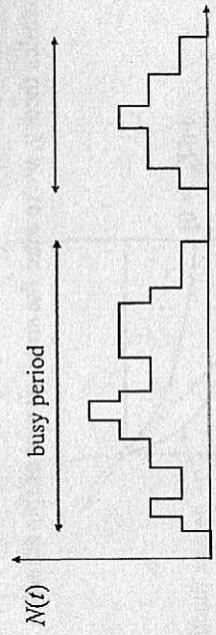
- M/D/1 ($c_b = 0$)
- M/M/1 ($c_b = 1$)
- M/Weibull($b_m, 0.5$)/1 ($c_b \approx 3.317$) !



J. Konarak, Operational Research/Queueing systems

Busy Period

Question that can't be answered experimentally within any definite time:



Will a busy period in progress ever terminate?

Not necessarily – with a probability dependent upon r :

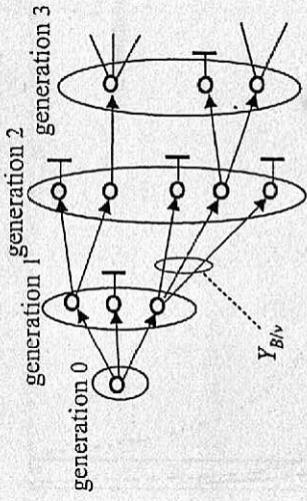
$$\omega = \Pr[\text{busy period eventually terminates}]$$

J. Konarak, Operational Research/Queueing systems

126

Busy Period (2)

$\omega = \Pr[\text{extinction of a population where each individual's number of offspring} = Y_{B\nu}]$



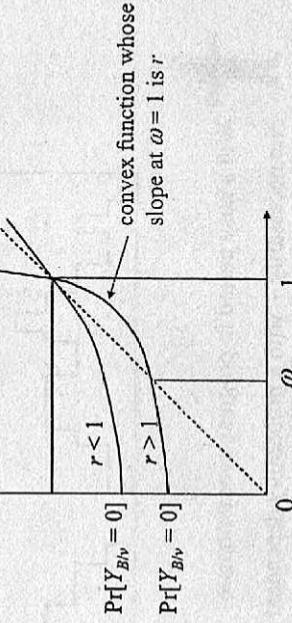
J. Konarski, Operational Research/Queuing systems

127

Busy Period (3)

$$\omega = \Pr[Y_{B\nu} = 0] \omega^0 + \Pr[Y_{B\nu} = 1] \omega + \Pr[Y_{B\nu} = 2] \omega^2 + \dots = B * \left(\frac{1-\omega}{a_m v} \right)$$

By population dynamics theory, we're after the smallest root of this equation in $(0, 1]$.



$$r > 1 \Rightarrow \omega < 1$$

$$r < 1 \Rightarrow \omega = 1$$

$r < 1 \Rightarrow \omega = 1$ surprising, given that $r = Y_{B\nu,m}$
Average offspring of 1 is a recipe for almost assured extinction!

J. Konarski, Operational Research/Queuing systems

128

Busy Period (4)

If $r < 1$, mean duration of busy period z_m is finite and easy to obtain:

z_m = service time of first request
+ sum of durations of busy periods triggered by requests arrived during
first request's service

$$z_m = \tau_m + Y_{B\lambda_m} z_m = \tau_m + (\tau_m/a_m) z_m = \tau_m + r z_m$$

$$\text{Hence, } z_m = \frac{\tau_m}{1-r}$$

M/G/1 Waiting Delay Distribution

Ingenious argument for FIFO: a departing request leaves behind a queue of requests arrived during its waiting time and service.

$$\text{So } N = Y_{W+B\nu}$$

Rewrite it as $N = Y_{(v\mu)\nu} + Y_{B\nu}$. Compare the corresponding polynomials in z for left-hand side (main result for M/G/1) and right-hand side:

$$(1-r) \frac{1-z}{1-z/B*(\frac{1-z}{a_m\nu})} = (v\mu)*(\frac{1-z}{a_m\nu}) \cdot B * (\frac{1-z}{a_m\nu})$$

Finally, from LT properties, $(c\lambda)^*(s) = X^*(cs)$ (verify this).

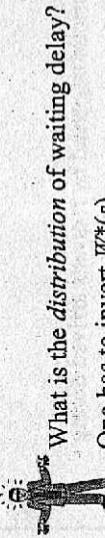
So, after simple algebra,

$$W^*(s) = \frac{1-r}{1-B*(\frac{s}{\nu})} - \frac{r}{a_m s}$$

M/G/1 Waiting Delay Distribution (2)

Mean and standard deviation follow by LT differentiation at $s = 0$:

$$w_m = -W^*(0), \quad \sigma_w = \sqrt{W^{**}(0) - w_m^2}$$



One has to invert $W^*(s)$.

Examples:

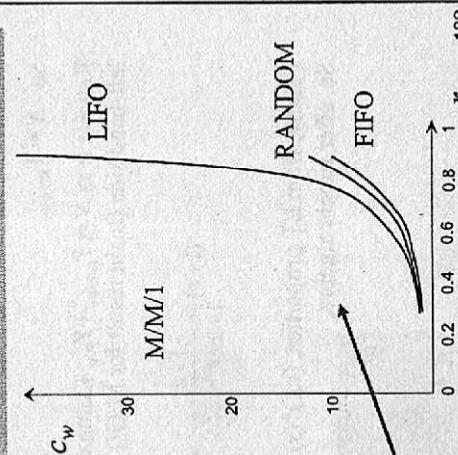
$$\text{- M/M/1: } B^*(s) = \frac{1}{b_m s + 1} \Rightarrow W^*(s) = 1 - r + r \frac{1}{v(1-r)} s + 1$$

$$\text{- M/D/1: } B^*(s) = e^{-b_m s} \Rightarrow W^*(s) = 1 - r + r \frac{\frac{(b_m/v)s}{1 - e^{-(b_m/v)s}}}{r} = r$$

J. Komorski, Operational Research/Queueing systems

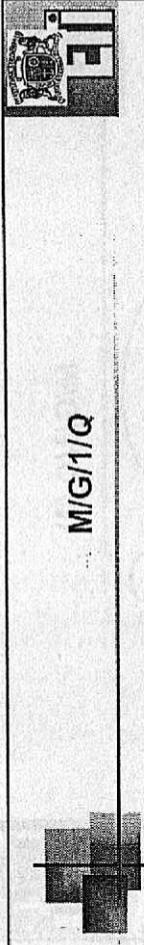
M/G/1 Waiting Delay Distribution (3)

Characteristics of $N(t)$ and busy period, as well as mean delay are QD invariant (provided QD is work-conserving and does not take advantage of information on request sizes).



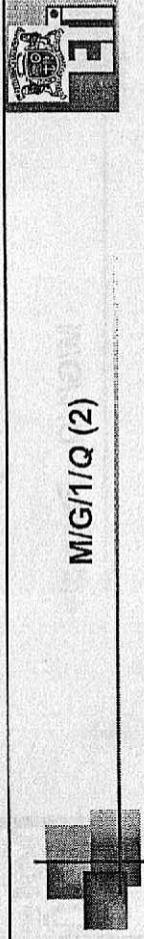
Delay distribution is not!
For various QDs, may differ dramatically
in the coefficient of variation.

J. Komorski, Operational Research/Queueing systems



- from practical viewpoint, the most interesting is determination of L
 - loss fraction due to buffer overflow
- no closed-form expression, have to resort to numerical calculation
- advantage over simulation when L is very small

133

J. Kouvatsi, Operational Research/Queueing Systems

Begin again with a stochastic equation:

$$N(t_n) = \min \{ \max[N(t_{n-1}) - 1, 0] + Y_{B_{t_n}}, Q \}$$

This time, "translation" into distributions will be done directly, without LT.

Distribution of $Y_{B_{t_n}}$ has been calculated before:

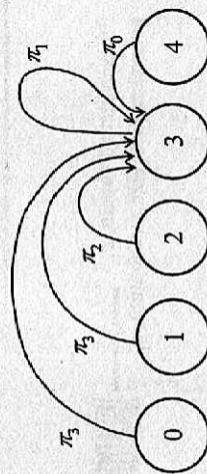
$$\pi_k = \Pr[Y_{B_{t_n}} = k] = \sum_0^{\infty} \frac{[(x/\nu)]/a_m]^k}{k!} e^{-[(x/\nu)]/a_m} (-dP(x))$$

and shown to follow from power series expansion of B^* $\left(\frac{1-z}{a_m \nu} \right)$.

J. Kouvatsi, Operational Research/Queueing Systems

134

M/G/1/Q (3)



Looking at the familiar picture, write down equations with unknowns p_k^- :

$$\begin{aligned} p_k^- &= p_0^- \pi_k + p_1^- \pi_k + p_2^- \pi_{k-1} + \dots + p_{k+1}^- \pi_0, \quad k = 0, \dots, Q-2 \\ p_0^- + \dots + p_{Q-1}^- &= 1 \end{aligned}$$

Furthermore, $p_k \equiv p_k^+$ (by PASTA) and $p_k^- \equiv p_k^+/(1-p_Q^+)$ (by unit jump argument).

Hence,

$$L = p_Q = 1 - \frac{1}{r + p_Q^-}$$

135

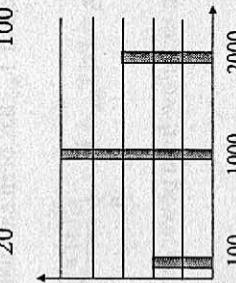
J. Konsztak, Operational Research & Queueing systems

M/G/1/Q – Example

Poisson arrival stream, $a_m = 4$ s

type	proportion [%]	size [s.u.]
http get	30	1000
telnet open	50	2000
snmp SetRequest	20	100

Processor speed = 500 s.u./s



Will $Q = 3$ be enough if $L \leq 10\%$ is required?

J. Konsztak, Operational Research & Queueing systems

136

M/G/1/Q – Example (2)

Distribution of $Y_{B,W}$:

$$\pi_i = 30\% \cdot \frac{(1000/500)/4)^i}{i!} e^{-(1000/500)/4}$$

$$+ 50\% \cdot \frac{(2000/500)/4)^i}{i!} e^{-(2000/500)/4}$$

$$+ 20\% \cdot \frac{(100/500)/4)^i}{i!} e^{-(100/500)/4}$$

$$\text{Offered load: } r = 30\% \cdot \frac{1000}{4 \cdot 500} + 50\% \cdot \frac{2000}{4 \cdot 500} + 20\% \cdot \frac{100}{4 \cdot 500} = 0.65$$

Equations:

$$p_0^- = p_0^- \pi_0 + p_1^- \pi_0$$

$$p_1^- = p_0^- \pi_1 + p_1^- \pi_1 + p_2^- \pi_0$$

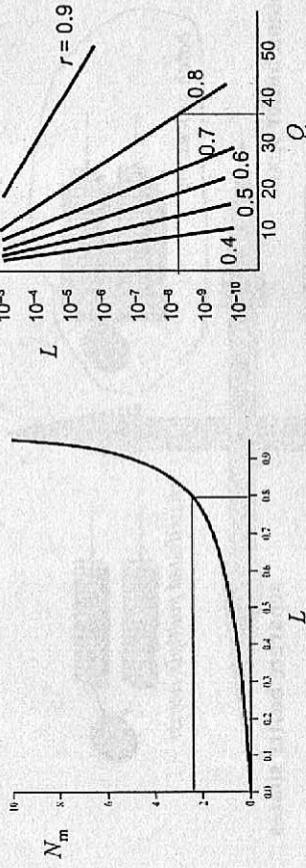
$$p_0^- + p_1^- + p_2^- = 1$$

J. Konior/R., Operational Research/Queueing systems

137

$$L = 1 - \frac{1}{r + p_0^-} = 7.6\%$$

M/G/1/Q Design



- fix tolerable N_m and L (e.g. 2.5 and 10^{-8})
- read maximum feasible r from left plot (drawn for M/D/1, assuming negligible L)
- read minimum required Q from right plot (drawn for M/D/1/ Q)

J. Konior/R., Operational Research/Queueing systems

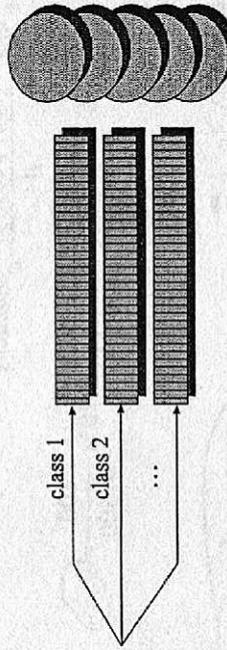
138

Priority Queuing

Prioritization is sometimes expected, meaning that requests have to be somehow classified ("we're all equal, only some of us more so than others").

Classification can be external (user-imposed) or internal (system-imposed).
Each class has its own logical queue.

(convention: lower class number designates higher priority)



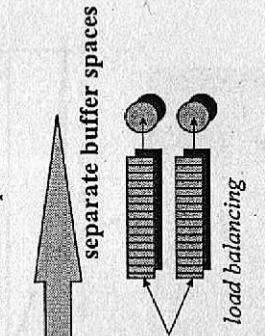
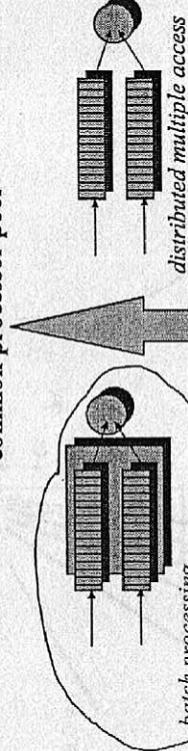
We still speak of *one* queuing system as long as some processors / buffer space / request sources are common to arrival streams of different classes.

J. Kowarski, Operational Research/Queuing systems

139

Priority Queuing (2)

common processor pool

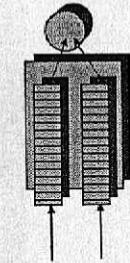


separate processor pools

J. Kowarski, Operational Research/Queuing systems

140

Priority Queuing (3)



Model:

- single processor
- infinite buffer capacity
- service mode: work-conserving, processor-bound, time-sharing

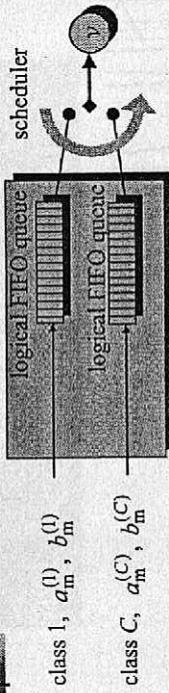
Prioritization objectives and tools:

- differentiation of waiting delays among request classes
QD – FIFO within class, what inter-class?
- differentiation of throughput among request classes
- limited access to buffer space / processor pool

J. Kemeristi, Operational Research/Queuing systems

141

Priority Queuing (4)



$$\text{offered load: } r = r^{(1)} + \dots + r^{(C)}, \text{ where } r^{(j)} = \frac{b_m^{(j)}}{a_m^{(j)}} \nu$$

At instant of service completion, scheduler selects logical queue, head request of which is to be served next.

Condition for selection of a higher-priority class can be arbitrary e.g., "select higher-priority class – unless lower-priority class queue is too long – unless higher-priority class has been waiting too long – unless several successive high-priority requests served – unless they arrived in a bulk – unless ..."

J. Kemeristi, Operational Research/Queuing systems

142

Priority Queuing: HOL

Head of Line (HOL) QD employs the weakest condition – just high-priority requests *being there*. Strongest priority, strongest delay differentiation.

Let $w_m^{(i)}$ – mean waiting delay of class i request

Component delays:

- residual service found in progress on arrival
- service of class $j \leq i$ requests found on arrival
- service of class $j < i$ requests arrived during $w_m^{(i)}$

$$w_m^{(i)} = w_m^{(\text{FIFO})} \Big|_{r^{(1)} + \dots + r^{(i)}} + \sum_{j < i} \frac{w_m^{(j)} b_m^{(j)}}{a_m^{(j)} v}$$

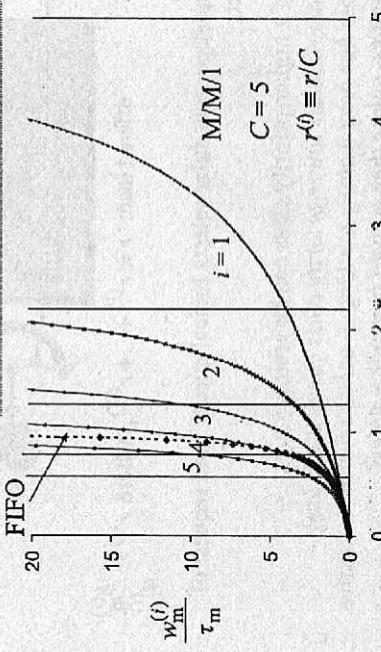
$r^{(i)}$

J. Komarst, Operational Research/Queuing systems

143

Priority Queuing: HOL (2)

$$\text{For M/G/1: } w_m^{(i)} = \frac{(1-r) w_m^{(\text{FIFO})}}{[1 - (r^{(1)} + \dots + r^{(i)})][1 - (r^{(1)} + \dots + r^{(i-1)})]}$$



- global stability \neq individual stability: starvation possible
- comparison with FIFO: tradeoff visible

J. Komarst, Operational Research/Queuing systems

144

Priority Queuing: SJF

HOL makes no use of information on request sizes.

If, however, request sizes are known on arrival, it is natural to favor small-size requests at the cost of others: **Shortest Job First (SJF)**.

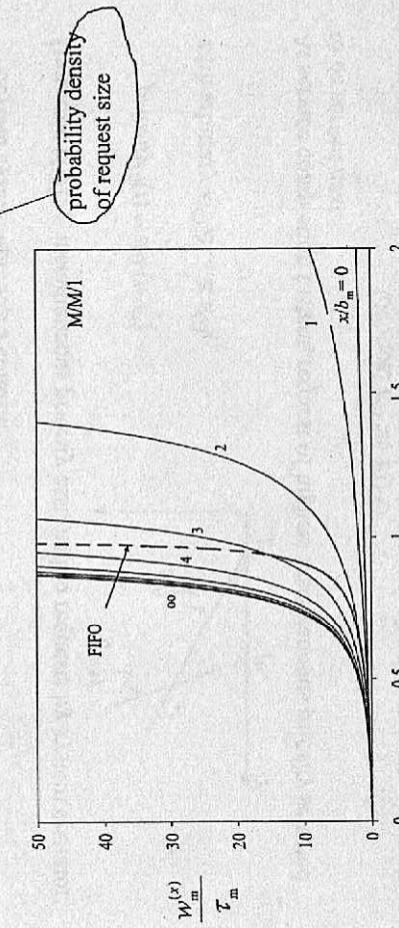
- special case of HOL – continuous class set
- class $x \equiv B \in (x, x+\Delta)$
- class x -induced offered load $\rho(x) = x/(\sigma_m v)$

145

J. Kemerki, Operational Research/Queuing Systems

Priority Queuing: SJF (2)

$$\bullet \text{ mean waiting delay of class } x \text{ requests: } w_m^{(x)} = \frac{(1-r)w_m^{(\text{FIFO})}}{\left[1 - \int_0^x \frac{y}{\sigma_m v} p(y) dy \right]^2}$$

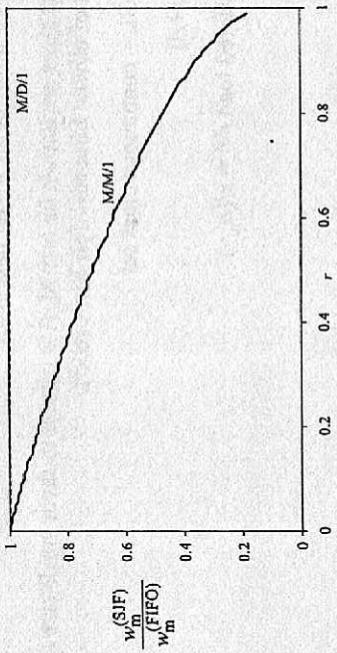


146

J. Kemerki, Operational Research/Queuing Systems

Priority Queuing: SJF (3)

- overall mean waiting delay : $w_m^{(SJF)} = \int_0^{\infty} w_m^{(x)} p(x) dx$



With respect to $w_m^{(SJF)}$, SJF always outperforms FIFO.

J. Konsztat, Operational Research/Queuing systems

147

Priority Queuing: Dynamic HOL

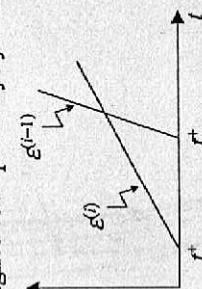
Problems with HOL and SJF:

- starvation
- delay differentiation decided by class-induced offered loads, beyond system operator's control!

Dynamic HOL – instantaneous priority assigned to requests by system operator:

$$\text{priority}^0(t) = \varepsilon^0 \cdot (t - t^*)$$

with arbitrary $\varepsilon^0 > \dots > \varepsilon^{(1)}$.



At service completion instant, request of highest instantaneous priority selected to be served next.

identical $\varepsilon^0 \Rightarrow$ FIFO
very large $\varepsilon^{(i-1)}/\varepsilon^0 \Rightarrow$ HOL

J. Konsztat, Operational Research/Queuing systems

148

Priority Queuing: Dynamic HOL (2)

Components of waiting delay of class i request:

- residual service found in progress on arrival
- service of class $j \leq i$ requests found on arrival
- service of class $j > i$ requests found on arrival it is unable to overtake
- service of later arrived class $j < i$ requests it is unable to escape from

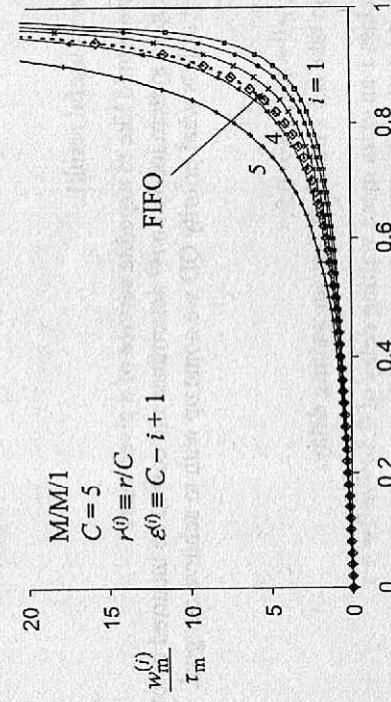
For M/G/1,

$$\frac{w_m^{(i)}}{w_m^{(1)}} = \frac{\frac{w(\text{FIFO}) - \sum_{j>i} r^{(j)} \left(1 - \frac{\varepsilon^{(j)}}{\varepsilon^{(i)}}\right) w_m^{(j)}}{1 - \sum_{j<i} r^{(j)} \left(1 - \frac{\varepsilon^{(j)}}{\varepsilon^{(i)}}\right)}} \quad \begin{array}{l} \cdot \text{recurrence relationship} \\ \cdot \text{only relative magnitudes of } \varepsilon^{(i)} \text{ matter} \end{array}$$

J. Konschmidt, Operational Research/Queuing systems

149

Priority Queuing: Dynamic HOL (3)



- global and individual stability coincide: starvation impossible
- adjustment of mean class delays through the $\varepsilon^{(i)}$
- comparison with FIFO: tradeoff visible again

J. Konschmidt, Operational Research/Queuing Systems

150

Priority Queuing: Conservation Law

Tradeoffs imply existence of invariants. As we change priority QD, what remains constant?

As long as service mode is work-conserving and time sharing, it is the shape of the $U(l)$ (unfinished work) process!



$$\sum_{i=1}^C r^{(i)} \psi_n^{(i)} = r \cdot w_n^{\text{FIFO}}$$

Conservation Law for time sharing

Note that r and $r^{(i)}$ reflect, respectively, total and class i -induced service demand.

Hence, the above may be expressed as follows:

Mean waiting delay of a service unit is invariant with respect to priority QD.

J. Kemerer, Operational Research/Queuing systems

151

Priority Queuing: Conservation Law (2)

This is a most useful result!

Suppose we would like to expedite service of a given class.

The breakup of system load into $r^{(i)}$ determines the cost to be incurred by other classes, regardless what priority QD we come up with to achieve our goal.

Example:

Let $C = 2$, $r^{(1)} = 0.1$, $r^{(2)} = 0.8$.

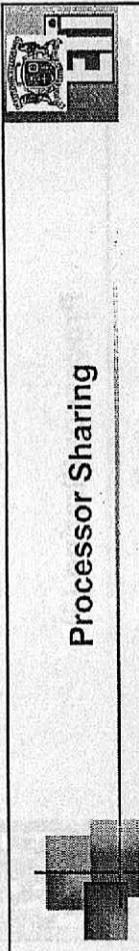
We want to take 0.5 s off class 2 mean waiting delay.

Will cost class 1 an extra mean waiting delay of $0.5 \cdot r^{(2)}/r^{(1)} = 4$ s.

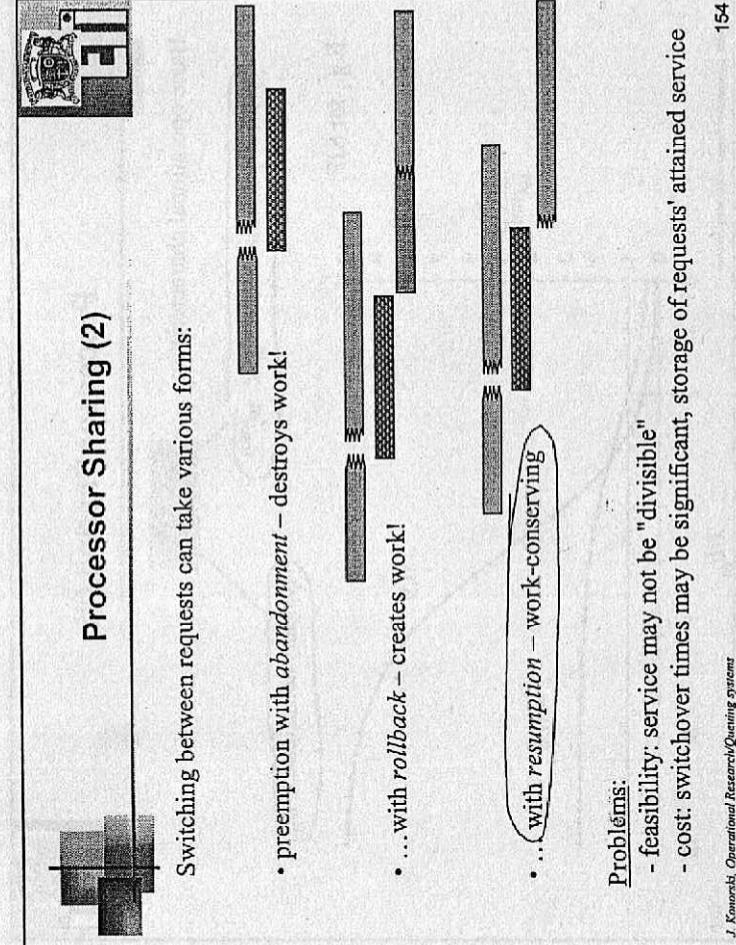
Is this tolerable?...

J. Kemerer, Operational Research/Queuing systems

152



153

J. Koenigst., Operational Research/Queueing Systems

154

J. Koenigst., Operational Research/Queueing Systems

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. Konsalik, Operational Research/Queueing systems

155

Processor Sharing (4)

Basic operational characteristic:

$w_m^{(x)}$

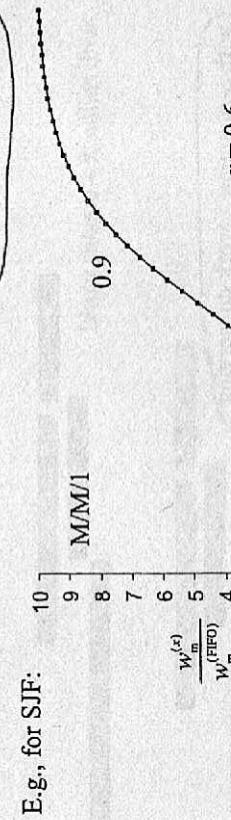
$w_m^{(FIFO)}$

$r = 0.6$

x/b_m

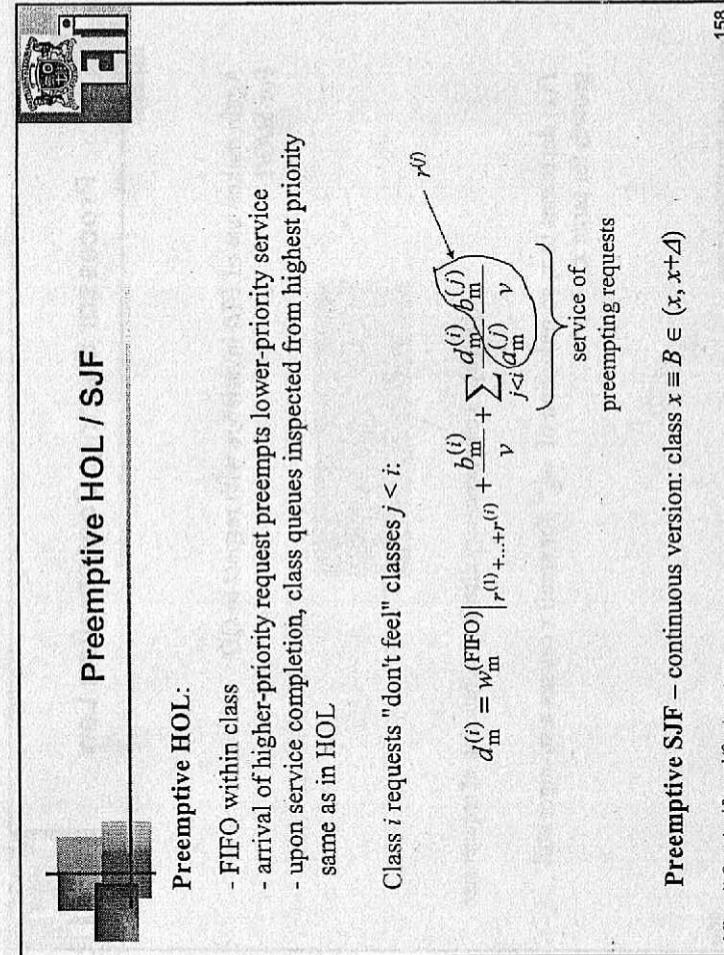
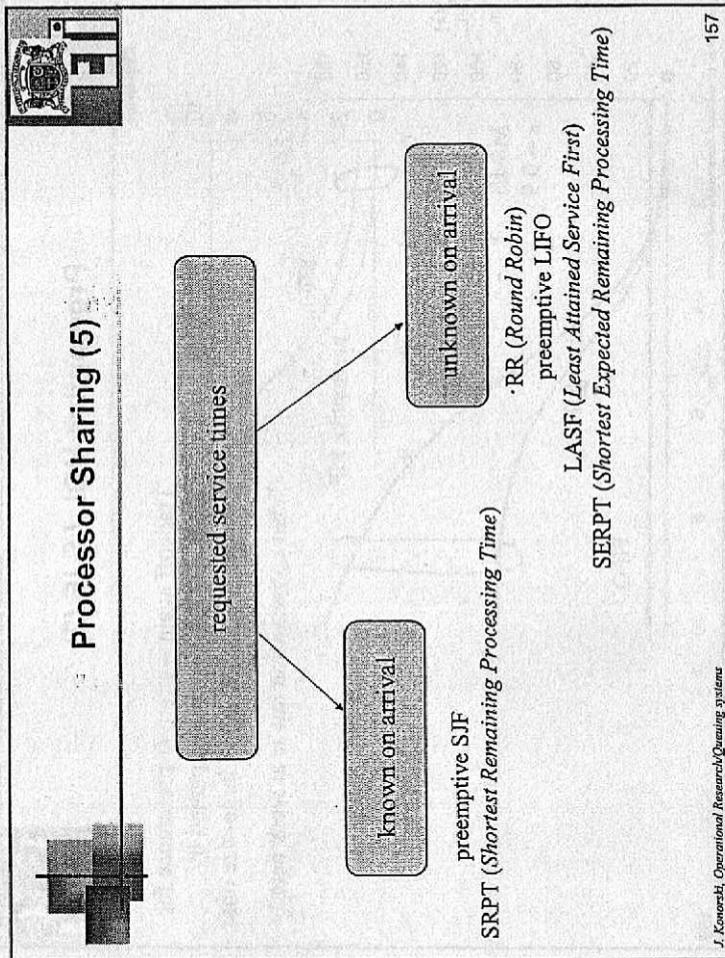
typically normalized to b_m

typically normalized to τ_m or $w_m^{(FIFO)}$



J. Konsalik, Operational Research/Queueing systems

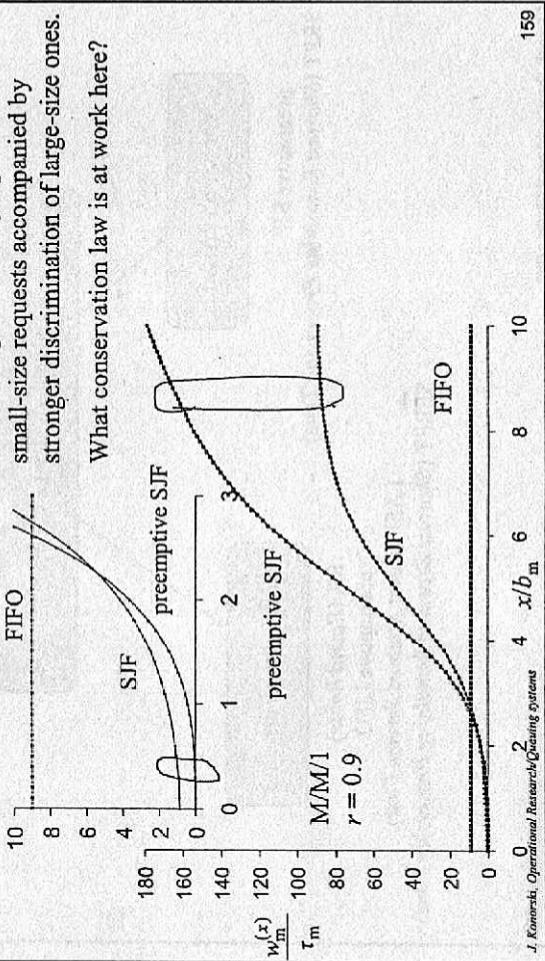
156



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

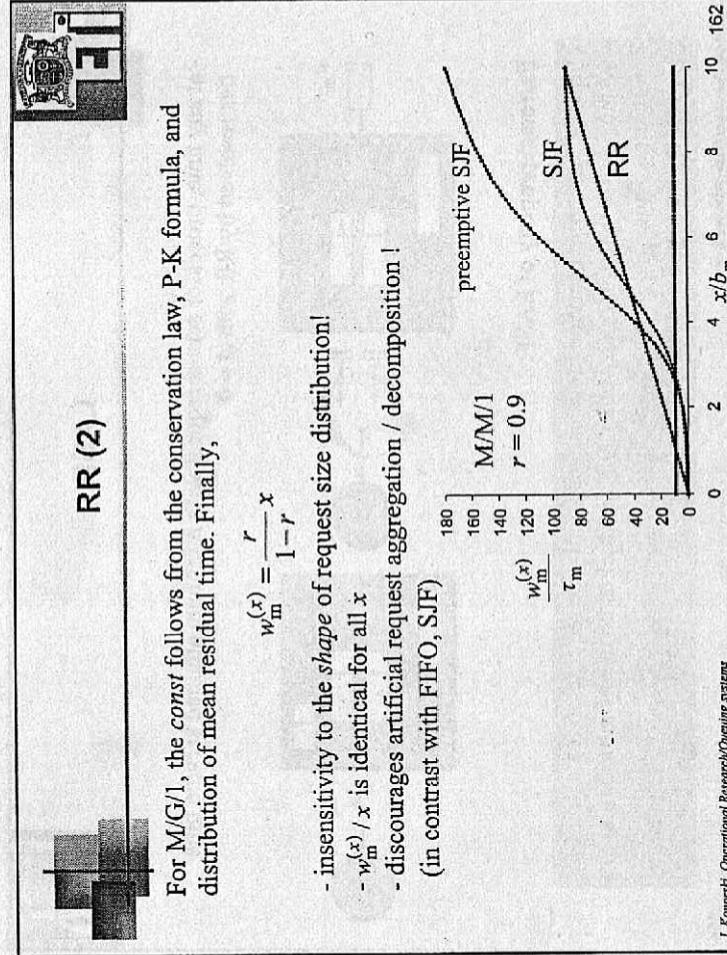
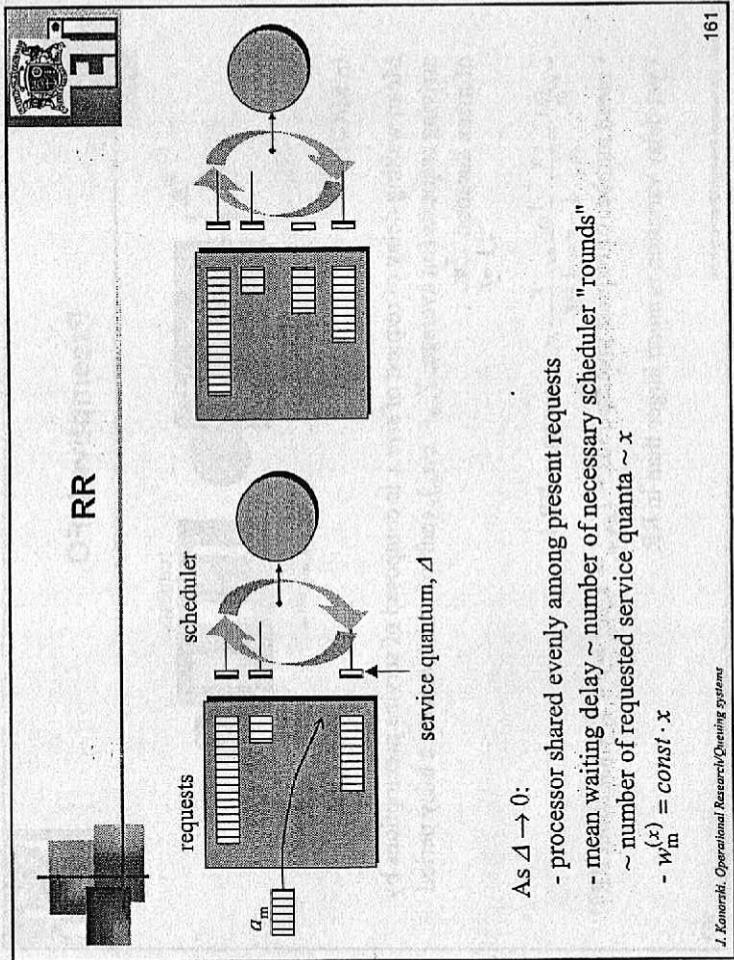
Again makes use of $U(t)$ invariance with respect to QD.

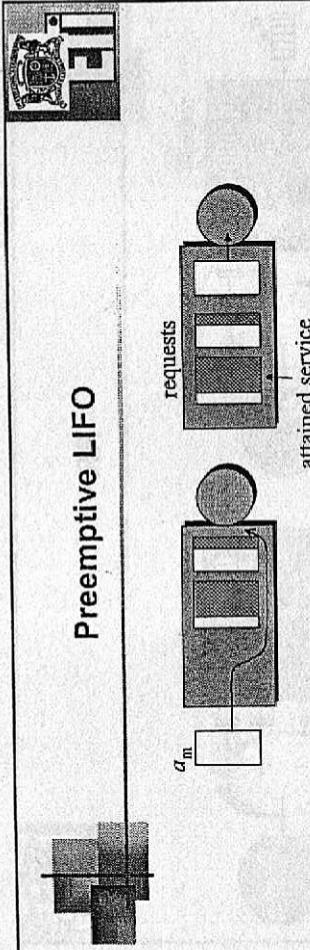
For M/G/1,

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

$\Pr[B \geq x]$ – 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:

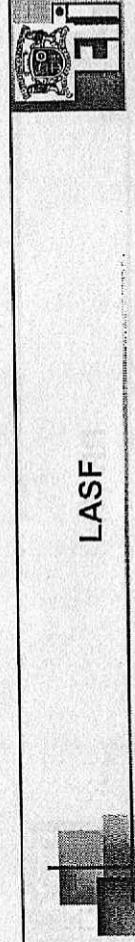
Mean waiting delay of request of size x is composed of service preemptions by arriving requests (on average, $Y_{x,m} = x/a_m$), each of which creates a busy period of mean duration $\frac{\tau_m}{1-r}$.

$$\bullet w_m^{(x)} = \frac{x}{a_m} \cdot \frac{\tau_m}{1-r} = \frac{r}{1-r} \cdot x, \text{ same as in RR!}$$

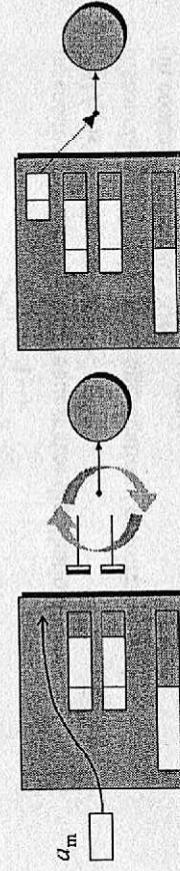
- mean number of request preemptions = $Y_{B\&R} = \tau_m/a_m = r < 1$, what about RR!
- yet delay variation is *much* larger than in RR

163

J. Konsztik, Operational Research/Queuing systems



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

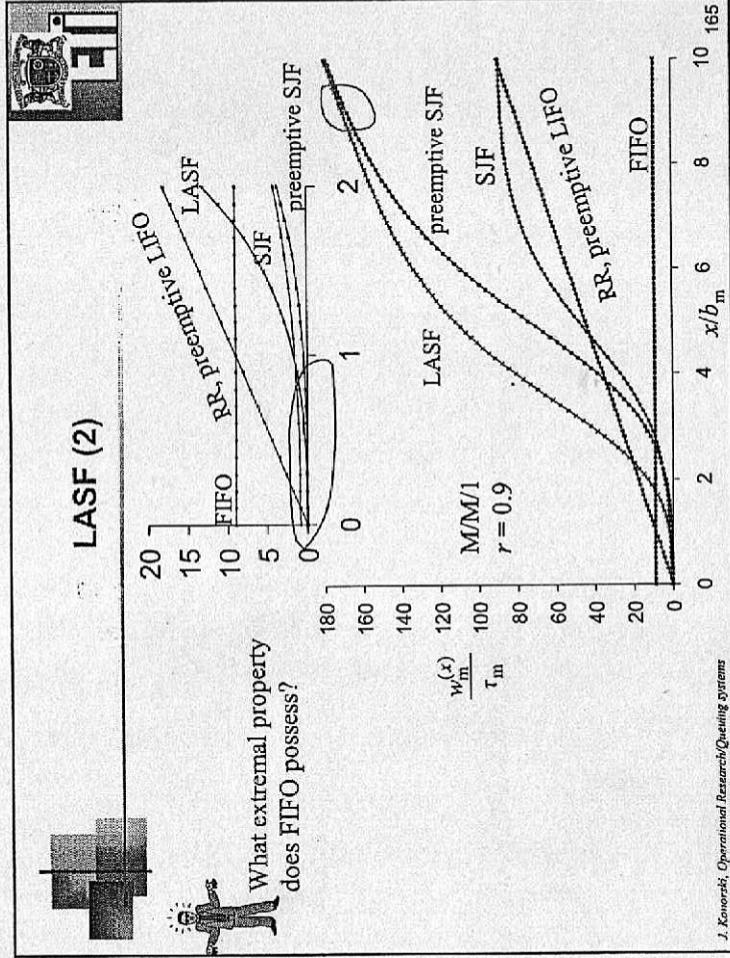


Extremal property of LASF:

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. Konsztik, Operational Research/Queuing systems

164



That's that done, then!
Thank you for your attention (for good).



József Konszak