



Operational Research

Queuing Systems 1: Description and Operation

Jerzy Konorski
Room 139 (old bldg)
office hours: tba
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

J. Konorski, Operational Research/Queuing systems

2

1

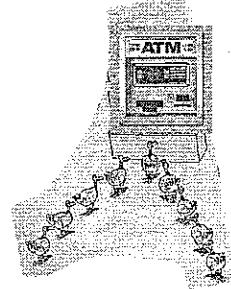
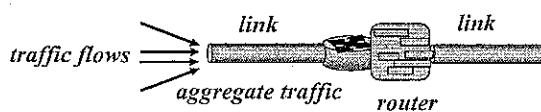
Queuing (Stochastic Service) Systems: Examples

- **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. Konorski, Operational Research/Queuing systems

3

Queuing Systems(2)



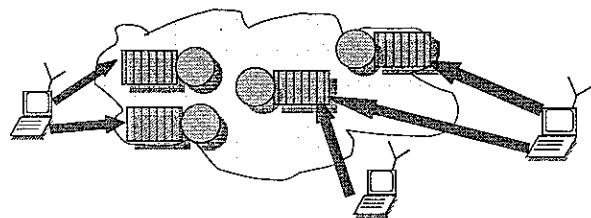
Queues everywhere!
Must be accounted for in systems design.

J. Konorski, Operational Research/Queuing systems

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. Konorski, 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. Konorski, Operational Research/Queuing systems

6

Queuing Theory: Mission (2)

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),

- 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. Kowarski, Operational Research/Queuing systems

7

Simplest Model

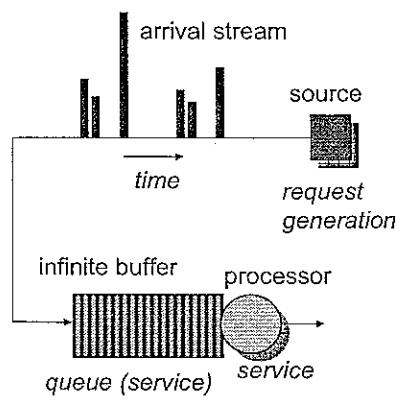
The simplest model of a queuing system consists of:

- a processor,
- a buffer, and
- an arrival stream.

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.

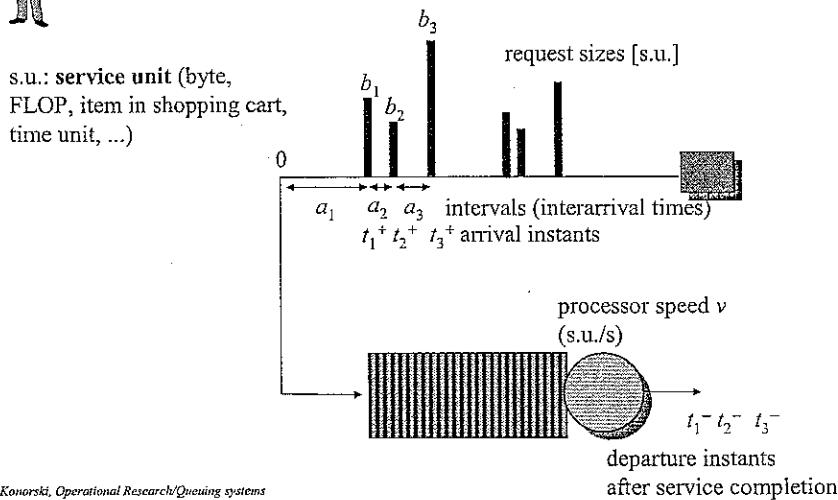


J. Kowarski, Operational Research/Queuing systems

8

Simplest Model (2)

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



Simplest Model (3)

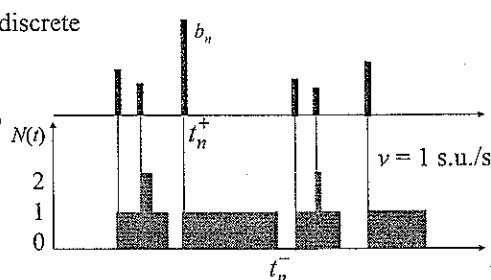
Remarks:

- $t_n^+ = a_1 + \dots + a_n$
- $b_n/v = \tau_n$ = request size / processor speed = requested service time n
- $d_n = t_n^- - t_n^+ =$ system delay of request n
- $w_n = d_n - b_n/v =$ 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

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

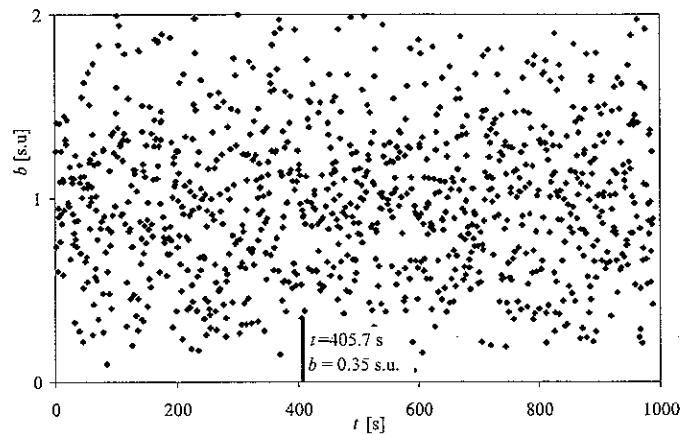
- processor speed v ,
- queuing discipline (order of service)

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



J. Konorski, Operational Research/Queuing systems

Arrival Stream and Service Process: Example

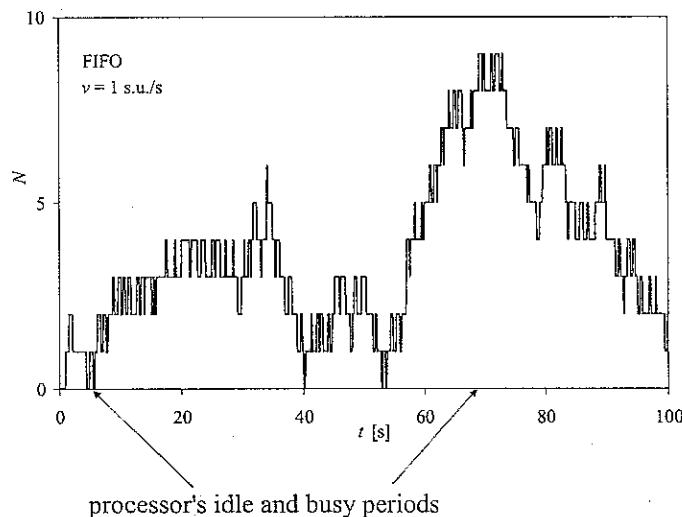


mean interval $a_m = 1$ s
mean request size $b_m = 1$ s.u.

J. Konorski, Operational Research/Queuing systems

11

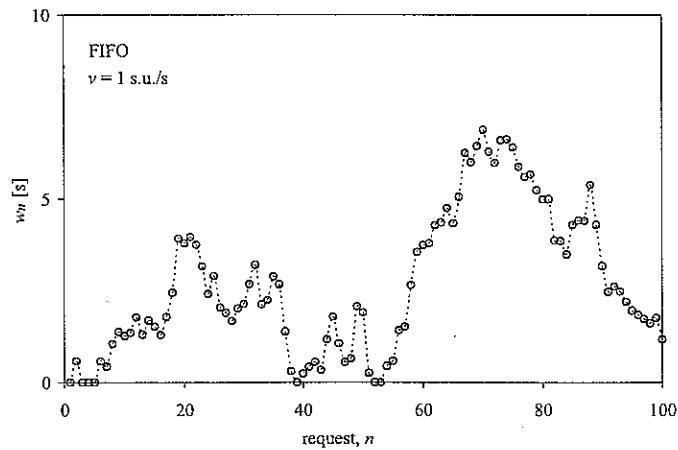
Arrival Stream and Service Process: Example (2)



J. Konorski, Operational Research/Queuing systems

12

Arrival Stream and Service Process: Example (3)



J. Konorski, Operational Research/Queuing systems

13

Properties of "Good" Service Process

- 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
- 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. Konorski, Operational Research/Queuing systems

14

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:

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

(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

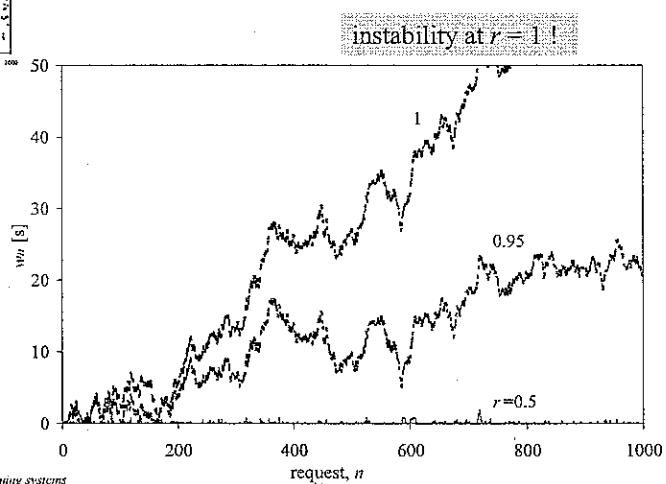
If $r > 1$ persists, processor "gets behind" – **instability**. If $r < 1$, processor "keeps pace" – system stable. What happens under $r = 1$?

J. Konorski, Operational Research/Queueing systems

15

Offered Load (2)

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

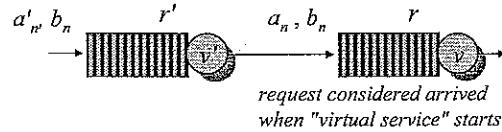


J. Konorski, Operational Research/Queueing systems

16

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 $r = (b_m/a_m)/v$.

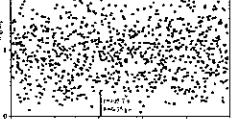
Clearly, $a_m = a'_m$, so $r' = (v/v)r$.

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

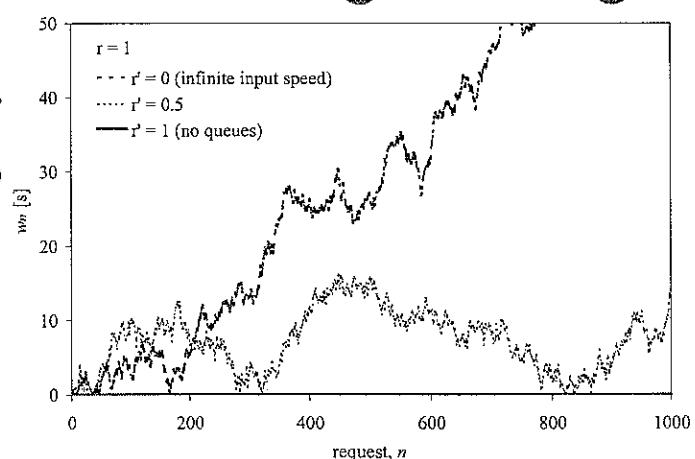
J. Konorski, Operational Research/Queuing systems

17

Impact of Input Speed (2)



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



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

J. Konorski, Operational Research/Queuing systems

18

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 Q ?
 - limited accessibility? when is a request rejected – *drop-tail*, ...?

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 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)?

Towards Richer Models (3)

- queuing discipline – what decides the order of service
 - arrival instant (FIFO, LIFO)?
 - pure chance (RANDOM)?
 - request size / current service advancement (SJF, LASF)?
 - predefined order (RR)?
 - special requirements of requests e.g., deadlines?
 - request classification for priority, fairness (HOL, WFQ)?
- service mode
 - work-conserving (busy period lasts exactly $\sum b_i/v$)
 - 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. Kowalski, Operational Research/Queuing systems

21

Steady State

In a **stationary** queuing system:

- mechanism governing request generation remains invariable in time (characteristics of the arrival stream do not change)
- processor speed v remains constant.

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).

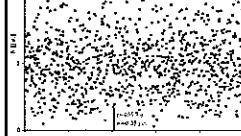
In particular,

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

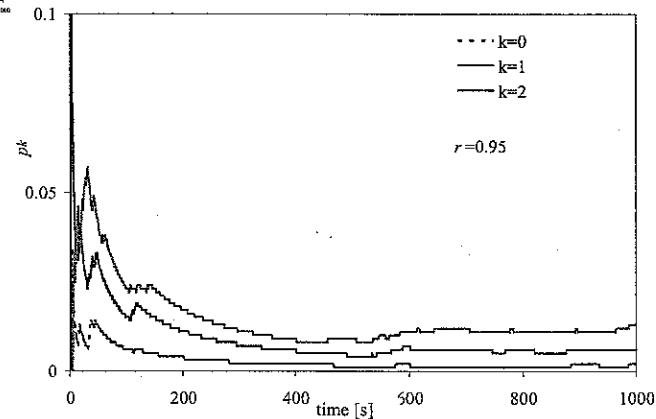
J. Kowalski, Operational Research/Queuing systems

22

Steady State (2)



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



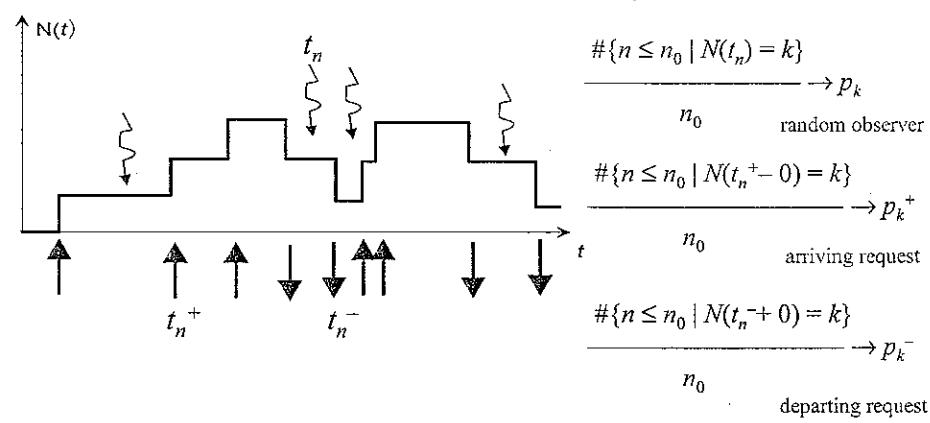
J. Konorski, Operational Research/Queuing systems

23

Steady State (3)

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

As $n_0 \rightarrow \infty$:

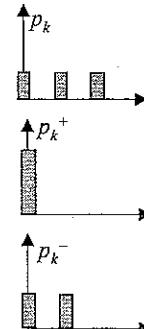
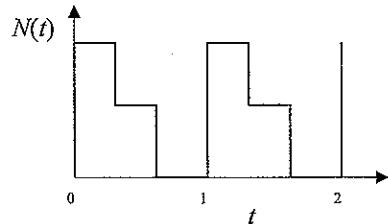


J. Konorski, Operational Research/Queuing systems

24

Steady State (4)

"Practitioners", beware!

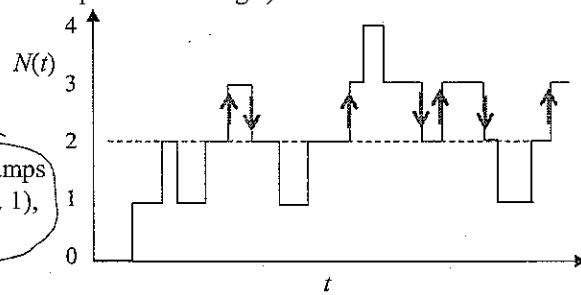


Steady State (5)

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

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

over any time period,
 $\#$ upward jumps - $\#$ downward jumps
= queue length (finite with prob. 1),
and $\#\uparrow - \#\downarrow \leq 1$



Remark: Only accepted requests count as arrivals.

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

$$\frac{p_k^+}{1-p_Q^+} = p_k^-, \quad k = 0, \dots, Q-1$$

Steady State (6)

Relevant evaluation criteria:

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

$$L = p_Q^+ \quad \begin{aligned} &\text{fraction of requests rejected (lost if no retries)} \\ &\text{due to buffer overflow} \end{aligned}$$

$$1 - p_0 \quad \begin{aligned} &\text{processor utilization} \\ &= 1 - \text{proportion of processor idle time} \end{aligned}$$

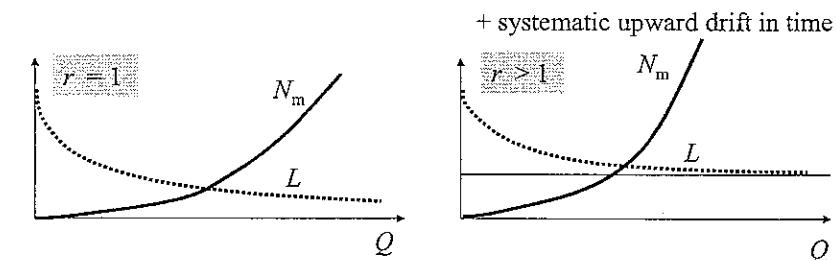
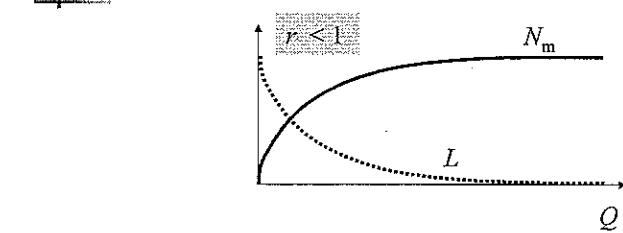
$$N_m = \lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T N(t) dt = \lim_{n_0 \rightarrow \infty} \frac{1}{n_0} \sum_{n=1}^{n_0} N(t_n) \quad \text{mean number of requests in system}$$

$$\left. \begin{aligned} \frac{w_m}{b_m} &= \frac{w_m}{\tau_m} \\ \frac{d_m}{\tau_m} &= \text{mean waiting delay} \\ &= \text{mean system delay} \end{aligned} \right\} \text{normalized to mean service time}$$

J. Konorski, Operational Research/Queuing systems

27

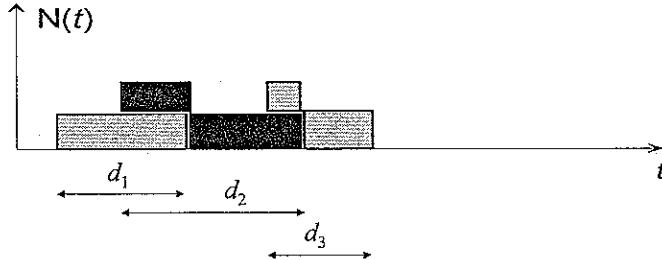
Steady State (7)



J. Konorski, Operational Research/Queuing systems

28

Little's Theorem



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

$$TN_m$$

$$(T/a_m)(1-L)d_m$$

J. Konorski, Operational Research/Queuing systems

29

Little's Theorem (2)

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

mean number of requests in system = (mean system thruput) x (mean system delay)

mean population = (mean circulation) x (mean lifetime)

Valid for any part of system:

- processor: $1 - p_0 = \frac{1-L}{a_m} \tau_m = (1-L)r$ flow conservation equation

- buffer: $N_m - (1 - p_0) = \frac{1-L}{a_m} w_m$

J. Konorski, Operational Research/Queuing systems

30

Is Queuing Theory Losing Momentum?

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

transatlantic transport	10 cruises 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. Konorski, 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. Konorski, Operational Research/Queuing systems

32

Is Queuing Theory Losing Momentum? (3)

Second, how do queues form anyway?

Not only because of $v < \infty$, 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 $v = \infty$.

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., $v > b_m/a_m$, but not much.

Meaning, allow queues at times.

Comments on Queuing Theory



Queuing theory is a mathematical analysis tool. When designing a queueuing 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

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

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)
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_m , 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_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**).

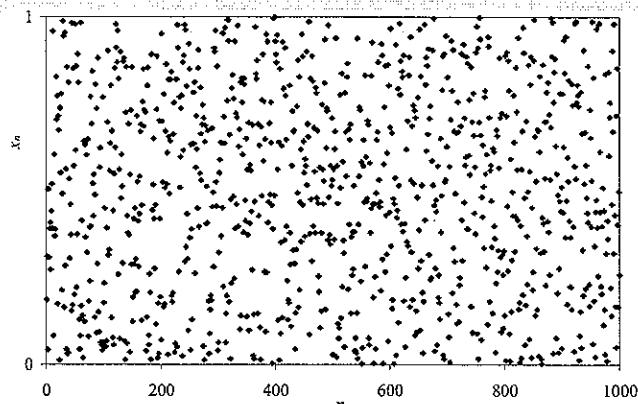
J. Konorski, Operational Research/Queuing systems

38

Random Variables and Stochastic Arrival Streams (2)

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



J. Konorski, Operational Research/Queueing systems

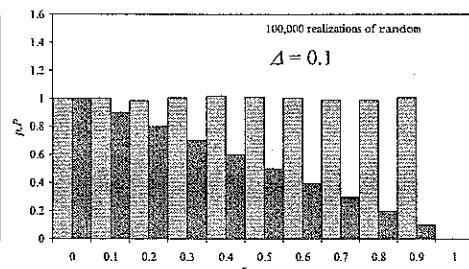
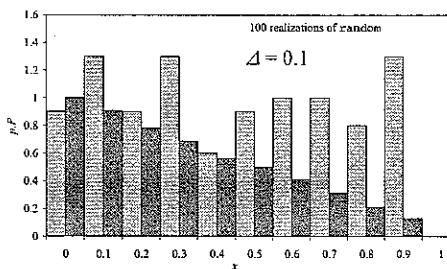
39

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 $C_i = \sum_{j \geq i} p_j$.

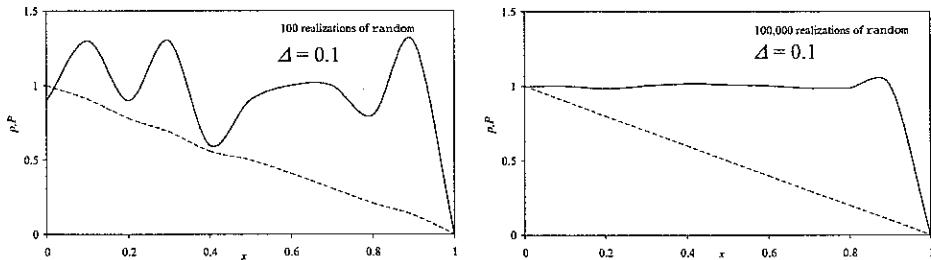


J. Konorski, Operational Research/Queueing systems

40

Empirical Distributions (2)

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



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

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

J. Konorski, Operational Research/Queuing systems

41

Theoretical Distributions

Probability density function and complementary distribution function
– histograms one would obtain taking $N \rightarrow \infty$ and $\Delta \rightarrow 0$.

Probability density function:
$$p(x) = \lim_{\Delta \rightarrow 0} \frac{\Pr[x \leq X < x + \Delta]}{\Delta}$$

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

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

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

J. Konorski, Operational Research/Queuing systems

42

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).

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

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

J. Konorski, Operational Research/Queuing systems

43

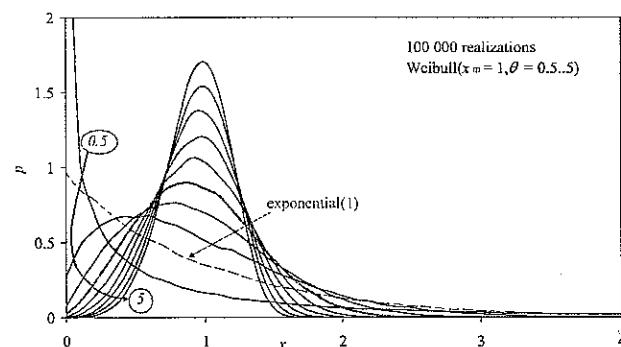
Computer Generation of Arrival Streams



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. Konorski, Operational Research/Queuing systems

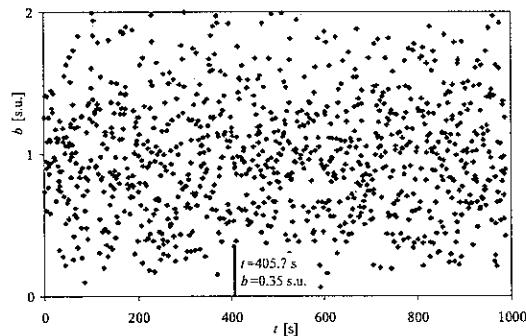
44

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. Konorski, Operational Research/Queuing systems

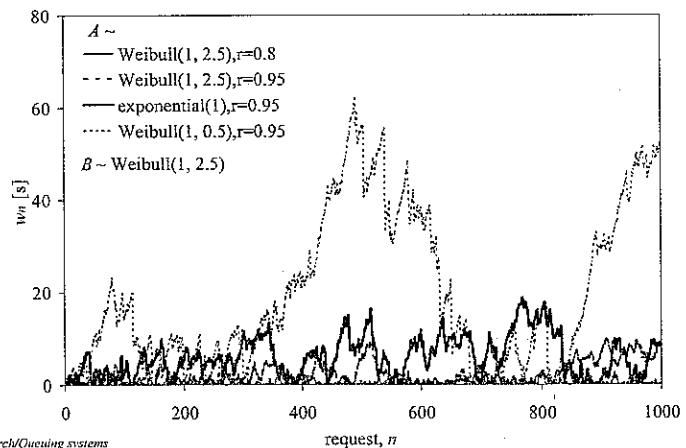
45

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 ?



J. Konorski, Operational Research/Queuing systems

46

Kendall Notation

A/B/S/Q/J ...

Examples:

M/M/1
M/D/1/10
E₃/H₂/5/10
M/M/5//20
M/G/S
G/G/S
M^X/G/5
MMPP/G/1
D+M/M/1//20
M/D/1/P*Q
M/G/1 LIFO, HOL, ...

S – # processors
Q – buffer capacity
J – request source population size

Types of distribution of A, B :

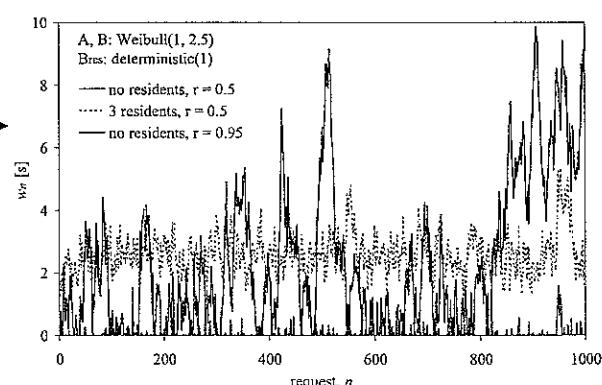
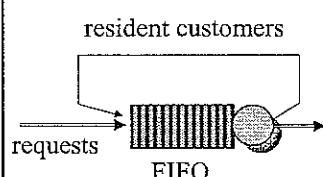
M – exponential (Markovian)
D – deterministic
E_k – Erlang of order k
H_k – hiperexponential of order k
G – general

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

J. Konarski, Operational Research/Queueing systems

47

Resident Customers



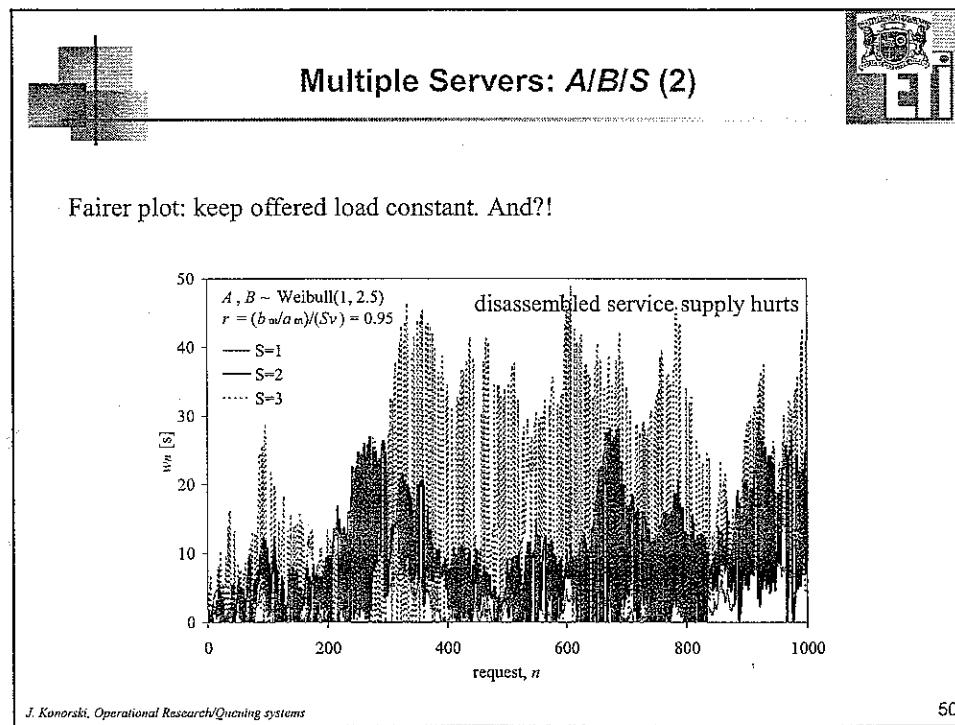
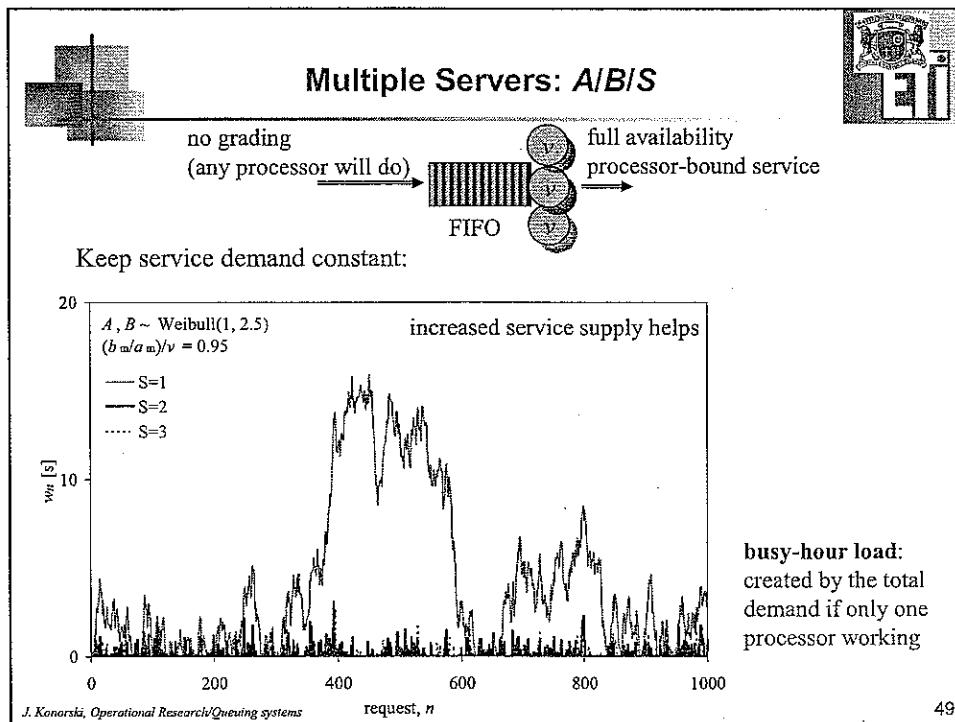
What worsens delays? Increased offered load?

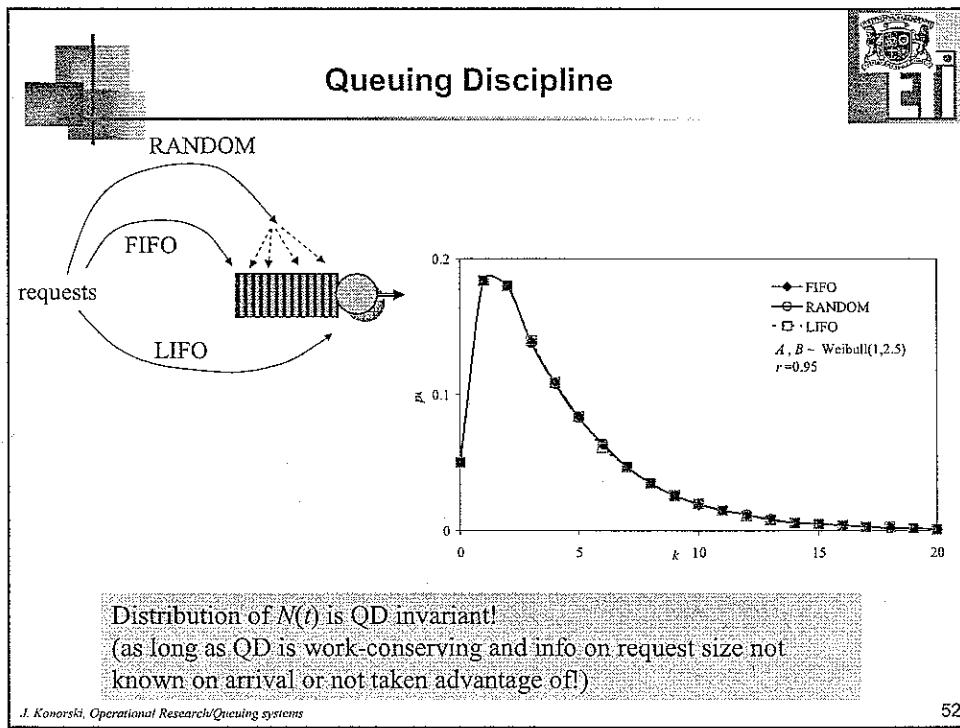
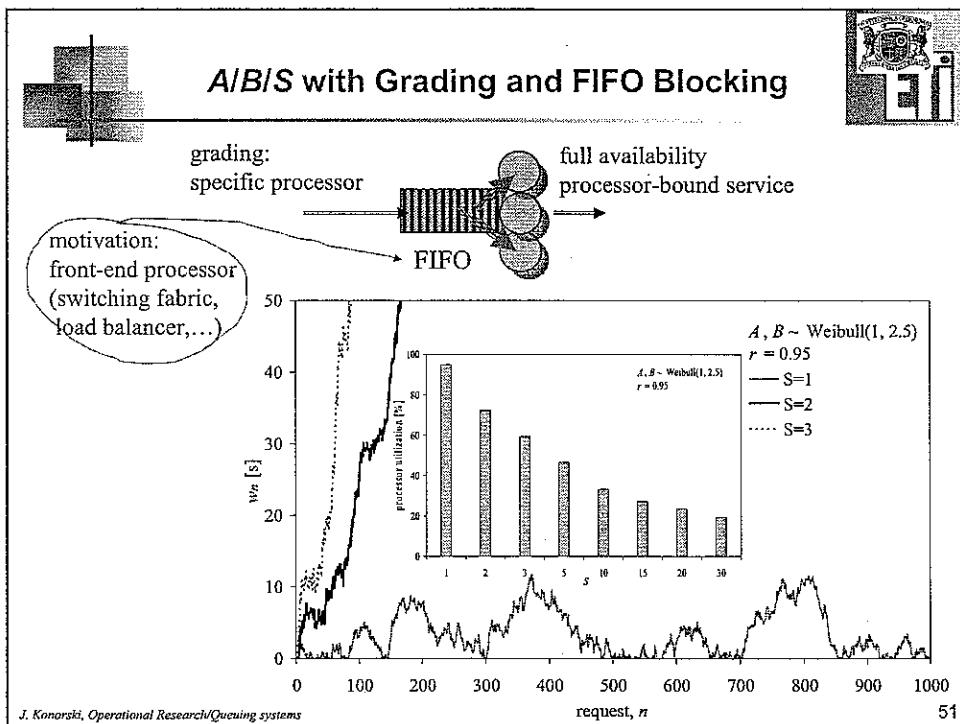
$$\frac{b_m}{a_m \nu} + \frac{\#residents \cdot b_m}{d_m \nu} \approx 0.5 + \frac{3 \cdot 1}{3.5 \cdot 2} = 0.95$$

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

J. Konarski, Operational Research/Queueing systems

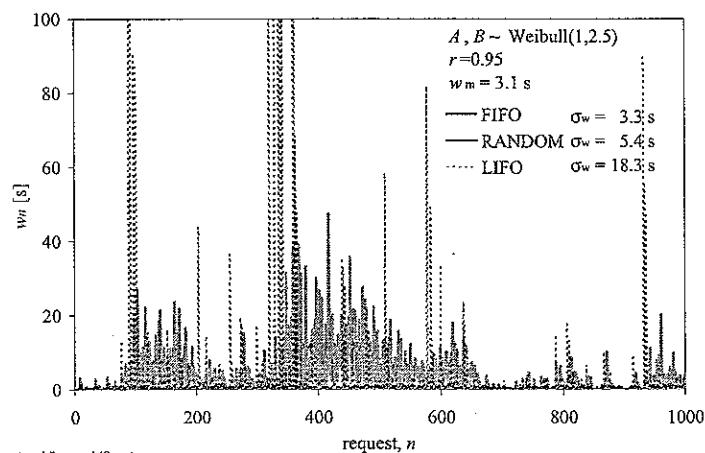
48





Queuing Discipline (2)

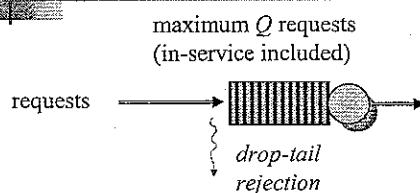
Which QD cause largest w_m ? largest σ_w ?



J. Konorski, Operational Research/Queuing systems

53

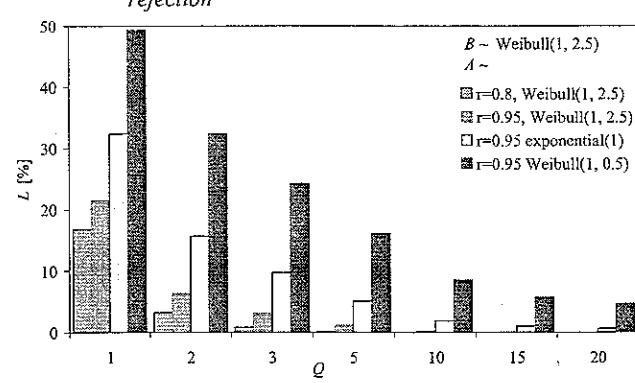
Finite Buffer Capacity: A/B/1/Q



Motivation:

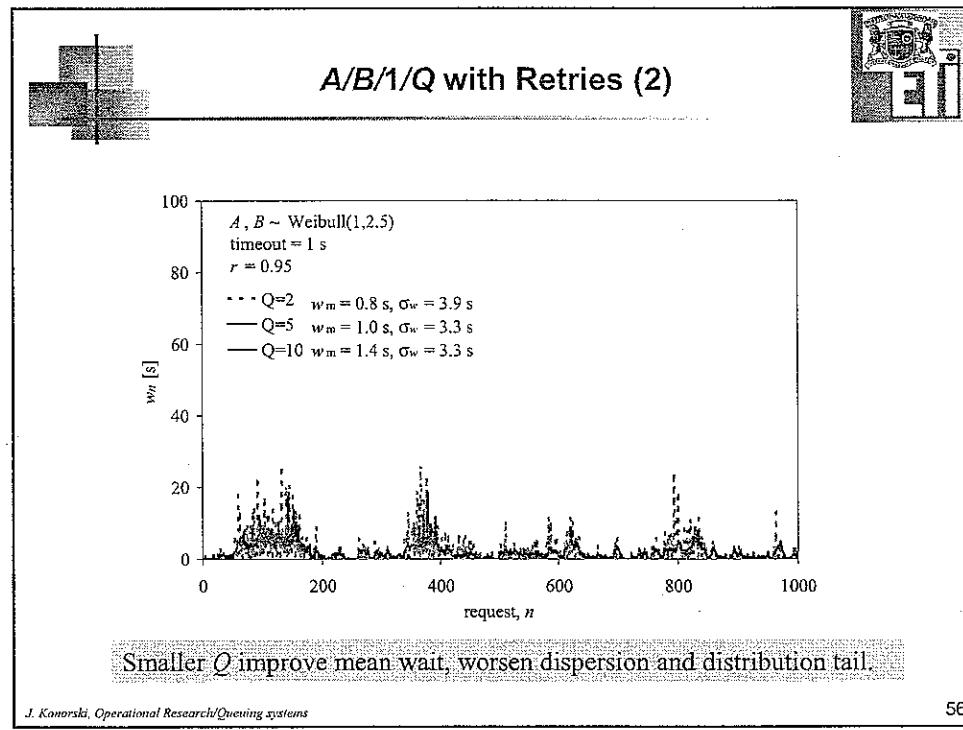
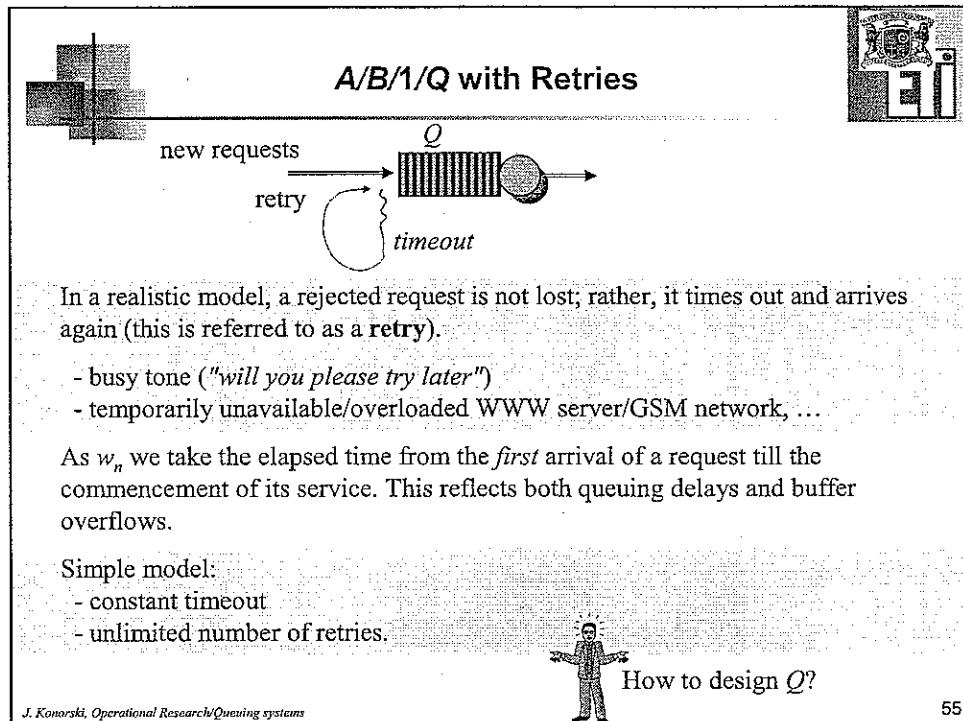
- memory cost?!
- delays, management cost!

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



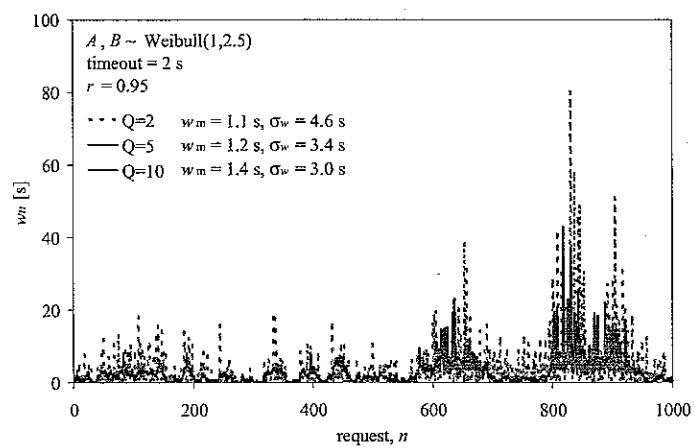
J. Konorski, Operational Research/Queuing systems

54



A/B/1/Q with Retries (3)

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

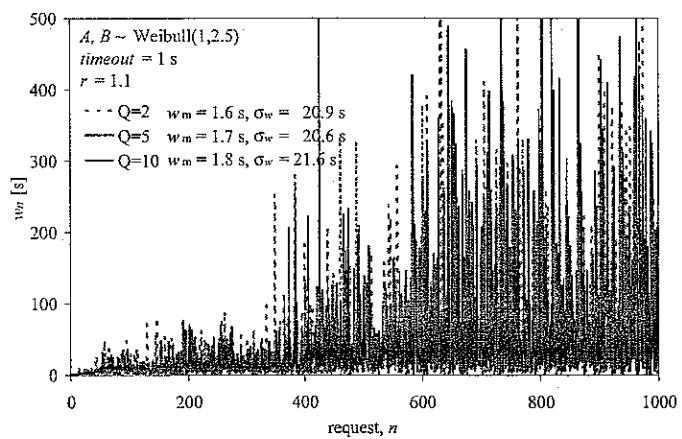


J. Konorski, Operational Research/Queueing systems

57

A/B/1/Q with Retries (4)

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



J. Konorski, Operational Research/Queueing systems

58

Common Models of Arrival Stream

- Weibull
- Bernouilli
- 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, ...

(a_n) are iid – independent, identically distributed renewal streams, often named after distribution of A

Arrival Stream: Impact of Autocorrelation

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):

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

Arrival Stream: Impact of Autocorrelation (2)

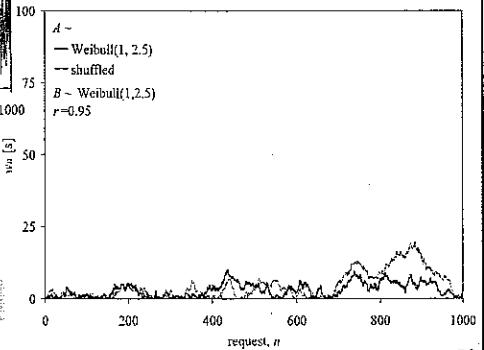
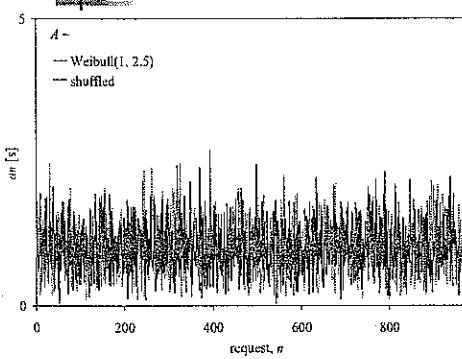
Experiment 1

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) .

Arrival Stream: Impact of Autocorrelation (3)



Shuffling makes almost no difference.

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.

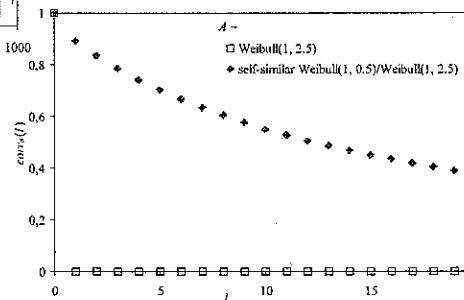
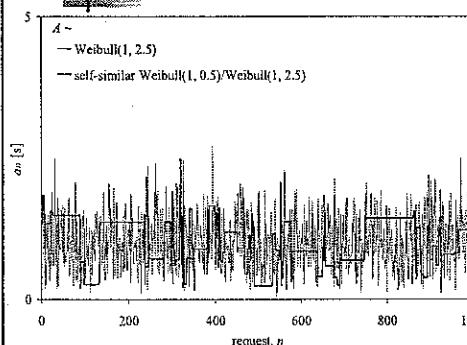
In both variants, distribution of A is the same.

However, the second variant yields (a_n) with long-range autocorrelation

– a self-similar arrival stream.

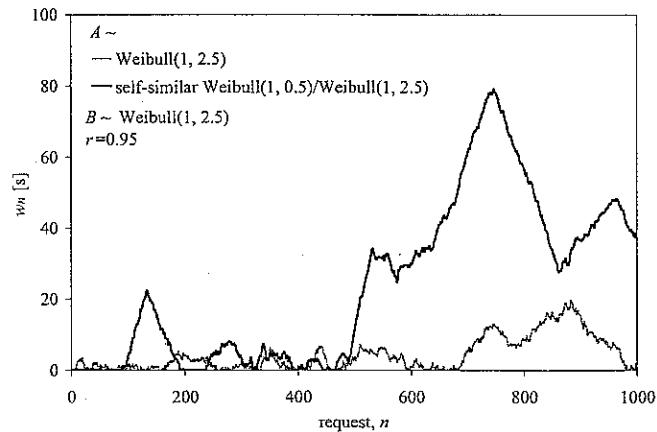
Using the same (b_n) as before, input the obtained arrival stream to a queuing system with $r = 0.95$.

Arrival Stream: Impact of Autocorrelation (5)



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!)



J. Konorski, Operational Research/Queuing systems

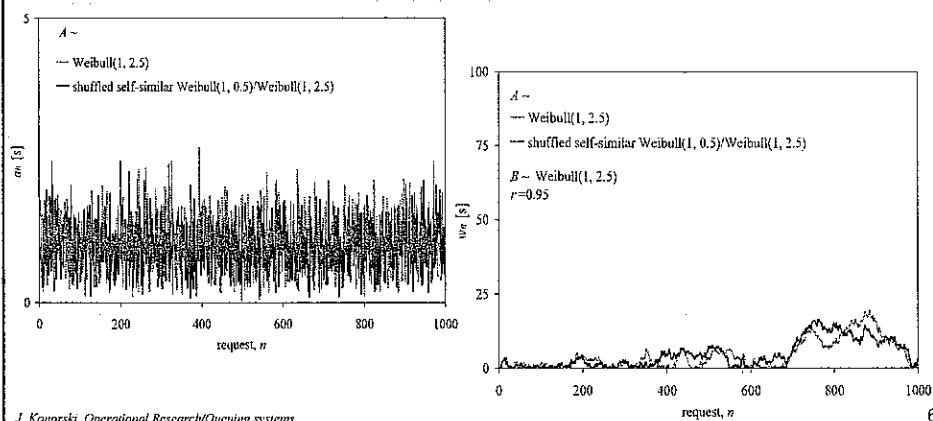
65

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...



J. Konorski, Operational Research/Queuing systems

66

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. Komorski, Operational Research/Queuing systems

67

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**.

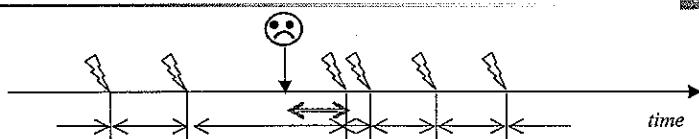
Events occur at random intervals. You arrive at a random instant.
How long do you wait for the next event?

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. Komorski, Operational Research/Queuing systems

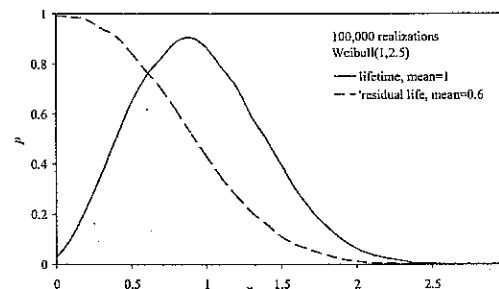
68

Renewal Arrival Stream: Residual Interval (2)



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

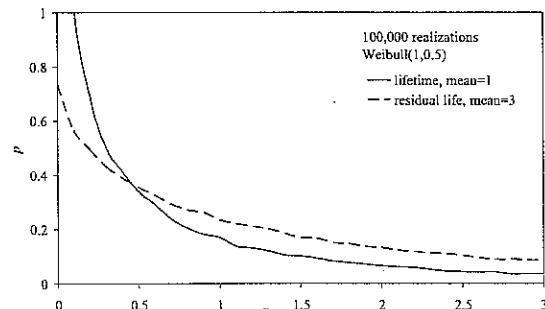
Given the distribution of inter-bus interval (lifetime), what distribution does the residual interval have? Shifted towards smaller realizations?



J. Konorski, Operational Research/Queuing systems

69

Renewal Arrival Stream: Residual Interval (3)



Paradox of residual life:

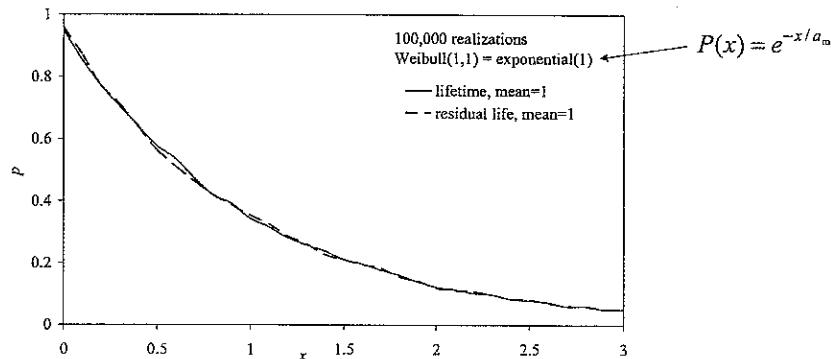
$$\bar{a}_m = \frac{a_m}{2} (1 + c_a^2) \quad c_a = \frac{\sigma_a}{a_m} - \text{coefficient of variation of } A$$

If $A \sim P(x)$ then $\tilde{A} \sim \tilde{P}(x) = P(x)/\bar{a}_m$

J. Konorski, Operational Research/Queuing systems

70

Renewal Arrival Stream: Residual Interval (4)



Exponential distribution is memoryless:

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

J. Konarski, Operational Research/Queuing systems

71

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/Queuing systems

72

Poisson Arrival Stream (2)

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

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

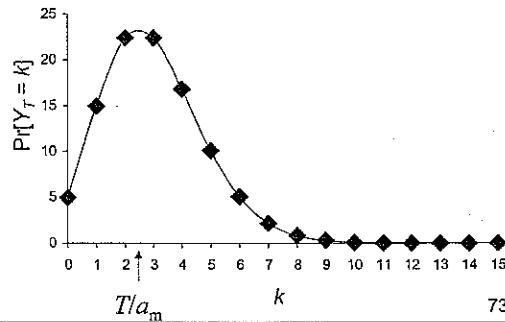


Y_T – number of requests arriving during interval of length T

$$Y_{Tm} = T/a_m$$

$$\Pr[Y_T = k] = \frac{(T/a_m)^k}{k!} e^{-T/a_m}$$

Poisson distribution



J. Konorski, Operational Research/Queueing systems

73

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.

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 distribution as does a random observer).

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

$$\text{Proof: } \Pr[N(t) = k \mid \text{request arrival in } (t, t+\Delta)] = \frac{\Pr[N(t) = k]}{\frac{\Delta}{a_m}} = \Pr[N(t) = k]$$

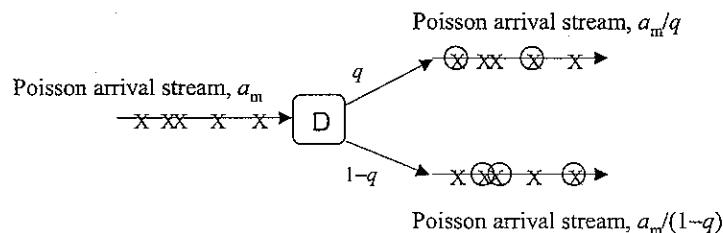
Hence, in M/G/S/Q: loss probability due to buffer overflow equals $L = p_Q^+ \equiv p_Q$.

J. Konorski, Operational Research/Queueing systems

74

Poisson Arrival Stream (4)

Random splitting:



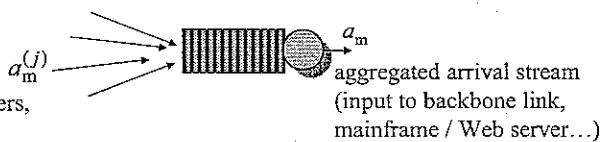
Random splitting preserves arrival stream's Poissonian nature!

Non-random splitting doesn't.

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

Aggregation of Renewal Arrival Streams

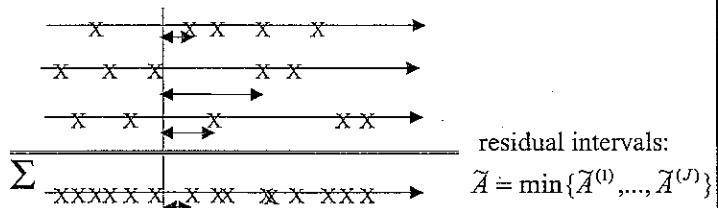
J independent components
– renewal arrival streams
(Web surfers, phone subscribers,
mobile terminals, ...)



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



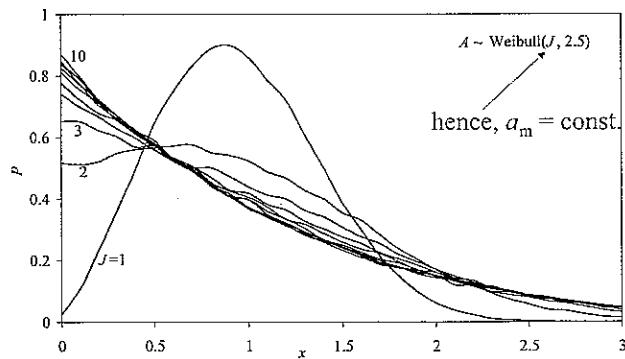
What interval distribution does the system "see"? At least a renewal stream?
Not necessarily. In general, analytic calculation difficult if not impossible.



Aggregation of Renewal Arrival Streams (2)

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

Aggregated arrival stream is Poisson (Palm theorem).



J. Konorski, Operational Research/Queuing systems

77

Aggregation of Renewal Arrival Streams (3)

Proof: omitted :)

$$\tilde{A} = \min\{\tilde{A}^{(0)}, \dots, \tilde{A}^{(J)}\}, \text{ so } \Pr[\tilde{A} \geq x] = \prod_{j=1}^J \Pr[\tilde{A}^{(j)} \geq x]$$

$J \rightarrow \infty$, 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,

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

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. Konorski, Operational Research/Queuing systems

78

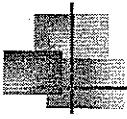



Operational Research
Queuing Systems 3: Markovian Models

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, ...)

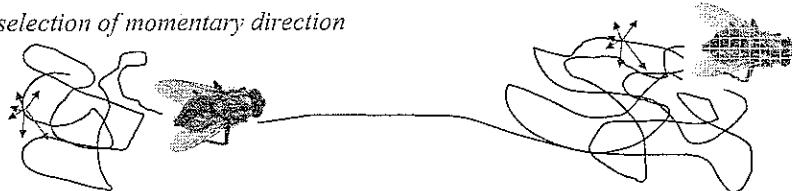
$$\text{state}(t + \Delta) = f[\text{state}(t), \phi]$$

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

Markovian Systems (3)

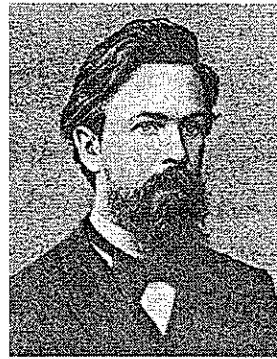
- The fly...

selection of momentary direction



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

Markovian Systems (4)



Andrei A. Markov (1856-1922)

Russian mathematician

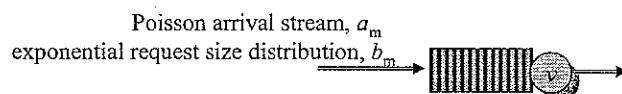
investigated stochastic processes of finite memory, now named after him,
that model many natural and man-made phenomena

J. Konorski, Operational Research/Queuing systems

83

Markovian Systems (5)

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



Practical impact stems from:

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

J. Konorski, Operational Research/Queuing systems

84

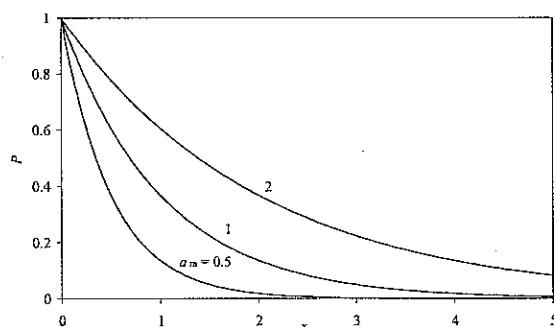
Markovian Systems (6)

v – (constant) processor speed

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

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

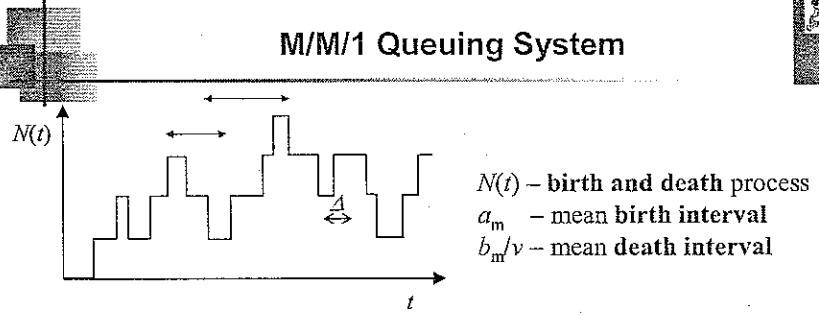
$$\Pr[\text{service time} \geq x] = e^{-x/(b_m/v)}$$



J. Konorski, Operational Research/Queuing systems

85

M/M/1 Queuing System



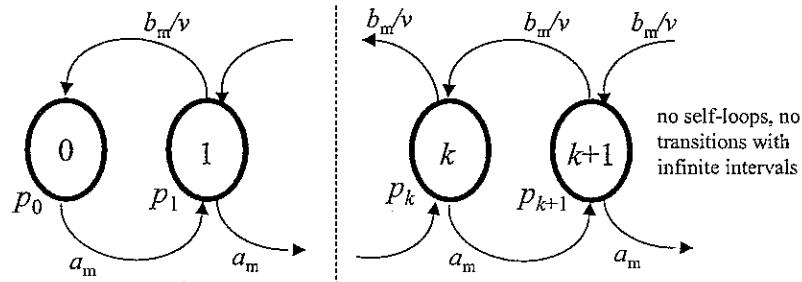
Let $p_k(t) = \Pr[N(t) = k]$. We know that $\lim_{t \rightarrow \infty} p_k(t) = p_k$ (statistical equilibrium). State distribution (p_k) can be derived from **birth and death equations**.

- Exponential distribution has no probability mass at 0
 \Rightarrow 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

J. Konorski, Operational Research/Queuing systems

86

M/M/1 Queuing System (2)



Suggestive "hydraulic" analogue :

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

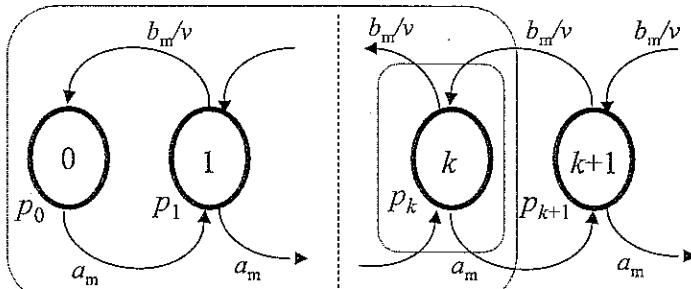
J. Konorski, Operational Research/Queuing systems

87

M/M/1 Queuing System (3)

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!



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

J. Konorski, Operational Research/Queuing systems

88

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:

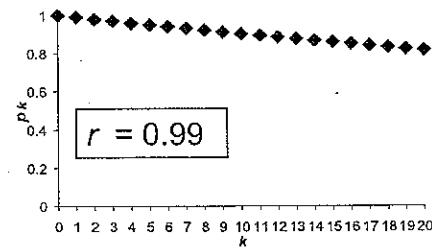
$$N_m = \frac{r}{1-r}$$

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

$$w_m = d_m - \tau_m$$

very suggestive!

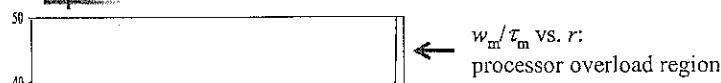
The most frequent queue length?!



J. Konorski, Operational Research/Queuing systems

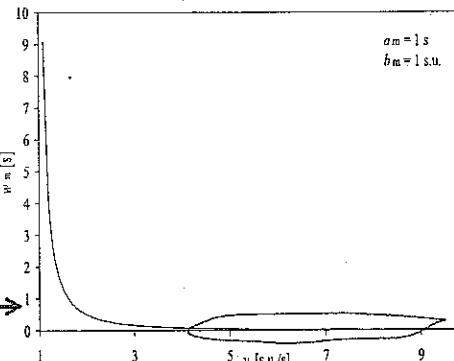
89

M/M/1 Queuing System (5)



w_m / t_m vs. r :
processor overload region

w_m vs. v :
processor overdimensioning region



Distribution of waiting time does depend on QD.

For FIFO, will be found with a little richer math... shortly.

J. Konorski, Operational Research/Queuing systems

90

M/M/... Systems

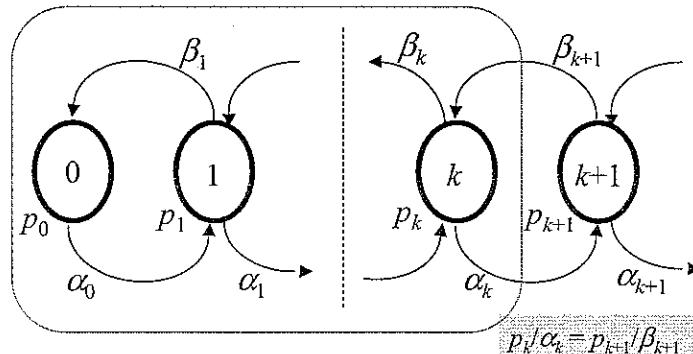
Birth-and-death processes help to analyze far richer and more realistic Markovian models of queuing systems featuring e.g.,

- 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!*

M/M/... Systems (2)

Make mean transition interval state-dependent:



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

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

M/M/1/Q System

$a_m, b_m \rightarrow Q \quad r = (b_m/a_m)/\nu$

Arrivals stop at $k=Q$:

$$\alpha_k = \begin{cases} a_m, & \text{if } k < Q \\ \infty, & \text{if } k \geq Q \end{cases} \quad \beta_k \equiv b_m/\nu$$

$$p_k = p_0 r^k, \quad k = 0, \dots, Q$$

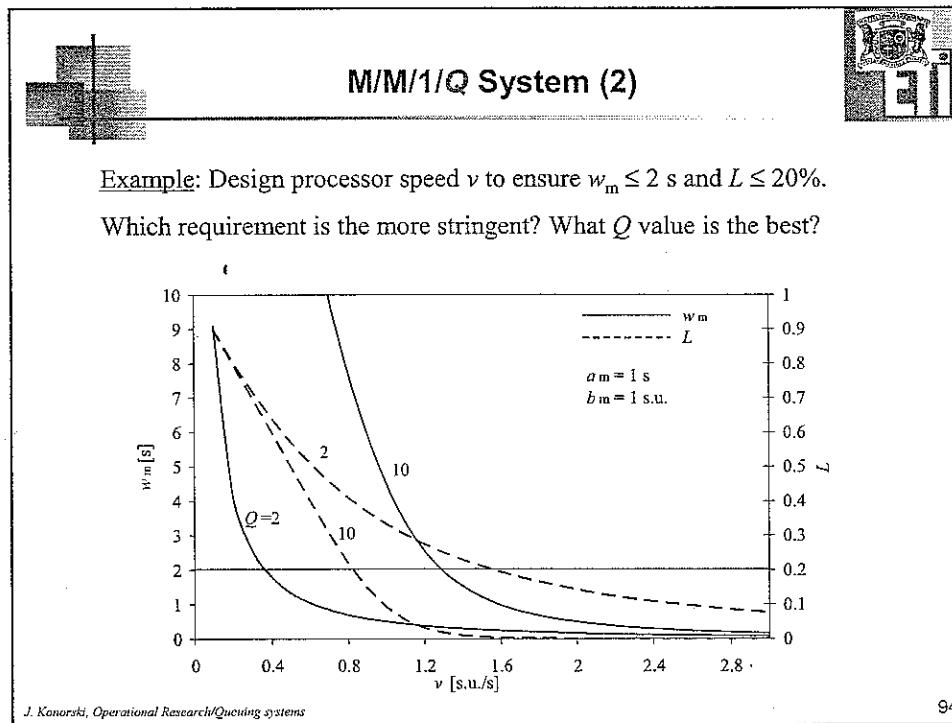
PASTA

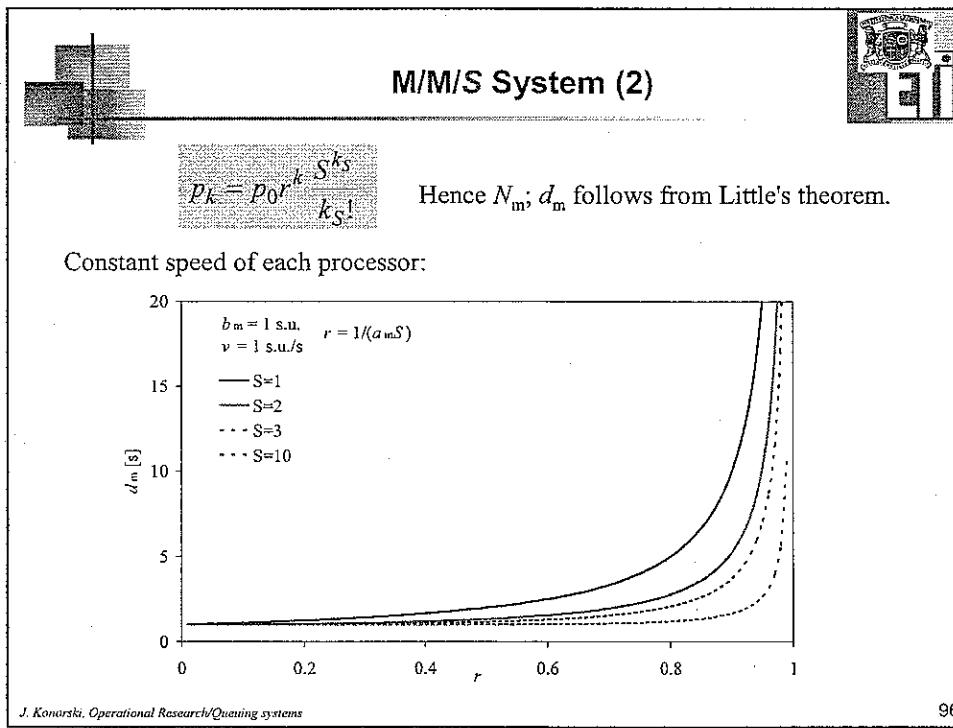
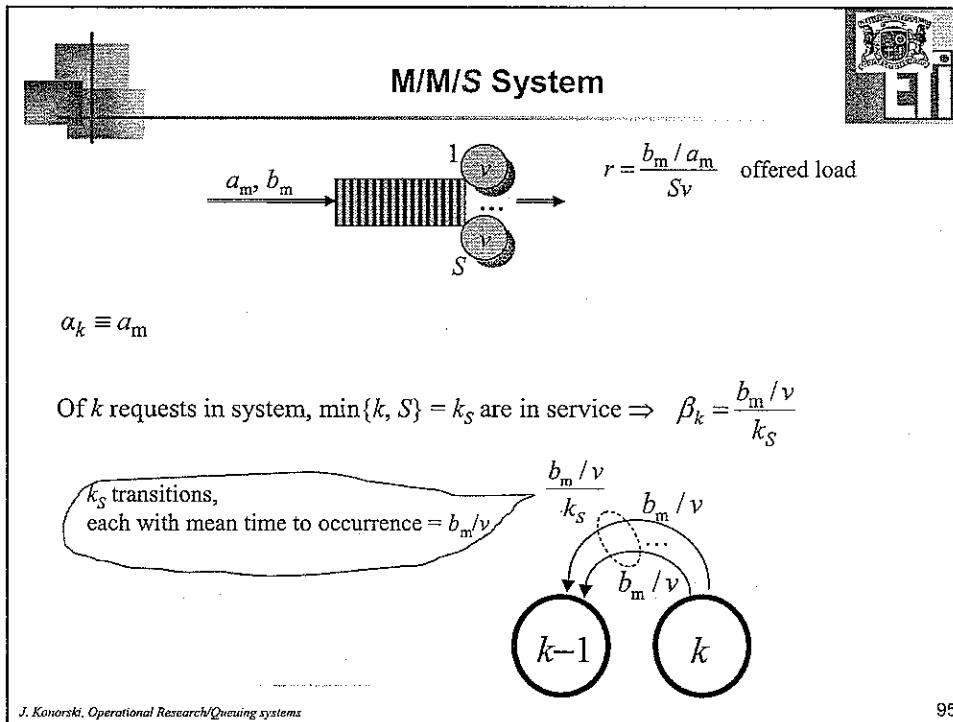
$$L = p_Q^{-1} = p_Q = \begin{cases} \frac{1-r}{1-r^{Q+1}} r^Q, & \text{if } r \neq 1 \\ \frac{1}{Q+1}, & \text{if } r = 1 \end{cases}$$

$$d_m = \frac{N_m}{1-L} \quad w_m = d_m - b_m/\nu$$

J. Konarski, Operational Research/Queuing systems

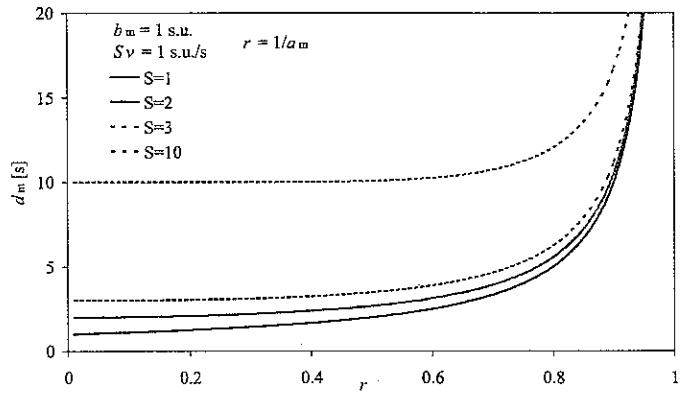
93





M/M/S System (3)

Constant speed of the processor pool:



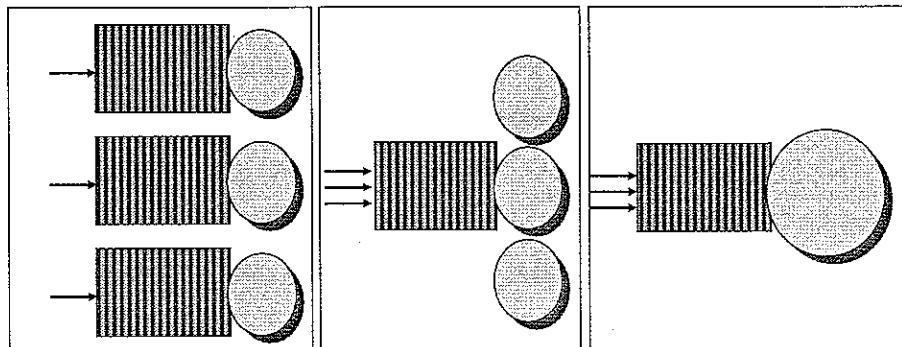
(We have seen already that disassembled service supply hurts!)

J. Konorski, Operational Research/Queuing systems

97

On Benefits of Buffer Sharing and Pooled Service Supply

$r = 0.667$



$3 \times M/M/1$

$$d_m = 1 \text{ s}$$

$M/M/3$

$$d_m = 0.48 \text{ s}$$

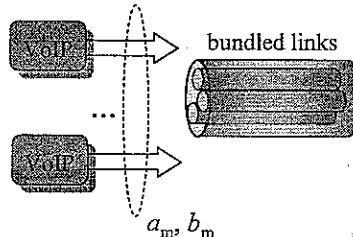
$M/M/1$

$$d_m = 0.33 \text{ s}$$

J. Konorski, Operational Research/Queuing systems

98

M/M/S/S System



no waiting room

v – single link capacity

$\rho = b_m / (a_m v)$ – busy-hour load

$\alpha_k, \beta_k, k = 0, \dots, S$ same as for M/M/S

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

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

- 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/)

J. Konorski, Operational Research/Queuing systems

99

M/M/ ∞ System

A huge hipermarket 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 hipermarket, 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).

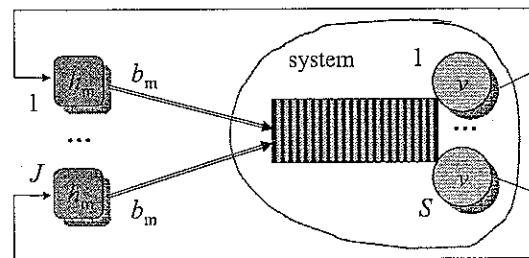
Special case of M/M/S/S with $S \rightarrow \infty$, therefore $p_0 = \frac{1}{\sum_{k=0}^{\infty} \frac{\rho^k}{k!}} = e^{-\rho}$

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

J. Konorski, Operational Research/Queuing systems

100

M/M/S//J System



Model:

- $\Pr[\text{think time} \geq x] = e^{-x/h_m}$
- population of terminals $J (> S)$

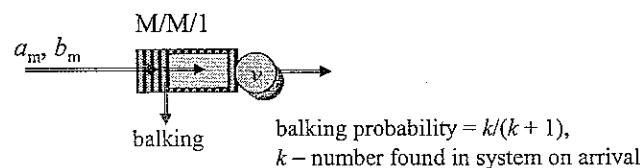
When k requests in system,

- $J - k$ 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. Konorski, Operational Research/Queuing systems

101

Impatient Requests: Balking



Assuming a Markovian system model, find the fraction of balking requests.

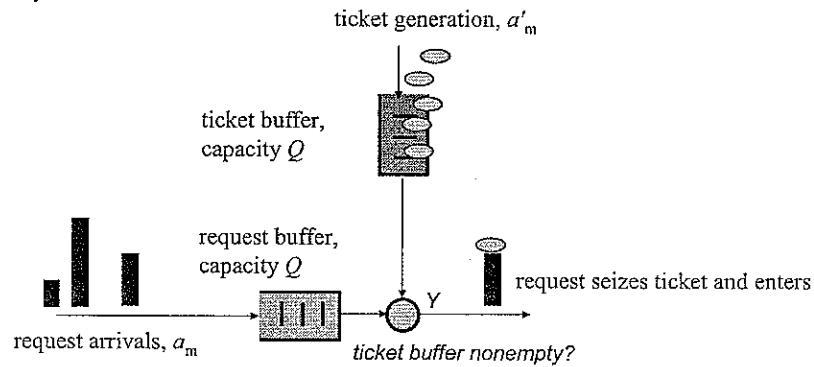
$$\alpha_k = \frac{a_m}{1 - \frac{k}{k+1}} = a_m \cdot (k+1), \quad \beta_k = b_m / v \Rightarrow (p_k) - \text{Poisson distribution (!)}$$

$$L = p_0 \cdot 0 + p_1 \cdot \frac{1}{2} + p_2 \cdot \frac{2}{3} + \dots + p_k \cdot \frac{k}{k+1} + \dots$$

J. Konorski, Operational Research/Queuing systems

102

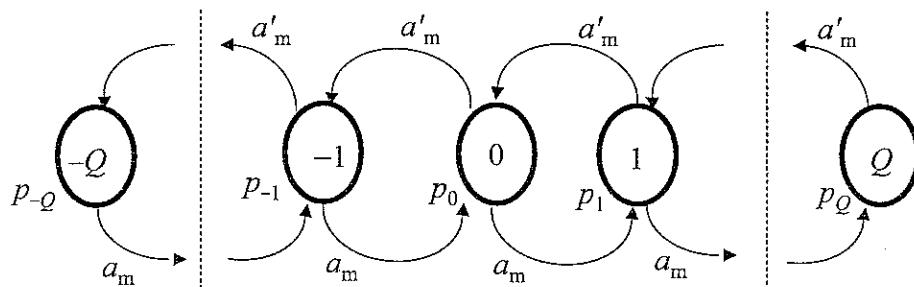
Ticket Counter



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

Ticket Counter (2)

Arrived request increments state, generated token decrements state:



Is statistical equilibrium possible with infinite Q ?

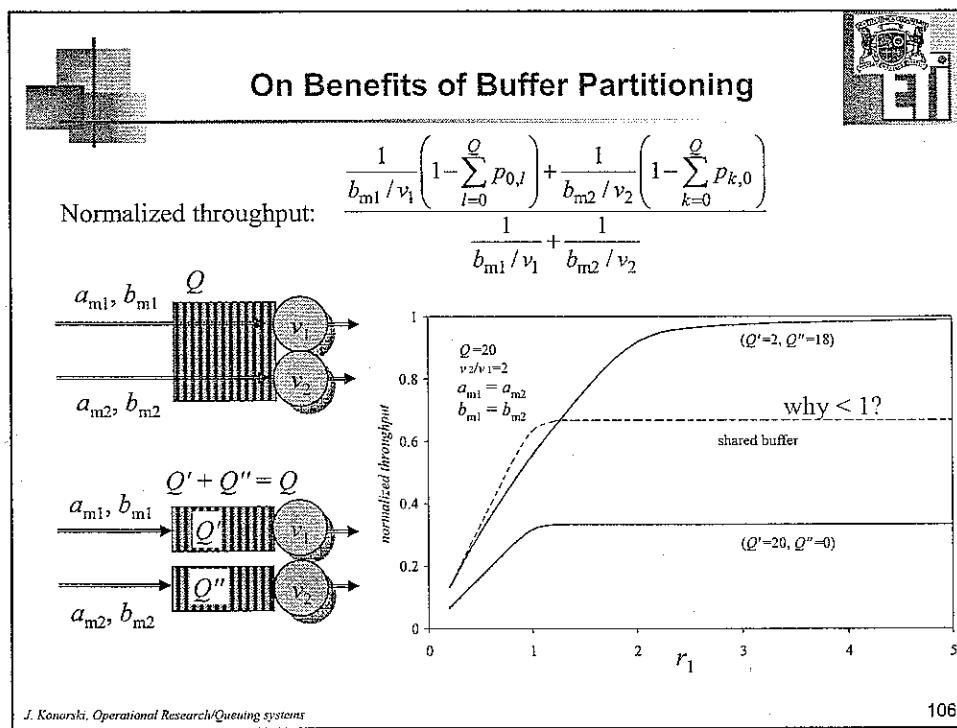
A/B/2/Q with Grading and No FIFO Blocking

Flow balance of probability mass:

$$p_{k,l}/a_{m1} + p_{k,l}/a_{m2} + p_{k,l}/(b_{m1}/v_1) + p_{k,l}/(b_{m2}/v_2) = p_{k-1,l}/a_{m1} + p_{k,l}/a_{m2} + p_{k+1,l}/(b_{m1}/v_1) + p_{k-1,l}/(b_{m1}/v_1)$$

Solution by guessing and inspection: $p_{k,l} = p_{0,0} r_1^k r_2^l$ where $r_1 = \frac{b_{m1}}{a_{m1} v_1}$, $r_2 = \frac{b_{m2}}{a_{m2} v_2}$

J. Konarski, Operational Research/Queueing systems 105



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. Konorski, Operational Research/Queuing systems

107

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/calculatorsalgebra.htm
- For some $X^*(s)$ all the above fail. Numerical algorithms exist, but due to inherent numerical instability, none is universal.

J. Konorski, Operational Research/Queuing systems

108

Laplace Transform (3)

[www.pe.tamu.edu/blasingame/data/P620_reference/P620_Lectures_\(pdf\)/P620_Mod1_Math/P620_Mod1_ML_05_LaplaceTrans.pdf](http://www.pe.tamu.edu/blasingame/data/P620_reference/P620_Lectures_(pdf)/P620_Mod1_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

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

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:
 $P(x) = e^{-x/(b_m/v)}$

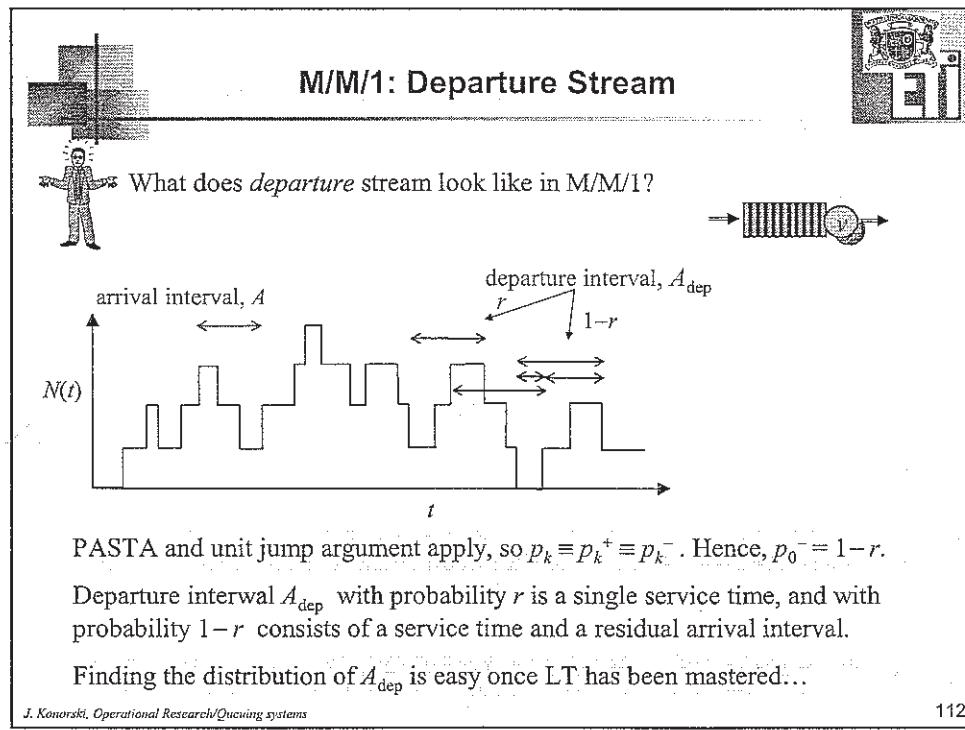
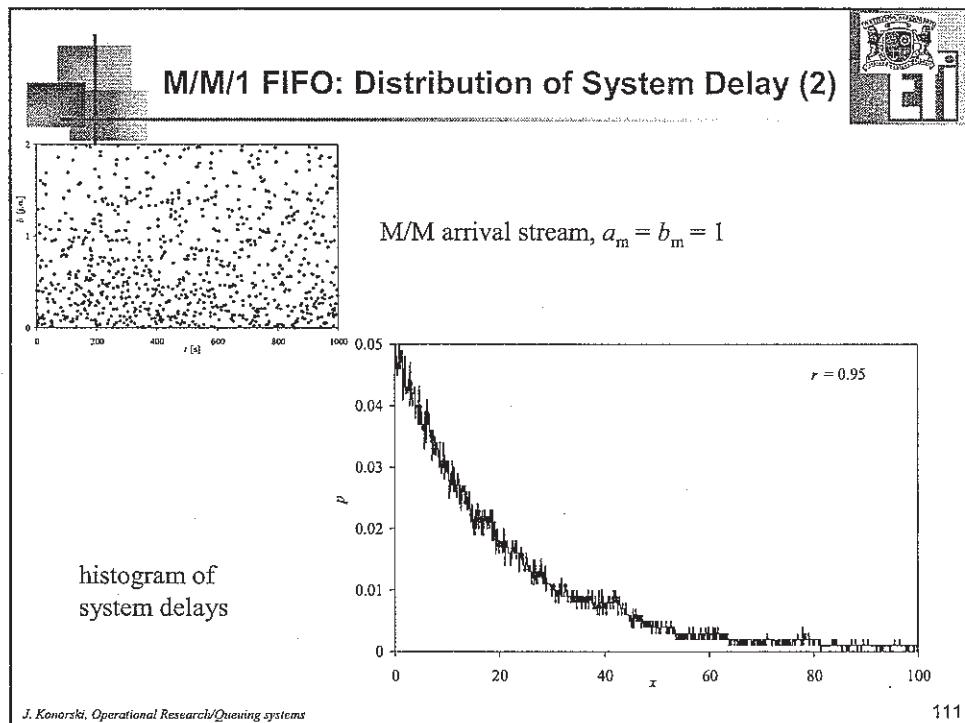
Hence, LT of system delay is $\left(\frac{1}{\frac{b_m}{v}s + 1}\right)^{k+1}$.

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

Averaging over k and using LT linearity gives:

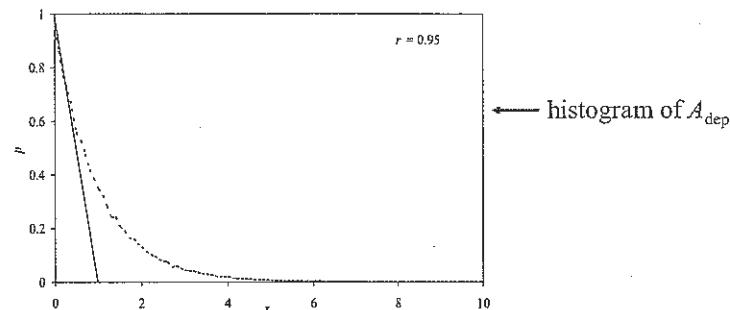
$$D^*(s) = \sum_{k=0}^{\infty} p_k \left(\frac{1}{\frac{b_m}{v}s + 1}\right)^{k+1} = \sum_{k=0}^{\infty} (1-r)r^k \left(\frac{1}{\frac{b_m}{v}s + 1}\right)^{k+1} = \frac{1}{\frac{r}{1-r} \frac{b_m}{v}s + 1}$$

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



M/M/1: Departure Stream (2)

$$A_{wyj}^*(s) = r \frac{1}{\frac{b_{sr}}{\nu} s + 1} + (1 - r) \frac{1}{a_{sr} s + 1} \cdot \frac{1}{\frac{b_{sr}}{\nu} s + 1} = \frac{1}{a_{sr} s + 1} = A^*(s)$$



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

Important in tandem configurations...



J. Konorski, Operational Research/Queuing systems

113

Operational Research

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

Jerzy Konorski
jekon@eti.pg.gda.pl

J. Konorski, Operational Research/Queuing systems

114

Motivation

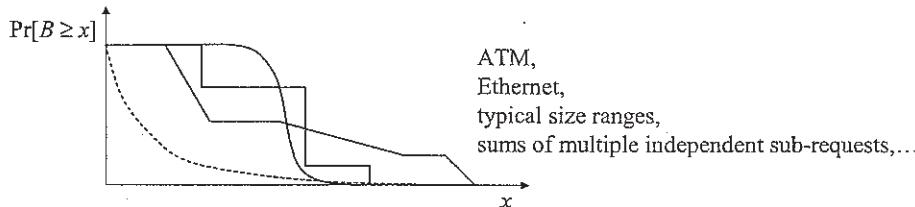
Poisson arrival stream still assumed: $A \sim \Pr[A \geq x] = e^{-x/\alpha_m}$

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!

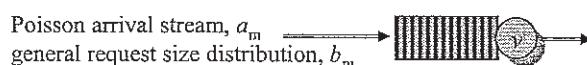
This leads to M/G/... models with general (arbitrary) B .



J. Konorski, Operational Research/Queueing systems

115

M/G/1 Analysis



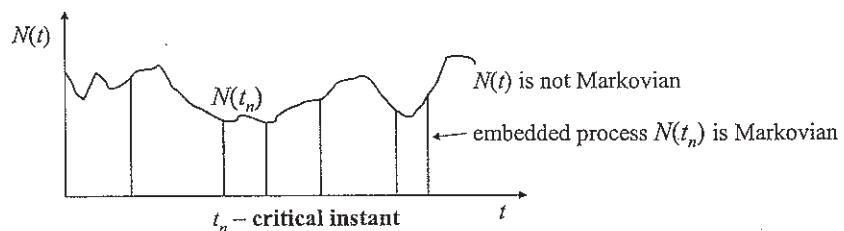
- M/G/1: 1 processor, infinite buffer capacity
- $B \sim P(x)$, arbitrary distribution with mean b_m
- $r = b_m / (a_m v)$, where v – processor speed
- $N(t)$ does not possess the Markov property
- two-dimensional process $[N(t), \text{residual service time}(t)]$ does possess it, but is analytically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

J. Konorski, 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)?

J. Konorski, Operational Research/Queuing systems

117

M/G/1 Analysis (3)

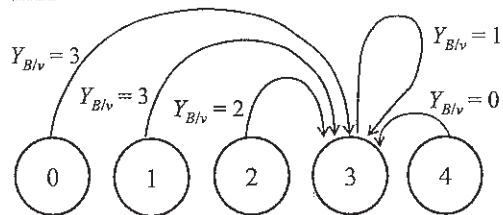
For M/G/1 the choice is easy:

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

Let $Y_{B/\nu}$ = number of request arrivals during one service time (B/ν).

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

$$N(t_n) = \max[N(t_{n-}) - 1, 0] + Y_{B/\nu} \quad (\text{a stochastic equation})$$



J. Konorski, Operational Research/Queuing systems

118

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

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

Q = ∞ & unit jump argument *PASTA*

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

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/...

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

- Main result for M/G/1:

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

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].

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:

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

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

power series expansion

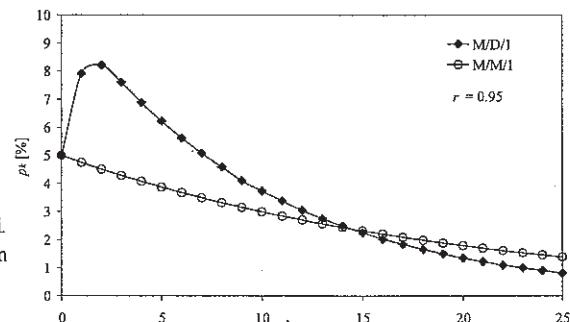
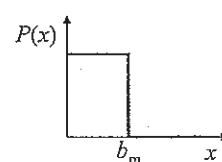
So $p_k = (1-r)r^k$, as yielded by Markovian analysis.

Special Case: M/D/1

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

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

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



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

J. Konorski, Operational Research/Queuing systems

123

P-K Formula

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

$$\frac{W_m}{N_m} = \frac{1 + c_b}{2} \frac{\tau_m}{1 - r}$$

$c_b = \sigma_b/b_m$ – coefficient of variation of request size distribution

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:

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

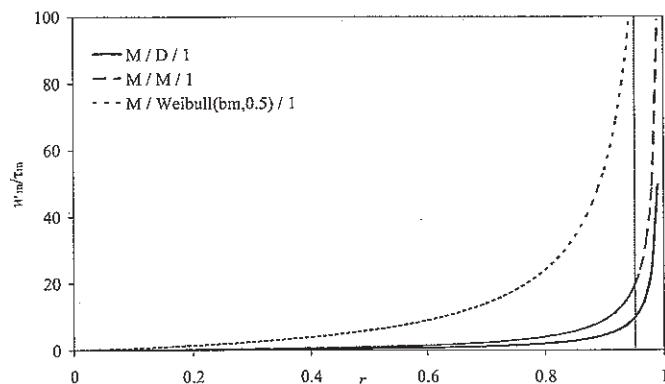
J. Konorski, Operational Research/Queuing systems

124

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 ($c_b \approx 3.317$) !

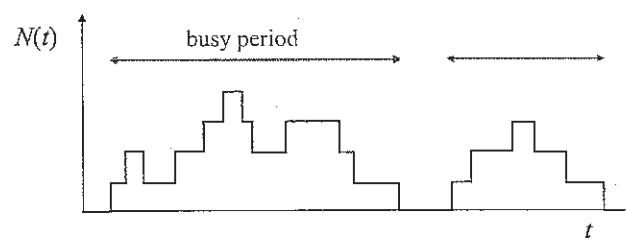


J. Konorski, Operational Research/Queueing systems

125

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 :

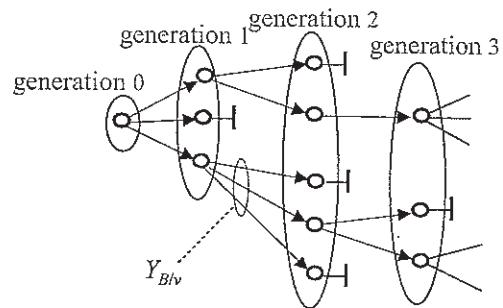
$$\omega = \Pr[\text{busy period eventually terminates}]$$

J. Konorski, Operational Research/Queueing systems

126

Busy Period (2)

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



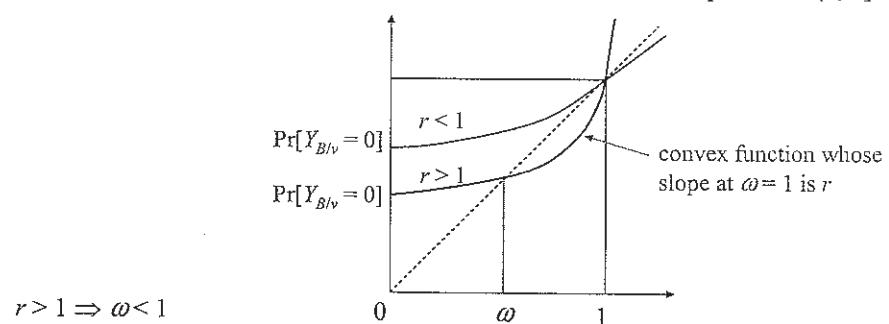
J. Konarski, Operational Research/Queuing systems

127

Busy Period (3)

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

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



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

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

$$r < 1 \Rightarrow \omega = 1 \text{ surprising, given that } r = Y_{B/v,m}$$

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

$$z_m = \tau_m + Y_{B/v,m} z_m = \tau_m + (\tau_m/a_m) z_m = \tau_m + r z_m$$

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

J. Konarski, Operational Research/Queuing systems

129

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.

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

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

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

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

So, after simple algebra,

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

J. Konarski, Operational Research/Queuing systems

130

M/G/1 Waiting Delay Distribution (2)

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

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

 What is the *distribution* of waiting delay?

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

Examples:

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

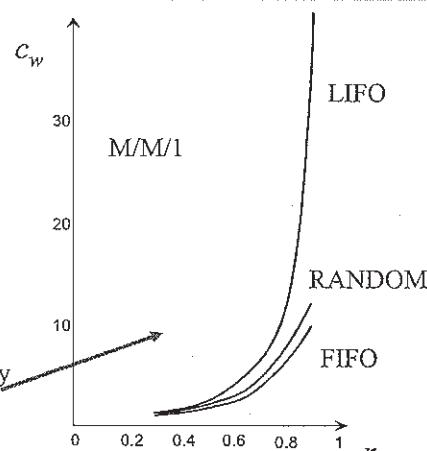
$$\text{- M/D/1: } B^*(s) = e^{-b_m s} \Rightarrow W^*(s) = 1 - r + r \frac{\frac{(b_m/v)s}{1 - e^{-(b_m/v)s}} - r}{1 - e^{-(b_m/v)s}}$$

J. Konorski, Operational Research/Queuing systems

131

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).



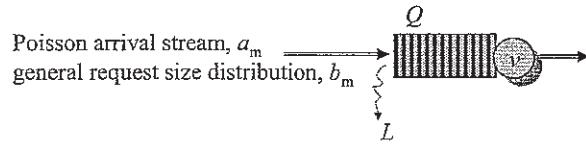
Delay distribution is not!

For various QDs, may differ dramatically in the coefficient of variation.

J. Konorski, Operational Research/Queuing systems

132

M/G/1/Q



- 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

M/G/1/Q (2)

Begin again with a stochastic equation:

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

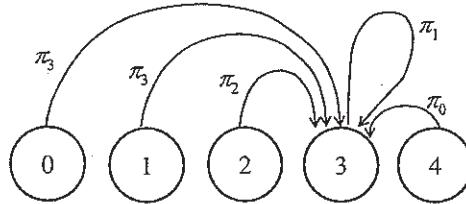
This time, "translation" into distributions will be done directly, without LT.

Distribution of $Y_{B/v}$ has been calculated before:

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

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

M/G/1/Q (3)



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

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

$$p_0^- + \dots + p_{Q-1}^- = 1$$

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

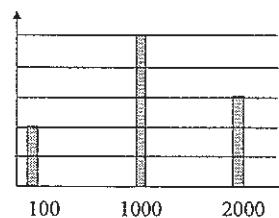
$$L = p_Q = 1 - \frac{1}{r + p_0}$$

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?

M/G/1/Q – Example (2)

Distribution of $Y_{B/N}$:

$$\begin{aligned}\pi_i &= 30\% \cdot \frac{(1000/500)/4}{i!} e^{-(1000/500)/4} \\ &+ 50\% \cdot \frac{(2000/500)/4}{i!} e^{-(2000/500)/4} \\ &+ 20\% \cdot \frac{(100/500)/4}{i!} e^{-(100/500)/4}\end{aligned}$$

$$\begin{aligned}\pi_0 &= 0.5561 \\ \pi_1 &= 0.2844\end{aligned}$$

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

Equations:

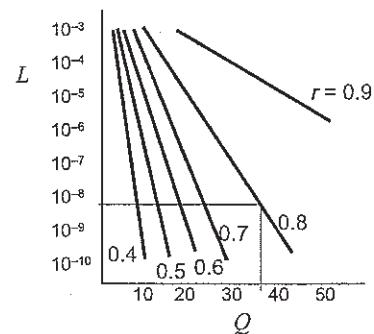
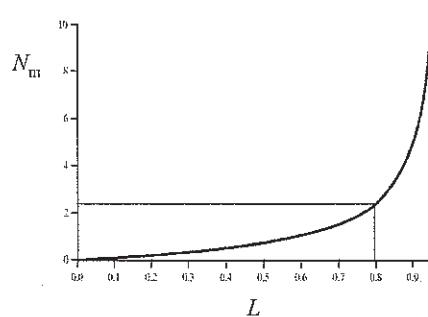
$$\begin{aligned}p_0^- &= p_0^- \pi_0 + p_1^- \pi_0 \\ p_1^- &= p_0^- \pi_1 + p_1^- \pi_1 + p_2^- \pi_0 \\ p_0^- + p_1^- + p_2^- &= 1\end{aligned}$$

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

J. Konorski, Operational Research/Queuing systems

137

M/G/1/Q Design



- fix tolerable N_m and L (e.g., 2.5 and 10^{-3})
- 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. Konorski, Operational Research/Queuing systems

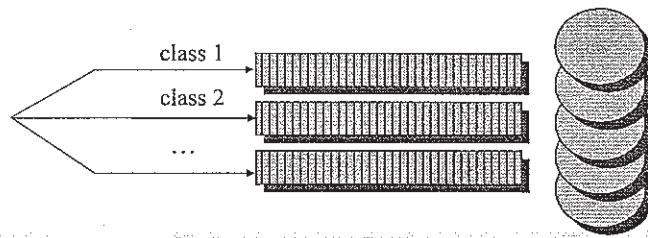
138

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").

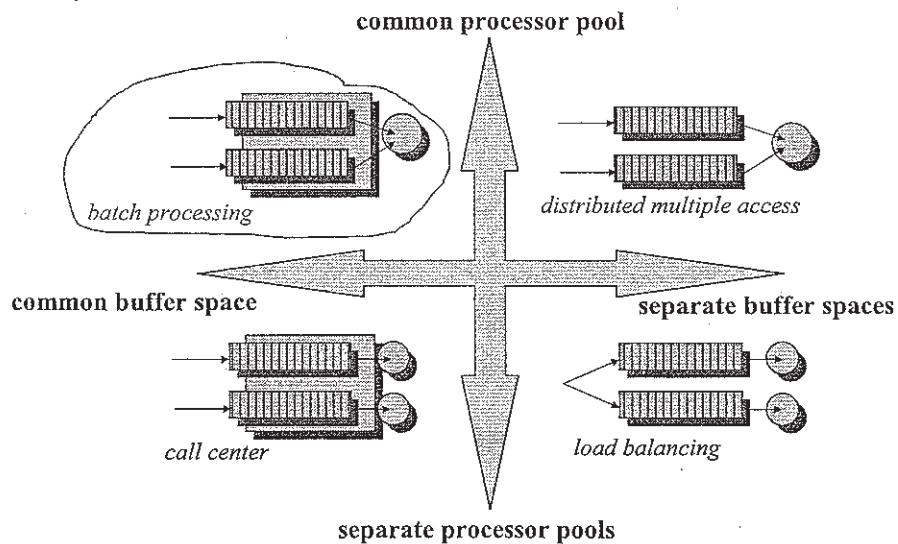
Classification can be external (user-imposed) or internal (system-imposed). Each class has its own **logical queue**.

(convention: lower class number designates higher priority)

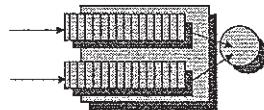


We still speak of *one* queuing system as long as some processors / buffer space / request sources are common to arrival streams of different classes.

Priority Queuing (2)



Priority Queuing (3)



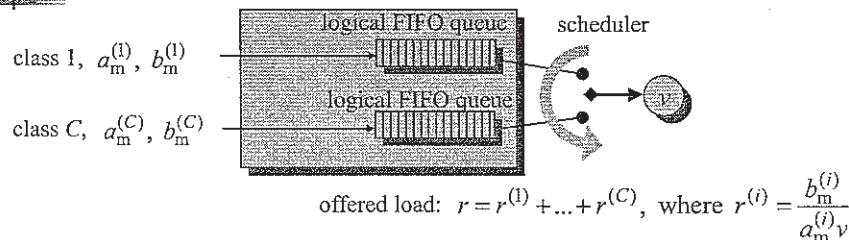
Model:

- 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

Priority Queuing (4)



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..."

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 $w_m^{(i)}$
- } FIFO with classes $j \leq i$ only

$$w_m^{(i)} = w_m^{(\text{FIFO})} \Big|_{r^{(1)} + \dots + r^{(i)}} + \sum_{j < i} \frac{w_m^{(i)} b_m^{(j)}}{a_m^{(j)}} \quad \text{linear equation, easy to solve for } w_m^{(i)}$$

$r^{(j)}$

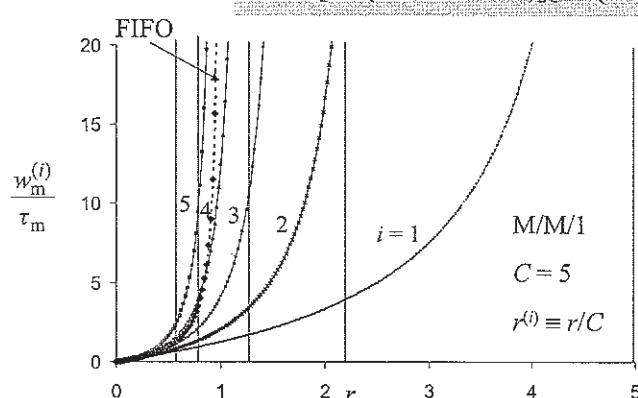
J. Konorski, Operational Research/Queuing systems

143

Priority Queuing: HOL (2)

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. Konorski, Operational Research/Queuing systems

144

Priority Queuing: SJF

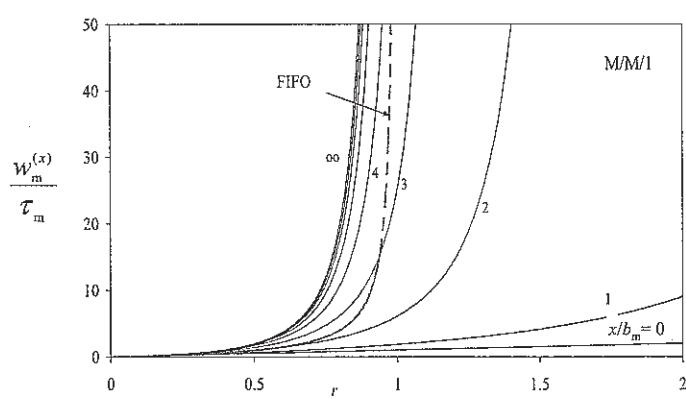
HOL makes no use of information on request sizes.

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 $r^{(x)} = x/(a_m v)$

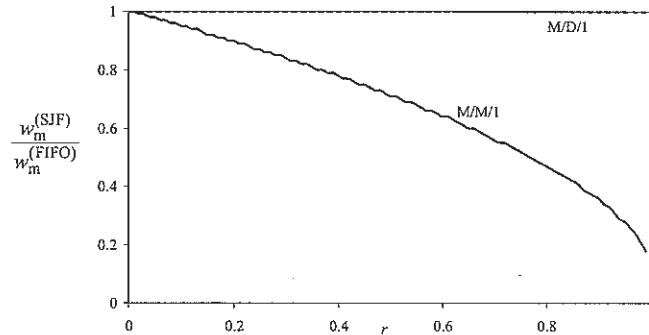
Priority Queuing: SJF (2)

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



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.

Priority Queuing: Dynamic HOL

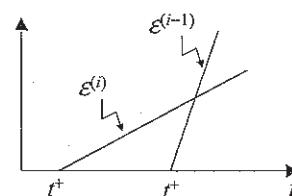
Problems with HOL and SJF:

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

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

$$\text{priority}^{(i)}(t) = \varepsilon^{(i)} \cdot (t - t^+)$$

with arbitrary $\varepsilon^{(C)} \geq \dots \geq \varepsilon^{(1)}$.



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

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

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 } unable to overtake
- 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,

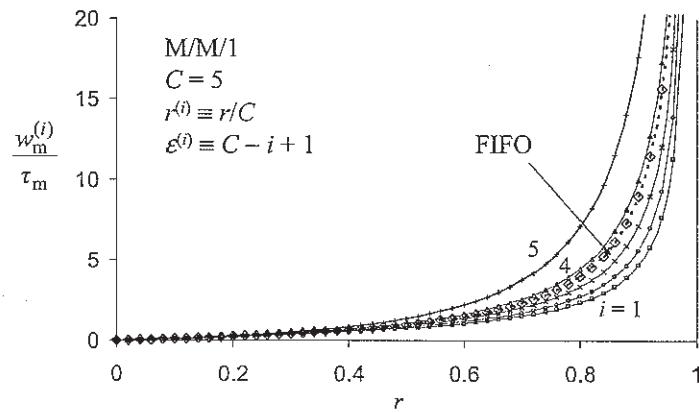
$$w_m^{(i)} = \frac{w_m^{(\text{FIFO})} - \sum_{j>i} r^{(j)} \left(1 - \frac{\varepsilon^{(j)}}{\varepsilon^{(i)}}\right) w_m^{(j)}}{1 - \sum_{j>i} r^{(j)} \left(1 - \frac{\varepsilon^{(i)}}{\varepsilon^{(j)}}\right)}$$

• recurrence relationship
• only *relative* magnitudes of $\varepsilon^{(i)}$ matter

J. Konorski, 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

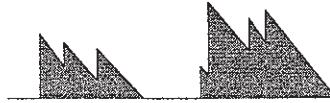
J. Konorski, Operational Research/Queuing systems

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 $U(t)$ (unfinished work) process!



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

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.

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:

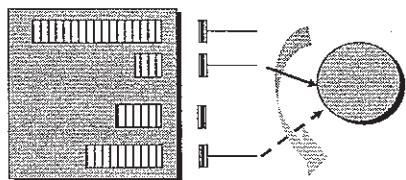
Let $C = 2$, $r^{(1)} = 0.1$, $r^{(2)} = 0.8$.

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?

Processor Sharing



Model

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

J. Konarski, Operational Research/Queuing systems

153

Processor Sharing (2)

Switching between requests can take various forms:

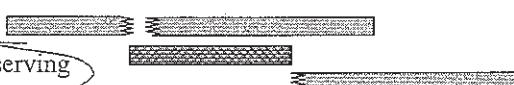
- preemption with *abandonment* – destroys work!



- ...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. Konarski, Operational Research/Queuing 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.

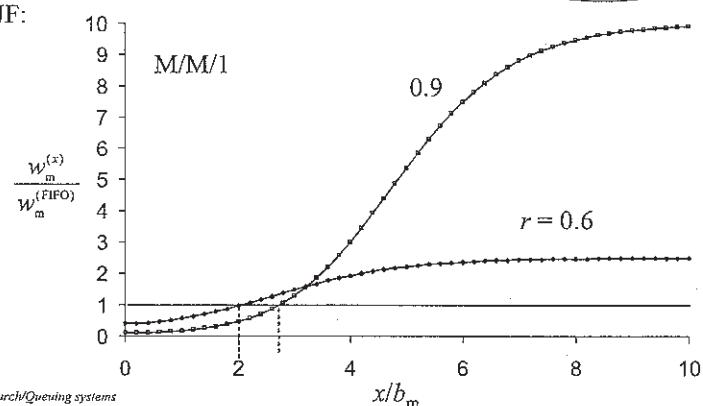
Processor Sharing (4)

Basic operational characteristic:

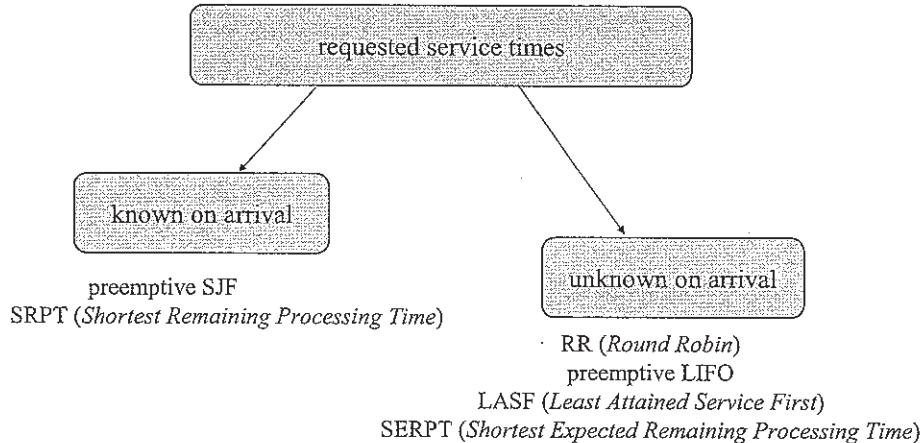
typically normalized to τ_m or $w_m^{(FIFO)}$

typically normalized to b_m

E.g., for SJF:



Processor Sharing (5)



J. Konorski, Operational Research/Queuing systems

157

Preemptive HOL / SJF

Preemptive HOL:

- FIFO within class
- arrival of higher-priority request preempts lower-priority service
- upon service completion, class queues inspected from highest priority same as in HOL

Class i requests "don't feel" classes $j < i$:

$$d_m^{(i)} = w_m^{(\text{FIFO})} \Big|_{r^{(1)} + \dots + r^{(i)}} + \frac{b_m^{(i)}}{\nu} + \sum_{j < i} \underbrace{\frac{d_m^{(i)} b_m^{(j)}}{a_m^{(j)}}}_{\text{service of preempting requests}} \cdot r^{(j)}$$

Preemptive SJF – continuous version: class $x \equiv B \in (x, x+\Delta)$

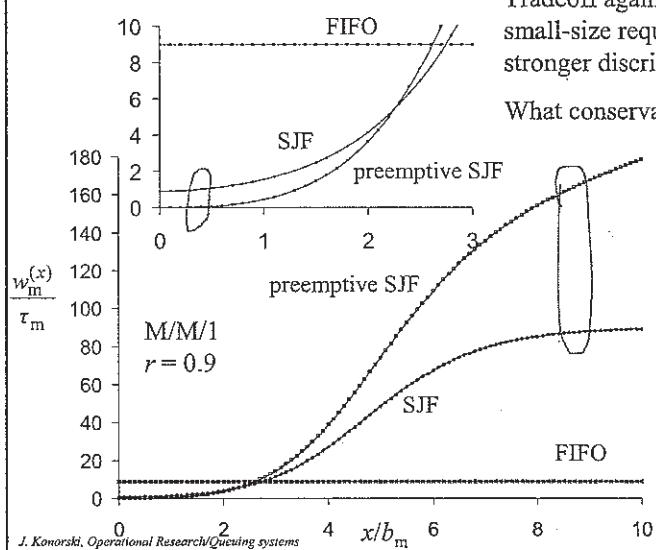
J. Konorski, Operational Research/Queuing systems

158

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?



159

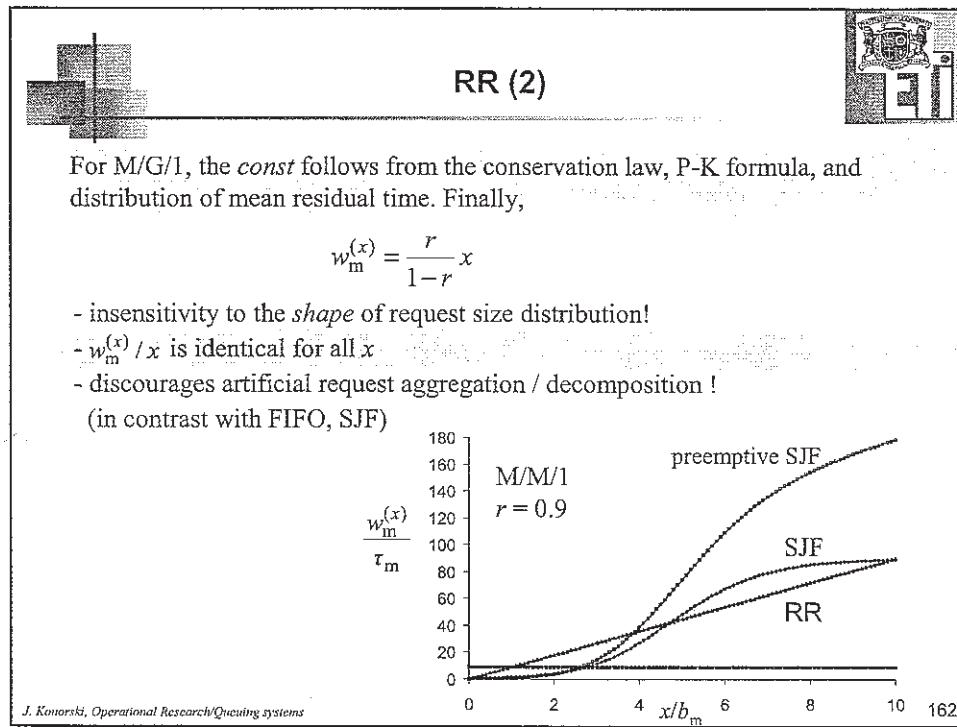
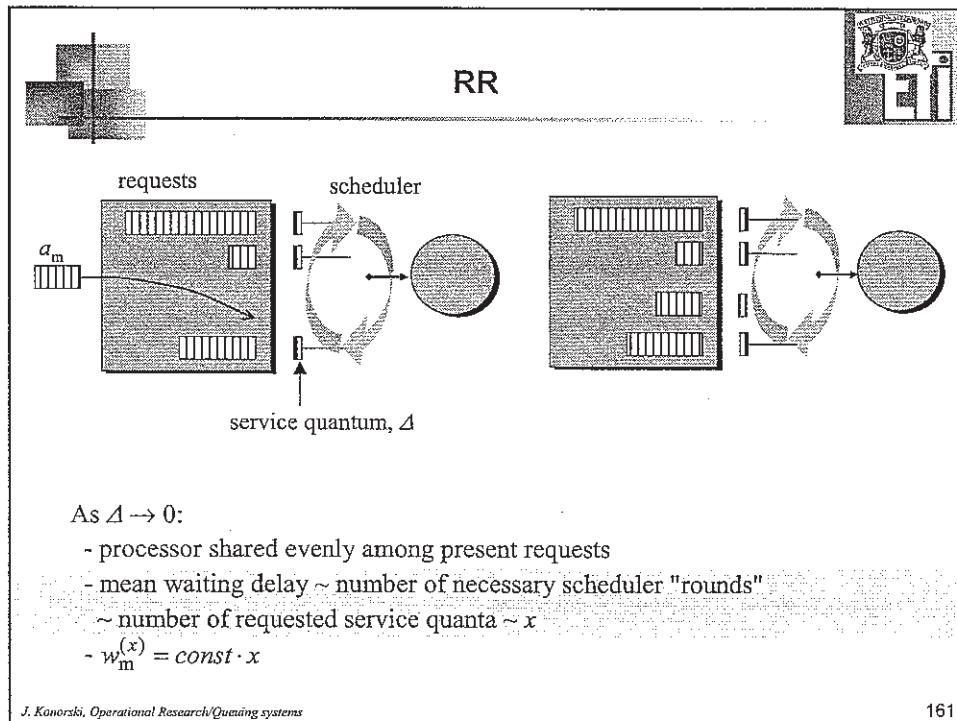
Processor Sharing: Conservation Law

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

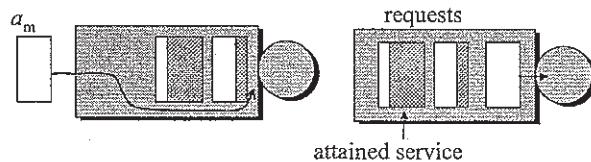
$$\int_0^{\infty} \frac{w_m(x)}{\tau_m} P(x) dx = w_m^{(\text{FIFO})}$$

$\Pr[B \geq x]$ – 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 .



Preemptive LIFO



In M/G/1:

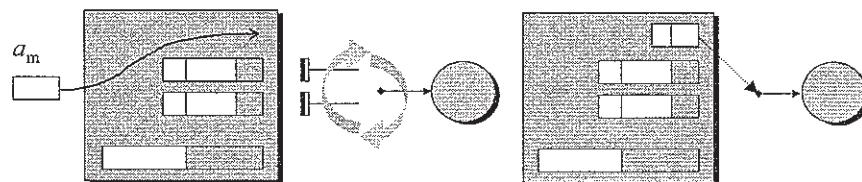
Mean waiting delay of request of size x is composed of service preemptions by arriving requests (on average, $Y_{x,m} = x/a_m$), each of which creates a busy period of mean duration $\frac{\tau_m}{1-r}$.

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

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

LASF

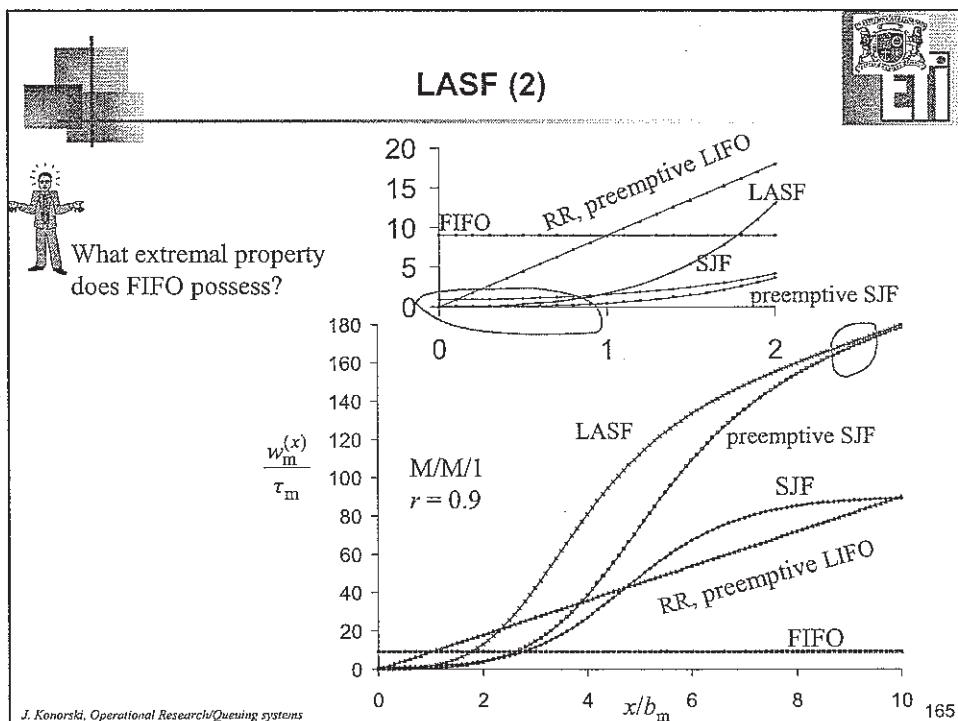
- at any time, processor serves request with minimum attained service time
- ties resolved by RR with $\Delta \approx 0$



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.



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



J. Konorski



Operational Research / Queuing Systems

Practice Set 1

Problem 1

Plot the queuing processes $N(t) = \text{queue length}$ and $U(t) = \text{unfinished work}$ for an arrival stream specified by $(t_n^+) = (0.5 \text{ s}, 2 \text{ s}, 5 \text{ s}, 6.5 \text{ s})$ and $(b_n) = (3 \text{ s.u.}, 2 \text{ s.u.}, 2 \text{ s.u.}, 3 \text{ s.u.})$, and for two cases:

- there are two processors each serving requests at the speed of $v = 1 \text{ s.u./s}$ (no grading, processor-bound and time-sharing service mode are assumed),
- there is one processor serving requests at the speed of $v = 2 \text{ s.u./s}$ (time-sharing service mode is assumed).

Compare the two cases from various viewpoints.

Be careful with the slopes of $U(t)$! What are the merits of case (b)? Is case (a) superior from some point of view?

Problem 2

Compare mean system delays in a single-processor queuing system under FIFO and RR with service quantum 2 s (partial use of assigned quantum causes earlier commencement of the next quantum). Processor speed is 1 s.u./s. Three requests, X, Y and Z, of sizes 7 s.u., 1 s.u. and 3 s.u., respectively, arrive simultaneously and queue up in the order (a) XYZ, (b) YZX, (c) XZY.

Can any general properties of FIFO and RR be inferred from such a limited set of scenarios?

Problem 3

A queuing system serves on average 800 transactions per second, each transaction on average requiring 5000 elementary operations to complete. An arriving transaction is immediately assigned a processor whose speed is 4,000,000 elementary operations/s. Find the mean number of transactions in system.

Use Little's theorem.

Operational Research / Queuing Systems

Practice Set 2

Problem 1

A single-processor infinite-buffer queuing system with processor speed $v = 1$ s.u./s serves a "dense" arrival stream of requests creating a 75% offered load. The total service demand in a one-second observation period is a random variable with standard deviation $\sigma = 0.1$ s. What are the chances that the processor can spare half of the second to deal with other (e.g., system) tasks without a backlog of requests forming at the end of the observation period?

Use the Laplace function: $\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_0^x e^{-y^2/2} dy$. Why Laplace function?

Tables of the Laplace function are widely available. Some characteristic values: $\Phi(0.5) = 0.191$, $\Phi(1) = 0.341$, $\Phi(1.5) = 0.433$, $\Phi(2) = 0.477$, $\Phi(2.5) = 0.4938$, $\Phi(3) = 0.4987$, $\Phi(\geq 4) \approx 0.5$. What to do if $x < 0$?

Problem 2

A single-processor queuing system with a finite buffer of capacity Q works under offered load $r > 1$. In such a system, the loss fraction L never drops below a certain level. What level is this?

Recall the flow conservation equation. How often is the processor idle if $Q \rightarrow \infty$?

Problem 3

A processor handles telemetric reports generated by a number of identical sources. To enable multiple access, the processor is equipped with a finite buffer that can accommodate up to Q reports. The table below contains the mean system delay of a report, normalized to the mean report processing time (in boldface), and report loss fraction due to buffer overflow, as dependent on Q and offered load r . Find the maximum number of sources that can be connected to the processor and necessary buffer capacity under the following assumptions:

- each source generates on average 20 reports per minute,
- a report contains on average 1800 records of data,
- the processor is capable of handling 12000 records per second,
- tolerable mean system delay of a report is 1.8 s,
- tolerable loss fraction due to buffer overflow is 4%.

$Q=$	20	21	22	23	24	25	
$r =$							
0.1	1.11	0 1.11	0 1.11	0 1.11	0 1.11	0 1.11	0
0.2	1.25	0 1.25	0 1.25	0 1.25	0 1.25	0 1.25	0
0.3	1.43	0 1.43	0 1.43	0 1.43	0 1.43	0 1.43	0
0.4	1.67	0 1.67	0 1.67	0 1.67	0 1.67	0 1.67	0
0.5	2	0 2	0 2	0 2	0 2	0 2	0
0.6	2.5	0 2.5	0 2.5	0 2.5	0 2.5	0 2.5	0
0.7	3.32	0 3.32	0 3.32	0 3.33	0 3.33	0 3.33	0
0.8	4.77	0 4.8	0 4.84	0 4.86	0 4.89	0 4.91	0
0.9	7.23	0.01 7.42	0.01 7.6	0.01 7.76	0.01 7.92	0.01 8.07	0.01
1	10.5	0.05 11	0.05 11.5	0.04 12	0.04 12.5	0.04 13	0.04
1.1	13.5	0.11 14.3	0.1 15.1	0.1 15.9	0.1 16.7	0.1 17.5	0.1
1.2	15.5	0.17 16.5	0.17 17.4	0.17 18.4	0.17 19.3	0.17 20.3	0.17
1.3	16.8	0.23 17.8	0.23 18.7	0.23 19.7	0.23 20.7	0.23 21.7	0.23
1.4	17.5	0.29 18.5	0.29 19.5	0.29 20.5	0.29 21.5	0.29 22.5	0.29
1.5	18	0.33 19	0.33 20	0.33 21	0.33 22	0.33 23	0.33
1.6	18.3	0.38 19.3	0.38 20.3	0.38 21.3	0.38 22.3	0.38 23.3	0.38
1.7	18.6	0.41 19.6	0.41 20.6	0.41 21.6	0.41 22.6	0.41 23.6	0.41
1.8	18.8	0.44 19.8	0.44 20.8	0.44 21.8	0.44 22.8	0.44 23.8	0.44
1.9	18.9	0.47 19.9	0.47 20.9	0.47 21.9	0.47 22.9	0.47 23.9	0.47
2	19	0.5 20	0.5 21	0.5 22	0.5 23	0.5 24	0.5

Operational Research / Queuing Systems

Practice Set 3

Problem 1

Each of 50 terminals connected to a common transceiver generates a request after a think time of average duration $\frac{2}{3}$ s. In 80% cases it is a message of average length 1000 bytes, and in 20% cases a control data report of average length 160 bytes. The transceiver works at 1 Mb/s in half-duplex; the average proportion of time it is switched to receive mode is 75% (during that time it is unavailable to the terminals). What is the resulting loss fraction?

Data in the problem permit to identify a_m , b_m and proportion of processor idle time. Now go back to the flow conservation equation...

Problem 2

In a single-processor queuing system with buffer capacity $Q = 2$ in statistical equilibrium and under offered load $r = 0.75$, we have $p_0 \geq p_1 \geq p_2$. Find the range of possible p_1 values.

Proceed as in the previous problem...

Problem 3 (simulation experiment)

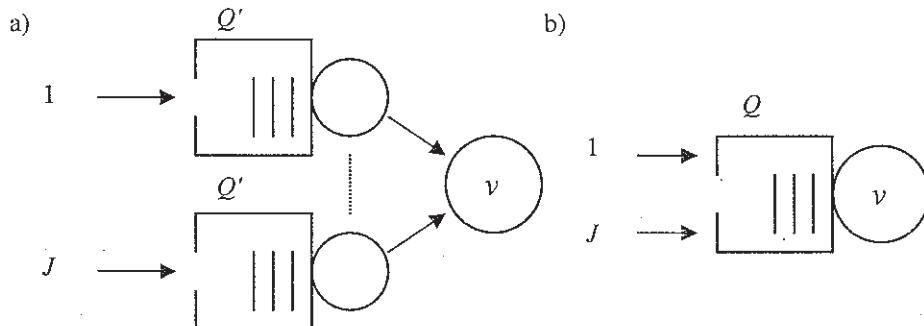
Each user of a single-processor queuing system generates an arrival stream of documents with mean interval a_m s. Documents have a mean length of b_m bytes. The processor handles documents at speed v bytes/s. Tolerable are:

- document loss fraction due to buffer overflow not exceeding L_{\max}
- mean system delay of a document not exceeding a given multiple c of b_m/v .

Subject to the above, compare the maximum number J_{\max} of users and required buffer capacity in two configurations:

- (a) *dedicated access* with a virtual processor and a separate buffer space assigned for each user, and
- (b) *common-buffer access* to the processor.

Perform simulations for $a_m = 6$ s, $b_m = 600$ bytes, $v = 24000$ bytes/s, $L_{\max} = 0.1\%$, $c = 5$.

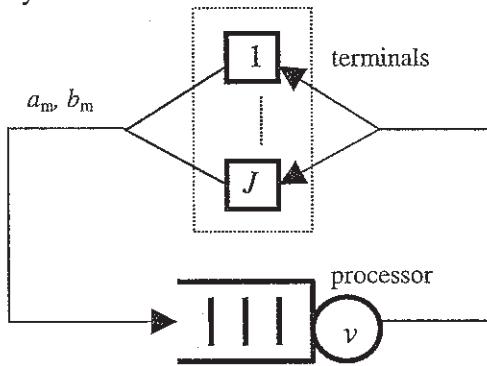


Operational Research / Queuing Systems

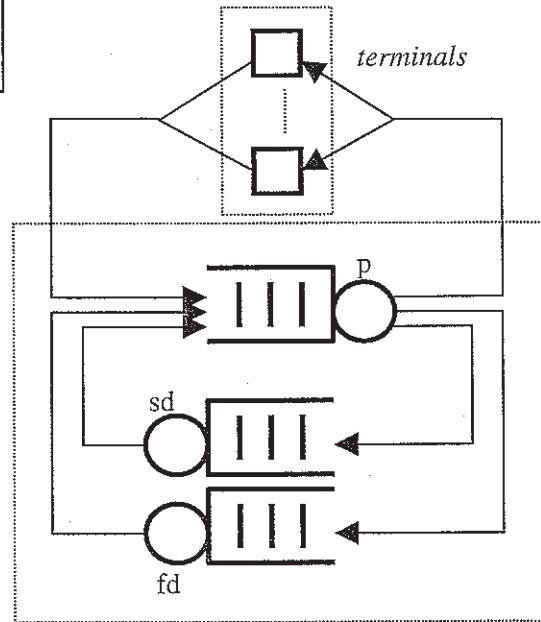
Practice Set 4

Problem 1

A single-processor queuing system interacts with $J = 10$ intelligent terminals in a query-response manner, as depicted below. Having received a response, a terminal generates a new query after a think time of average duration $h_m = 4$ s. The average number of elementary operations needed to generate a response is $b_m = 15000$, and processor is $v = 5000$ elementary operations/s. Find the relationship between the proportion of processor idle time and mean waiting delay.



The flow conservation equation again comes in handy.



Problem 2

Each of $J = 30$ computer terminals generates requests that require sequential service at a processor, slow-disk controller and fast-disk controller, as depicted below. On average, a request has to visit these devices $l_p = 21$, $l_{sd} = 12$, $l_{fd} = 8$ times, respectively, whereas average service times there equal $\tau_p = 0.05$ s, $\tau_{sd} = 0.07$ s and $\tau_{fd} = 0.02$ s. Upon notification of service completion for its request, a terminal enters a think time of average $h_m = 15$ s, and subsequently generates another request.

- Which device is the bottleneck, and which one is the most overdimensioned? How will his change if the processor is tuned up so that $\tau_p = 0.03$ s?
- What processor speedup do we need in order for the mean system delay (request time within the system) to become $d^*_m = 12$ s, and what speedup would ensure $d^*_m = 9$ s?

Assume a certain request arrival interval at the terminal-system interface. Use it to express the offered load at each device. For (b), will Little's theorem be of any use?

Operational Research / Queuing Systems

Practice Set 5

Problem 1

Answer the following questions related to finite-buffer queuing systems.

- An average of 40 requests per second arrive in an M/M/1/2 queuing system, each requesting 20 ms of processor time. How many requests on average are lost per day due to buffer overflow?
- In an M/M/1/5 queuing system, two requests arrive on average during the service of a request of average length. What is the loss fraction?
- An M/M/1/ Q queuing system serves an arrival stream whose mean characteristics are a_m and b_m . For sustained operation, the processor needs a $p\%$ proportion of idle time (used for maintenance), which in turn requires a minimum processor speed v_{min} . How does v_{min} change with increasing Q ?

In (b), infer the offered load from the information given.

Problem 2

Consider an M/M/S/S system (also called a loss system since it has no waiting room in buffer) with 150 users. Each user generates on average 10 transactions per second and tolerates $L \leq 3\%$ rejections (transactions lost on arrival due to lack of available processors). An average transaction request the processor to execute 800 elementary operations. The system operator faces a choice between leasing a number of processors of speed 0.5 million elementary operations/s, or four times faster processors leased at a 2.5 times higher monthly fee. Which type and how many processors should the operator lease to minimize the total fee while meeting the users' requirements?

In the calculations, you'll find useful the following table obtained from an Erlang-B calculator, where L values (in %) are given for $S = 1$ to 10 processors, and busy-hour load ranging from $\rho = 0.2$ to 0.6 erlangs.

$S =$	$\rho =$	0.2	0.4	0.6	0.8	1	1.2	1.4	1.6	1.8	2	2.2	2.4	2.6	2.8	3
1	16.67	28.57	37.5	44.44	50	54.55	58.33	61.54	64.29	66.67	68.75	70.59	72.22	73.68	75	
2	1.64	5.41	10.11	15.09	20	24.66	28.99	32.99	36.65	40	43.06	45.86	48.42	50.78	52.94	
3	0.11	0.72	1.98	3.87	6.25	8.98	11.92	14.96	18.03	21.05	24	26.84	29.56	32.15	34.61	
4	0.01	0.07	0.3	0.77	1.54	2.62	4	5.65	7.5	9.52	11.56	13.87	16.12	18.37	20.61	
5	0	0.01	0.04	0.12	0.31	0.63	1.11	1.77	2.63	3.67	4.88	6.24	7.73	9.33	11	
6	0	0	0	0.02	0.05	0.12	0.26	0.47	0.78	1.21	1.76	2.44	3.24	4.17	5.21	
7	0	0	0	0	0.01	0.02	0.05	0.11	0.2	0.34	0.55	0.83	1.19	1.64	2.19	
8	0	0	0	0	0	0	0.01	0.02	0.05	0.09	0.15	0.25	0.39	0.57	0.81	
9	0	0	0	0	0	0	0	0	0.01	0.02	0.04	0.07	0.11	0.18	0.27	
10	0	0	0	0	0	0	0	0	0	0	0.01	0.02	0.03	0.05	0.08	
$S =$	$\rho =$	3.2	3.4	3.6	3.8	4	4.2	4.4	4.6	4.8	5	5.2	5.4	5.6	5.8	6
1	76.19	77.27	78.26	79.17	80	80.77	81.48	82.14	82.76	83.33	83.87	84.37	84.85	85.29	85.71	
2	54.94	56.78	58.48	60.07	61.54	62.91	64.19	65.39	66.51	67.57	68.56	69.49	70.38	71.21	72	
3	36.95	39.15	41.24	43.21	45.07	46.83	48.49	50.06	51.55	52.97	54.3	55.57	56.78	57.92	59.02	
4	22.81	24.97	27.07	29.1	31.07	32.96	34.78	36.54	38.22	39.83	41.38	42.86	44.29	45.65	46.96	
5	12.74	14.51	16.31	18.11	19.91	21.68	23.44	25.16	26.84	28.49	30.09	31.64	33.15	34.62	36.04	
6	6.36	7.6	8.91	10.29	11.71	13.18	14.66	16.17	17.68	19.18	20.68	22.16	23.63	25.07	26.49	
7	2.83	3.56	4.38	5.29	6.27	7.33	8.44	9.6	10.81	12.05	13.32	14.6	15.9	17.2	18.5	
8	1.12	1.49	1.93	2.45	3.04	3.7	4.44	5.23	6.09	7	7.97	8.97	10.01	11.09	12.19	
9	0.4	0.56	0.77	1.02	1.33	1.7	2.12	2.6	3.15	3.74	4.4	5.11	5.86	6.67	7.51	
10	0.13	0.19	0.28	0.39	0.53	0.71	0.93	1.18	1.49	1.84	2.24	2.68	3.18	3.72	4.31	

What is the maximum tolerable busy-hour load in either of the two options?

Problem 3

The arrival stream to an M/M/1 queuing system has mean interarrival interval $a_m = 10$ s and mean request size $b_m = 10$ s.u. Processor speed is v s.u./s. Draw a state transition diagram corresponding to the underlying birth-and-death process for the following model specifications:

- a) with probability 25% a request whose service has been completed immediately returns to the queue instead of departing from the system,
- b) at three or more requests in system, the processor speeds up by 50%,
- c) upon termination of a busy period, the processor "goes on vacation" i.e., ignores arriving requests, and only resumes operation when three requests are queued,
- d) the processor "goes on vacation" after completion of each request's service; "vacation" duration is exponentially distributed with mean h_m ,
- e) the processor occasionally breaks down and comes up after a while, whereupon the interrupted service is resumed (all requests arriving during breakdown are queued); the down- and up-times are exponentially distributed with mean f_m i g_m , respectively,
- f) requests are admitted in pairs – the first request from a pair is held back until the second one arrives, then the pair is regarded as an arriving request of size equal to the size of the second request.

For each of these models carefully define a system state. Examine all events that can any moment occur at a given state, and the resulting state transitions.

Problem 3

Consider impatient requests in an M/M/1 system. The arrival stream has mean interval $a_m = 10$ s and mean request size $b_m = 10$ s.u., and processor speed is $v = 1$ s.u./s. On arrival, each request draws its individual "patience threshold" from exponential distribution with mean c_m , and upon its expiration escapes from the queue if still waiting, otherwise departs upon service completion. For ease of calculation let $c_m = b_m/v$. Find the distribution (p_k) of the number of requests in system, as well as the fraction L of escaped requests.

Use the general birth-and-death solution, also apply the flow conservation equation.

