



Operational Research

Queuing Systems 1: Description and Operation

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

to complete this course part, pass test (min 17 marks out of 34)

J. Konorski, Operational Research/Queuing system

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

I. Konorski, Operational Research/Queuing system

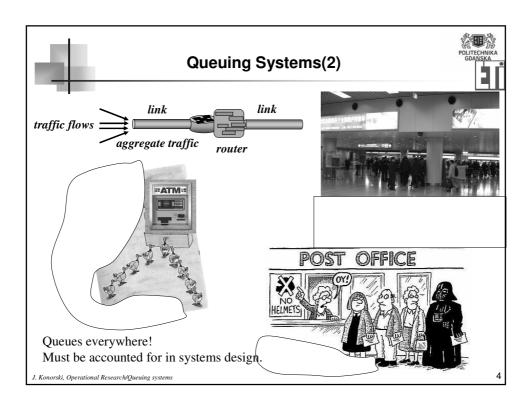


Queuing (Stochastic Service) Systems: Examples



- 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/Oueuing system

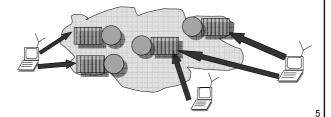




Queuing Systems(3)



- 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 system:



QueuingTheory: 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



QueuingTheory: 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. Konorski, Operational Research/Queuing system

7



Simplest Model



The simplest model of a queuing system consists of:

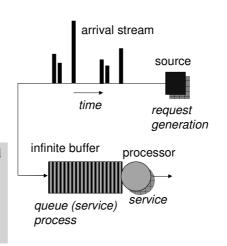
- a processor,
- a buffer, and
- an arrival stream.

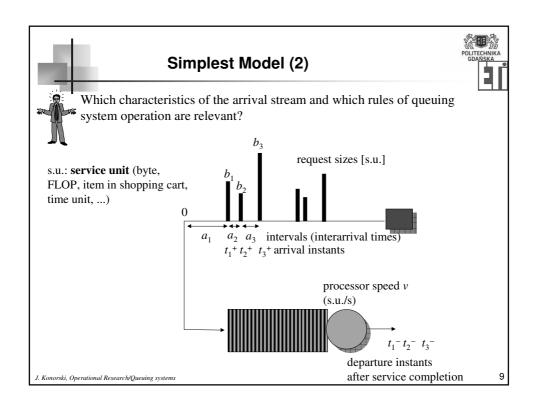
Characteristics of service process depend on those of arrival stream and the way buffer & processor system operates.

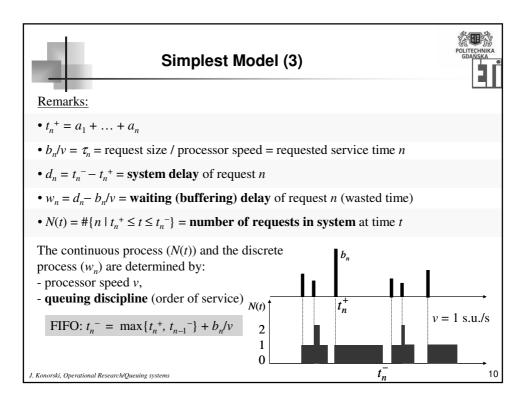
How?

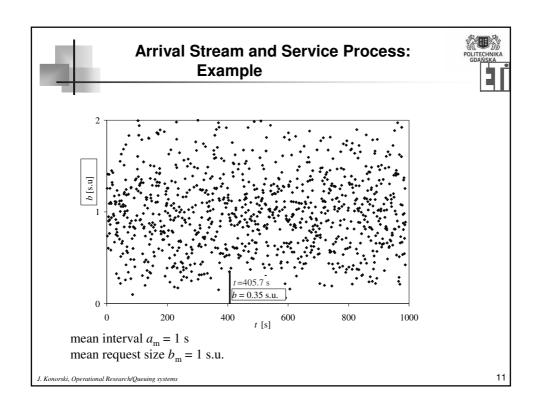
This is what queuing theory is about.

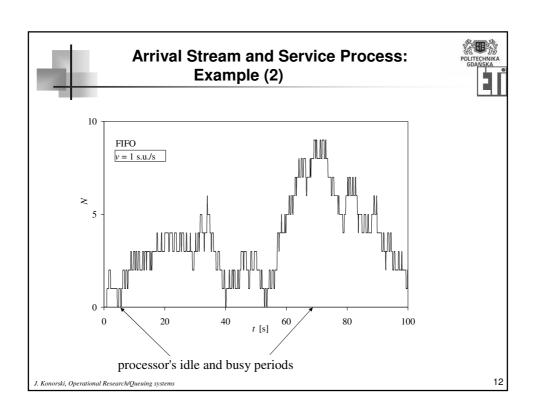
J. Konorski, Operational Research/Queuing systems

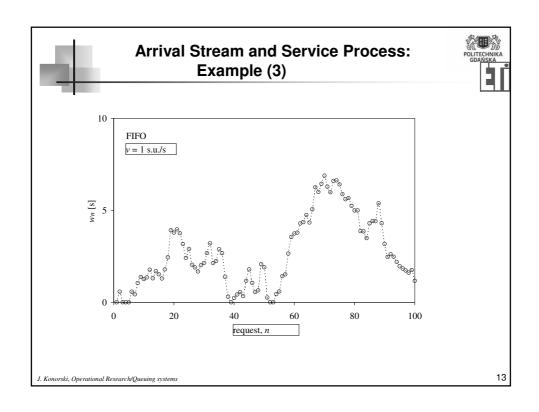














Properties of "Good" Service Process



- 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



Offered Load



Consider a long observation period T. Within it,

- approximately $T/a_{\rm m}$ arrivals occur, in total creating mean **service demand** = $b_{\rm m}(T/a_{\rm m})$ s.u.
- service supply is vT.

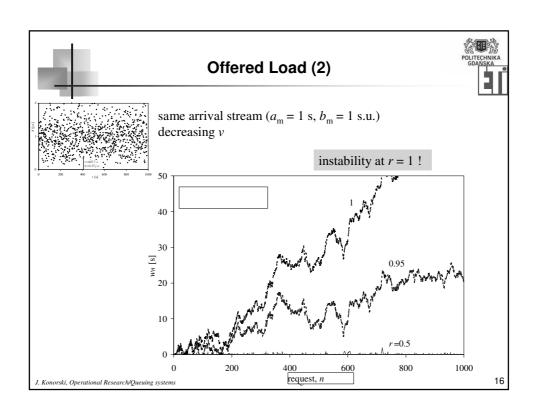
Offered load = ratio of mean service demand and service supply:

$$r = \frac{b_{\rm m}(T/a_{\rm m})}{vT} = \frac{b_{\rm m}/v}{a_{\rm m}} = \frac{b_{\rm m}/a_{\rm m}}{v} \quad ({\rm dimensionless})$$

- = ratio of mean service time $b_{\rm m}/v$ and mean request interval $a_{\rm m}$,
- = ratio of mean service demand rate $b_{\rm m}/a_{\rm m}$ and service supply rate ν

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/Oueuing system





Impact of Input Speed



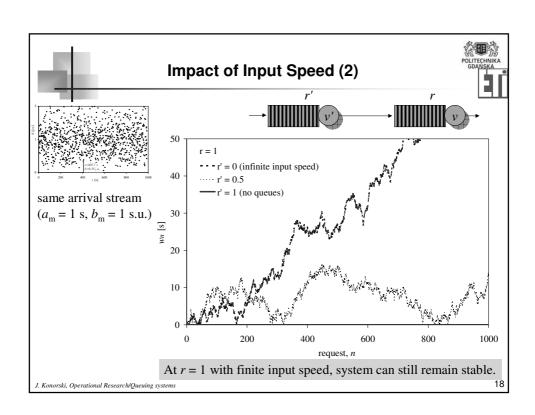
So far, immediate input assumed of requests from source to system (arbitrary a_n). In reality, request transfer from source occurs at finite speed.

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

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

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

I. Konorski, Operational Research/Oueuing system





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, ...?

J. Konorski, Operational Research/Oueuing system.

19



Towards Richer Models (2)



• request behavior in case of buffer overflow / long queue found on arrival loss?

retry upon timeout? pushout of some queued requests? impatience of the $1^{\rm st}$ kind (balking), $2^{\rm nd}$ 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)?

J. Konorski, Operational Research/Queuing system



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 $\Sigma b_n / v$) 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?

Operational Research/Queuina systems

0.



Steady State



In a **stationary** queuing system:

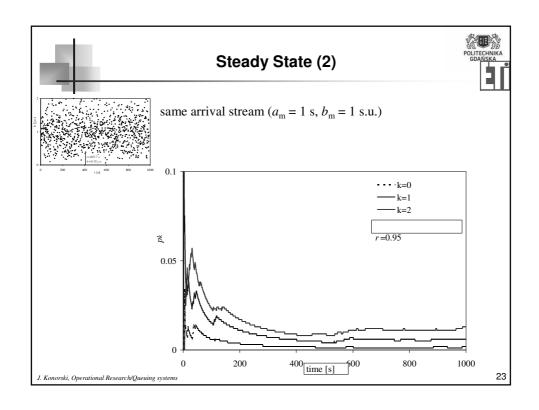
- mechanism governing request generation remains invariable in time (characteristics of the arrival stream do not change)
- processor speed v remains constant.

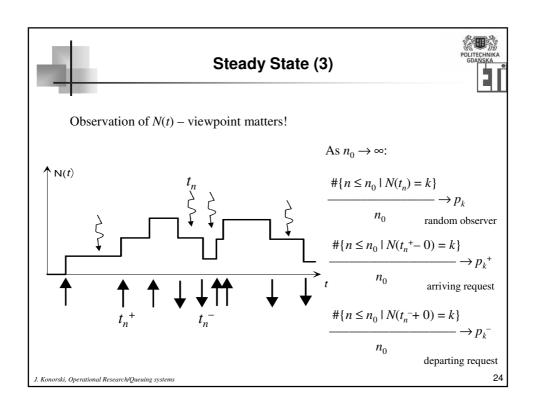
A *well-designed* stationary queuing system (where r < 1) tends to a **steady state** (characteristics of the queuing process exhibit asymptotic behavior as the observation period lengthens).

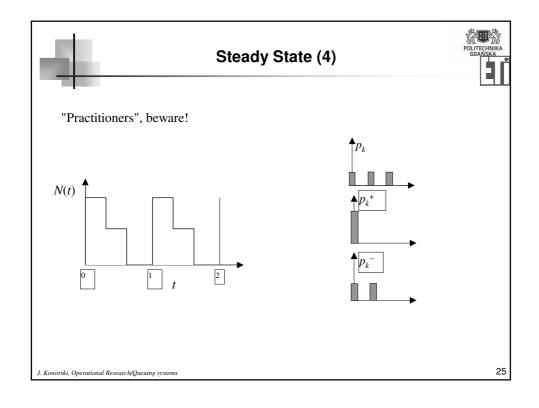
In particular,

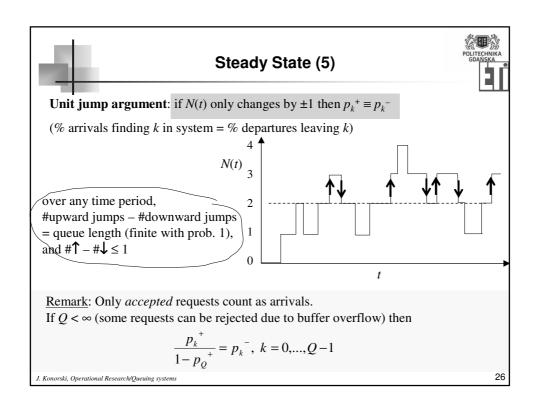
$$\frac{|\{t \le T \mid N(t) = k\}|}{T} \to p_k \text{ as } T \to \infty$$

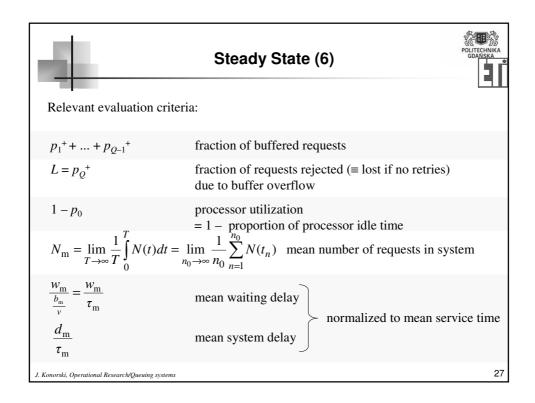
J. Konorski, Operational Research/Queuing systems

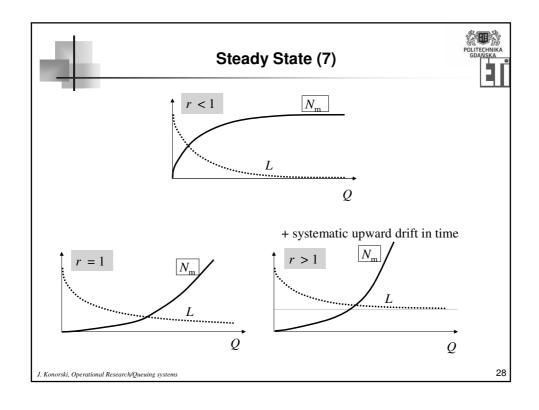


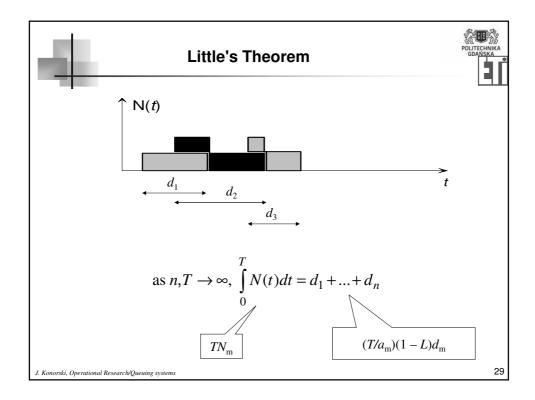














Little's Theorem (2)



$$N_{\rm m} = \frac{1 - L}{a_{\rm m}} d_{\rm m}$$

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: $1 - p_0 = \frac{1 - L}{a_{\rm m}} \tau_{\rm m} = (1 - L)r$ flow conservation equation
• buffer: $N_{\rm m} - (1 - p_0) = \frac{1 - L}{a_{\rm m}} w_{\rm m}$



Is Queuing Theory Losing Momentum?



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

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

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 $b_{\rm m}/a_{\rm m}$ grows in step with service supply rate v, and so r isn't dropping any!

J. Konorski, Operational Research/Queuing systen



Is Queuing Theory Losing Momentum? (3)



Second, why do queues arise 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 $b_{\rm m}/a_{\rm m} < \infty$ 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., $v > b_{\rm m}/a_{\rm m}$, but not much. Meaning, allow queues at times.

J. Konorski, Operational Research/Queuing system

33



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.

J. Konorski, Operational Research/Queuing system



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 system

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 system





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 system

37



Random Variables and Stochastic Arrival Streams



For an arrival stream observed over a finite time, $a_{\rm m}$, $b_{\rm m}$ 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; either is known or can be derived somehow (the **Bayesian approach**).

J. Konorski, Operational Research/Queuing systems

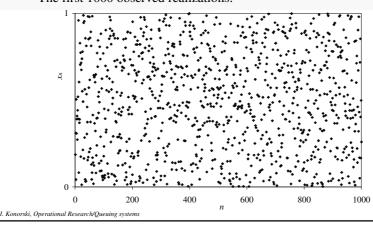


Random Variables and Stochastic Arrival Streams (2)



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:



39



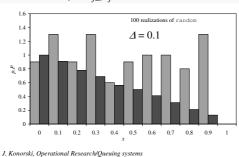
Empirical Distributions

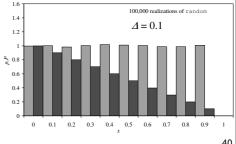


Having realizations $x_1,...,x_N$ one constructs a **histogram**:

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

Complementary cumulative histogram constructed analogously, with bars of height $C_i = \sum_{j \ge i} p_j$.



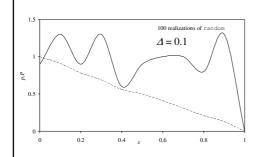


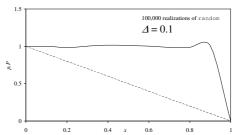


Empirical Distributions (2)



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





Mean value:

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



Theoretical Distributions



Probability density function and complementary distribution function

– histograms one would obtain taking $N \to \infty$ and $\Delta \to 0$.

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

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

Complementary cumulative distribution function: $P(x) = \Pr[X \ge x] = \int p(y)dy$

$$P(x) = \Pr[X \ge x] = \int_{0}^{\infty} p(y) dy$$

$$x_{\rm m} = \int_{-\infty}^{\infty} x p(x) dx, \quad \sigma_x = \sqrt{\int_{-\infty}^{\infty} (x - x_{\rm m})^2 p(x) dx}$$

J. Konorski, Operational Research/Queuing systems



Theoretical Distributions (2)



Modeling for engineering applications often uses Weibull distribution:

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

 λ – scale parameter,

 θ – shape parameter (= 1: exponential distribution).

$$x_{\rm m} = \int_{0}^{\infty} \theta \sqrt{\frac{y}{\lambda}} e^{-y} dy$$

$$\sigma_{x} = \sqrt{\int_{0}^{\infty} \theta \sqrt{\frac{y^{2}}{\lambda^{2}}} e^{-y} dy - x_{\mathrm{m}}^{2}}$$

J. Konorski, Operational Research/Queuing system

43



Computer Generation of Arrival Streams

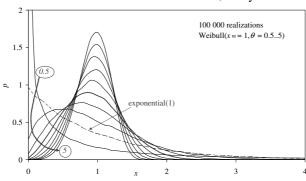




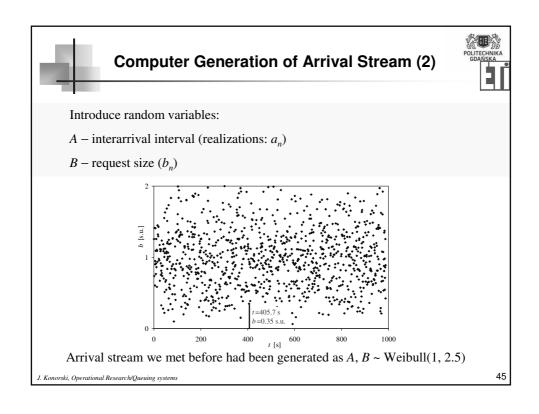
Given random generator that returns values $(z_n)_{n=1,2,...}$ how to generate pseudorandom numbers $(x_n)_{n=1,2,...}$ with arbitrary P(x)?

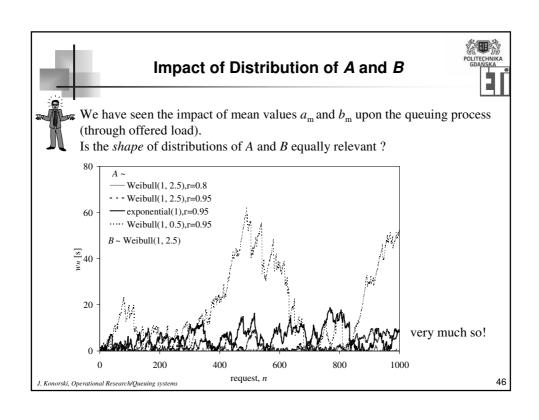
$$x_n$$
 solves $P(x) = z_n$ e.g., for Weibull distribution: $x_n = \theta - \frac{\ln z_n}{\lambda}$

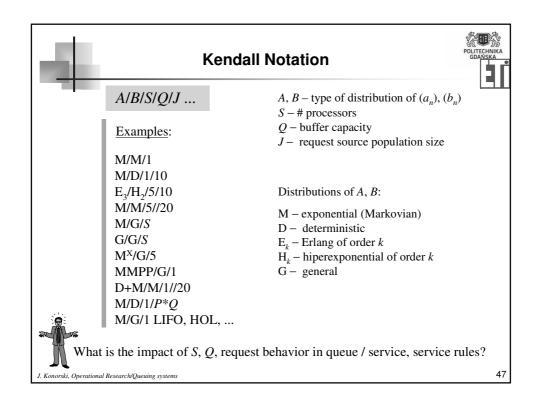
(method of inverted distribution function; many others exist).

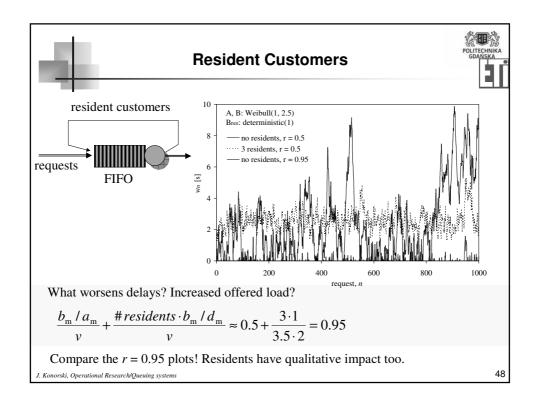


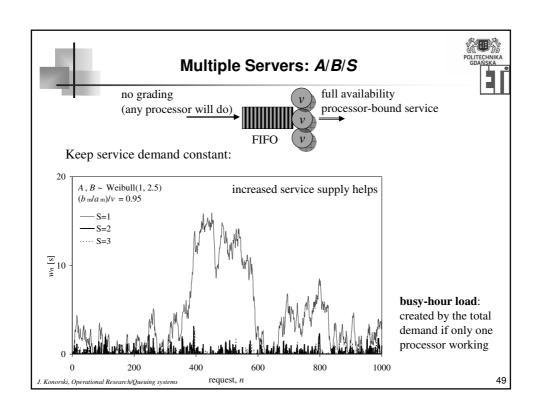
J. Konorski, Operational Research/Queuing systems

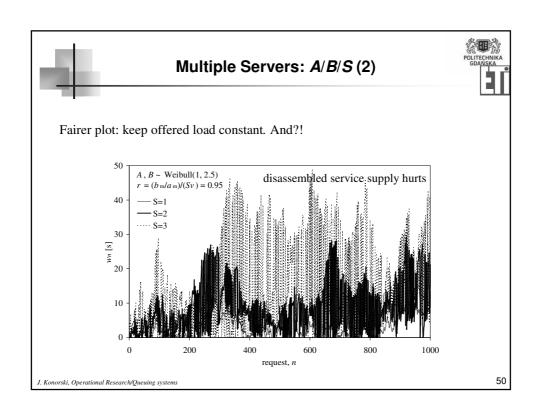


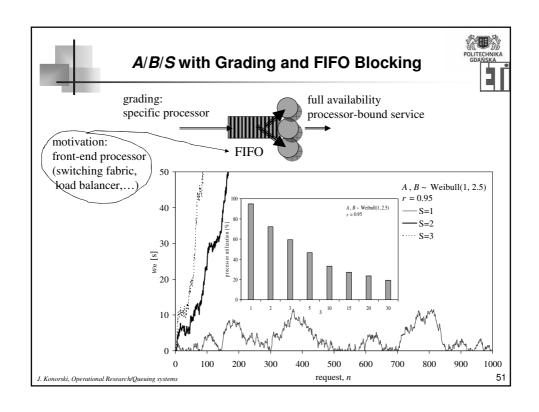


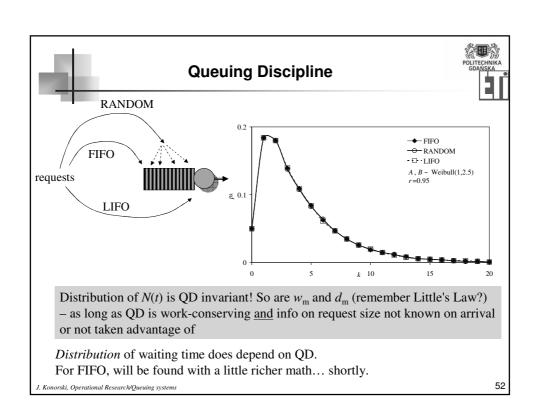


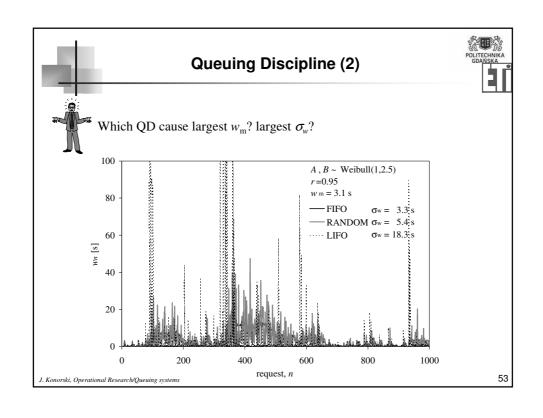


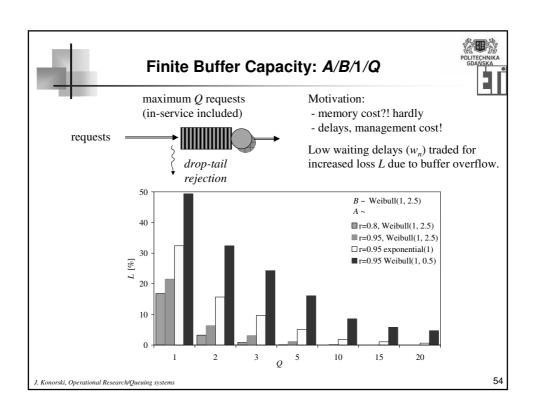


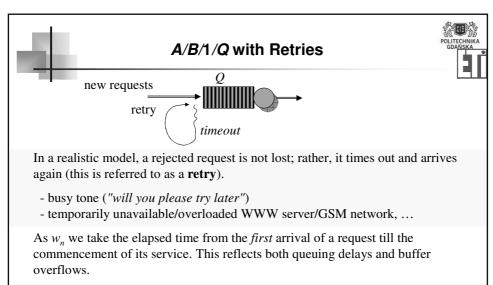






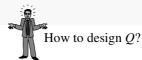






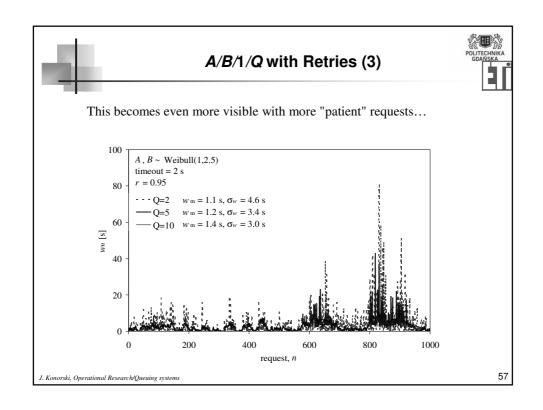
Simple model:

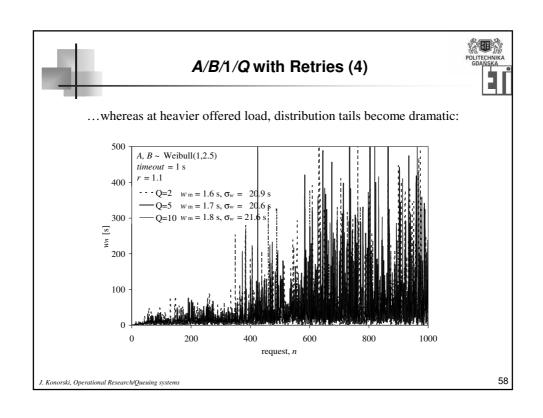
- constant timeout
- unlimited number of retries.



. Konorski, Operational Research/Queuing systems

A/B/1/Q with Retries (2) 100 $A, B \sim \text{Weibull}(1,2.5)$ timeout = 1 sr=0.9580 -Q=2 $w = 0.8 \text{ s}, \sigma_w = 3.9 \text{ s}$ Q=5 $w_m = 1.0 \text{ s}, \sigma_w = 3.3 \text{ s}$ Q=10 $w_m = 1.4 \text{ s}, \sigma_w = 3.3 \text{ s}$ w_n [s] 40 20 1000 200 400 600 800 request, n Larger Q worsen mean wait, reduce std deviation and distribution tail. 56 J. Konorski, Operational Research/Queuing system.







Common Models of Arrival Stream



- Weibull
- Bernouilli
- Erlang
- gamma
- (a_n) are **iid** *independent*, *identically distributed* **renewal streams**, often named after distribution of A
- more complex e.g.,

time-varying distributions of A and B, Markov Modulated Poisson Process, Batch Markovian Arrival Process, fractal Brownian process, ...

 model nonstationarity, dependence on the queuing process, bulk arrivals, internal correlation, ...

J. Konorski, Operational Research/Queuing system

59



Arrival Stream: Impact of Autocorrelation





Are distributions of A and B enough to determine the queuing process (given fixed v), or do we need information on the internal correlation in (a_n) ?

$$corr_a(l) = \frac{1}{\sigma_a^2} \frac{1}{M} \sum_{n=1}^{M} (a_n - a_m)(a_{n+l} - a_m), \quad M \to \infty, l = 0, 1, 2, ...$$

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

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

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

J. Konorski, Operational Research/Queuing systems



Arrival Stream: Impact of Autocorrelation (2)



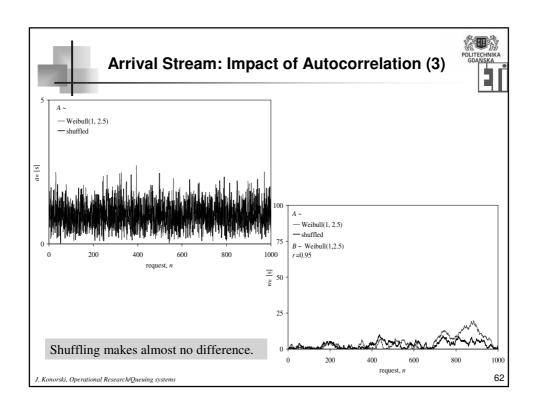
Experiment 1

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) .

J. Konorski. Operational Research/Queuing system





Arrival Stream: Impact of Autocorrelation (4)



Experiment 2

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

- as iid intervals from successive random numbers renewal stream,
- by repeating each successive interval a number of times drawn from Weibull(1, 0.5) distribution, rounded to nearest integer.

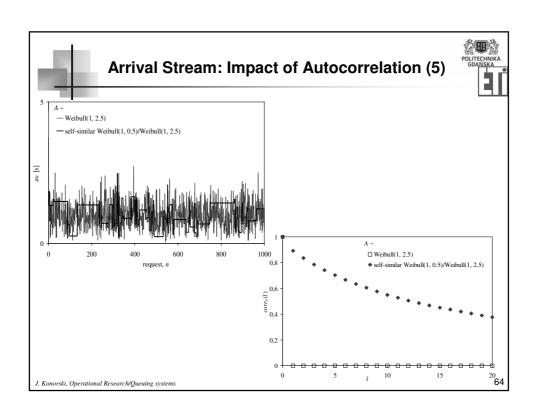
In both variants, distribution of *A* is the same.

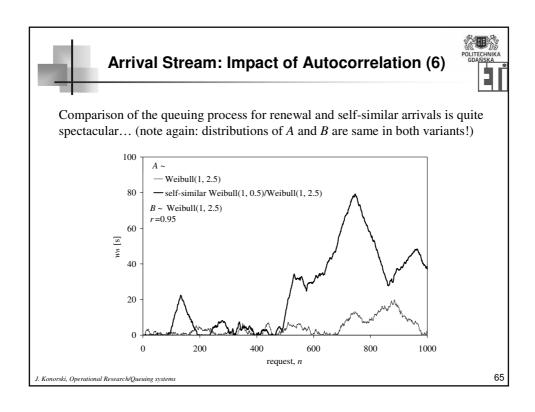
However, the second variant yields (a_n) with long-range autocorrelation

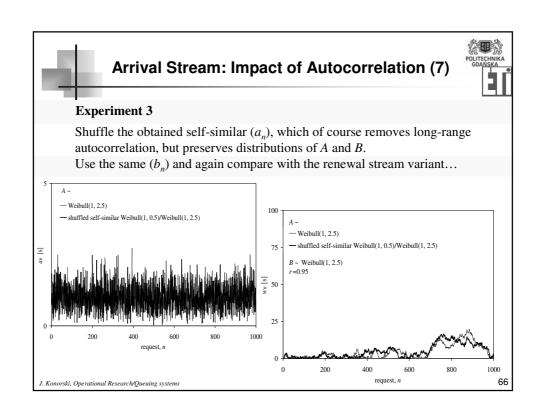
- a **self-similar** arrival stream.

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

J. Konorski, Operational Research/Queuing systems









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. Konorski, Operational Research/Queuing system

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.*



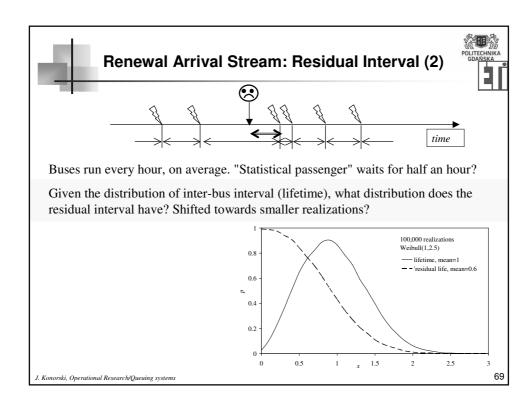
What types of distribution of *A* can one encounter most often?

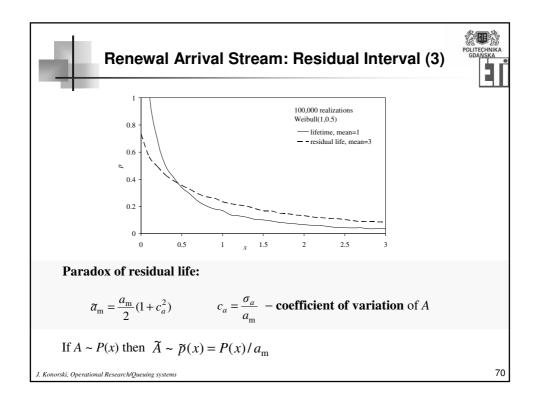
Clearly, very diverse. However, one of them has a suggestive meaning and a unique, "magic" property. To understand it, consider **residual interval**.

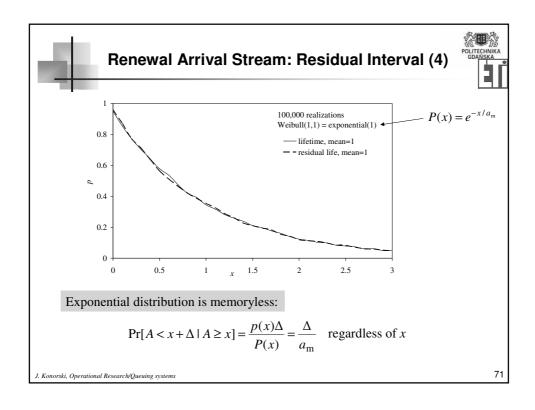
Events occur at random intervals. You arrive at a random instant. How long do you wait for the next event?

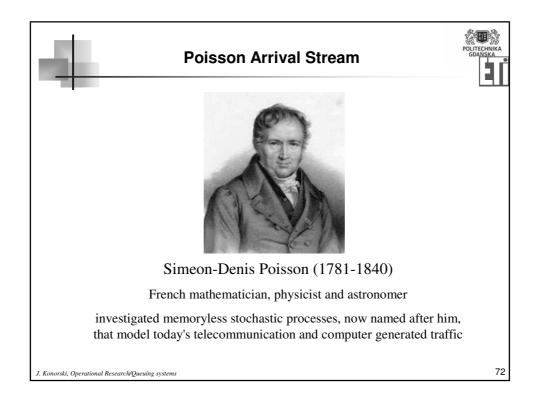
This is what may preoccupy a passenger at a bus stop, a subscriber trying to get through to a busy number, a VIP yet nonpreemptive customer urgently seeking access to a server etc.

J. Konorski, Operational Research/Queuing system.











Poisson Arrival Stream (2)

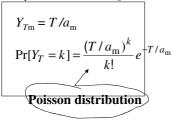


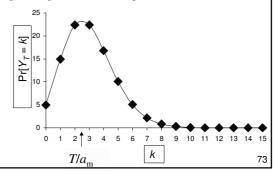
At any instant of time, new request arrival occurs with constant probability. In Kendall notation: M/... systems.

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



 Y_T – number of requests arriving during interval of length T





J. Konorski, Operational Research/Queuing system.

4

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).

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

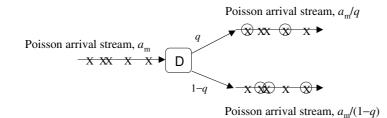
Hence, in M/G/S/Q: loss probability due to buffer overflow equals $L = p_Q^+ \equiv p_Q$.



Poisson Arrival Stream (4)



Random splitting:



Random splitting preseves arrival stream's Poissonian nature!

Non-random splitting doesn't.

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

J. Konorski. Operational Research/Queuing system

75

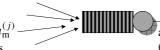


Aggregation of Renewal Arrival Streams



J independent components

- renewal arrival streams $a_{rr}^{(l)}$ (Web surfers, phone subscribers, mobile terminals, ...)

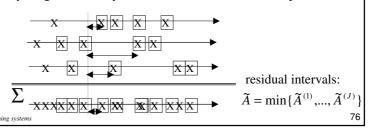


aggregated arrival stream (input to backbone link, mainframe / Web server...)

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



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



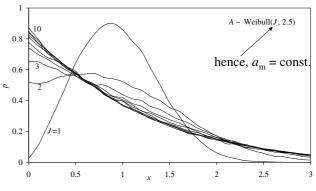


Aggregation of Renewal Arrival Streams (2)



With $J \to \infty$, but $a_{\rm m} = \frac{a_{\rm m}^{(j)}}{J} > 0$ (a practical model) and independent component streams we have:

Aggregated arrival stream is Poisson (Palm theorem).





Aggregation of Renewal Arrival Streams (3)



Proof: omitted :)

$$\widetilde{A} = \min{\{\widetilde{A}^{(1)}, ..., \widetilde{A}^{(J)}\}}, \text{ so } \Pr{[\widetilde{A} \ge x]} = \prod_{j=1}^{J} \Pr{[\widetilde{A}^{(j)} \ge x]}$$

 $J \to \infty$, but $a_{\rm m} > 0$. That is, $a_{\rm m}^{(j)} \to \infty$.

For any finite x, $x/a_{\rm m}^{(j)} \to 0$ and we can neglect $Y_x^{(j)} > 1$.

$$\Pr[\widetilde{A}^{(j)} \geq x] = 1 - \Pr[Y_x^{(j)} > 0] \approx 1 - \Pr[Y_x^{(j)} = 1] \approx 1 - \left(Y_x^{(j)}\right)_{\mathrm{m}} \approx 1 - \frac{x}{a_{\mathrm{m}}^{(j)}}$$

Primary,

$$\Pr[\widetilde{A} \ge x] \approx \prod_{j=1}^{J} \left(1 - \frac{x}{a_{\text{m}}^{(j)}} \right) \approx \exp\left(-\sum_{j=1}^{J} \frac{x}{a_{\text{m}}^{(j)}} \right) = e^{-x/a_{\text{m}}}$$

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





Operational Research

Queuing Systems 3: Markovian Models

Jerzy Konorski jekon@eti.pg.gda.pl

J. Konorski. Operational Research/Queuing system

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/... 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

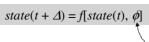


Markovian Systems (2)



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, ...)



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

J. Konorski, Operational Research/Queuing system

81

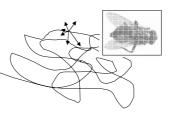


Markovian Systems (3)



- 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)
- $trend(t + \Delta) = (1 c) \cdot trend(t) + c \cdot \phi(t)$ (current observation)
- $market_share(t + \Delta) = \phi \cdot market_share(t)[1 market_share(t)]$ (current management performance)
- Internet topology (number and points of attachment of new networks)

J. Konorski, Operational Research/Queuing systems



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 system

83



Markovian Systems (5)



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

Poisson arrival stream, $a_{\rm m}$ exponential request size distribution, $b_{\rm m}$

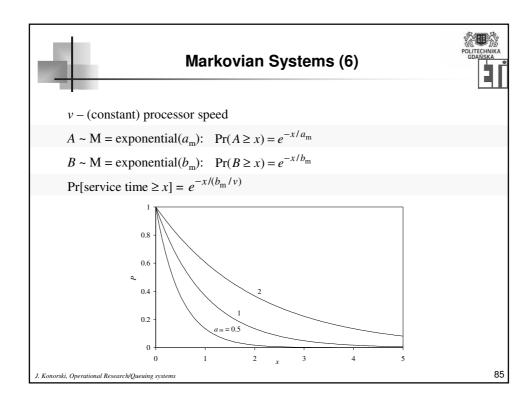


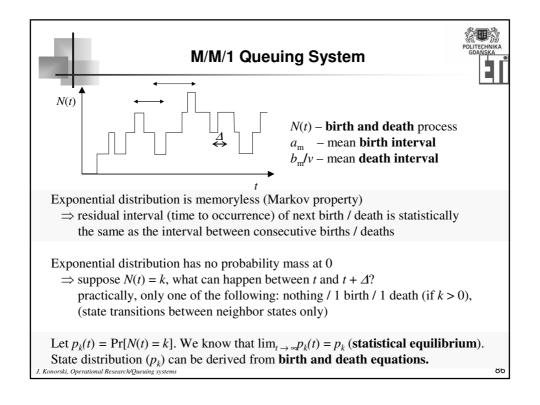
Practical impact stems from:

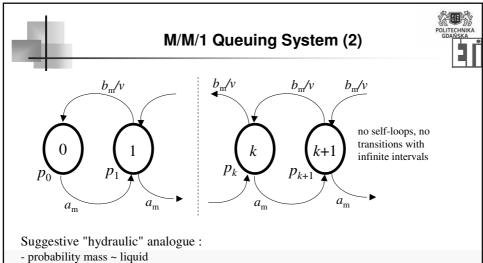
- Poisson arrival stream
- Palm theorem
- random splitting
- PASTA
- pessimistic (meaning: fail-safe) performance characteristics
- exponential request size distribution: crude approximation of
- call holding time, Web / P2P file transfer
- batch processing time

- ...

J. Konorski, Operational Research/Queuing systems







- state ~ vessel
- p_k ~ liquid pressure in vessel k
- state transition ~ pipe
- mean residual interval (time to occurrence) of event triggering state transition
- ~ pipe's flow resistance (=1/cross-sectional area)

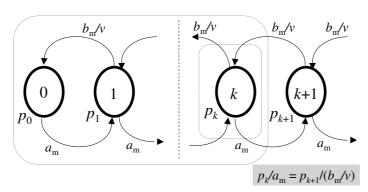


M/M/1 Queuing System (3)



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!



J. Konorski, Operational Research/Queuing systems



M/M/1 Queuing System (4)



$$p_k = p_0 \cdot r^k$$

$$p_0 + p_1 + p_2 + \dots = 1 \Rightarrow p_0 = 1 - r$$

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

$$N_{\rm m} = \frac{r}{1 - r}$$

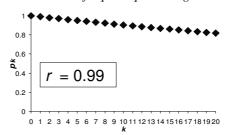
$$d_{\rm m} = a_{\rm m} N_{\rm m} = \tau_{\rm m} \left(\frac{1}{1 - r} \right)$$

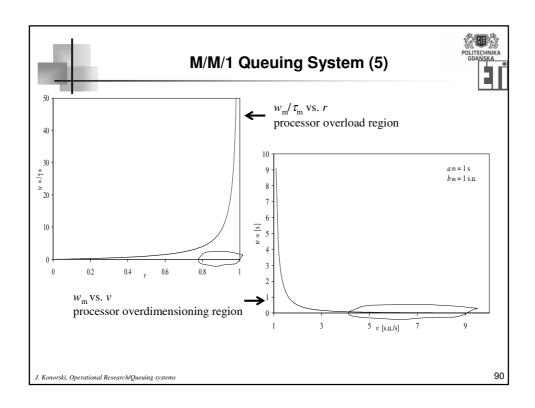
$$w_{\rm m} = d_{\rm m} - \tau_{\rm m}$$

$$w_{\rm m} = d_{\rm m} - \tau_{\rm m}$$

very suggestive!

The most frequent queue length?!







M/M/... Systems



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

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

J. Konorski, Operational Research/Queuing system

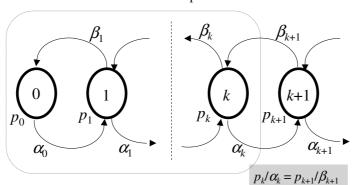
91



M/M/... Systems (2)



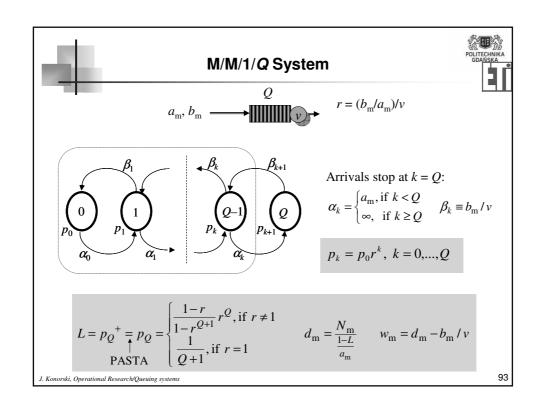
Make mean transition interval state-dependent:

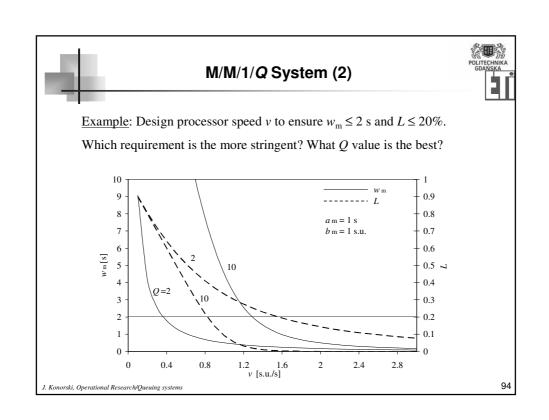


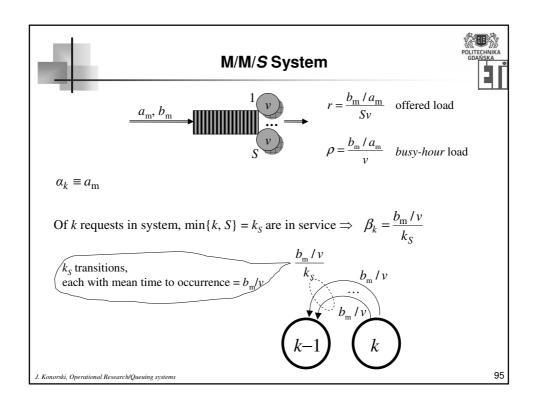
Solution now becomes : $p_k = p_0 \frac{\beta_1 ... \beta_k}{\alpha_0 ... \alpha_{k-1}}, \ k = 0,1,2,...$

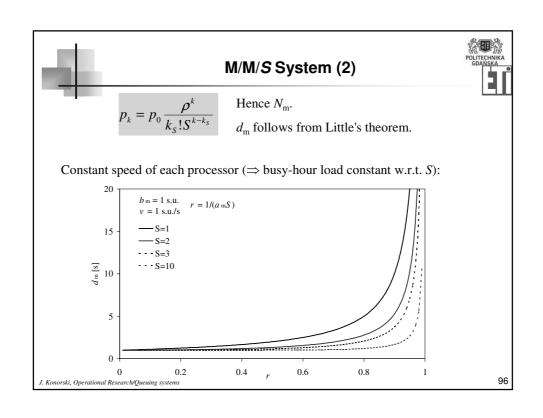
whence all the interesting characteristics: L, p_0 , $N_{\rm m}$, $d_{\rm m}$, $w_{\rm m}$, ...

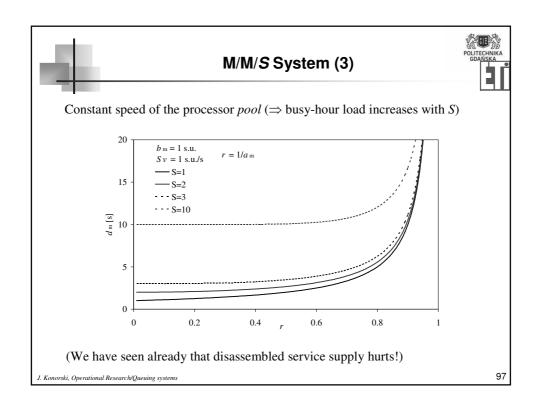
I. Konorski, Operational Research/Queuing systems

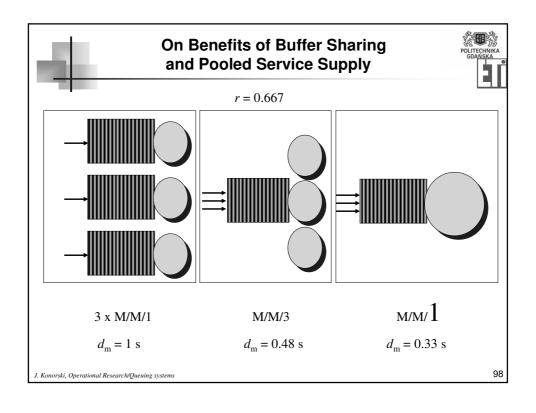








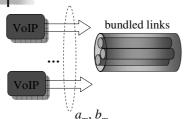






M/M/S/S System





no waiting room

v - single link capacity

 $\rho = b_{\rm m}/(a_{\rm m}v)$ – busy-hour load [erlangs] α_k , β_k , k = 0,...,S same as for M/M/S

$$p_k = p_0 \frac{\rho^k}{k!}$$
 , where $p_0 = \frac{1}{\sum_{k=0}^{S} \frac{\rho^k}{k!}}$

• famous Erlang B formula:

$$p_{S} = L = \frac{\frac{\rho^{S}}{S!}}{\sum\limits_{k=0}^{S} \frac{\rho^{k}}{k!}}$$

- magic: holds for any request size distribution i.e., for M/G/S/S (!)
- online calculators exist (www.voip-calculator.com/calculator/)



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.

$$a_{\rm m} = 3 \text{ s}, \ \tau_{\rm m} = 900 \text{ s}, \ \rho = 300 \text{ erlangs}$$

 $N_{\rm m} = \rho = 300$ (this we know from Little's theorem).

Special case of M/M/S/S with $S \to \infty$, therefore $p_0 = \frac{1}{\sum_{k=0}^{\infty} \frac{\rho^k}{k!}} = e^{-\rho}$ $p_k = e^{-\rho} \frac{\rho^k}{k!} - \text{Poisson distribution (!)}$

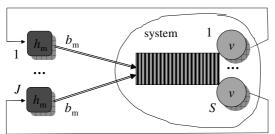
$$p_k = e^{-\rho} \frac{\rho^k}{k!}$$
 – Poisson distribution (!)

J. Konorski, Operational Research/Queuing system.



M/M/S//J System





Model:

- Pr[think time $\ge x$] = $e^{-x/h_{\text{m}}}$
- population of terminals J > S

When k requests in system,

- J k terminals during think time $\Rightarrow \alpha_k = \frac{h_{\rm m}}{J k}$
- k_S in service $\Rightarrow \beta_k$ same as for M/M/S
- as $J \to \infty, \, h_{\rm m} \to \infty, \, h_{\rm m}/J \to a_{\rm m}$, the system becomes M/M/S

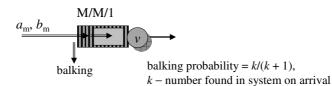
J. Konorski. Operational Research/Queuing system

101



Impatient Requests: Balking



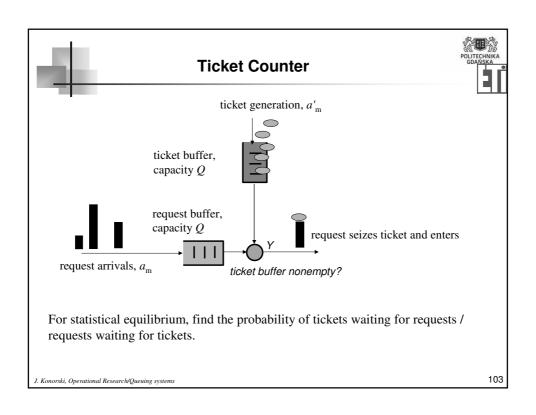


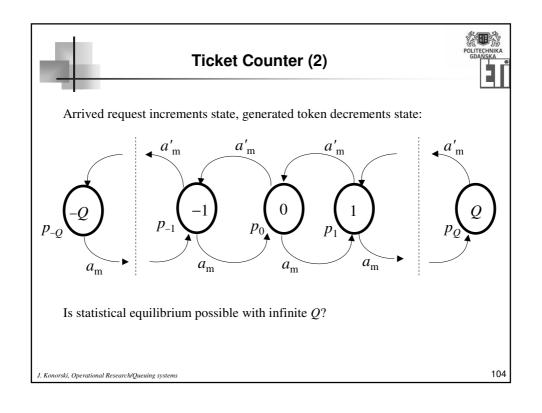
Assuming a Markovian system model, find the fraction of balking requests.

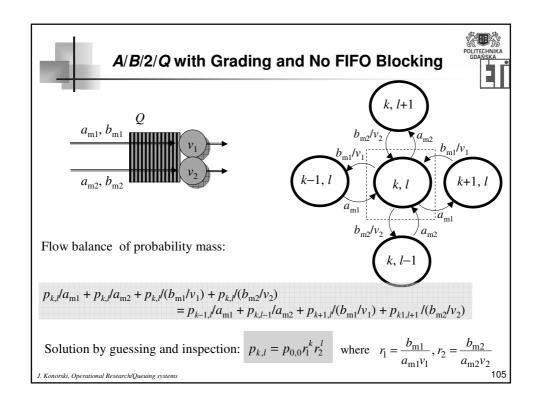
$$a_k = \frac{a_{\rm m}}{1 - \frac{k}{k+1}} = a_{\rm m} \cdot (k+1), \quad \beta_k \equiv b_{\rm m} / v \implies (p_k) - \text{Poisson distribution (!)}$$

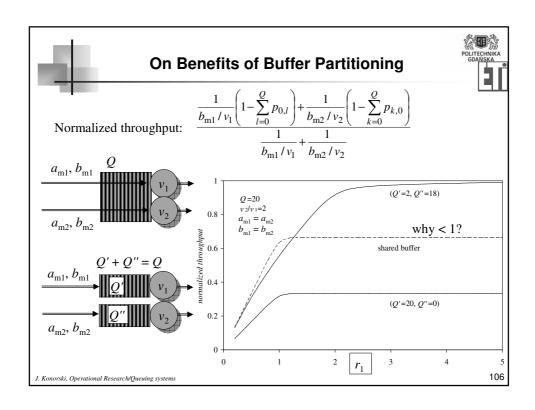
$$L = p_0 \cdot 0 + p_1 \cdot \frac{1}{2} + p_2 \cdot \frac{2}{3} + \dots + p_k \cdot \frac{k}{k+1} + \dots$$

J. Konorski, Operational Research/Queuing system











Laplace Transform



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

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

$$X * (s) = \int_{0}^{\infty} e^{-sx} (-dP(x)) = (e^{-sX})_{m}$$

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

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

J. Konorski, Operational Research/Queuing system

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 $X*(s) = \left(\frac{1}{cs+1}\right)^K$ (with integer K) then

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

- Direct application of inverse LT troublesome, involves complex numbers arithmetic and the Bromwitch 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 system



Laplace Transform (3)



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

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

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

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

and the Stehfest extrapolation coefficients are given

$$V_{i} = (-1)_{2}^{\frac{n}{L}+i} \sum_{k=\left[\frac{i+1}{2}\right]}^{Min\left[i,\frac{n}{2}\right]} \frac{k_{2}^{n}(2k)!}{\left[\frac{n}{2}-k\right]!k!(k-1)!(i-k)!(2k-i)!}$$

J. Konorski, Operational Research/Queuing system

109



M/M/1 FIFO: Distribution of System Delay



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

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

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

Averaging over k and using LT linearity gives:

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

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

J. Konorski, Operational Research/Queuing systems

