# Internet Performance

An introduction to performance evaluation with application to the Internet

Alain Jean-Marie

LIRMM, University of Montpellier

161 Rue Ada, 34392 Montpellier Cedex 5, France

ajm@lirmm.fr

# Contents of the course

**Introduction:** Motivation and approach

- Networks, contention, delays and losses
- Methodology

## Part VI: Deterministic models _____ 120

- Traffic envelopes and $(\sigma, \rho)$-bounds
- Bounds on the delay and on the buffer sizes
- Traffic shapers
- Service curves
- Network calculus

## Part VII: Application: Traffic management _____ 137

- Capacity planning
- Route planning, routing
- Window-based congestion control
- Models of TCP

## Part VIII: Application: Forward Error Correction _____ 153

## Bibliography _____ 160

# Introduction

In a communication network using routing/switching (Internet, ATM, Frame Relay...), *queues* form along the communication path. Management of statistical multiplexing, contention.

These queues create *delay* and *losses*.

The problem is to know how to quantify these.

The approach is *stochastic*, given the uncertain nature of traffic.

Queueing Theory: a set of concepts, tools, general and particular results for approaching these problems.

Research for results permitting to define, calculate and guarantee the celebrated *quality of service* (QoS).

## Methodology

How to obtain performance measures?

**Real System:** Define objectives

Instrument the system: place control points, place measurement points (not easy! intrusive)

Perform measurements

Change parameters

Do it again

**Simulated System:** Define objectives

Program a sufficient representation of the system, elements and behavior

Perform measurements

Change parameters

Do it again

## Mathematical analysis: Define objectives

Establish a sufficient mathematical representation of the system, elements and behavior

Calculate measures

For both Simulation and Analysis, one needs <span style="color:red">Models</span>.

# Modeling Issues

- Uncertainty and randomness

- Definition of the "performance measures", "Quality of Service"

- Parameters, controllable $(x_1, \ldots, x_n)$, uncontrollable (input) $(y_1, \ldots, y_n)$

- Tractability of models

  $$\mathrm{QoS} \;=\; f(x_1, \ldots, x_n; y_1, \ldots, y_m)$$

  − Analysis, exact: formulas, numerical methods.

  − Analysis, approximate.

  − Simulation.

- Validation of assumptions

- Optimization, dimensioning, capacity assignment.

  $$\max_{x_1, \ldots, x_n} \; f(x_1, \ldots, x_n; y_1, \ldots, y_m)$$

- Optimization, design choices (protocols, architecture, topology).

$$f(x_1, \ldots, x_n; y_1, \ldots, y_m) \overset{?}{<>} g(x_1, \ldots, x_n; y_1, \ldots, y_m)$$

- Statistics, measures for the input parameters (workload characterization).

# On the use of simulation

Quite common use of simulation

- new idea for a protocol
- implementation in a simulator
- run with various experimental conditions
- it works! → publish

Use of simulation in conjunction with modeling

- imagine a reasonable model
- solve it
- use simulation to validate the solution (esp. if approximations involved)
- vary assumptions to show robustness
- if it works, publish! if not, try to revise model...

# Uncertainty and randomness

Unknown quantities: arrival times of "events", amount of resources claimed on the system.

☐ Stochastic models.

Unknown quantities are random variables.

Random in, Random out ⟹ performance measures are random in nature

⟹ compute or measure their statistics (mean, variance, distribution...).

⟹ necessity to define performance these measures rigorously

⟹ understand the stochastic issues: stationarity, transience, ergodicity

⟹ necessity to perform measurements, statistics, estimators

? Collect the statistics on unknown quantities. Validate stochastic assumptions against real data.

□  Deterministic models.

Unknown quantities have bounds.

Analysis reveals the worst case scenarios ⟹ guaranteed performance.

? Accuracy of the bounds. How frequent are those bad cases?

Difficulty: worst case quantities do not always imply worst performance measures...

# Part I: Stochastic Processes

- Random variables
- Random processes
- Stationarity, ergodicity
- Covariance, autocorrelation
- Markov Chains

# Random Variables

Probability space: $\Omega$ set of *trajectories* or *realizations*.

Random variable $X$: function from the space of trajectories $\Omega$ into a space of values.

Distribution:

$$\mathbb{P}\{X \leq x\} = \mathbb{P}\{\omega \mid X(\omega) \leq x\}.$$

Expectation (mean), variance:

$$\mathbb{E}X = \int x \, d\mathbb{P}\{X \leq x\}$$

$$\mathrm{Var}(X) = \int x^2 \, d\mathbb{P}\{X \leq x\} - \mathbb{E}X^2$$

If the variable is *discrete*:

$$\mathbb{E}X = \sum_n n\, \mathbb{P}\{X = n\}$$

$$\text{Var}(X) = \sum_n n^2\, \mathbb{P}\{X = x\} - \mathbb{E}X^2$$

Variance: measure of the variability of $X$ around its mean.

Covariance of two r.v.:

$$\text{Cov}(X, Y) = \mathbb{E}(XY) - \mathbb{E}X\,\mathbb{E}Y.$$

Measure of the dependence between $X$ and $Y$. If $X$ and $Y$ are *independent*, $\text{Cov}(X, Y) = 0$.

Laplace transform (Laplace-Stiltjes) of $X$:

$$X^*(s) = \int_0^\infty e^{-st}\, d\mathbb{P}\{X \leq t\} = \mathbb{E}(e^{-sX}).$$

Generating function of a discrete random variable:

$$X^*(z) = \sum_{n=0}^{\infty} z^n \, \mathbb{P}\{X = n\} = \mathbb{E}(z^X).$$

Addition law: if $X \perp\!\!\!\perp Y$ then,

$$(X + Y)^*(s) = X^*(s) \, Y^*(s).$$

# Stochastic processes

A stochastic process "lives" in a state space $\mathcal{E}$.

Two categories:

discrete time $\quad \{X_n, n \in \mathbb{Z}\}$

continuous time $\quad \{X(t), t \in \mathbb{R}\}$

Discrete time: a sequence of random variables.

Continuous time: a family of random functions $\omega \mapsto X(t; \omega)$.

Classical examples:

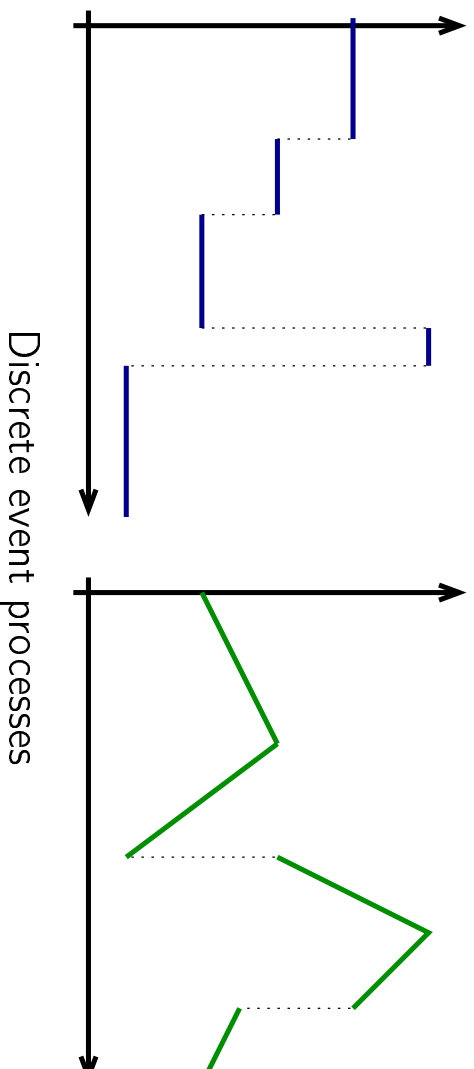- sequence of independent Bernoulli (Heads/Tails) tosses:

$$X_n \;=\; 0 \quad \text{with proba } 1/2, \quad X_n \;=\; 1 \quad \text{with proba } 1/2.$$

- Brownian motion: $\{X(t), t \in \mathbb{R}\}$ such that

$$X(s+t) - X(t) \;\sim\; \mathcal{N}(0, \sigma t).$$

# Discrete event systems

In the domain of information systems (computers, networks), one works with *discrete event systems* such that $X(t)$ or $\dot{X}(t)$ is piecewise constant.



Discrete event processes

Mathematical models for this situation:

- Event arrival processes: Point processes (Baccelli, Bremaud).

- More generally: deterministic dynamics + random jumps in space and time

  $\Rightarrow$ PDP = Piecewise Deterministic (Markov) Processes (Davis).

Frameworks for studying stationarity, distribution, optimal control.

# Stationarity

Stationarity in the strict sense: $X(\cdot) = X(\cdot + s)$ in distribution.

In particular, $\mathbb{E}f(X(t_1)) = \mathbb{E}f(X(t_2))$

Stationarity in the mean: $\mathbb{E}X(t_1) = \mathbb{E}X(t_2)$.
Stationarity in covariance: $\mathbb{E}X(t_1)X(t_1 + s) = \mathbb{E}X(t_2)X(t_2 + s)$ for all $t_1$, $t_2$, $s$.
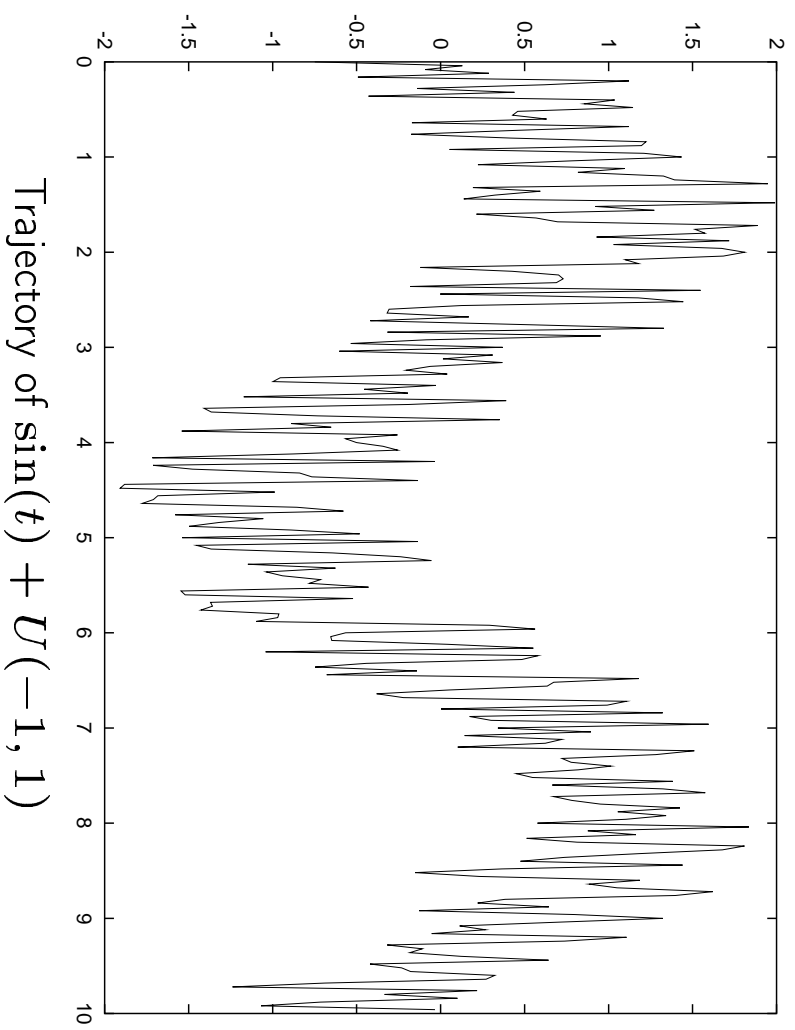
Stationarity excludes periodicity. Example:

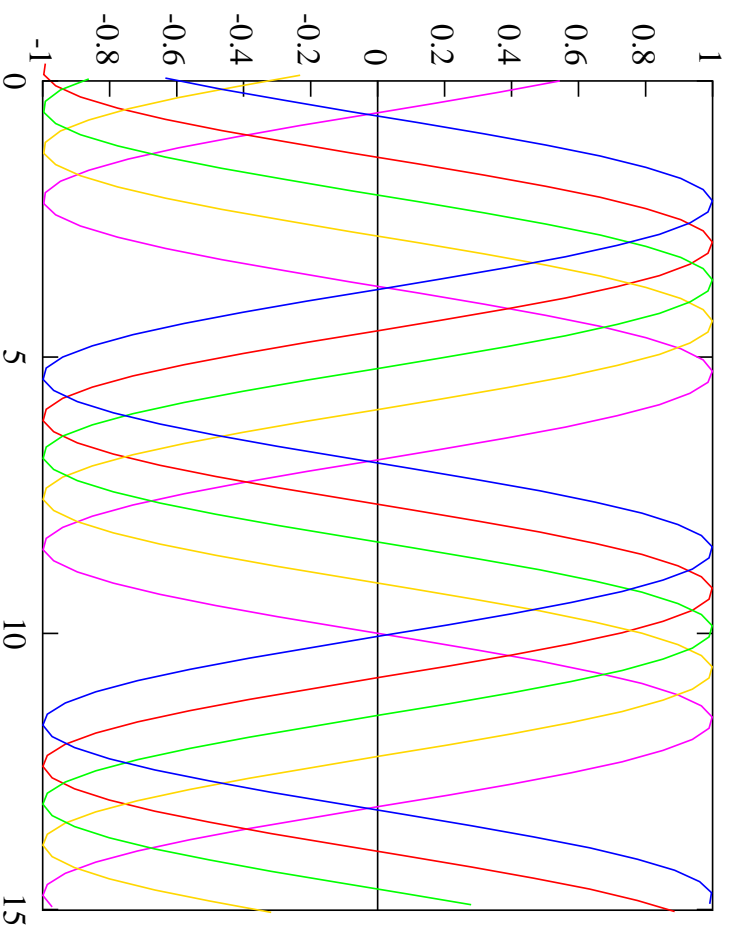$$X(t) = \sin(t) + \xi_t$$

with $\xi_t$ random and small.

Then

$$\mathbb{P}\{X(t + \pi) > 0\} \neq \mathbb{P}\{X(t) > 0\}.$$

Trajectory of $\sin(t) + U(-1, 1)$

But there exist processes essentially periodic and stationary:

$$X(t) = \sin(t + \xi), \quad \xi \sim U(0, \pi).$$

Trajectories of $\sin(t + \xi(\omega))$:

# Convergence

A process is in general not stationary, but it can become so:

$$X(t) \rightarrow X, \quad t \rightarrow \infty$$

$$X_n \rightarrow X, \quad n \rightarrow \infty$$

in distribution (or otherwise).

If for any $s$,

$$X[t, t+s] \rightarrow \hat{X}[0, s], \quad \text{in distribution} \quad t \rightarrow \infty.$$

The process converges to a *steady state*.

If convergence is fast enough, one can use the distribution of $X$ as an approximation for that of of $X(t)$.

# Ergodicity

Ergodicity: coincidence of spatial and temporal averages:

$$\mathbb{E}f(X(s)) = \lim_{T \to \infty} \frac{1}{T} \int_0^T f(X(t))\,\mathrm{d}t\,,$$

$$\mathbb{E}f(X(n)) = \lim_{N \to \infty} \frac{1}{N} \sum_{n=1}^N f(X_n)\,.$$

Application: the law of large numbers for *statistical estimators* of quantities.

There exist processes that are stationary but not ergodic.

# Autocorrelations

Autocorrelation:

$$R(s, s+t) = \mathbb{E}[X(t) X(t+s)].$$

Autocovariance: dependence of the state of $X$ a instant $t+s$ with respect to instant $t$.

$$h(t, t+s) = \operatorname{Cov}(X(t) X(t+s)) = \mathbb{E}[X(t) X(t+s)] - \mathbb{E}X(t) \mathbb{E}X(t+s).$$

If $X(t) \perp\!\!\!\perp X(t+s)$, then $h(t, t+s) = 0$.

Definition: $X$ stationary in the large sense (or at the second order): for any $t$:

$$h(t, t+s) = h(s) = \mathbb{E}X(0) \mathbb{E}X(s) - (\mathbb{E}X)^2.$$

Note: $h(0) = \operatorname{Var}(X)$.

## Memory

Total Autocorrelation:

continuous time $\qquad \int_0^\infty |h(s)|\,ds \qquad$ discrete time $\qquad \sum_{n=0}^\infty |h(n)|$ .

A process has a *short memory* if

$$\int_0^\infty |h(s)|\,ds < \infty.$$

Otherwise, it has a *long memory*.

Long memory $\Rightarrow$ a $\boxed{\text{slow decrease}}$ of the dependence of $X(t+s)$ et $X(t)$.

## Markov chains

$\{X(n),\, n \in \mathbb{N}\}$ is a homogeneous **discrete time Markov chain** if:

i/ (Markov property) $\forall t \in \mathbb{N}$, et $\forall$ $(j_0, j_1, \ldots, j_t, j_{t+1}) \in \mathcal{E}^{t+2}$:

$$\mathbb{P}\{X(t+1) = j_{t+1} | X(t) = j_t, \ldots, X(0) = j_0\} = \mathbb{P}\{X(t+1) = j_{t+1} | X(t) = j_t\}\, ;$$

ii/ (homogeneity) $\forall t \in \mathbb{N}$, et $(i, j) \in \mathcal{E} \times \mathcal{E}$,

$$\mathbb{P}\{X(t+1) = j | X(t) = i\} = P_{i,j}\, .$$

$P_{i,j},\, (i, j) \in \mathcal{E} \times \mathcal{E}$: transition probabilities

**P** *transition matrix.*

# Dynamics of probabilities

One looks for *transition probabilities at $n$ steps*:

$$p(i, j; n) = \mathbb{P}\{X(n) = j \mid X(0) = i\},$$

Let $P(n)$ be the matrix of $p(i, j; n)$. Then:

$$\boxed{P(n) = \mathbf{P}^n.}$$

Let now, for $n \in \mathbb{N}$ and $j \in \mathcal{E}$,
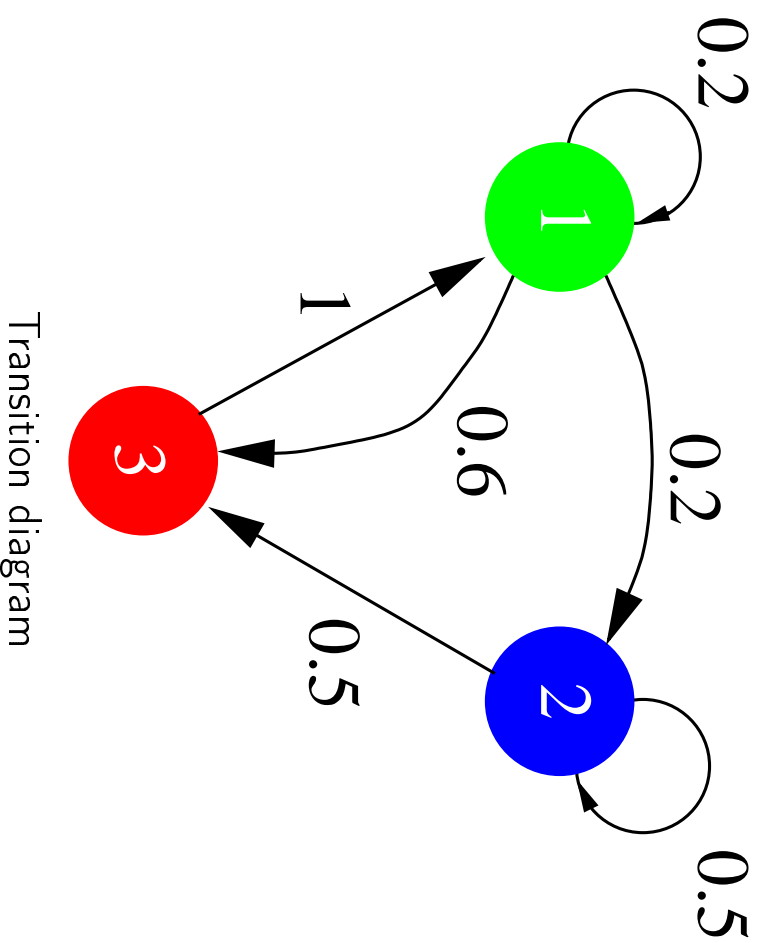
$$\pi_n(j) = \mathbb{P}\{X(n) = j\}.$$

Then:

Algebraic form: for any $n \in \mathbb{N}$:

$$\pi_n(j) = \sum_{i \in \mathcal{E}} \pi_0(i) \, p(i, j; n) \, .$$

$$\boxed{\pi_n = \pi_0 \, \mathsf{P}^n \, .}$$

# Example of Markov chain



Transition diagram

Transition Matrix:

$$P = \begin{pmatrix} 0.2 & 0.2 & 0.6 \\ 0 & 0.5 & 0.5 \\ 1 & 0 & 0 \end{pmatrix}.$$

Probability vectors:

$$
\begin{aligned}
\mathbf{p}_0 &= (1, 0, 0) \\
\mathbf{p}_1 &= (0.2 \ \ 0.2 \ \ 0.6) \\
\mathbf{p}_2 &= (0.64 \ \ 0.14 \ \ 0.22) \\
\mathbf{p}_3 &= (0.348 \ \ 0.198 \ \ 0.454) \\
\mathbf{p}_4 &= (0.5236 \ \ 0.1686 \ \ 0.3078) \\
&\vdots \quad \cdots \\
\mathbf{p}_\infty &= (5/11, 2/11, 4/11)
\end{aligned}
$$

# Equilibrium equations

If $\lim_n \boldsymbol{\pi}_n = \boldsymbol{\pi}$ exists, then:

$$\boldsymbol{\pi} = \boldsymbol{\pi} \, \mathsf{P} .$$

These equilibrium equations are written: $\forall i \in \mathcal{E}$,

$$\pi(i) = \sum_{j \in \mathcal{E}} \pi(j) \, P_{j,i} .$$

They define the *stationary probability*.

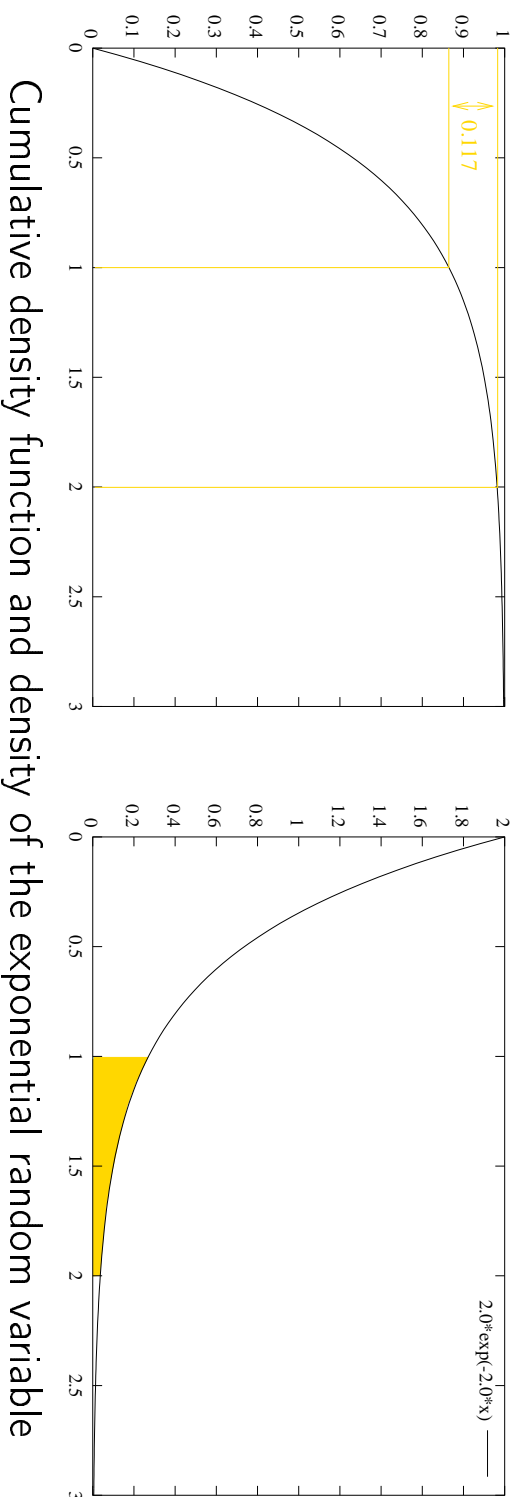The computation of stationary probabilities is reduced to the solution of a linear system!

Problem: it is often very large, and even infinite.

# The continuous time

A random variable $X$ has an *exponential* distribution of parameter $\lambda > 0$ ($X \sim \mathrm{Exp}(\lambda)$) if:

$$F_X(x) = \mathbb{P}\{X \leq x\} = 1 - e^{-\lambda x}.$$

Cumulative density function and density of the exponential random variable

The exponential distribution is *memoryless*: $\forall s, t > 0$,

$$\mathbb{P}\{X > s + t \mid X > s\} = \mathbb{P}\{X > t\}.$$

The family of exponential distributions is stable under minimization:

- If $X_1 \sim \text{Exp}(\lambda_1)$, $X_2 \sim \text{Exp}(\lambda_2)$ and $X_1$ and $X_2$ are independent: then

$$\min\{X_1, X_2\} \sim \text{Exp}(\lambda_1 + \lambda_2)$$

- Moreover:

$$\mathbb{P}\{\min\{X_1, X_2\} = X_i\} = \frac{\lambda_i}{\lambda_1 + \lambda_2}.$$

# The Poisson process

Consider a random sequence $T_0 \leq T_1 \leq \ldots \leq T_n \leq T_{n+1} \leq \ldots$. The counting process:

$$N(a, b) = \#\{n \mid a \leq T_n < b\} = \sum_{n=0}^{\infty} \mathbf{1}_{\{a \leq T_n < b\}}$$

is a **Poisson process** of parameter $\lambda$ if $\{T_{n+1} - T_n\}$ is a i.i.d. sequence of variables $\mathrm{Exp}(\lambda)$.

For all $u$:

$$\mathbb{P}\{N(x, x+u) = k\} = \frac{(\lambda u)^k}{k!} e^{-\lambda u} .$$

In particular, $\mathbb{E}N(x, x+u) = \lambda u$: $\lambda$ is the *arrival rate* of the process.

Limit Theorem: if one superposes a large number of "rare" processes, the resulting process is asymptotically Poisson.

# Continuous time Markov chains

Let $\{X(t),\, t \in \mathbb{R}^+\}$, having the following properties. When $X$ enters state $i$:

- $X$ stays in state $i$ a random time, exponentially distributed with parameter $\tau_i$, independent of the past; then

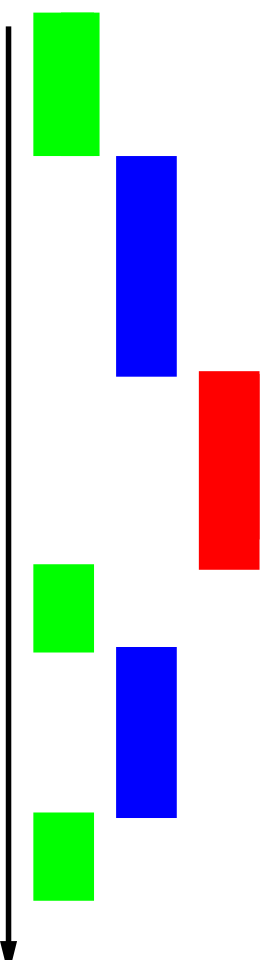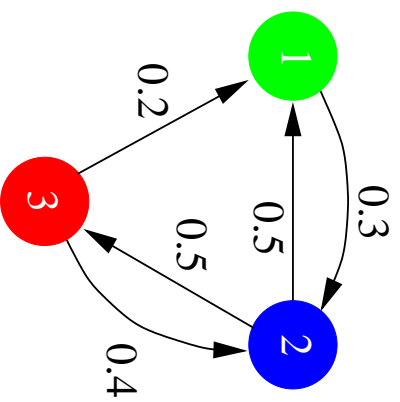- $X$ jumps instantly in state $j$ with probability $p_{ij}$. We have $p_{ij} \in [0, 1]$, $p_{ii} = 0$ and

$$\sum_j p_{ij} = 1.$$

The exponential distribution being *memoryless*, we obtain a process which has the property that:

$$\mathbb{P}\{X(t_{n+1}) = j_{n+1}|X(t_n) = j_n, \ldots, X(t_0) = j_0\}$$

$$= \mathbb{P}\{X(t_{n+1}) = j_{n+1}|X(t_n) = j_n\}.$$

## Example

$$\tau = \begin{pmatrix} 0.3 \\ 1 \\ 0.6 \end{pmatrix} \qquad P = \begin{pmatrix} 0 & 1 & 0 \\ \frac{1}{2} & 0 & \frac{1}{2} \\ \frac{1}{3} & \frac{2}{3} & 0 \end{pmatrix} \qquad Q = \begin{pmatrix} -0.3 & 0.3 & 0 \\ 0.5 & -1.0 & 0.5 \\ 0.2 & 0.4 & -0.6 \end{pmatrix} .$$

# The Definition

A process $\{X(t), t \in \mathbb{R}^+\}$ is a homogeneous continuous time Markov chain (or: *Markov process*) iff:

i/ (Markov property) For all $n \in \mathbb{N}$, every $n + 2$-uple of reals $t_0 < t_1, < \cdots < t_n < t_{n+1}$ and every $n + 2$-uple $(j_0, j_1, \ldots, j_n, j_{n+1})$ of elements of $\mathcal{E}$:

$$\mathbb{P}\{X(t_{n+1}) = j_{n+1}|X(t_n) = j_n, \ldots, X(t_0) = j_0\}$$
$$= \mathbb{P}\{X(t_{n+1}) = j_{n+1}|X(t_n) = j_n\} ;$$

ii/ (homogeneity) For all reals $s, t$ and $u$, and every pair $(i, j)$ of $\mathcal{E}$, independently of $t$ we have:

$$\mathbb{P}\{X(t + u) = j|X(s + u) = i\} = \mathbb{P}\{X(t) = j|X(s) = i\} = P_{t-s}(i, j) .$$

# Dynamics of probabilities

Chapman-Kolmogorov equations:

$$P_{t+s}(i,j) = \sum_{k \in \mathcal{E}} P_t(i,k)\, P_s(k,j) \,,$$

or, in algebraic form:

$$\mathsf{P}_{t+s} = \mathsf{P}_t\, \mathsf{P}_s \,,$$

If the process $\{X(t)\}$ is "regular" enough, then there exists a matrix $\mathsf{Q} = \mathsf{P}'(t)$ such that:

$$\frac{d\mathsf{P}_t}{dt} = \mathsf{Q}\mathsf{P}_t = \mathsf{P}_t\mathsf{Q} \,.$$

This is *infinitesimal generator*.

Then:

$$\mathbf{P}_t \;=\; e^{t\mathbf{Q}}$$

Theoretically, computation of transient probabilities, of the speed of convergence etc.

$$\mathbf{P}_t \;=\; \mathbf{p}_0 \mathbf{P}_t \;=\; \mathbf{p}_0 e^{t\mathbf{Q}} \,.$$

# Construction of generators

Under the evolution assumptions above, the process $\{X(t), t \in \mathbb{R}^+\}$ is a CTMC of infinitesimal generator:

$$
\begin{cases}
q(i, j) & = \tau_i p_{ij} \quad \text{if } i \neq j \\
q(i, i) & = -\tau_i \,.
\end{cases}
$$

## Construction #2.

Consider a stochastic process in continuous time, $\{X(t), t \in \mathbb{R}^+\}$, having the following properties. When $X$ enters state $i$:

- For each state $j \neq i$, a random variable $Y_{ij}$ with exponential distribution of parameter $\mu_{ij}$, is drawn, independently between them and of the past. It is possible that $\mu_{ij} = 0$, in which case $Y_{ij} = +\infty$.

- The minimum of the $Y_{ij}$ is one of them: $Y_{ik}$. At time $Y_{ik}$, $X$ jumps instantly in state $k$.

Then $\{X(t), t \in \mathbb{R}^+\}$ is a CTMC of infinitesimal generator $Q$:

$$
\begin{cases}
q(i, j) = \mu_{ij} & \text{if } i \neq j \\
q(i, i) = -\sum_{j \neq i} \mu_{ij}.
\end{cases}
$$

# Equilibrium equations

If $\lim_t \boldsymbol{\pi}_n = \boldsymbol{\pi}$ exists, then:

$$\mathbf{0} = \boldsymbol{\pi} \, Q \, .$$

These equilibrium equations can be written: $\forall i \in \mathcal{E}$,

$$\left( \sum_{j \neq i} q_{i,j} \right) \pi(i) = \sum_{j \neq i} \pi(j) q_{j,i} \, .$$
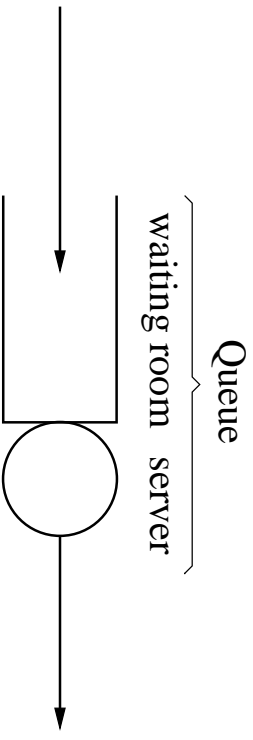
Interpretation: entering flow = outgoing flow.

Generalization: global equilibrium equations. For $S \subset \mathcal{E}$:

$$\sum_{i \in S, j \in \overline{S}} \pi(i) \, q_{i,j} = \sum_{i \in \overline{S}, j \in S} \pi(i) q_{i,j} \, .$$

# Part II: Queuing Theory

- Discrete queues, fluid queues

- Arrival process, service process; Kendall's notation

- Performance measures: number of customers, waiting/response time, loss probability, jitter

- Dynamics of the queue; workload curves; evolution equations

- Capacity: finite or infinite?

- Simple queues or networks of queues?

- Stochastic models of traffic
  - I.i.d processes
  - Poisson process
  - Markov modulated processes

# Queues

waiting room   server

Queue

arrivals        wait      service     departure

$a_1, a_2 \ldots$              $\sigma_1 \, \sigma_2 \ldots$

Usual representation of a queue

The elements that compose a queue are:

- one or several servers

- a waiting room

- (possibly) classes of customers

- an arrival process per class

- a process of service durations

- a service discipline

# Kendall's notation

This notation allows to identify certain queues among the variety of possibilities.

A queuing model is denoted by:

$$A/S/P/K/D$$

**A** the inter-arrival distribution

**S** the service time distribution

**P** the number of servers

**K** the size of the waiting room (by default: $\infty$)

**D** the discipline of service (by default: FIFO)

Examples: the queue M/M/1, M/GI/1/K, etc.

## Performance measures

**Stability condition** Under which conditions the queue admits a stationary behavior? $X(t)$ is a dynamic quantity:

$$\lim_{t \to \infty} \mathbb{P}\{X(t) \leq x\} = ?$$

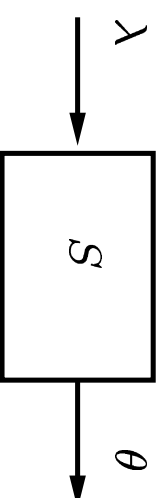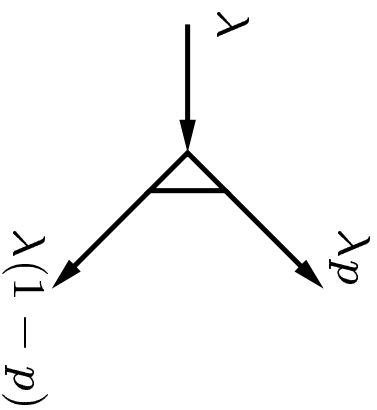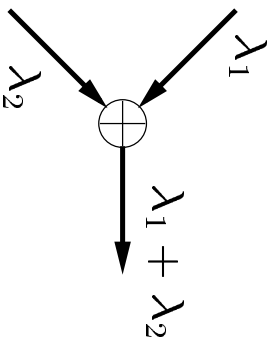**Throughput** If $N(a, b)$ counts the number of arrivals in $[a, b[$, the throughput of arrivals:

$$\lambda = \limsup_{t \to \infty} \frac{N(0, t)}{t} = \mathbb{E}N(0, 1) = \limsup_{n \to \infty} \frac{n}{a_n}.$$

If the departure instants of customers are $d_1, \ldots, d_n, \ldots$, the throughput of outputs is:

$$\theta = \limsup_{n \to \infty} \frac{n}{d_n}.$$

The throughputs are conserved:

If stability:

$\lambda_1$

$\lambda_2$

$\lambda_1 + \lambda_2$

$\lambda$

$\lambda(1-p)$

$\lambda p$

$\lambda$

$S$

$\theta$

Laws of conservation of throughputs

$$\boxed{\lambda = \theta.}$$

**Utilization**  Fraction of the time some resource is used:

$$\rho = \limsup_{T \to \infty} \frac{U(0, T)}{T} \, .$$

**Response time**  $R_n = d_n - a_n$.

Also: waiting time $W_n$ and service time $\sigma_n$:

$$R_n = W_n + \sigma_n \, .$$

**Loss rate/probability**  Fraction of customers "lost". Ratio of "effective" throughput to the "offered" throughput.

**Cycle time**  For cyclic systems.

**Jitter**  Measure of the variability of the network's response:

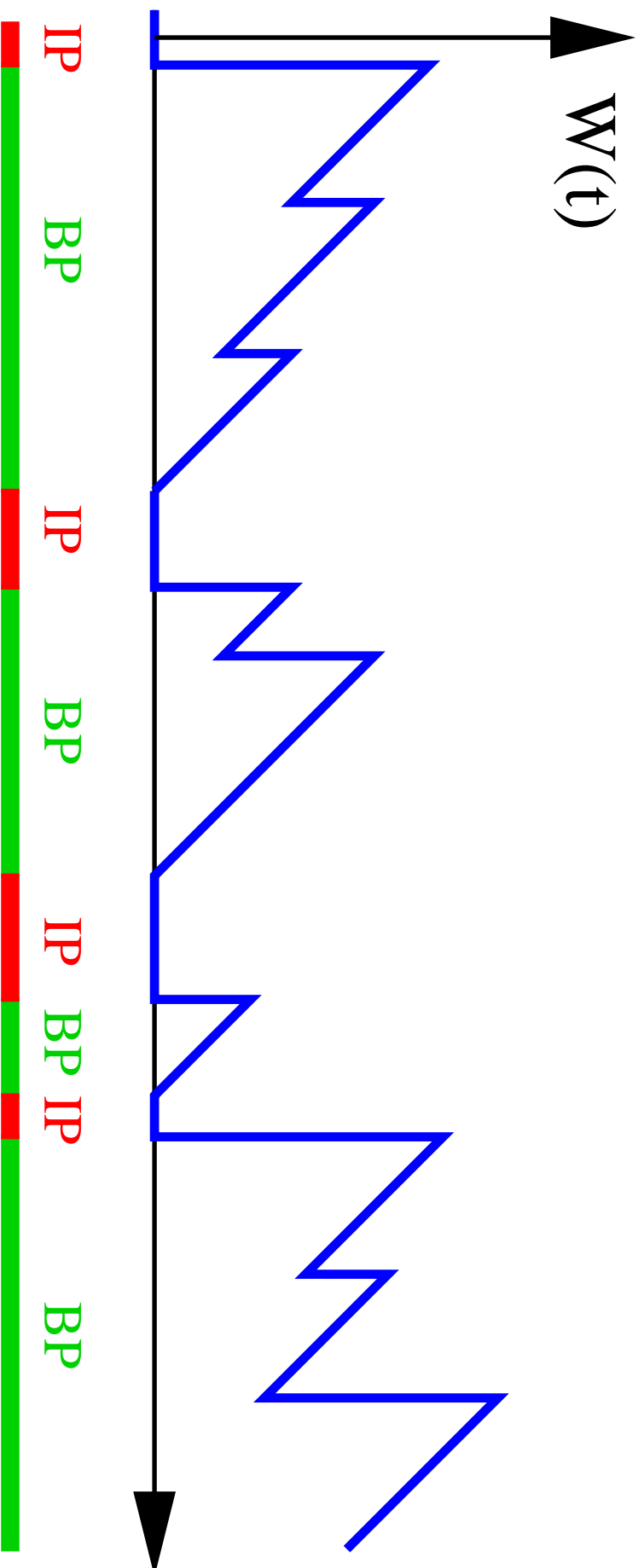$$J_n = |(d_{n+1} - d_n) - (a_{n+1} - a_n)| = |R_n - R_{n+1}| \, .$$

# Dynamics of a queue

Fundamental quantities:

$N(t)$   number of customers present in the system at time $t$;

$W(t)$   quantity of work (workload, backlog) present in the queue at time $t$

Evolution of $W(t)$: the workload curve.
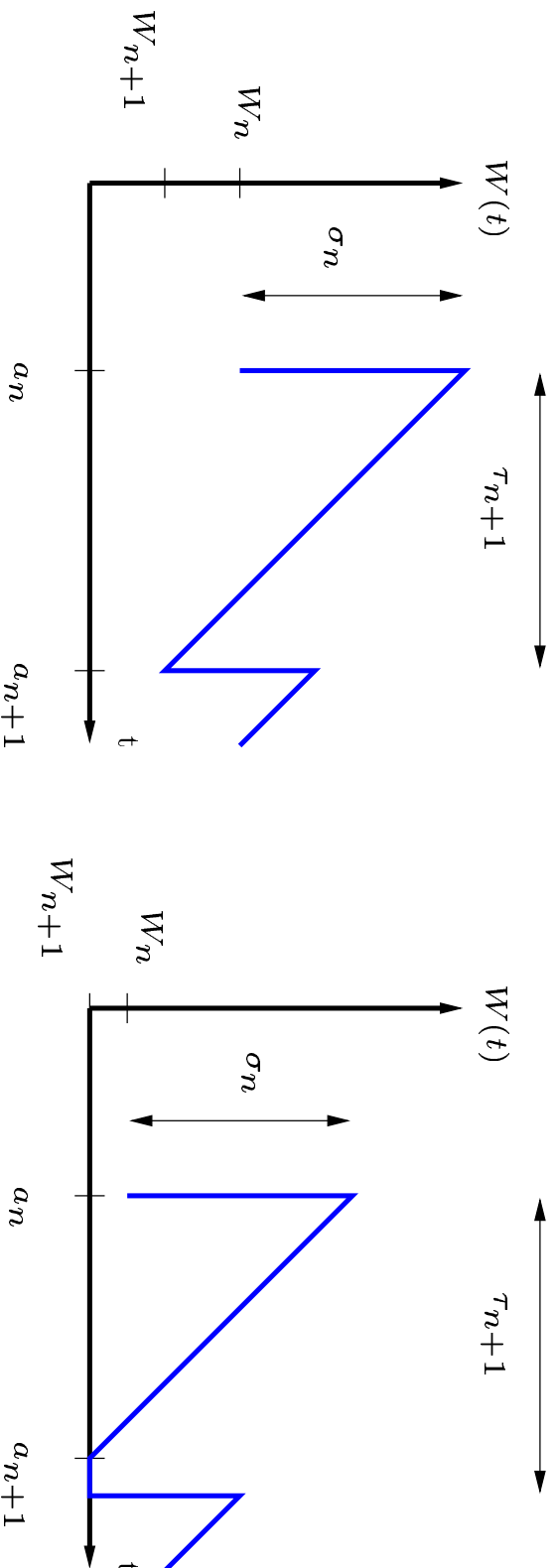
# W(t)

IP · BP · IP · BP · IP · BP · IP · BP

A workload curve

Busy (or Activity) periods (AP) and Idle periods (IP).
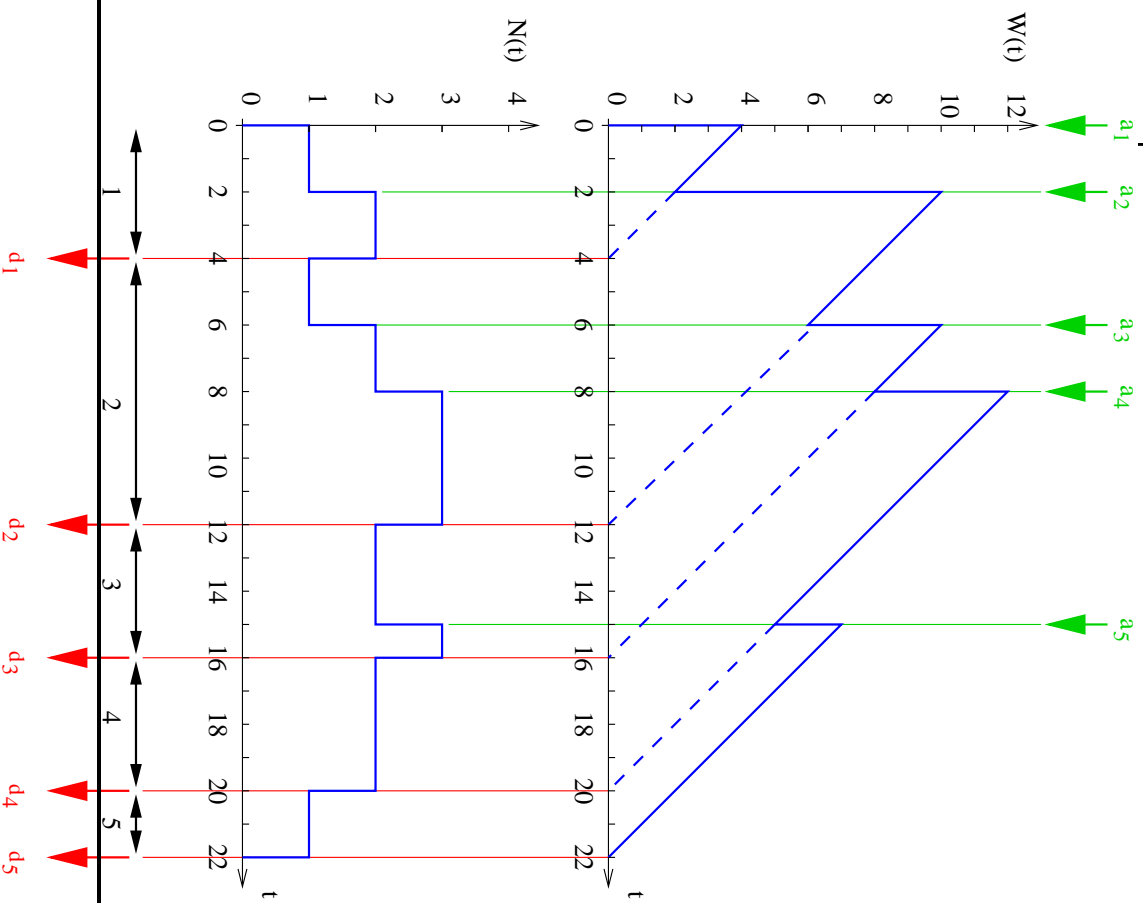
# Waiting times – the FIFO case

$W_n$: waiting time of customer $n$ before service. $\quad \sigma_n$: duration of the service of customer $n$.

**Lindley's Equation**:

$$W_{n+1} = [W_n + \sigma_n - \tau_{n+1}]^+.$$

# Relationship between number of customers and workload

# Virtual waiting time, real waiting time

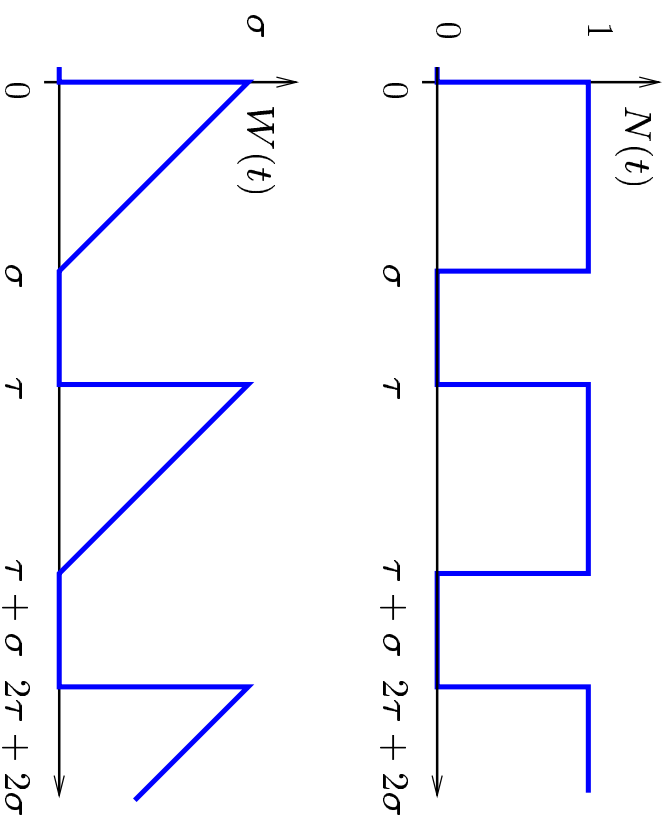If $a_n$ is the arrival time of customer $n$ and if FIFO:

$$W_n = W(a_n^-) .$$

$\Rightarrow W(t)$ is also called: virtual waiting time.

**Warning!** $W_n$ and $W(t)$ do not necessarily have the same distribution!

Example: the D/D/1 queue.
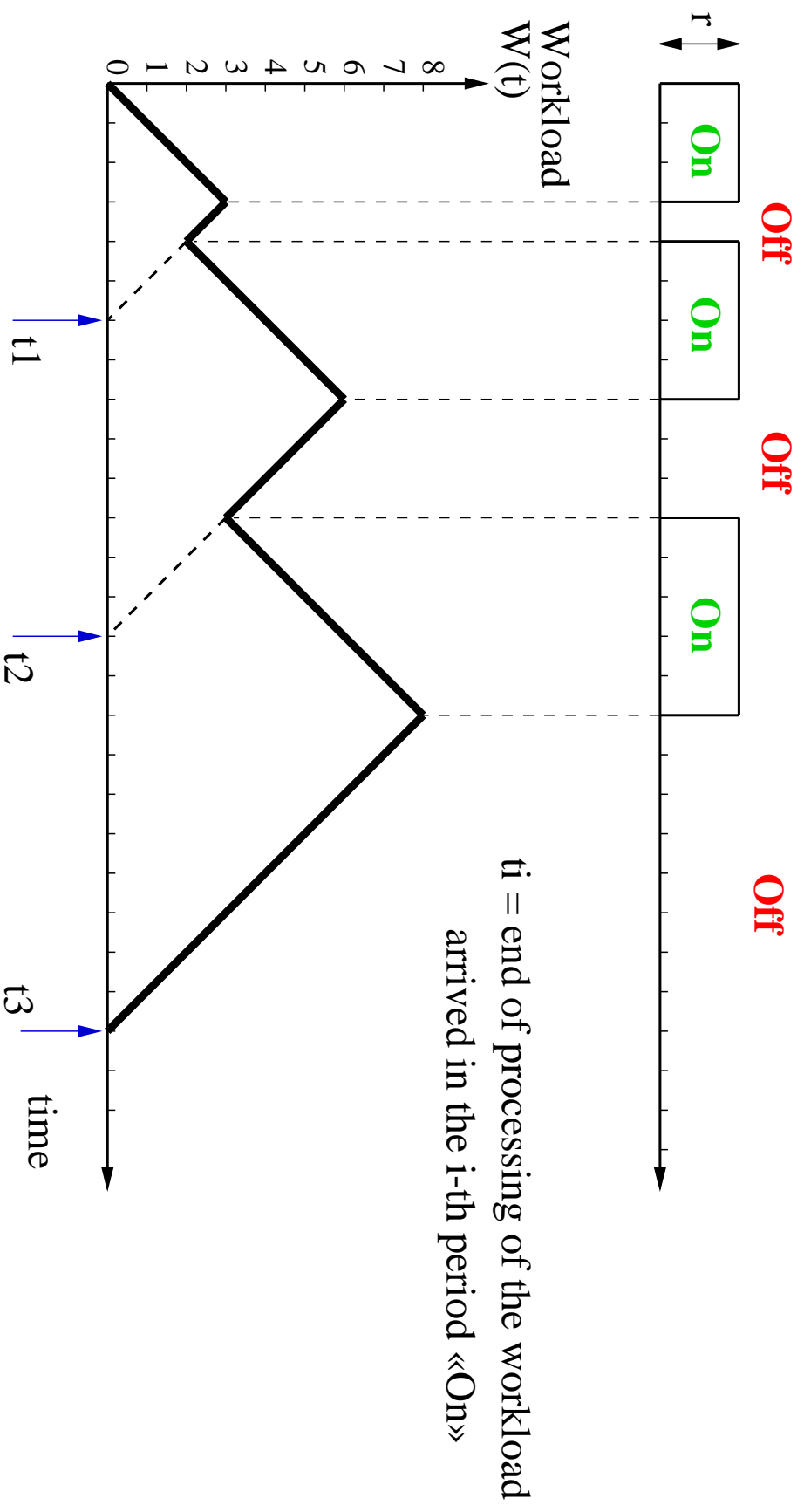


**Property PASTA** (Poisson Arrivals See Time Averages)

arrivals are Poisson, the stationary distributions of $W_n$ and $W(t)$ <span style="color:red">do</span> coincide.

## Fluid queues

No more customers, but some "fluid" arriving with a certain rate $r(t)$ (variable) and served at a certain speed $C$ (possibly variable too).

Example: arrivals according to an "on/off" process (typical of digitized voice, video, etc.):

# General Results

## Stability

Stability: $W_n$ admits a stationary regime.

Result:

The G/G/1 queue is stable if and almost if

$$\boxed{\mathbb{E}\sigma \; < \; \mathbb{E}\tau}$$

## Little's formula

The average response time $R$ and the average number of customers $N$ are linked by the formula:

$$\boxed{\lambda T \; = \; N}$$

# Traffic models

The traffic is described by:

- the arrival process $\{a_n\}_{n\in\mathbb{N}}$ or the distribution of inter-arrivals $\{\tau_n\}_{n\in\mathbb{N}}$.

- the service process $\{\sigma_n\}_{n\in\mathbb{N}}$.

"iid" models. Distribution of the inter-arrival time is fixed + independence. Idem for services.

Classical cases: deterministic, exponential distribution:

$$\mathbb{P}\{\tau > x\} = e^{-\lambda x}$$

Gamma/Erlang distribution (sums of exponentials).

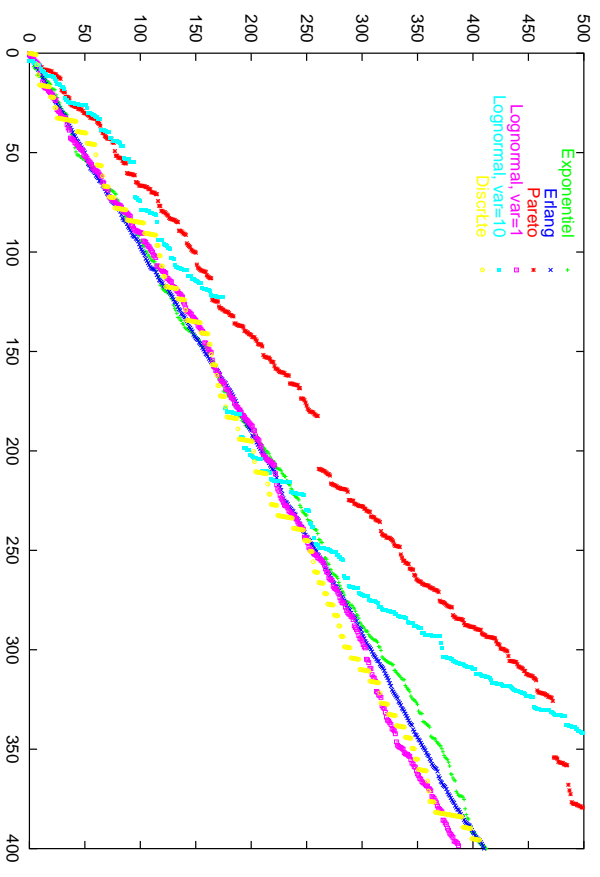New trends: laws with a "heavy tail": Pareto

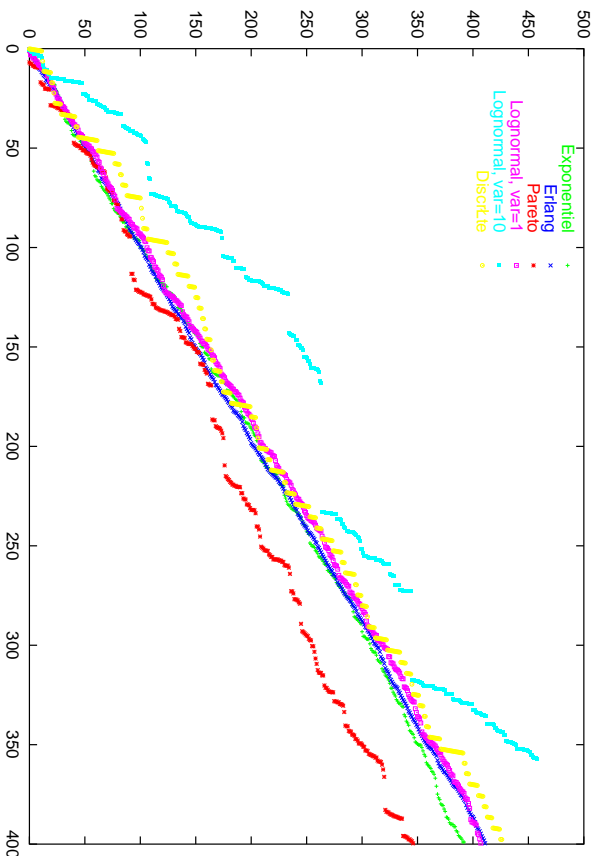$$\mathbb{P}\{\tau > x\} = \left(\frac{a}{a+x}\right)^{\alpha}.$$

Weibull, LogNormal.

$$\mathbb{P}\{\tau > x\} = \mathbb{P}\{X > \log(x)\}, \quad X \sim \mathcal{N}(m, \sigma).$$

# Example

Comparison of irregularities in arrival times for various laws of $\tau$.



x-axis: $a_n$, y-axis: $n$.

Markov modulated arrivals

**MAP: Markov Arrival Process**: The arrivals occur at the times where a Markov chain changes state.

# MMPP: Markov Modulated Poisson Process: Arrivals according to Poisson processes with an intensity depending on the state of a Markov chain (or a semi-Markov process).

In particular, the IPP processes: Interrupted Poisson Processes.

# MMRP: Markov Modulated Rate Process: according to a fluid process of rate depending on the state of a Markov chain.

# Superposition of sources

If several sources of traffic are superimposed, the resulting process is still modulated by Markov.

$$
\begin{pmatrix}
-\lambda & \lambda & 0 \\
0 & -\mu & \mu \\
\nu & 0 & -\nu
\end{pmatrix}
\oplus
\begin{pmatrix}
-\alpha & \alpha \\
\beta & -\beta
\end{pmatrix}
=
\begin{pmatrix}
- & \lambda & 0 & \alpha & 0 & 0 \\
\beta & - & 0 & 0 & \alpha & 0 \\
\nu & 0 & - & \mu & 0 & \alpha \\
0 & \nu & \beta & 0 & - & \lambda \\
0 & 0 & \nu & \beta & 0 & -\mu \\
0 & 0 & 0 & \nu & \beta & -
\end{pmatrix}
$$

# Modeling hypotheses

- Finite or infinite capacity?

  Queues with infinite capacity are easier to analyze: they can be used as approximations.

$$\mathbb{P}\{\text{loss}\} \leftrightarrow \mathbb{P}\{N = K\} \leftrightarrow \mathbb{P}\{W > K\}$$

- Modeling networks?

  Few results exist on networks of queues. Analysis relies on the single bottleneck assumption (often valid).

- What traffic models? Compromise between what can be calculated and what is reasonable in practice.

# Part III: Exact analysis

- Exact analysis in the case of infinite buffers
  - the M/M/1 queue, the M/GI/1 queue, the GI/M/1 queue.
  - the MMPP/GI/1 queue
- Exact analysis in the case of finite buffers
  - the M/M/1/K queue.
  - QBDs (Quasi Birth-Death processes)
- Queueing networks

# The M/M/1 queue

Characteristics: Infinite waiting room, 1 server.

**inter-arrivals:** exponential distribution with parameter $\lambda$:

**services:** exponential distribution with parameter $\mu$:

$$\mathbb{P}\{\tau \leq x\} = 1 - e^{-\lambda x}, \qquad \mathbb{P}\{\sigma \leq x\} = 1 - e^{-\mu x}.$$

Stability: $\lambda < \mu$.

$\{N(t)\}$ is a Markov Chain: a *birth and death process*

Performances:

$$\mathbb{P}\{W > x\} = \frac{\lambda}{\mu} e^{-(\mu - \lambda)x}$$

$$\mathbb{P}\{N \geq n\} = \left(\frac{\lambda}{\mu}\right)^n$$

# The M/GI/1 queue

Arrivals: exponentials, rate $\lambda$,

Services: arbitrary distribution, Laplace transform $S^*(s)$.

The Laplace transform of the waiting time and of the number of customers (*Pollaczek-Khinchine* formula):

$$W^*(s) = \frac{1 - \rho}{s - \lambda(1 - S^*(s))}$$

$$N^*(z) = S^*(\lambda(1 - z)) \frac{(1 - \lambda/\mu)(1 - z)}{S^*(\lambda(1 - z)) - z}.$$

In particular, the averages are:

$$\mathbb{E}W = \frac{\lambda \mathbb{E}\sigma^2}{2(1-\rho)} \qquad \mathbb{E}N = \rho + \rho^2 \frac{\mu^2 \mathbb{E}\sigma^2}{2(1-\rho)}.$$

$\rho = \lambda/\mu$: is the utilization.

# The GI/M/1 queue

Arrivals: arbitrary law, Laplace transform $A^*(s)$,

Services: exponential, average $1/\mu$.

Distribution of the waiting time:

$$\boxed{\; \mathbb{P}\{W > x\} \;=\; \theta \, e^{-\mu(1-\theta)x} \;, \;}$$

with:

$$\theta \;=\; A^*(\mu(1-\theta)) \;.$$

$\Rightarrow$ exponential distribution!

## Equivalent service rate

$$\widehat{\lambda} \;=\; \theta\mu \;.$$

$\Rightarrow$ equivalent Bandwidth for networks.

# The MMPP/GI/1 queue

Arrivals: MMPP with $N$ states, generator $\mathbf{Q}$ and matrix of rates $\mathbf{\Lambda}$;

Services: independent with a general distribution $H(x)$, of Laplace transform $H^*(s)$.

Distribution of the workload $W$:

$$W^*(s) = s(1-\rho)\,\mathbf{g}\,[s\mathbf{I} + \mathbf{Q} - (1 - H^*(s))\mathbf{\Lambda}]^{-1}\,\mathbf{1},$$

$\mathbf{g}$ vector to be determined.

If $\sigma \sim \mathrm{Exp}(\mu)$ (the MMPP/M/1 queue), then:

$$\mathbb{P}\{W > x\} = \sum_{k=0}^{N} a_k e^{-\theta_k x} \sim a_1 e^{-\theta_1 x},$$

with $\theta_k$ such that:

$$det\{-\theta_k s \mathbf{I} + \mathbf{Q} - (1 - H^*(-\theta_k))\mathbf{\Lambda}\} = 0.$$

$\Rightarrow$ again: tail of distribution asymptotically exponential.

# The M/M/1/K queue

As the M/M/1 but with a finite capacity $K$.

Markov Chain: it is finite

Performances: let $\rho = \lambda/\mu$.

$$\mathbb{P}\{N = K\} = \rho^K \frac{1 - \rho}{1 - \rho^{K+1}}.$$

Probability of losing a customer: it is precisely $\mathbb{P}\{N = K\}$ (PASTA).

Note: with the approximation of an infinite buffer, one would have had:

$$\mathbb{P}\{N = K\} \simeq \rho^K (1 - \rho).$$

# QBD processes

QBD: "quasi birth-death".

This is a structure of Markov chain obtained when the arrival process or the service process has "phases".

For example:



They are solved using (for instance) Neuts' method.

## Neuts' method for QBD

The state space $\mathcal{E}$ is partitioned in finite blocks of the same size

$$\mathcal{E}_k = \{(k, 1), (k, 2), \ldots, (k, N)\}$$

such that the transition matrix of the Markov chain has the $N \times N$ *block tridiagonal* structure:

$$P = \begin{pmatrix} S_0 & L & & & \\ M & S & L & & \\ & M & S & L & \\ & & \ddots & \ddots & \ddots \\ & & & M & S & L \\ & & & & M & S_K \end{pmatrix}$$

The stationary probabilities are also grouped in blocks:

$$\pi_k = (\pi_{k,1}, \dots, \pi_{k,N})$$

The equilibrium equation $\pi P = \pi$ becomes:

$$\begin{cases} \pi_{k-1} L + \pi_0 S_0 + \pi_1 M = \pi_0 \\ \pi_{k-1} L + \pi_k S + \pi_{k+1} M = \pi_k \quad 0 < k < K \\ \pi_{K-1} L + \pi_K S_K = \pi_K \end{cases}$$

$\Rightarrow$ numerical resolution of the recurrence by iterative methods.

The same analysis for continuous time Markov chains.

# Product form solutions: Jackson Networks

$N$ queues (stations) which services are $\sim$ Exp:

- the vector $\boldsymbol{\lambda}_0 = (\lambda_{0,1}, \ldots, \lambda_{0,N})$ of external arrivals rates in each queue,
- the vector $(\mu_1, \ldots, \mu_N)$ of service rates,
- square matrix $N \times N$ of internal routing $\mathbf{R}$: $r_{i,j} = \mathbb{P}\{$a customer going out of $i$ goes to $j\}$.

Entering flow in stations: vector $\boldsymbol{\lambda} = (\lambda_1, \ldots, \lambda_N)$ solution of:

$$\boldsymbol{\lambda} = \boldsymbol{\lambda}_0 + \boldsymbol{\lambda} \mathbf{R}.$$

The stability condition of the system:

$$\forall 1 \leq i \leq N, \quad \lambda_i < \mu_i.$$

If stability, the stationary probability distribution is:

$$p(n_1, \ldots, n_N) = \prod_{i=1}^{N} \left( 1 - \frac{\lambda_i}{\mu_i} \right) \left( \frac{\lambda_i}{\mu_i} \right)^{n_i},$$

$\Rightarrow$ as if queues were $M/M/1$ in isolation, and independent

$\Rightarrow$ justification of the end-to-end response time formula:

$$T = \sum_{i=1}^{N} \frac{1}{C_i - L_i}$$

$C_i$: capacity of the link/switch, $L_i$: entering traffic.

## Kelly networks

Customers belong to several classes.

To each class $k$ corresponds a *route* in the network:

$$r_k = \left( r_k^1, \ldots, r_k^{n_k} \right).$$

Customers arrive according to Poisson processes, and servers deliver service times with an exponential distribution.

The flow (throughput) $\hat{\lambda}_{ik}$ of customers of class $k$ entering in queue $i$:

$$\hat{\lambda}_{ik} = \lambda_k \times (\text{number of } i \text{ in } r_k).$$

Total throughput, all classes aggregated:

$$\hat{\lambda}_i = \sum_k \hat{\lambda}_{ik}.$$

A Kelly network



Routes

| | |
|---|---|
| (red) | $(1,2,3)$ |
| (blue) | $(1,2,5)$ |
| (green) | $(1,4)$ |
| (black) | $(1,4,5)$ |

## Stationary probabilities

Let $M = ((m_{ik}))$ be a matrix of populations per queue and per class. The stationary probability that the network is in the state $M$ is:

$$\mathbb{P}\{M\} = \prod_{i=1}^{N} \left(1 - \frac{\hat{\lambda}_i}{\mu_i}\right) \left(\sum_{k=1}^{K} m_{ik}\right)! \prod_{k=1}^{K} \frac{1}{m_{ik}!} \left(\frac{\hat{\lambda}_{ik}}{\mu_i}\right)^{m_{ik}} .$$

Marginal probabilities: if $\mathbf{n} = (n_1, n_2, \ldots, n_K)$ is a possible population at queue $i$:

$$\mathbb{P}\{\mathbf{n}\} = \left(1 - \frac{\hat{\lambda}_i}{\mu_i}\right) \left(\sum_{k=1}^{K} m_{ik}\right)! \prod_{k=1}^{K} \frac{1}{m_{ik}!} \left(\frac{\hat{\lambda}_{ik}}{\mu_i}\right)^{m_{ik}} .$$

If $\mathbf{m} = (m_1, m_2, \ldots, m_N)$ is a possible network population, then:

$$\mathbb{P}\{\mathbf{m}\} = \prod_{i=1}^{N} \left(1 - \frac{\hat{\lambda}_i}{\mu_i}\right) \left(\frac{\hat{\lambda}_i}{\mu_i}\right)^{m_i} .$$

# Statistics

Average number of customers in queue $i$,

$$\overline{N}_i = \frac{\hat{\lambda}_i}{\mu_i - \hat{\lambda}_i}$$

End-to-end response time on route $k$:

$$\overline{T}_k = \frac{1}{\lambda_k} \sum_{j=1}^{n_k} \frac{\hat{\lambda}_{r_k^j, k}}{\mu_{r_k^j} - \hat{\lambda}_{r_k^j}},$$

Average response time, all classes aggregated:

$$\overline{T} = \frac{1}{\sum_{k=1}^{K} \lambda_k} \sum_{i=1}^{N} \frac{\hat{\lambda}_i}{\mu_i - \hat{\lambda}_i}.$$

# Questions for networks

There exist extensions to the product forms of Jackson and Kelly networks: multiclass networks, mixt open/closed, and with various service disciplines: The BCMP theorem.

Numerous questions stay open. For example:

- less restrictive assumptions on traffic models: non-Poisson arrival processes, non-exponential services
- finite capacities, losses, feedback
- service disciplines and stability
- distributions of end-to-end response times

# Part IV: Asymptotic Analysis

- Principle

- Bounds and exponential asymptotics
  - Chernoff Bounds and Kingman's bound
  - Markov Additive Processes
  - Equivalent Bandwidth

- Long memory, autosimilarity, sub-exponentiality
  - Autosimilar Processes in Nature
  - Sub-exponentiality and Asymptotic Dominance
  - Long Memory and Finite Capacity

## Principle

Direct asymptotic analysis: find an equivalent to:

$$\mathbb{P}\{W > x\}, \quad x \to \infty$$

of which one hopes to find an approximation.

Typically: an *exponential asymptotic equivalent*:

$$\mathbb{P}\{W > x\} \sim C\, e^{-\theta x}, \quad x \to \infty.$$

Bounds: one tries to find bounds of this nature (for $x$ "large", or for all $x$)

$$B(\theta)\, e^{-\theta x} \leq \mathbb{P}\{W > x\} \leq C(\theta)\, e^{-\theta x},$$

# Chernoff's bound

Let $t$ be a fixed real number, and $X$ a random variable.

Laplace-Stiltjes transform of $X$:

$$X^*(s) = \mathbb{E}(e^{-sX}).$$

The following holds:

$$\mathbf{1}_{\{x \leq t\}} \leq e^{\theta(x-t)} \qquad \forall x, \theta$$

$$\mathbb{E}\mathbf{1}_{\{X \leq t\}} \leq \mathbb{E}e^{\theta(X-t)} \qquad \forall \theta$$

$$\mathbb{P}\{X \leq t\} \leq X^*(-\theta)\, e^{-\theta t} \qquad \forall \theta$$

$$\mathbb{P}\{X \leq t\} \leq \inf_{\theta}\{X^*(-\theta)\, e^{-\theta t}\}$$

## Kingman's bound

We consider the GI/GI/1 queue

For all number $\theta \geq 0$, such that:
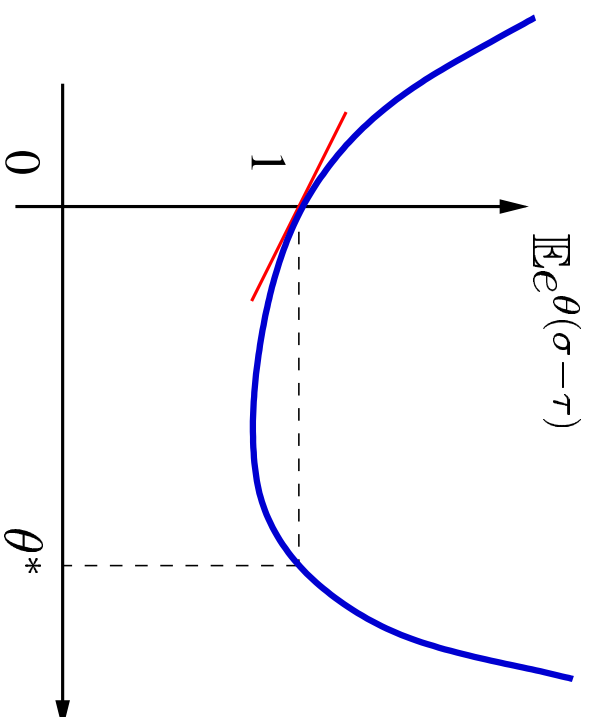
$$\mathbb{P}\{W > x\} \leq e^{-\theta x}$$

$$\mathbb{E}\left(e^{\theta(\sigma-\tau)}\right) \leq 1.$$

Hence, taking the largest possible $\theta$:

$$\theta^* = \sup\{\theta \geq 0 \mid \mathbb{E}(e^{\theta(\sigma-\tau)}) \leq 1\}$$

Conclusion: exponential decrease in the case $\theta^* > 0$.

Illustration: value of $\theta$.

# Large deviations

This result is generalized to arrival/service processes less simple:

If the process $\{U_n\} = \{\sigma_n - \tau_n\}$, stationary and ergodic, satisfies:

$$\Phi(\theta) = \lim_{t \to \infty} \frac{1}{t} \log \mathbb{E}[e^{\theta(U_0 + \cdots + U_{n-1})}],$$

then:

$$\lim_{x \to \infty} \frac{1}{x} \mathbb{P}\{W > x\} = -\theta^*.$$

with

$$\theta^* = \sup\{\theta \geq 0 \mid \Phi(\theta) = 0\}.$$

# Long memory and autosimilarity

Measures have shown that process of arrival of information exhibits a certain <span style="color:red">autosimilarity</span> and <span style="color:red">long term correlation</span>.

But the "classical" models do not have this property.

From where does this phenomenon come from?

What is the influence of this long term memory on the loss probabilities? Is it necessary to throw the known models away?

What new models are analyzable? Models with arrivals/services "heavy tailed".

Notion of sub-exponentiality of probability distributions.

# Autosimilarity

Let $\mathbf{X} = \{X(n)\}_n$ be a stationary process in the large sense.

$\mathbf{X}$ is autosimilar if:

$$\mathbf{X} =_d \frac{1}{m^H} \left( X_{t(m-1)+1} + \ldots + X_{tm} \right)$$

for all $m$.

$H$: Hurst parameter.
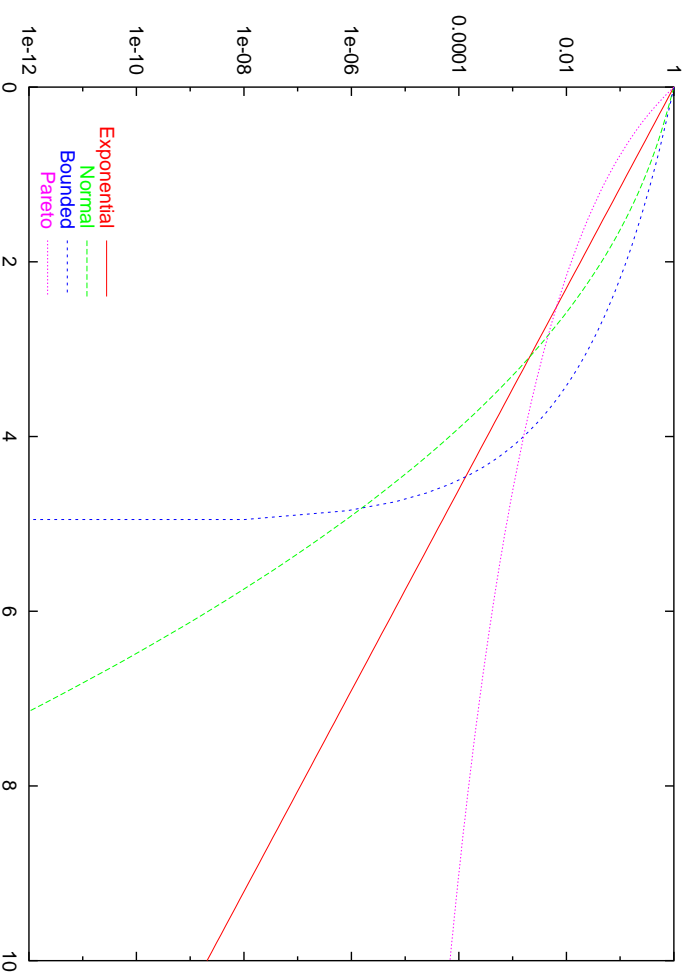
Example: the fractional Brownian motion is autosimilar.

It is also a long memory process.

# Sub-Exponentiality

The problem, graphically: consider the plot of several distributions

A queueing example: two independent On/Off sources.

- $A_1$ duration of "On" *heavy tailed*
- $A_2$ duration of "On" is arbitrary, of average throughput $\rho$
- $C$: capacity of the server.

Compare the stationary workload of two queueing systems:

$W^{1+2}$ : workload when $A_1$ and $A_2$ are superimposed.

$W^1$ : workload with $A_1$ alone but with capacity $C - \rho$.

Then:

$$\boxed{\mathbb{P}\{W^{1+2} > x\} \sim \mathbb{P}\{W^1 > x\}}$$

Conclusion: $A_1$ "dictates" the asymptotic behavior $W$.

# Application 1: Differentiated Services

- Position of the problem
- Model for throughputs
  – Analysis
  – Results
  – Validation
- Model for delays
  – Analysis
  – Results
  – Validation
- Conclusions

# The idea of Differentiated Services

**Objective** Improve the "quality of service" of the Internet by adding some kind of service definition and guarantees.

Move away from "best effort" and its lack of response time/throughput guarantees.

**Means** Use the TOS (half) octet in IP headers

Specify mechanisms at routing nodes that use this information

**Problems** What1 information? What mechanisms.

Progression of the idea:

- Clark '95 – Service discrimination

- Crowcroft '96 – All you need is just one bit

- Bolot et al. – 1-bit schemes for service Discrimination: INRIA report '97

- IETF Diffserv working group $\rightarrow$ RFC 2475 '98

- May et al. – Simple performance models: INFOCOM' 99

- Martin May's PhD Thesis, oct' 99

Two types of differentiation: using 1 bit, define 2 classes with one of them having:

- better throughput

or

- better delay characteristics

In Diffserv:

- Assured Forwarding
- Expedited Forwarding (aka: Premium Service of V. Jacobson)

This is per class Qos and NOT per flow Qos.

# Model for "Assured Forwarding"

A simple markovian model:

$\lambda_{out}$ OUT

$\lambda_{in}$ IN

Buffer size K

**Router** = single server queue

**Input traffic** = two classes, IN (high priority, tagged) and OUT, low priority.

Arrivals according to Poisson processes

**Services** Exponential distribution.

**Buffer** Finite capacity $K$

**Service discipline** = FIFO

# Buffer management = RIO

RED = Random Early Detection / Discard

RIO = RED on IN and OUT

# Analysis

Superposition of two independent Poisson processes: a dual view:

P($\lambda_1$)

P($\lambda_2$)

P($\lambda_1$)

P($\lambda_2$)

P($\lambda_1 + \lambda_2$)

P($\lambda$)

P($\lambda p$)

P($\lambda(1 - p)$)

Probability that a given packet is In = proportion of In packets:

$$p = \frac{\lambda_{\mathsf{in}}}{\lambda_{\mathsf{in}} + \lambda_{\mathsf{out}}} .$$

Probability of accepting a packet:

$$\alpha(n) = \frac{\lambda_{\text{in}}}{\lambda} \alpha_{\text{in}}(n) + \frac{\lambda_{\text{out}}}{\lambda} \alpha_{\text{out}}(n)$$

with $\lambda = \lambda_{\text{in}} + \lambda_{\text{out}}$.

Evolution of the number of customers $N(t)$:

$$n \rightarrow n+1 \quad \text{with rate } \lambda_{\text{in}} \times \alpha_{\text{in}}(n) \quad n < K$$

$$n \rightarrow n+1 \quad \text{with rate } \lambda_{\text{out}} \times \alpha_{\text{out}}(n) \quad n < K$$

$$n \rightarrow n-1 \quad \text{with rate } \mu \quad n > 0$$

From constructions 1 and 2, this is a Markov chain. Actually: a Birth and Death process

## Analysis: solution

Equilibrium equations for the stationary probabilities $\pi(n)$:

$$(\lambda \alpha(n) + \mu)\pi(n) \;=\; \lambda \alpha(n-1)\pi(n-1) \;+\; \mu \pi(n+1) \,.$$

provided that $n \geq 1$.

Global balance equations:

$$\lambda \alpha(n)\pi(n) \;=\; \mu \pi(n+1) \,.$$

Solution of the recurrence:

$$\pi(n) \;=\; \pi(0)\left(\frac{\lambda}{\mu}\right)^{n}\prod_{i=0}^{n-1}\alpha(i)$$

$$\pi(0) \;=\; \left[\sum_{n=0}^{K}\left(\frac{\lambda}{\mu}\right)^{n}\prod_{i=0}^{n-1}\alpha(i)\right]^{-1}$$

Computation of performance measures:

**throughputs**

$$\lambda_{\text{in}}^{\text{eff}} \;=\; \lambda_{\text{in}}\sum_{n=0}^{K-1}\alpha_{\text{in}}(n)\pi(n)$$

$$\lambda_{\text{out}}^{\text{eff}} \;=\; \lambda_{\text{out}}\sum_{n=0}^{K-1}\alpha_{\text{out}}(n)\pi(n)$$

## average queue length

$$N = \sum_{n=0}^{K-1} n \, \pi(n)$$

## response times of accepted packets

$$R = \sum_{n=0}^{K-1} \frac{n+1}{\mu} \, \pi(n)$$

$$= \frac{N}{\lambda_{\text{in}} + \lambda_{\text{out}}}$$

Effective throughputs, global and per class, p=0.5, K=100



**Effective throughput** vs **Offered Load**

Legend:
- HP/LP tail drop — black
- In/HP RED — green
- HP Threshold — blue
- Out/LP RED — red
- LP Threshold — cyan

# Validation against simulation

Is the result robust? Does it depend on the arrival process?

Testing the Poisson assumption using simulation of the system with several input traffic characteristics:

- Poisson processes (just checking the analytical formulas)

- A superposition of 32 On/Off process with *constant* inter-arrivals and *exponentially distributed* on and off periods

- A superposition of 32 On/Off processes with *constant* inter-arrivals and *Pareto distributed* on and off periods ($\alpha = 1.4$: Hurst parameter = 0.8).

Per class loss probabilities, analytical and simulation, p=0.9, K=100

| | |
|---|---|
| In/HP Analytic | — |
| Out/LP Analytic | — |
| HP loss simul. (Poisson) | ◇ |
| LP loss simul. (Poisson) | + |
| HP loss simul. (On/Off Expo) | □ |
| LP loss simul. (On/Off Expo) | × |
| HP loss simul. (On/Off Pareto) | ▷ |
| LP loss simul. (On/Off Pareto) | ∗ |

offered load

## Asymptotic sharing of the processor.

The existence of formulas allow to compute the asymptotic expansion when the load increases.

Define

$$\phi = \frac{\lambda_{\text{in}} \alpha_{\text{in}}(K-1)}{\lambda_{\text{in}} \alpha_{\text{in}}(K-1) + \lambda_{\text{out}} \alpha_{\text{out}}(K-1)}.$$

Then:

$$\lambda_{\text{in}}{}^{\text{eff}} = \mu \phi + \frac{1}{\rho} \frac{p\mu}{\alpha(K-1)} \left( \frac{\alpha_{\text{in}}(K-2)}{\alpha(K-2)} - \frac{\alpha_{\text{in}}(K-1)}{\alpha(K-1)} \right) + O\left(\frac{1}{\rho^2}\right).$$

Proposal: take the term $(\ldots)=0$ so that predictible sharing occurs as soon as possible.

# The impact of bursts

A model with Poisson batch arrivals: each arrival brings $B$ packets at the same time.
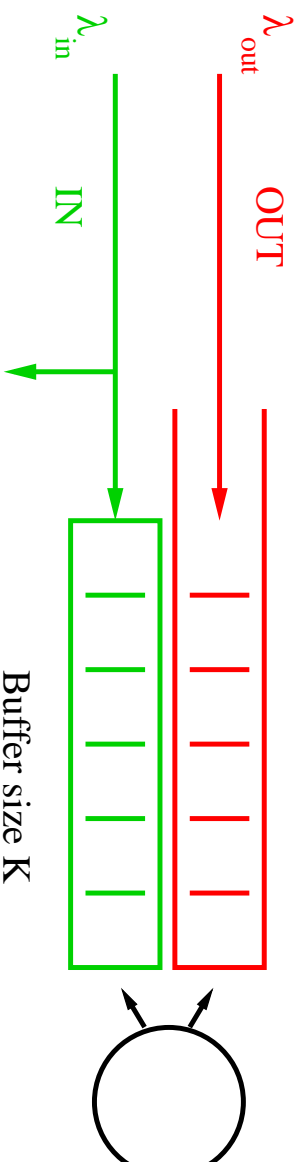
Possible transitions for $N(t)$:

$$
\begin{array}{llll}
n & \to & n + B & \text{with rate } \lambda \times \alpha(n) \quad n \leq K - B \\
n & \to & K & \text{with rate } \lambda \times \alpha(n) \quad K - B < n \leq K - 1 \\
n & \to & n - 1 & \text{with rate } \mu \quad n > 0
\end{array}
$$

This Markov chain can be solved numerically $\to \pi(n)$.

Results for $B$ large can be compared to that of the normal Poisson process ($B = 1$), and the Tail Drop mechanism ($\alpha(n) = 1$ if $n < K$).

May and Bonald conclude that RED does not discriminate between bursty and smooth arrival processes, whereas Tail Drop does.

# The model for Expedited Forwarding

$\lambda_{out}$    OUT

$\lambda_{in}$    IN

Buffer size K

Service discipline: (preemptive) priority of In over Out.

Analysis:

- For In: a $M/M/1/K$ queue $\Rightarrow$ use known results
- For Out: a $M/M/1$ queue with priorities (not quite) $\Rightarrow$ compute stochastic bounds, use results of Miller and Takacs

## Analysis of preemption

For a lower priority customer: the response time is

$$R = X + \sum_{j=1}^{A(X)} B_j^{(K)}$$

with

$X$: low priority workload $W$ found upon arrival

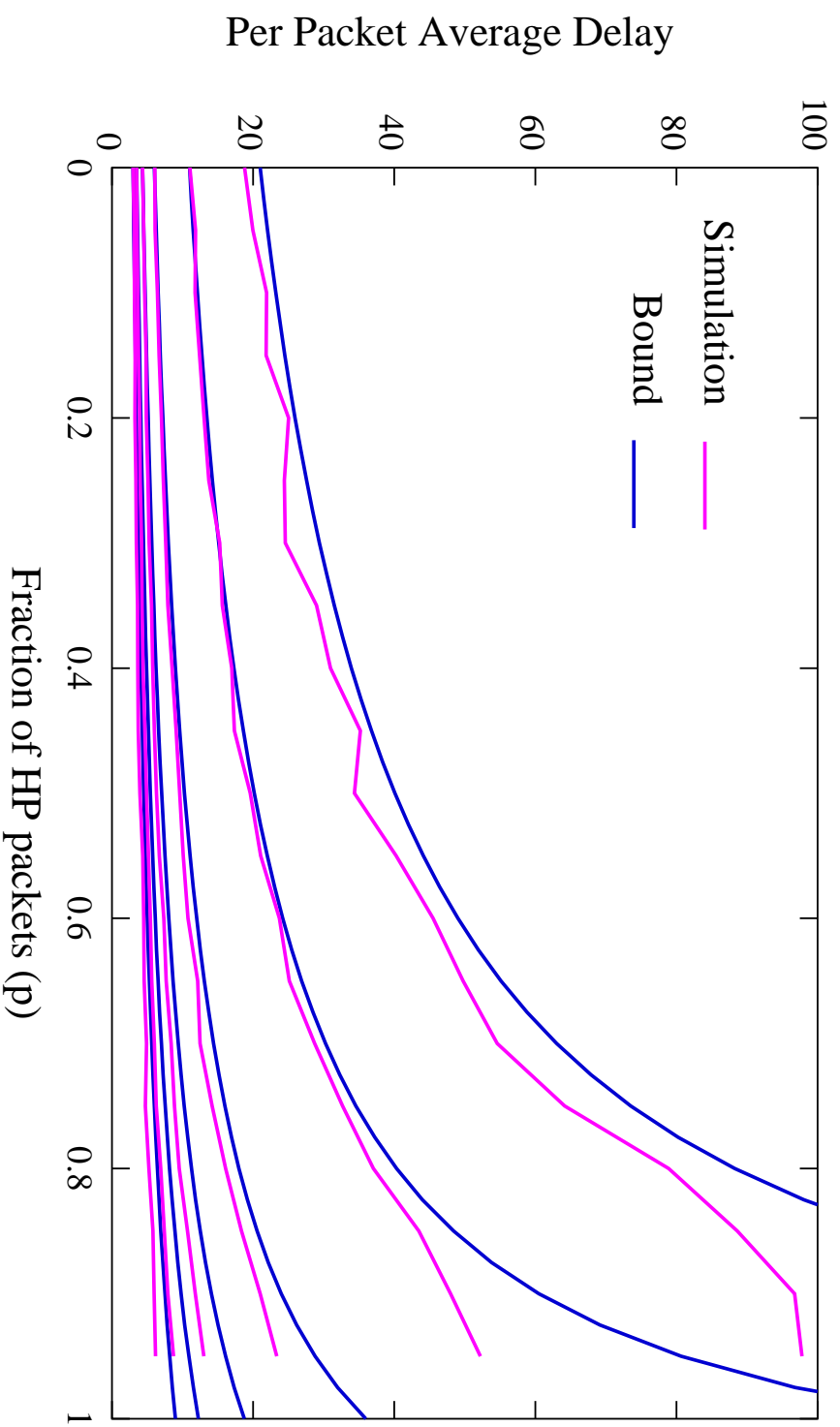$B_j^{(K)}$: length of a *busy* period of high priority: BP in a $M/M/1/K$ (known!)

Because of $K$: $\chi$ and $A(\chi)$ are difficult. However:

$$\chi \quad \leq_{\text{st}} \quad \text{Response}(M/M/1) \quad =_d \quad \text{Exp}(\mu - \lambda_{\text{in}})$$

$$A(\chi) \quad \leq_{\text{st}} \quad A(\text{Exp}(\mu - \lambda_{\text{in}})) \quad =_d \quad \text{Geom}(\rho_{\text{in}})$$

Finally, with $b_n = \mathbb{E}B^n$:

$$b_n \;=\; \left( b_{n-1} \;-\; \sum_{j=1}^{n-1} b_j \frac{\rho_1}{1 + \rho_1}^{\,n-j} \right) \frac{1 + \rho_1}{\rho_1}.$$

# Conclusion

## Conclusions

- A fairly efficient yet simple scheme
- A fairly good model
- Insight on the behavior of RED/RIO at high loads

## Research issues

- Investigate average queue length measurements:

$$\hat{q}_{n+1} = \alpha q_n + (1 - \alpha)\hat{q}_n$$

- Investigate RIO based on the queue length of tagged packets instead of total queue length

# Part VI: Deterministic Models

- Traffic envelopes and $(\sigma, \rho)$ bounds

- Bounds on the delay and on the buffer sizes

- Traffic shapers

- Service curves

- Network calculus

# Principle

*Work arrival function*:

$$S(a, b) = \sum_{a \leq a_n < b} \sigma_n \quad \text{(discrete)}$$

$$= \int_a^b r(t) \, dt \quad \text{(fluid)}$$

Envelop of the arrival of work: the worst of situations for intervals of a certain length $t$:

$$\alpha(t) = \sup_s S(s, s + t) \, .$$

Example: bounds "$(\sigma, \rho)$":

$$S(s, s + t) \leq \sigma + \rho \, t, \quad \forall s.$$

# Results

If a process with an envelop bounded by an affine "$(\sigma, \rho)$" function, feeds a queue, then:

$$W(t) \leq \sigma,$$

in addition, the response time is bounded: for all customer $n$

$$R_n \leq \frac{\sigma}{1 - \rho},$$

whatever the work conserving service policy.

$\implies$ dimensioning of buffers.

Proof: for all $a < b$

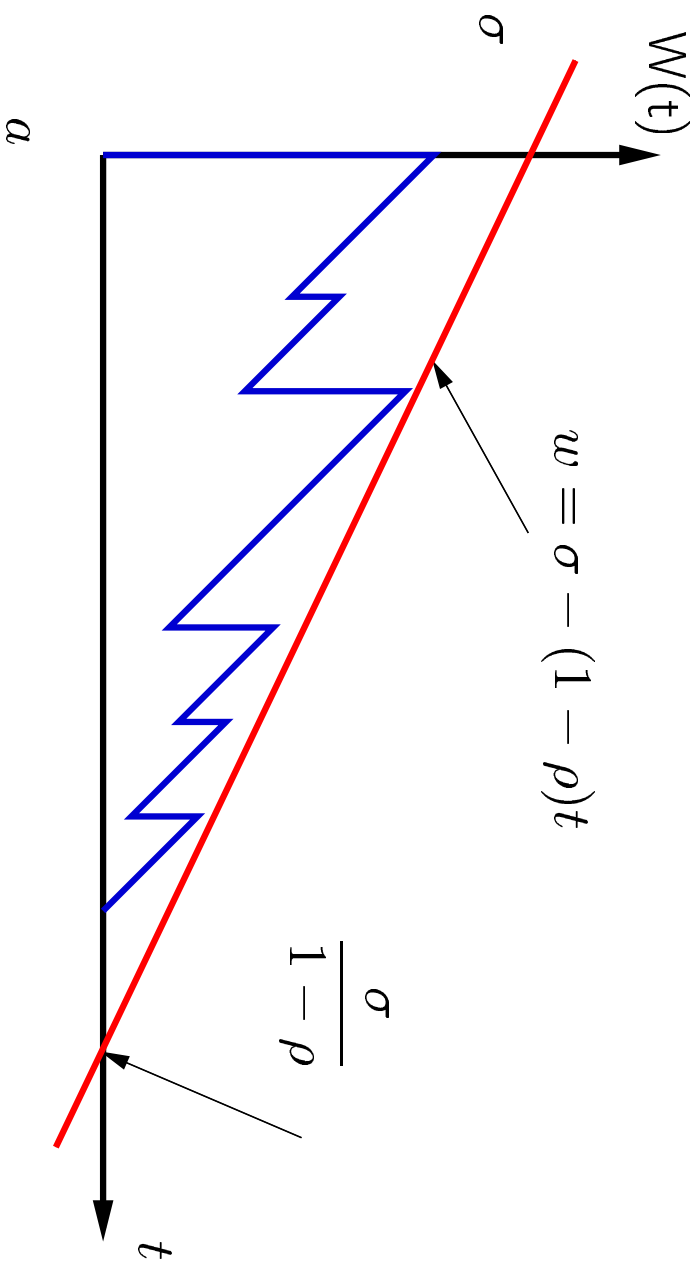$$W(b) = W(a) + S(a, b) - \int_a^b \mathbf{1}_{\{\text{served at } u\}} du.$$

If $a$ starts a Busy Period, then

- $W(a) = 0$
- customers are served on $[a, t]$ as long as $W(t) > 0$ (work conserving).

Therefore:

$$
\begin{aligned}
0 \leq W(t) &= S(a, t) - (t - a) \\
&\leq \sigma + \rho(t - a) - (t - a) \\
\Rightarrow (t - a)(1 - \rho) &\leq \sigma
\end{aligned}
$$

## Networks: addition of burstiness

Consider arrival processes with $(\sigma, \rho)$ envelops:
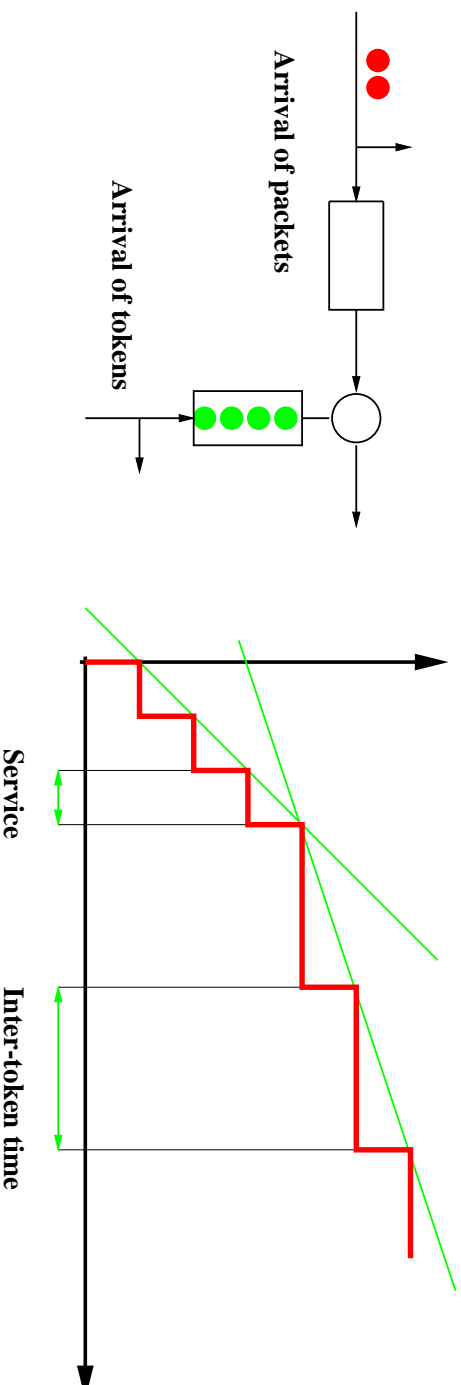
$$S_i(a, a + t) \leq \sigma_i + \rho_i\, t .$$

Then the superposition of the processes is also $(\sigma, \rho)$ with:

$$\sigma = \sum_i \sigma_i \quad \rho_i = \sum_i \rho_i .$$

# Traffic shapers

Elements of a network in charge of reducing the impact of bursts: shaping, *smoothing* of traffic.

Example: the Token Bucket (also: Leaky Bucket).

**Arrival of packets**

**Arrival of tokens**

**Service**

**Inter-token time**

Possible applications: definition of *traffic contracts* based on:

(peak rate,sustained rate,packet length,burst length)

# Service curves

Input/Output view of a (lossless) system:

A(t) → **Service System** → D(t)

With:

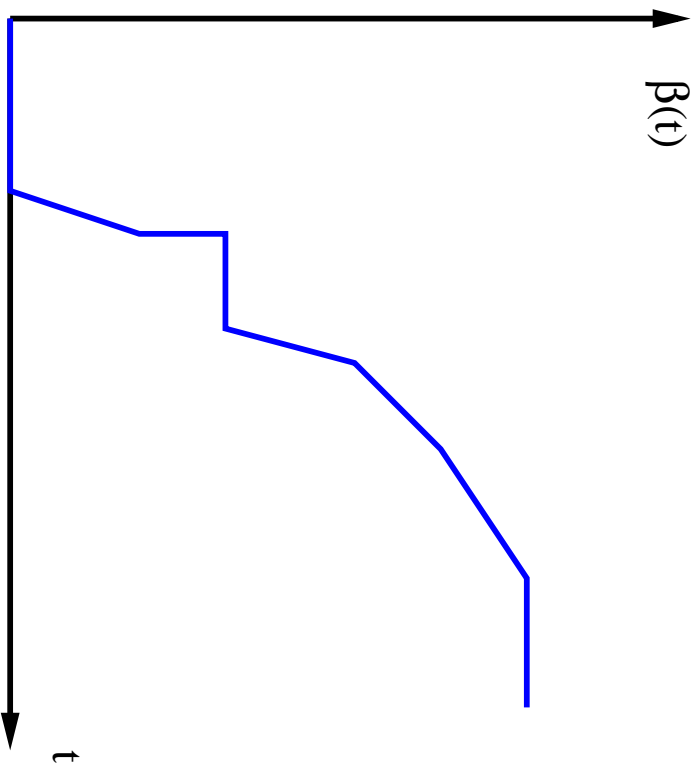- $A(t)$ quantity of information arrived at time $t$
- $D(t)$ quantity of information departed at time $t$
- $W(t) = A(t) - D(t)$ backlog of information at time $t$

The service system offers the service curve $\beta(t)$ if for some $t_0$:

$$D(t) \geq A(t - t_0) + \beta(t).$$

# Backlog and delay bounds



quantity of
information

q

t

time

Backlog at time t

Response time of element
of information "q"

β(t)

α(t)

Formulas for the bounds:

$$W(t) \leq \sup_{s \geq 0}\{\alpha(s) - \beta(s)\}$$

$$R(q) \leq \sup_{s \geq 0}\inf_{\tau \geq 0}\{\alpha(s) \leq \beta(s + \tau)\}$$

# Example: video/voice playout



quantity of
information

Original stream

Response time of the network

Maximum backlog

Smoothed stream

time

# The importance of service disciplines

Service curve of <span style="color:red">FIFO</span>?

- fluid: $\beta(t) = C \times t$

- discrete, with packet service time $\leq T$: $\beta(t) = C \times T \times \left\lfloor \dfrac{t}{T} \right\rfloor$



Ideal fluid

Actual

T

t

# Application: IETF's IntServ

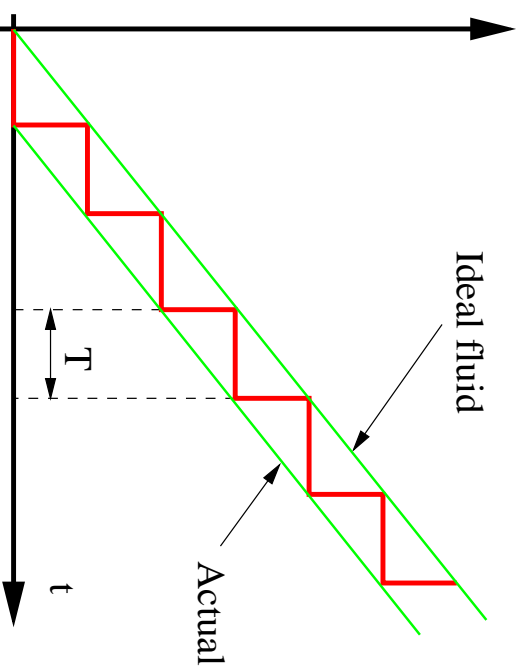Integrated Services: assume a "leaky-bucket-like" shaped input traffic,

$$\alpha(t) = \min\{M + pt, b + rt\},$$

and "RSVP" nodes with a service curve of the form

$$\beta_{C,T}(t) = C \times (t - T)^+.$$

Bounds:

$$W_{\max} = b + r \max\left(\frac{b - M}{p - r}, T\right)$$

$$R_{\max} = \frac{1}{C}\left(M + \frac{b - M}{p - r}(p - C)^+\right) + T.$$

# Network calculus

The output flow of a network element with service curve $\beta(\cdot)$ is bounded:

$$D(t+s) - D(t) \leq \alpha_{\text{out}}(s) = \sup_{u \geq 0}\{\alpha(s+u) - \beta(u)\} .$$

$\Rightarrow$ propagation of boundedness.

Nodes in series:

$$\beta(t) = \inf_{0 \leq s \leq t}\{\beta_1(s) + \beta_2(t-s)\}$$

$$\boxed{\beta_1} \longrightarrow \boxed{\beta_2} \longrightarrow$$

"Pay bursts only once": $D \leq D_1 + D_2$.

# Problems to solve

- superposition of flows $\Rrightarrow$ increase of burstiness

$\rightarrow$ Reshaping

$\rightarrow$ per-flow service disciplines, such as Generalized Processor Sharing (GPS), or Earliest Deadline First (EDF)

- loose bounds $\Rrightarrow$ over-reservation $\Rrightarrow$ waste of bandwidth and buffer

$\rightarrow$ improve accuracy

$\rightarrow$ combined stochastic and deterministic analysis

- losses

# Application 2: Traffic Management

Queueing analysis gives an insight into several issues of Traffic Management in networks (included the Internet), among which:

- capacity planning
- route planning, routing
- window-based congestion control and TCP

# Traffic Matrices

The network is formed of

- pairs of origin/destination pairs, generating traffic with rates $\lambda_{(O,D)}$

- routers with capacities $C_n$

- links between routers with capacities $C_{i,j}$

- routes representing the path followed by the information (cells, frames, datagrams...) between some O and D

Each location where queuing takes place is modeled by a queuing "node".

Offered load at node $n$:

$$\hat{\lambda}_n = \sum_{r \text{ route going through node } n} \lambda_r .$$

Kelly/Jackson's theorem: the average response time, all classes aggregated:

$$\overline{T} = \frac{1}{\Lambda} \sum_{n=1}^{N} \frac{\hat{\lambda}_n}{\mu_n - \hat{\lambda}_n} .$$

with $\Lambda = \sum_{k=1}^{K} \lambda_k$ the total offered traffic.

Taking into consideration the propagation delay $d_n$ associated with node $n$, a reasonable formula is:

$$\boxed{\overline{T} = \frac{1}{\Lambda} \sum_{n=1}^{N} \frac{\hat{\lambda}_n}{\mu_n - \hat{\lambda}_n} + d_n .}$$

# Capacity Planning

Assuming known: traffic rates and routes.

Problem: allocate link/node capacity so as to minimize collective average.

$$\min_{(\mu_1,\ldots,\mu_N)\in\mathcal{M}} \overline{T}(\mu_1,\ldots,\mu_N)$$

where $\mathcal{M}$ = set of feasible allocations (economical/technical/ethical constraints).

# Route Planning

Assuming known: node and link capacities, O/D traffics.

Problem: allocate routes so as to minimize collective average.

Decision variables: $x_{O,D,r}$ = quantity of traffic sent on route $r$ between O and D.

$\mathbf{x}$: vector of all such variables.

$$\min_{\mathbf{x} \in \mathcal{R}} \quad \overline{T}(\mathbf{x})$$

where $\mathcal{R}$ = set of feasible route allocations.

Typical constraints:

- if traffic can be shared between routes (e.g. datagrams):

$$x_{O,D,r} \in [0, \lambda_{(O,D)}] \qquad \sum_r x_{O,D,r} = \lambda_{(O,D)}$$

- if traffic cannot be shared (e.g. virtual circuits)

$$x_{O,D,r} \in \{0, \lambda_{(O,D)}\} \qquad \exists! \, r, \quad x_{O,D,r} = \lambda_{(O,D)}$$

- capacity constraints: for all queueing node $n$

$$\hat{\lambda}_n < C_n.$$

# Routing

Consider a network with distributed routing based on distance vector tables.

According to the response time formula at nodes for Kelly networks (the $M/M/1$ formula) plus the propagation delay, a reasonable metric for link $n = (i \rightarrow j)$ is:

$$D_{i,j} = \frac{\hat{\lambda}_n}{\mu_n - \hat{\lambda}_n} + d_{i,j} \, .$$

# Flow Control

If at some node $\hat{\lambda}_n > C_n$, then the buffers fill up and losses are unavoidable.

Even if the sustained rate $\hat{\lambda}_n < C_n$, temporary bursts may cause losses.

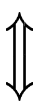Flow control uses feedback at the source to limit the offered load.

Among the possibilities: window flow control.

Application of Little's law

- $W$ the average size of the window
- $T$ the average round trip time
- $\theta$ the throughput

$$\boxed{W \;=\; \theta\, T}$$

$$\Longleftrightarrow$$

$$\boxed{\theta \;=\; \frac{W}{T}.}$$

The larger the window, the better... until the queues inside the network overflow.

# Analytical approach to flow control

Consider the closed Jackson network with $W$ customers, each node having capacity $\mu$



The throughput and RTT of customers are

$$\theta \;=\; \frac{W\,\mu}{W + 2N - 1}$$
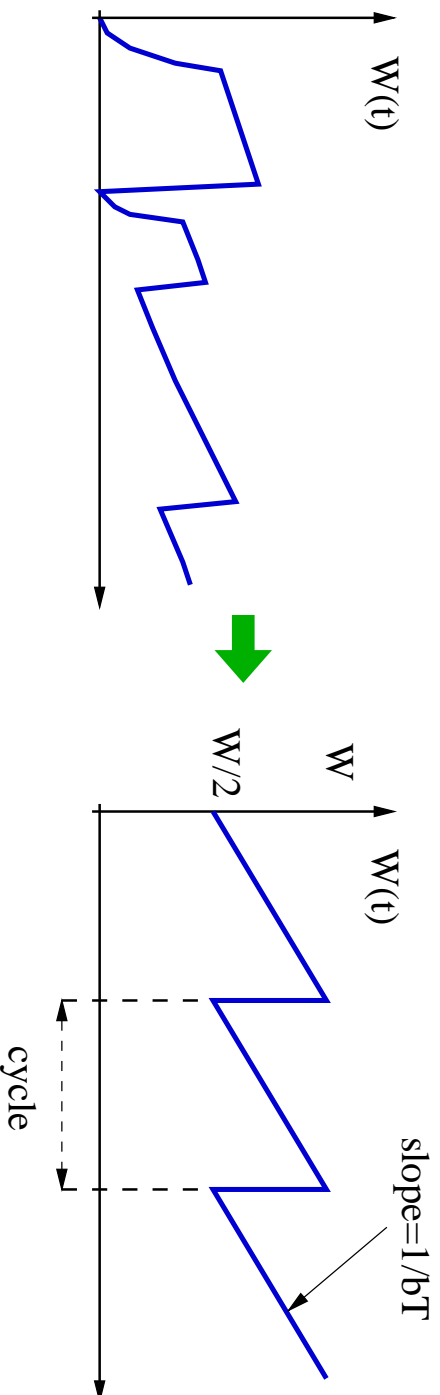
$$RTT \;=\; \frac{W + 2N - 1}{\mu}$$

# Models of TCP

## Principles of TCP Reno

- TCP sources use a window $W$ (bytes) of unacknowledged packets

- packets are acknowledged by the receiver (by groups of $b$: delayed ACKs)

- packet losses are detected

  – because multiple ACKs of a packet are received $\Rightarrow$ W $\leftarrow$ W/2

  – after a time-out has expired $\Rightarrow$ W $\leftarrow$ 1

- the window grows

  – W $\leftarrow$ W+1 each time an ACK is received in the slow start or fast recovery modes

  – W $\leftarrow$ W+1/W each time an ACK is received in the congestion avoidance

Observe: after each RTT, $W$ increases of $(W/b) \times 1/W = 1/b$.

Cycle analysis:



- length of the cycle: $T_0 = T \times W/2b$:

- number of packets transmitted: $N = \dfrac{1}{T}\displaystyle\int_0^{T_0} W(t)\,dt = \dfrac{3bW^2}{8}$:

- number of packets lost: 1, a proportion $p = \dfrac{8}{3bW^2}$

Finally, the effective throughput is:

$$\theta = \frac{N-1}{T_0} \sim \frac{1}{T}\sqrt{\frac{3}{2pb}}$$

This "square root formula" (S. Floyd, T. Ott) gives a relationship between the loss probability $p$ and the throughput $\theta$.

$\Rightarrow$ concept of "TCP friendly" services.

## A Stochastic Model of TCP

Using the same principle but with random times $S_n$ between losses (Altman *et al.*).



The process $\{W_n\}_{n \in \mathbb{N}}$ is a discrete-time Markov chain. It evolves as (with $\alpha = 1/bT^2$):

$$W_{n+1} = \frac{1}{2}W_n + \alpha S_n$$

Solving for the recurrence:

$$W_n = \alpha \sum_{k=0}^{\infty} \frac{1}{2^k} S_{n-1-k}.$$

The chain admits a stationary behavior.

Computation of the moments: introduce

- $\lambda =$ intensity ("throughput") of the loss process
- $R(k) = \mathbb{E}(S_0 S_k)$ the autocorrelation of inter-loss times

$$\mathbb{E}W = \frac{2\alpha}{\lambda}$$

$$\mathbb{E}W^2 = \frac{4\alpha^2}{3} \left( R(0) + 2 \sum_{k \geq 1} \frac{1}{2^k} R(k) \right)$$

The loss probability and the throughput are related by:

$$\theta = \frac{1}{T} \frac{1}{\sqrt{2pb}} \sqrt{R(0) + 2\sum_{k \geq 1} \frac{1}{2^k} R(k)}.$$

Finally, assuming $\{S_n\}$ form a sequence of *independent and identically distributed* random variables with average $d$, variance $\sigma_S^2$ and coefficient of variation $c^2 = \sigma_S^2/d^2 - 1$, then:

$$R(0) = \sigma_S^2 + d^2 \quad R(k) = d^2, \quad k \geq 1$$

$$\boxed{\theta = \frac{1}{T}\sqrt{\frac{3+c^2}{2pb}}}$$

# Models of FEC

Forward Error Correction consists in adding redundancy to data so that it can cope with loss.

Assume a stream of packets of the same size, grouped in blocks of size $n$.

It is possible to add $k$ packets to each block so that any $k$ losses in the super-block of $n + k$ packets can be recovered.

→ improved loss recovery for the group

→ increased load, increased loss rate for individual packets

Does the compromise bring a global benefit? What is the optimal value of $k$?

# A Model based on the M/M/1/K queue
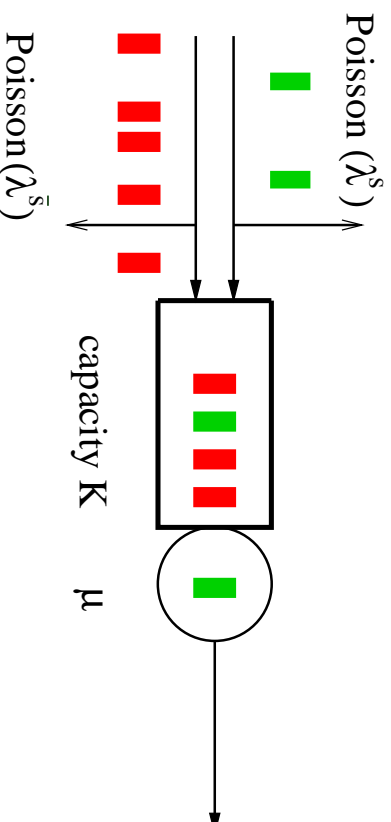
Assume packets arrive according to

- for the *tagged source*, a Poisson process with rate $\lambda^s$.

- for the other sources, a Poisson process with rate $\lambda^{\bar{s}}$.

The problem is to compute

$$P(j, n) = \mathbb{P}\{j \text{ packets are lost in a block of } n \text{ consecutive ones}\}.$$

Poisson ($\lambda^s$)

Poisson($\lambda^{\bar{s}}$)

capacity K    $\mu$

The analysis proceeds in several steps:

- Consider the distribution of the number of packets found by the <span style="color:green">first</span> of a block

$$\Pi(i) = \rho^i / (\sum_{\ell=0}^{K} \rho^\ell).$$

- Compute $Q_i(k)$, the probability that $k$ packets out of $i$ leave the system during an inter-arrival epoch

$$Q_i(k) = \rho\alpha^{k+1} \quad 0 \leq k \leq i-1$$

$$Q_i(i) = \alpha^i,$$

(1)

where $\alpha := (1+\rho)^{-1}$.

• Write down recurrent equations for $P^{s,a}(j, n)$

$$= \mathbb{P}\{j \text{ packets of } s \text{ are lost in a block of } n \text{ consecutive ones, given that the first finds } i\}$$

$$P_i^{s,a}(j, 1) = \begin{cases} 1 & j = 0 \\ 0 & j \geq 1, \end{cases} \qquad i = 0, 1, ..., K - 1$$

$$P_K^{s,a}(j, 1) = \begin{cases} 1 & j = 1 \\ 0 & j = 0, j \geq 2. \end{cases}$$

For $n \geq 2$, we have for $0 \leq i \leq K - 1$ and for $i = K$, respectively:

$$P_i^{s,a}(j, n) = \sum_{k=0}^{i+1} Q_{i+1}(k) \left[ p_s P_{i+1-k}^{s,a}(j, n - 1) + p_{\bar{s}} P_{i+1-k}^{s,\bar{a}}(j, n - 1) \right]$$

$$P_K^{s,a}(j,n) = \sum_{k=0}^{K} Q_K(k) \left[ p_s P_{K-k}^{s,a}(j-1,n-1) + p_{\bar{s}} P_{K-k}^{\bar{s},a}(j-1,n-1) \right],$$

where $P_i^{\bar{s},a}(j,n)$ for $n \geq 1$ is given by:

$$P_i^{\bar{s},a}(j,n) = \sum_{k=0}^{i+1} Q_{i+1}(k) \left[ p_s P_{i+1-k}^{s,a}(j,n) + p_{\bar{s}} P_{i+1-k}^{\bar{s},a}(j,n) \right], \quad 0 \leq i \leq K-1$$

$$P_K^{\bar{s},a}(j,n) = P_{K-1}^{\bar{s},a}(j,n)$$

This is a set of linear equations which can be solved numerically.

Another approach is to compute the generating functions.

Define

$$q_s(y,z) \triangleq \sum_{j=0}^{\infty} \sum_{n=1}^{\infty} y^j z^{n-1} P^s(j,n).$$

Then:

$$q_s(y,z) = \frac{R_K}{(1-z)} \left( \frac{\rho^{K-1}(1-z)[\delta_{K+1}]^2}{z\phi_K} \left[ \frac{1}{z\delta_K - \delta_{K+1} - \rho z\phi_K y} \right] + R_K^{-1} + \frac{\rho^{K-1}\delta_{K+1}}{z\phi_K} \right).$$

with

- $x_1(z)$ and $x_2(z)$ be the solutions in $x$ of $x^2 - (1+\rho)x + \rho(p_{\bar{s}} + p_s z) = 0$
- $\delta_k = x_1^k - x_2^k$, $\phi_k = (p_{\bar{s}} + p_s z)\delta_{k-1} - \delta_k$,
- $R_K = (\sum_{l=0}^{K} \rho^l)^{-1}$.

Complicated at it may seem, this expression allows:

- a faster computation of the quantities
- asymptotic expansions such as: for $\rho$ fixed, and $1 \leq n < K$, $\tilde{P}^s_\rho(> 0, n) =$

$$
\overbrace{
\begin{cases}
\dfrac{R_K \rho^K}{1 - \rho_{\bar{s}}} \left[ (1-\rho)n + \left( \dfrac{\rho_s}{1-\rho} - \dfrac{\rho_s \rho_{\bar{s}}}{1-\rho_{\bar{s}}} \right) + \theta^n O\left( n^{-3/2} \right) \right] & \text{if } \rho < 1 \\[3ex]
\dfrac{1}{K+1} \dfrac{1}{p_s} \left[ \left( 2\sqrt{p_s} + \dfrac{p_{\bar{s}}}{\sqrt{p_s n}} \right) \dfrac{\sqrt{n}}{\sqrt{\pi}} \left( 1 + O\left( \dfrac{1}{n} \right) \right) - p_{\bar{s}} \right] & \text{if } \rho = 1 \\[3ex]
1 - \left( \dfrac{4\rho_s^2}{\rho(\rho-1)^3} - \dfrac{\rho_s \rho_{\bar{s}}}{\rho(1-\rho_{\bar{s}})} \left( \dfrac{1}{\rho-1} - \dfrac{\rho-1}{(1+\rho_s - \rho_{\bar{s}})^2} \right) \right) \\
\quad \theta^{n-1} \dfrac{\beta\, n^{-3/2}}{(1-\rho_{\bar{s}})\sqrt{\pi}} (1 + o(1)) & \text{if } \rho > 1
\end{cases}
}
$$

# Brief Bibliography and sources

Stochastic Processes, Discrete Event Systems:

- Ross, *Stochastic Processes*

- Cinlar, *Introduction to Stochastic Processes*, Prentice Hall, 1975.

- Baccelli, Bremaud, *Elements of Queueing Theory, Applications of Mathematics*, vol. 26, Springer-Verlag, 1994.

- Davis, *Markov Models and Optimisation*, Prentice Hall, 1993.

- F. Baccelli, G. Cohen, G.J. Olsder, and J.-P. Quadrat. *Synchronization and Linearity*. Wiley, 1992.

Queueing theory and its application to networks:

- Kleinrock, *Queueing Networks* (2 volumes), Wiley, 1975.

- D. Bertsekas et R. Gallager, *Data Networks*, Prentice Hall, 1987.

- J. Walrand, *Introduction to Queuing Networks*, Prentice-Hall, 1989.

- E. Gelenbe, G. Pujolle, Wiley, new edition 1999.

Numerical aspects, QBD, MMPP:

- Neuts, *Matrix-Geometric Solutions in Stochastic Models - - An Algorithmic Approach*. Johns Hopkins University Press, Baltimore, 1981.

- Tijms, *Stochastic Models, an algorithmic approach*, Wiley, 1994.

- Mitra *et. al*, numerous papers.

- A. Jean-Marie, Z. Liu, Ph. Nain, D. Towsley, "Computational aspects of the Workload Distribution in the MMPP/GI/1 queue", *JSAC 99*.

Long memory (very brief)

- Ph. Nain: "Impact of Bursty Traffic on Queues", to appear in *Statistical Inference in Stochastic Processes*. `http://www-sop.inria.fr/mistral/personnel/Philippe.Nain/PAPERS/LRD/impact_bursty.pdf`.

- Bolot and Grossglauser: "On the relevance of long-range dependence in network traffic", INRIA research report RR-2830, 1996.

## Deterministic Models/Network Calculus

- R. L. Cruz, "A calculus for network delay, Parts I and II", *IEEE Trans. on Information Theory*, vol. 37, no 1, Jan. 1999, pp. 114-141.

- C.S. Chang, numerous papers.

- J.-Y. Le Boudec, several texts. E.g. "Application of Network Calculus to Guaranteed Service Networks", *IEEE Trans. on Information Theory*, vol. 44, no 3, May 1998, pp. 1087–1095.

- R. Agrawal, F. Baccelli and R. Rajan, "An Algebra for Queueing Networks with Time Varying Service", INRIA research report RR-3435, May 1998.

## Differential Services

- M. May, J.-C. Bolot, C. Diot and A. Jean-Marie, "1-Bit Schemes for Service Discrimination in the Internet: Analysis and Evaluation", INRIA research report RR-3238, aug. 1997.

- M. May, J.-C. Bolot, C. Diot and A. Jean-Marie, "Simple performance models of Differentiated Service Schemes in the Internet", INFOCOM'99, New-York, march 1999.

- M. May, T. Bonald and J.-C. Bolot, "Analytic Evaluation of RED Performance", INRIA research report, may 1999.

- M. May, PhD Thesis, oct. 1999. `http://www-sop.inria.fr/rodeo`

## TCP

- J. Padhye, V. Firoiu, D. Towsley and J. Kurose, "Modeling TCP throughput: A simple model and its empirical validation", ACM SIGCOMM, Sep. 1998.

- E. Altman, K. Avrachenkov, C. Barakat, "A Stochastic Model of TCP/IP with Stationary Random Losses", INFOCOM'2000.