

On overloaded queues

Alain Jean-Marie

INRIA et LIRMM, University of Montpellier 2
161 Rue Ada, 34392 Montpellier Cedex 5, France
ajm@lirmm.fr

25 Ans du Gérard
13 May 2005

Based on results obtained with Philippe Robert
A longer version was presented at the Lunteren 2005 Conference

Plan of the talk

Introduction	2
General properties of overloaded queues	7
The FIFO Case	11
The Overloaded Processor Sharing Queue	12
Final word	25

The Single-Server Queue

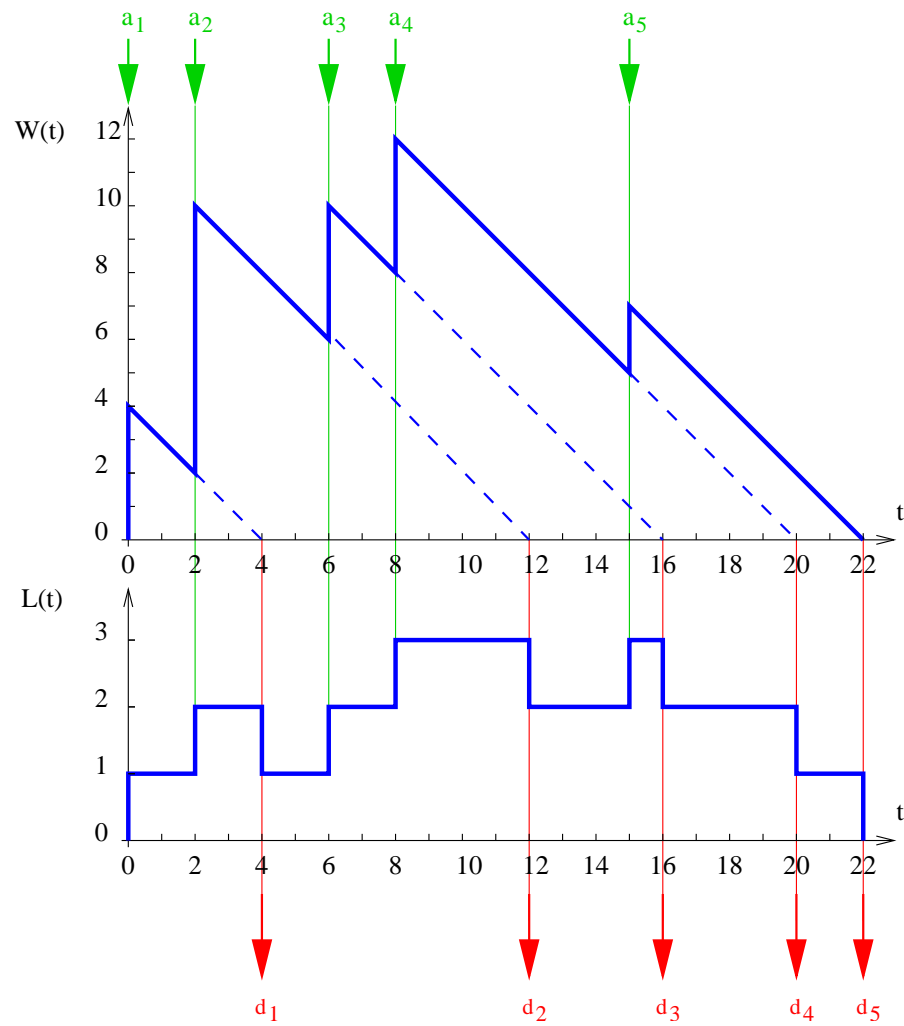
The basic $G/G/1/\infty/D$ queue

- Arrival of customers, according to a stationary process with inter-arrival times $\tau_1, \tau_2, \dots, \tau_n, \dots$
- Service requirements $\sigma_1, \sigma_2, \dots, \sigma_n, \dots$
- One single server
- Infinite waiting room for customers
- Service discipline (scheduling) D , assumed to be work conserving.

State description of a queue

The state of the queue can be described by several evolving quantities:

- $A(t)$: number of arrivals up to time t
- $D(t)$: number of departures up to time t
- $L(t)$: number of customers at time t (“length” of the queue)
- $W(t)$: workload at time t , number of units of work left to do for the server
- $R(t) = (r_1(t), r_2(t), \dots, r_{L(t)}(t))$, vector of the residual service times of customers that are in the queue.



Stability

The case usually considered in Queueing Theory is when

$$L(t) \rightarrow L(\infty), \quad W(t) \rightarrow W(\infty)$$

in distribution as $t \rightarrow \infty$. Such a queue is called **stable**.

When stability occurs:

- the **response times** T_n of customers also have a stationary distribution,
- the queue empties (\iff the server becomes idle) infinitely often.

Stability condition

When does this happen? If **inter-arrival** times τ_n and **service times** σ_n are stationary sequences, there is stability if and only if

$$\mathbb{E}(\sigma_0) < \mathbb{E}(\tau_0) .$$

Equivalently,

$$\lambda \quad := \quad \frac{1}{\mathbb{E}(\tau_0)} \quad = \quad \lim_{t \rightarrow \infty} \frac{A(t)}{t} \quad < \quad \frac{1}{\mathbb{E}(\sigma_0)} \quad =: \quad \mu$$

input rate service capacity

And consequently:

$$\text{output rate} \quad := \quad \lim_{t \rightarrow \infty} \frac{D(t)}{t} \quad = \quad \lambda \quad \text{input rate}$$

Plan of the talk

Introduction

General properties of overloaded queues

The FIFO Case

The Overloaded Processor Sharing Queue

Final word

Unstable queues

What happens when $\lambda > \mu$?

The queue is **overloaded**: too much work arrives. Also called unstable, transient, . . .

The number of customers waiting grows, the queue “explodes”.

The waiting time of customers tends to grow with time.

...

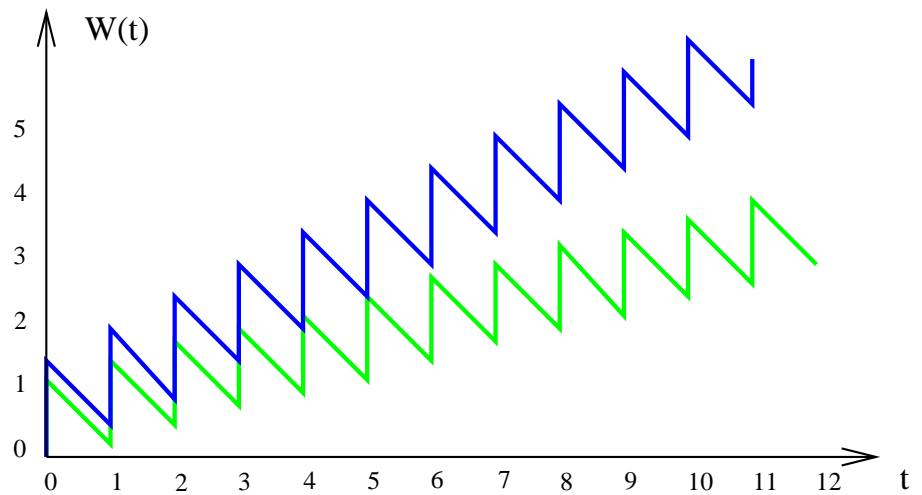
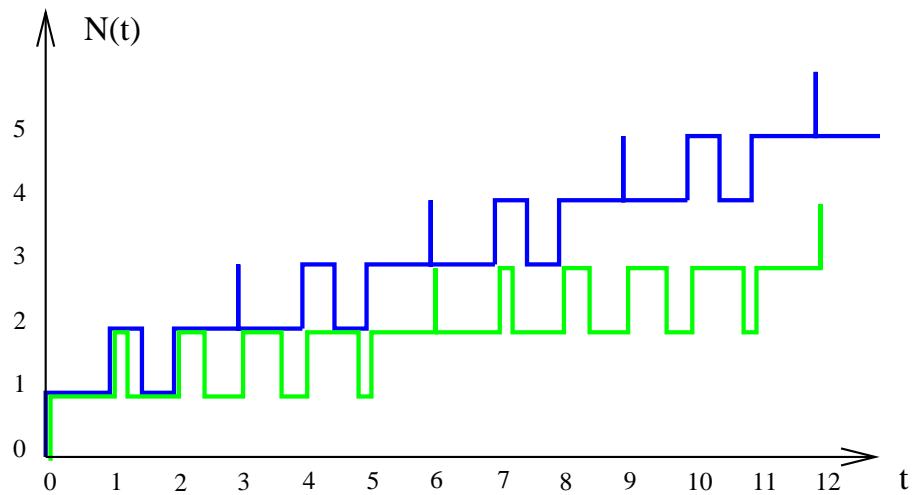
Does it really?

How fast does it grow? How bad is it?

The answers turn out to depend (only) on the service discipline

An example

The $D/D/1/\infty/FIFO$ queue.



General properties

Some general properties can be stated for an overloaded queue:

Properties

- the workload $W(t)$ goes to infinity almost surely,
- its growth rate is

$$\lim_{t \rightarrow \infty} \frac{W(t)}{t} = \frac{\lambda - \mu}{\mu} = \frac{\lambda}{\mu} - 1 ,$$

- there exists almost surely a time t_0 such that the server is always busy after t_0 :

$$W(t) > 0 \quad \forall t > t_0 .$$

Unstable queues (cdt)

The situation for the number of customers $L(t)$ is not so clear:

- does $L(t) \rightarrow \infty$?
- if it does, is there a **growth rate**

$$\alpha := \lim_{t \rightarrow \infty} \frac{L(t)}{t} ?$$

- what is the output rate $\theta = \lim_t D(t)/t$? According to the conservation law of customers:

$$\lambda = \alpha + \theta .$$

Plan of the talk

Introduction

General properties of overloaded queues

The FIFO Case

The Overloaded Processor Sharing Queue

Final word

The FIFO case

For a FIFO queue, we have:

Properties

- the growth rate of $L(t)$ is $\alpha = \lambda - \mu$
- the output rate of customers is $\theta = \mu$
- the response time of customers grows linearly with time:

$$\lim_{n \rightarrow \infty} \frac{T_n}{n} = \frac{1}{\lambda} - \frac{1}{\mu} .$$

Plan of the talk

Introduction

General properties of overloaded queues

The FIFO Case

The Overloaded Processor Sharing Queue

Final word

The overloaded Processor Sharing queue

Under the **Processor Sharing** discipline, the server serves each of the $L(t)$ customers at rate $1/L(t)$.

Recall: vector of residual service times

$$R(t) = (r_1(t), r_2(t), \dots, r_{L(t)}(t)) .$$

As long as no arrival occurs and all $r_i(t)$ remain positive:

$$\frac{dr_i(t)}{dt} = - \frac{1}{L(t)} .$$

Growth rates

Properties

- the growth rate of $L(t)$ is α , unique positive solution of:

$$x = \lambda (1 - \mathbb{E}(e^{-x\sigma_0})) ,$$

- the response time of the n -th customer, grows linearly with n : given its service time,

$$\frac{T_n}{n} \xrightarrow{n \rightarrow +\infty} \frac{(e^{\alpha\sigma_0} - 1)}{\lambda} ,$$

- the output rate θ is solution of:

$$y = \lambda \mathbb{E}E(e^{-(\lambda-y)\sigma_0}) .$$

A “proof”

Idea of the proof: consider a customer with service time σ_n arriving at time a_n . Its response time T_n is such that:

$$\begin{aligned}\sigma_0 &= \int_{a_n}^{a_n+T_n} \frac{1}{L(u)} du \\ &\simeq \int_{a_n}^{a_n+T_n} \frac{1}{\alpha u} du \\ &= \frac{1}{\alpha} \log \left(\frac{a_n + T_n}{a_n} \right) \\ \implies T_n &\simeq a_n (e^{\alpha \sigma_0} - 1) .\end{aligned}$$

Consider now the number of customers **still present** at time t .

Customer n with $a_n \leq t$ and service time σ , is still there if

$$\begin{aligned} a_n + T_n \geq t &\stackrel{\sim}{\iff} a_n e^{\alpha\sigma} \geq t \\ &\iff a_n \geq t e^{-\alpha\sigma} . \end{aligned}$$

Therefore, since $a_n \cong \lambda n$, there are approximately

$$\lambda t - \lambda t e^{-\alpha\sigma} = \lambda t (1 - e^{-\alpha\sigma}) \text{ of these.}$$

De-conditioning on σ , we get:

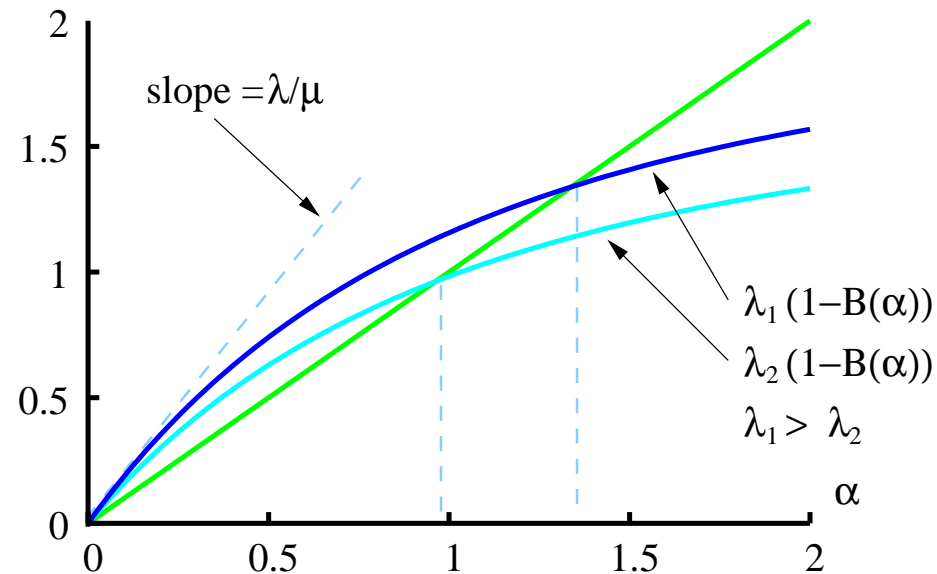
$$\begin{aligned} L(t) &\cong \lambda t \mathbb{E}(1 - e^{-\alpha\sigma}) \\ \implies \alpha &= \lambda (1 - \mathbb{E}(e^{-\alpha\sigma})) . \end{aligned}$$

Input/Output rate relations

How does the output rate θ vary with λ ?

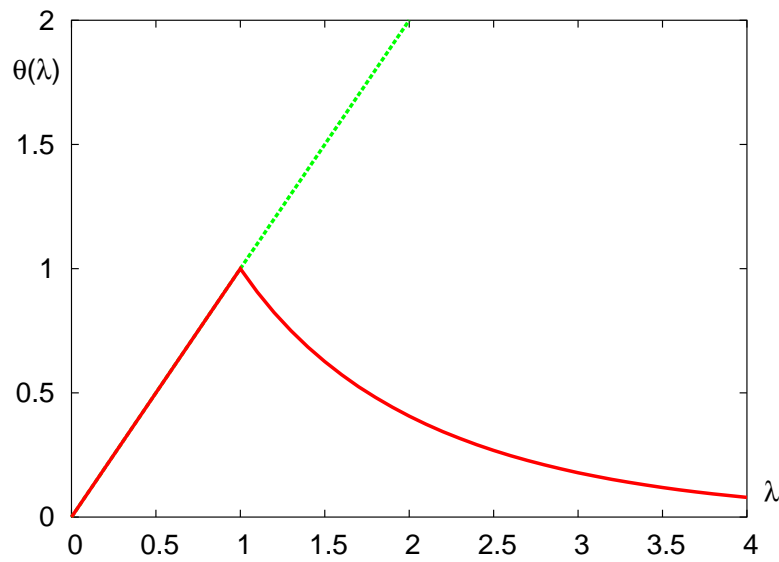
The growth rate α of the queue is increasing with respect to λ :

$$\alpha = \lambda(1 - B(\alpha))$$

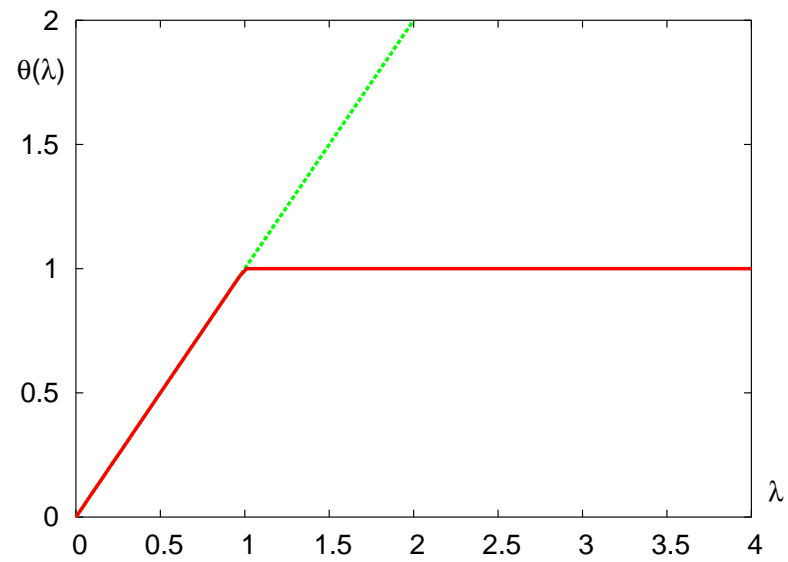


Input/Output rate relations (ctd)

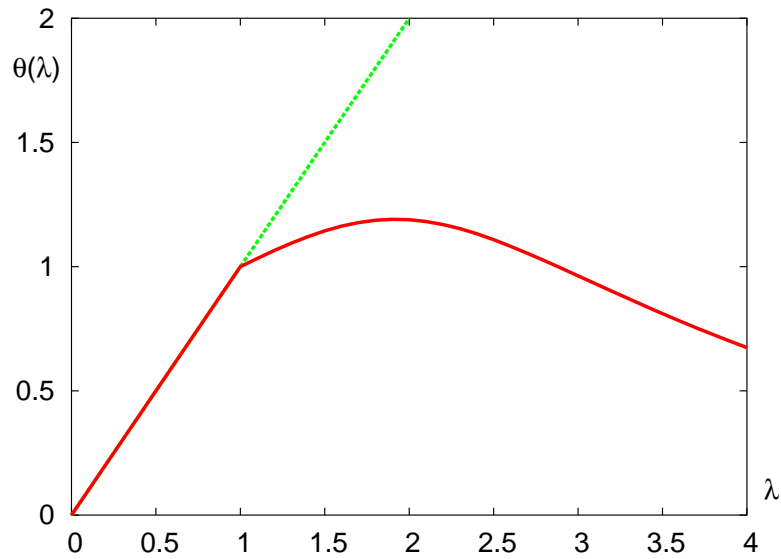
It is also true that $\alpha(\lambda)/\lambda$ is increasing, and $\theta(\lambda)/\lambda$ decreasing.
However: the output rate $\theta(\lambda)$ is **not** monotone, nor convex.



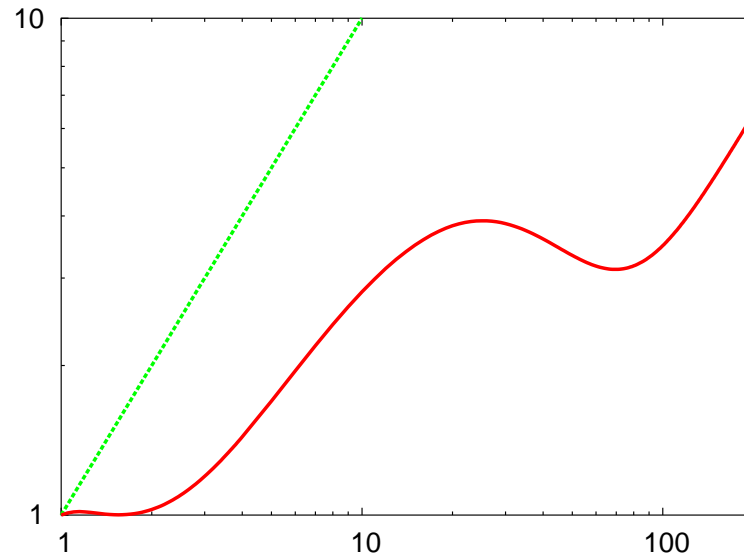
$. / D / 1 / \infty / PS$



$. / M / 1 / \infty / PS$



$$\sigma = \begin{cases} 1/2 & \text{wp } 8/9 \\ 5 & \text{wp } 1/9 \end{cases}$$



$$\sigma = \begin{cases} 0 & \text{wp } 4/125 \\ 18 & \text{wp } 1/250 \\ 3/2 & \text{wp } 4399/7250 \\ 1/20 & \text{wp } 259/725 \end{cases}$$

On Elephants and Mice

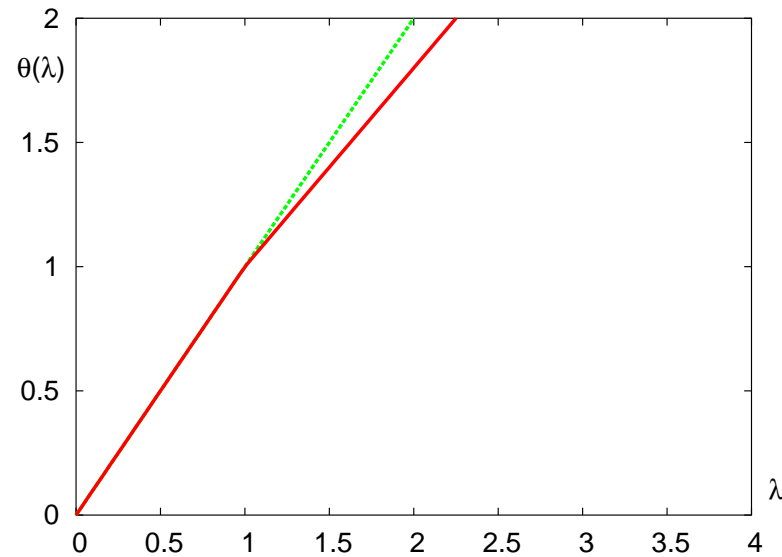
Consider the service time distribution:

$$\sigma = \begin{cases} 0 & \text{wp } 1 - \varepsilon & \text{lots of mice} \\ \text{Exp}(\varepsilon\mu) & \text{wp } \varepsilon & \text{few elephants} \end{cases}$$

It has mean $\frac{1}{\mu}$ and variance $\left(\frac{2}{\varepsilon} - 1\right) \frac{1}{\mu^2}$.

In this case:

$$\theta(\lambda) = \lambda - \varepsilon(\lambda - \mu) \quad \alpha(\lambda) = \varepsilon(\lambda - \mu) .$$



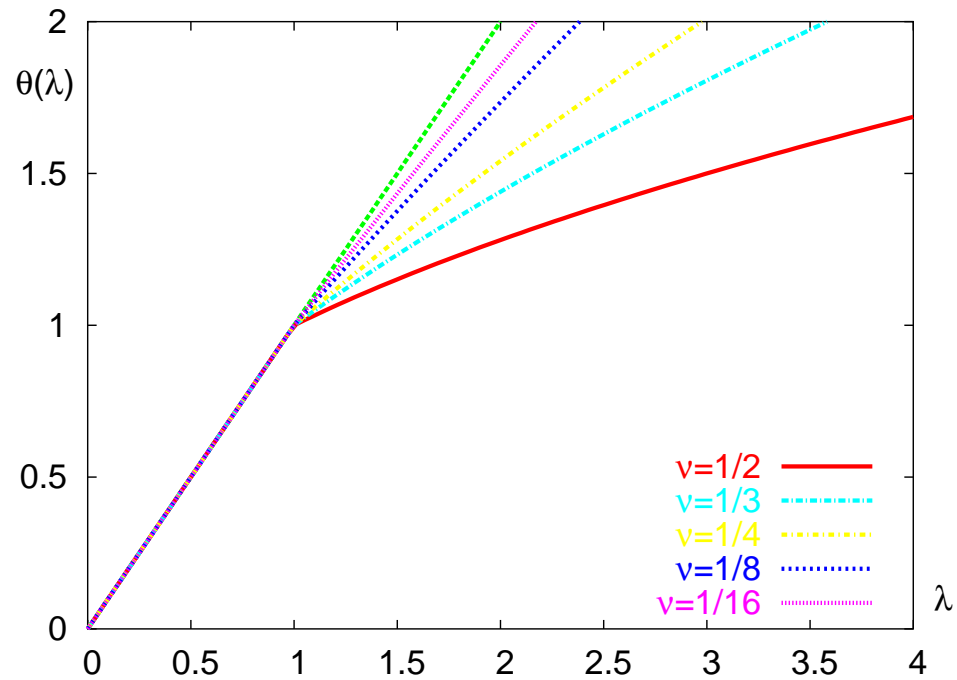
When $\varepsilon \rightarrow 0$, the output rate remains close to λ . The customers accumulate, but at a very small rate! The mice do not experience a lower throughput.

\implies application to the TCP protocol, see [Bonald and Roberts \(2003\)](#).

Asymptotic Behavior

The assumption that $\sigma = 0$ is not essential. What is relevant is the [density close to 0](#) of the service time distribution.

$\frac{dP(\sigma \leq x)/dx}{x \rightarrow 0}$	$\frac{B^*(s)}{s \rightarrow \infty}$	$\frac{\theta(\lambda)}{\lambda \rightarrow \infty}$
$o(x)$	$o(s^{-1})$	0
Ax	$\frac{A}{s}$	A
$> O(x)$	$< o(s^{-1})$	$+\infty$



$$\sigma \sim \text{Gamma}(\mu; \nu) \quad \mathbb{E}(e^{-\sigma s}) = \left(\frac{\nu\mu}{\nu\mu + s} \right)^\nu .$$

A stronger result on Residual Service Times

Residual service times also converge:

Property For any continuous function $f : \mathbb{R}_+ \rightarrow \mathbb{R}$, almost surely

$$\begin{aligned} \lim_{t \rightarrow +\infty} \frac{1}{t} \sum_{i=1}^{L(t)} f(r_i(t)) &= \lambda \mathbb{E} \left(\int_0^{\sigma_0} f(x) \alpha e^{-\alpha(\sigma_0 - x)} dx \right) \\ &= \lambda \int_0^{\infty} \int_0^u f(x) \alpha e^{-\alpha(u-x)} dx d\mathbb{P}\{\sigma_0 \leq u\} , \end{aligned}$$

provided that $\mathbb{E} \left(\sup_{x \leq \sigma_0} |f(x)| \right) < +\infty$. Moreover the result is valid for all the indicator functions of intervals.

Other Oddities: Response Times

We have seen that the distribution of T_n behaves as:

$$\frac{T_n}{n} \xrightarrow{n \rightarrow +\infty} \frac{(e^{\alpha\sigma_0} - 1)}{\lambda},$$

In expectation,

$$\mathbb{E}(T_n) \simeq \frac{n}{\lambda} (\mathbb{E}(e^{\alpha\sigma_n}) - 1).$$

For instance, for service times $\sigma_n \sim \text{Exp}(\mu)$: $\alpha = \lambda - \mu$ and

$$\mathbb{E}(T_n) \simeq \frac{n}{\lambda} \frac{\lambda - \mu}{2\mu - \lambda}.$$

This is infinite if $\lambda \geq 2\mu$!! A consequence of results by [Coffman, Muntz and Trotter \(1970\)](#).

Plan of the talk

Introduction

General properties of overloaded queues

The FIFO Case

The Overloaded Processor Sharing Queue

Final word

Open questions

Various more or less open issues:

- what about networks of queues?
- what about weighted processor sharing and variants (head-of-the line PS, Fair Queuing, . . .)?
- what about threshold-based disciplines, Foreground-Background, Earliest-Deadline-First, . . . ?
- . . .

Bibliography

JEAN-MARIE, A. AND ROBERT, PH. On the transient behavior of the processor sharing queue. *QUESTA*, 17 (1994), 129–136.

COFFMAN, E. G. JR, MUNTZ, R. AND TROTTER, H. Waiting time distributions for processor-sharing systems. *JACM*, 17 (1970), 123–130.

YASHKOV, S. On a heavy traffic limit theorem for the M/G/1 processor sharing queue. *Commun. Statist. - Stochastic Models* 9, 3 (1993), 467–471.

BONALD, T. AND ROBERTS, J.W. Congestion at flow level and the impact of user behaviour. *Computer Networks*, 42 (2003), 521–536.

BANSAL, N. AND HARCHOL-BALTER, M. Scheduling Solutions for Coping with Transient Overload. Technical Report #CMU-CS-01-134, School of Computer Science, Carnegie Mellon University, May 2001.