



# Stochastic Network Modeling (SNM)

## Queuing Theory

Introduction

Kendal  
Notation

Little Theorem

PASTA  
Theorem

The M/M/1  
Queue

M/G/1 Queue

M/G/1/K  
Queue

Queues in  
Tandem

Networks of  
Queues

M/G/1 Busy  
Period

M/G/1 Delays

## Stochastic Network Modeling (SNM)

Llorenç Cerdà-Alabern

Universitat Politècnica de Catalunya

Departament d'Arquitectura de Computadors

llorenc@ac.upc.edu

### Parts

- I Introduction
- II Discrete Time Markov Chains (DTMC)
- III Continuous Time Markov Chains (CTMC)
- IV **Queuing Theory**



# Stochastic Network Modeling (SNM)

## Queuing Theory

Introduction

Kendal  
Notation

Little Theorem

PASTA  
Theorem

The M/M/1  
Queue

M/G/1 Queue

M/G/1/K  
Queue

Queues in  
Tandem

Networks of  
Queues

M/G/1 Busy  
Period

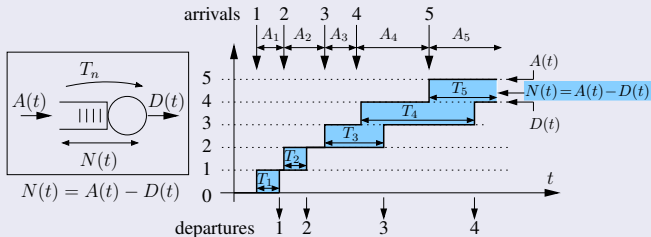
M/G/1 Delays

## Part IV

# Queuing Theory

### Outline

- Introduction
- Kendal Notation
- Little Theorem
- PASTA Theorem
- The M/M/1 Queue
- M/G/1 Queue
- M/G/1/K Queue
- Queues in Tandem
- Networks of Queues
- M/G/1 Busy Period
- M/G/1 Delays
- Matrix Geometric Method



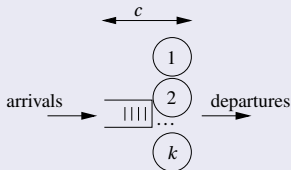
- Queueing theory is the mathematical study of waiting lines, or queues.
- Common notation:
  - $A(t)$ : number of arrivals  $[0, t]$ .
  - $A_n$ : interarrival time between customers  $n$  and  $n + 1$ .
  - $T_n$ : time in the system (response time) for customer  $n$ .
  - $N(t)$ : number in the system at time  $t$ .

## Kendal Notation

$$A/S/k[/c/p]$$

- **A**: arrival process,
- **S**: service process,
- **k**: number of servers,
- **c**: maximum number in the system (number of servers + queue size). Note: some authors use the queue size.
- **p**: population.

If “c” or “p” are missing, they are assumed to be **infinite**.





# Kendal Notation

## Queuing Theory

Introduction

Kendal Notation

Little Theorem

PASTA Theorem

The M/M/1 Queue

M/G/1 Queue

M/G/1/K Queue

Queues in Tandem

Networks of Queues

M/G/1 Busy Period

M/G/1 Delays

## Common arrivals/service processes

- **G**: general (non specific process is assumed),
- **M**: Markovian (exponentially or geometrically distributed),
- **D**: deterministic,
- **P**: Poisson (discrete RV,  $N$ , equal to the number of arrivals exponentially dist. in a time  $t$ ):

$$P_p(N = n, t) = \frac{(\lambda t)^n e^{-\lambda t}}{n!}, n \geq 0, t \geq 0.$$

- **Er**: Erlang (continuous RV equal to the time  $t$  that last  $n$  arrivals exponentially dist.):

$$f_e(t) = \lambda P_p(N = n - 1, t) = \frac{\lambda^n t^{n-1} e^{-\lambda t}}{(n-1)!}, t \geq 0, n \geq 1$$

## Examples

- **M/M/1**: M. arr. / M. serv. / 1 server,  $\infty$  queue and population.
- **M/G/1**: M. arr. / Gen. serv. / 1 server,  $\infty$  queue and population.



## Part IV

# Queuing Theory

### Outline

- Introduction
- Kendal Notation
- **Little Theorem**
- PASTA Theorem
- The M/M/1 Queue
- M/G/1 Queue
- M/G/1/K Queue
- Queues in Tandem
- Networks of Queues
- M/G/1 Busy Period
- M/G/1 Delays
- Matrix Geometric Method



# Little Theorem

## Queuing Theory

### Introduction

### Kendal Notation

### Little Theorem

#### Graphical proof

#### Application to the waiting line and the server

#### Mean number in the Server

### PASTA Theorem

### The M/M/1 Queue

### M/G/1 Queue

### M/G/1/K Queue

### Queues in Tandem

### Networks of

## Little Theorem

- Define the stochastic processes:
  - $A(t)$ : number of arrivals  $[0, t]$ .
  - $T_n$ : time in the system (response time) for customer  $n$ .
  - $N(t)$ : number in the system at time  $t$ .

- And the mean values:

- Mean number of customers in the system:

$$N = \lim_{t \rightarrow \infty} \frac{1}{t} \int_0^t N(s) ds$$

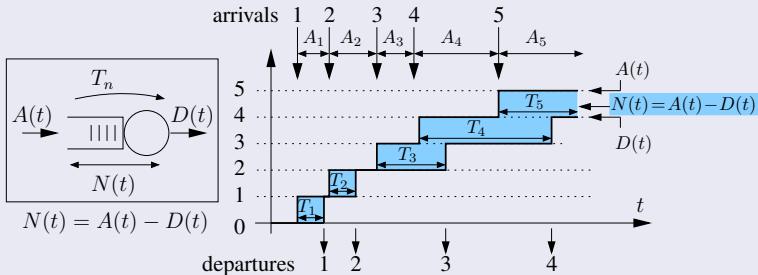
- Arrival rate:  $\lambda = \lim_{t \rightarrow \infty} A(t) / t$
- Mean time in the system:  $T = \lim_{t \rightarrow \infty} (\sum_n T_n) / A(t)$

- The following relation follows:

$$N = \lambda T$$

**Mnemonic: NAT** (Number = Arrivals x Time).

## Graphical proof



- From the graph we have:

$$\frac{1}{t} \int_0^t N(s) ds = \frac{1}{t} \sum_{i=1}^{A(t)} T_i = \frac{A(t)}{t} \frac{\sum_{i=1}^{A(t)} T_i}{A(t)}$$

- Taking the limit  $t \rightarrow \infty$ :  $N = \lambda T$

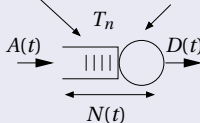


## Application to the waiting line and the server

- We can apply the Little theorem to the **waiting line** and the **server**:

Waiting time in the queue  
of customer  $n$ :  $W_n$

Service time:  $S_n$



Time in the system:

$$T_n = W_n + S_n$$

Expected value:

$$T = W + S$$

where

$$T = E[T_n], W = E[W_n],$$

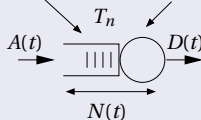
$$S = E[S_n]$$

- Mean number of customers in the queue:**  $N_Q = \lambda W$ .
- Mean number of customers in the server:**  $N_S = \rho = \lambda S$ .

## Mean number in the Server

Waiting time in the queue  
of customer  $n$ :  $W_n$

Service time:  $S_n$



Time in the system:

$$T_n = W_n + S_n$$

Expected value:

$$T = W + S$$

where

$$T = E[T_n], W = E[W_n],$$

$$S = E[S_n]$$

- In a **single server queue** (even if not Markovian):

$$\rho = N_S = E[N_S(t)] = \lambda E[S]$$

$$E[N_S(t)] = 0 \times \pi_0 + 1 \times (1 - \pi_0) = 1 - \pi_0 \Rightarrow \pi_0 = 1 - \rho$$

- $\rho = N_S = \lambda E[S] = 1 - \pi_0$  is the proportion of time the system is busy, in other words, is the **server utilization or load**.



Master in Innovation and Research in Informatics (MIRI)  
Computer Networks and Distributed Systems  
**Stochastic Network Modeling (SNM)**

Queuing  
Theory

Introduction

Kendal  
Notation

Little Theorem

PASTA  
Theorem

Example of PASTA

The M/M/1  
Queue

M/G/1 Queue

M/G/1/K  
Queue

Queues in  
Tandem

Networks of  
Queues

M/G/1 Busy  
Period

## Part IV

# Queuing Theory

### Outline

- Introduction
- Kendal Notation
- Little Theorem
- **PASTA Theorem**
- The M/M/1 Queue
- M/G/1 Queue
- M/G/1/K Queue
- Queues in Tandem
- Networks of Queues
- M/G/1 Busy Period
- M/G/1 Delays
- Matrix Geometric Method



# PASTA Theorem

## Queuing Theory

### Introduction

### Kendal Notation

### Little Theorem

### PASTA Theorem

#### Example of PASTA

### The M/M/1 Queue

### M/G/1 Queue

### M/G/1/K Queue

### Queues in Tandem

### Networks of Queues

### M/G/1 Busy Period

## PASTA Theorem: Poisson Arrivals See Time Averages

- The mean time the chain is in state  $i$  is  $\pi_i \Rightarrow$  using **PASTA**, the **probability that a Markovian arrival see the system in state  $i$  is  $\pi_i$**  (proof: see [1]).
- The equivalent theorem in **discrete time** is the **arrival theorem, RASTA**: Random Arrivals See Time Averages: the **probability that a random arrival see the system in state  $i$  is  $\pi_i$** .

[1] Ronald W Wolff. “**Poisson arrivals see time averages**”. In: *Operations Research* 30.2 (1982), pp. 223–231.



# Little Theorem

## Queuing Theory

Introduction

Kendal Notation

Little Theorem

PASTA Theorem

Example of PASTA

The M/M/1 Queue

M/G/1 Queue

M/G/1/K Queue

Queues in Tandem

Networks of Queues

M/G/1 Busy Period

## Example of PASTA

- Assume that a system can have, at most,  $N$  customers (e.g.  $N - 1$  in the queue and 1 in service).
- Assume that an arrival is **lost** when the system is full.
- By **PASTA** the proportion of Poisson arrivals that see the system full, and are lost, is equal to the proportion of time the system has  $N$  in the system,  $\pi_N$ .
- Thus, **the loss probability is  $\pi_N$** .



# Stochastic Network Modeling (SNM)

## Queuing Theory

Introduction

Kendal  
Notation

Little Theorem

PASTA  
Theorem

The M/M/1  
Queue

Q-matrix

Stationary  
Distribution

Properties

Stability

Example: Loss  
probability in a  
telephone switching  
center

M/G/1 Queue

M/G/1/K  
Queue

Queues in

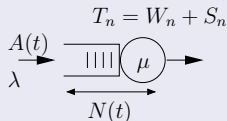
## Part IV

# Queuing Theory

## Outline

- Introduction
- Kendal Notation
- Little Theorem
- PASTA Theorem
- **The M/M/1 Queue**
- M/G/1 Queue
- M/G/1/K Queue
- Queues in Tandem
- Networks of Queues
- M/G/1 Busy Period
- M/G/1 Delays
- Matrix Geometric Method

## The M/M/1 Queue



- Markovian **arrivals** with rate  $\lambda \Rightarrow$  the **interarrival time** is exponentially distributed with mean  $1/\lambda$ :

$$P\{A_n \leq x\} = 1 - e^{-\lambda x}, x \geq 0$$

$\Rightarrow A(t)$  is a **Poisson process**:

$$P(A(t) = i) = \frac{(\lambda t)^i}{i!} e^{-\lambda t}, i \geq 0, t \geq 0$$

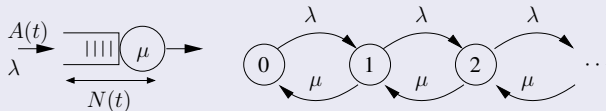
- **Markovian Services** with rate  $\mu \Rightarrow$  **service time** exponentially distributed with mean  $1/\mu$ :

$$P\{S_n \leq x\} = 1 - e^{-\mu x}, x \geq 0$$

## Q-matrix

- The process  $N(t) = \{\text{number in the system at time } t \geq 0\}$  is a CTMC.

OBSERVATION: for a non Markovian service, the process  $N(t)$  would not be a MC! State transition diagram:



- Q-matrix:

$$Q = \begin{bmatrix} -\lambda & \lambda & 0 & 0 & \dots \\ \mu & -(\mu + \lambda) & \lambda & 0 & \dots \\ 0 & \mu & -(\mu + \lambda) & \lambda & \dots \\ \dots & \dots & \dots & \dots & \dots \end{bmatrix}$$





# The M/M/1 Queue

## Queuing Theory

Introduction

Kendal Notation

Little Theorem

PASTA Theorem

The M/M/1 Queue

Q-matrix

Stationary Distribution

Properties

Stability

Example: Loss probability in a telephone switching center

M/G/1 Queue

M/G/1/K Queue

Queues in

## Stationary Distribution

- Solving the M/M/1 queue using flux balancing (or the general solution of a reversible chain):

$$\pi_i = (1 - \rho) \rho^i, i = 0, \dots, \infty$$

$$\text{where } \rho = \frac{\lambda}{\mu}$$



# The M/M/1 Queue

## Queuing Theory

Introduction

Kendal Notation

Little Theorem

PASTA Theorem

The M/M/1 Queue

Q-matrix

Stationary Distribution

Properties

Stability

Example: Loss probability in a telephone switching center

M/G/1 Queue

M/G/1/K Queue

Queues in

## Properties

- Mean **customers in the system**:

$$N = \lim_{t \rightarrow \infty} \frac{1}{t} \int_0^t N(s) ds = \sum_{i=0}^{\infty} i \pi_i = \sum_{i=0}^{\infty} i (1 - \rho) \rho^i = \frac{\rho}{1 - \rho}$$

- Mean **time in the system** (response time):

$$\text{Little: } N = \lambda T \Rightarrow T = \frac{N}{\lambda} = \frac{\rho}{\lambda(1 - \rho)} = \frac{1}{\mu - \lambda}$$

- Mean **time in the queue**:  $W = T - \frac{1}{\mu} = \frac{\rho}{\mu - \lambda}$

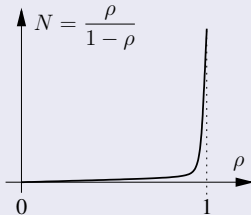
- Mean **Number in the queue**:  $N_Q = \lambda W = \frac{\rho^2}{1 - \rho}$

- Mean **number in the server**:  $N_s = N - N_Q = \rho$

NOTE:  $\pi_0 = 1 - \rho$

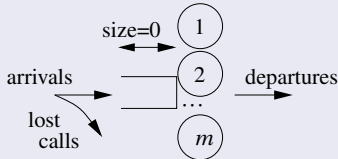
## Stability

- $N$  and  $T$  are proportional to  $1/(1 - \rho) \Rightarrow$  when  $\rho \rightarrow 1 \Rightarrow N, T \rightarrow \infty$ .
- The process  $N(t)$  is **positive recurrent**, **null recurrent** or **transient** according to whether  $\rho = \lambda/\mu$  is below, equal or greater than 1, respectively.



## Example: Loss probability in a telephone switching center

- Hypothesis: Switching center with  $m$  circuits and “lost call”, infinite population, Markovian arrivals with rate  $\lambda$  and exponentially distributed call duration with mean  $1/\mu \Rightarrow$  **M/M/m/m** queue.





# The M/M/1 Queue

## Queuing Theory

### Introduction

### Kendal Notation

### Little Theorem

### PASTA Theorem

### The M/M/1 Queue

### Q-matrix

### Stationary Distribution

### Properties

### Stability

### Example: Loss probability in a telephone switching center

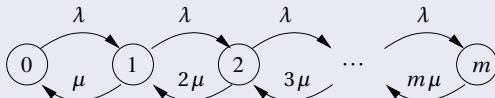
### M/G/1 Queue

### M/G/1/K Queue

### Queues in

## Example: Loss probability in a telephone switching center

- Since the minimum of  $i$  independent and identically exponentially distributed RV with parameter **service time** is exponentially distributed with parameter  $i\mu$ :





# The M/M/1 Queue

## Queuing Theory

### Introduction

### Kendal Notation

### Little Theorem

### PASTA Theorem

### The M/M/1 Queue

### Q-matrix

### Stationary Distribution

### Properties

### Stability

### Example: Loss probability in a telephone switching center

### M/G/1 Queue

### M/G/1/K Queue

### Queues in

## Example: Loss probability in a telephone switching center

- Stationary Distribution of the queue M/M/m/m:
- Solving using the **general solution of a reversible chain**:

$$\text{Define } \rho_k = \frac{\lambda}{(k+1)\mu}, k = 0, \dots, m-1$$

$$\pi_0 = \frac{1}{G}, \pi_i = \frac{1}{G} \prod_{k=0}^{i-1} \rho_k = \frac{1}{G} \left(\frac{\lambda}{\mu}\right)^i \frac{1}{i!}, 0 < i \leq m \Rightarrow$$

$$\pi_i = \frac{1}{G} \left(\frac{\lambda}{\mu}\right)^i \frac{1}{i!}, 0 \leq i \leq m. G = \sum_{k=0}^m \left(\frac{\lambda}{\mu}\right)^k \frac{1}{k!}.$$

- Using **PASTA** Theorem (Poisson Arrivals See Time Average): the **loss call probability** is the probability that the queue is in state  $m$ :  $\pi_m$ , “Erlang B Formula”.