# Performance Modeling
# of Computer Systems and Networks

*Prof. Vittoria de Nitto Personè*

## Analytical models
### (single resource)

Università degli studi di Roma Tor Vergata
Department of Civil Engineering and Computer Science Engineering

1

---

# Performance evaluation techniques

Computational and mathematical techniques to *model*, *simulate* and *analyze* the performance of *stochastic systems*

Modeling: conceptual framework describing a system

Simulate: perform experiments using computer implementation of the model
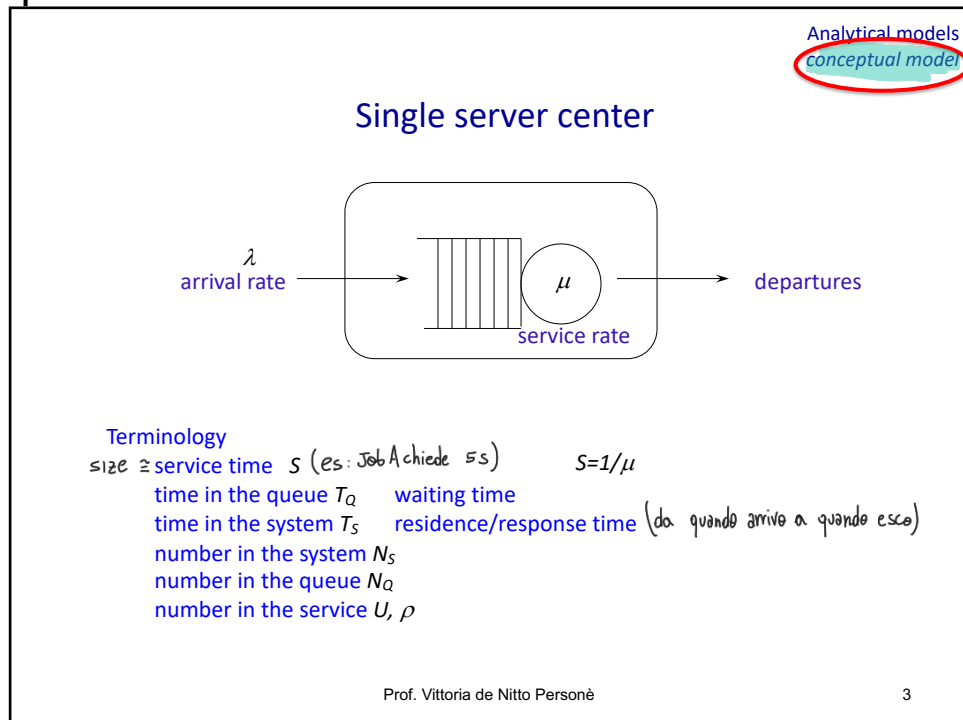
Analyze: draw conclusions from output

*Simulation models*

*Analytical models*

2

1

## p 22 discrete

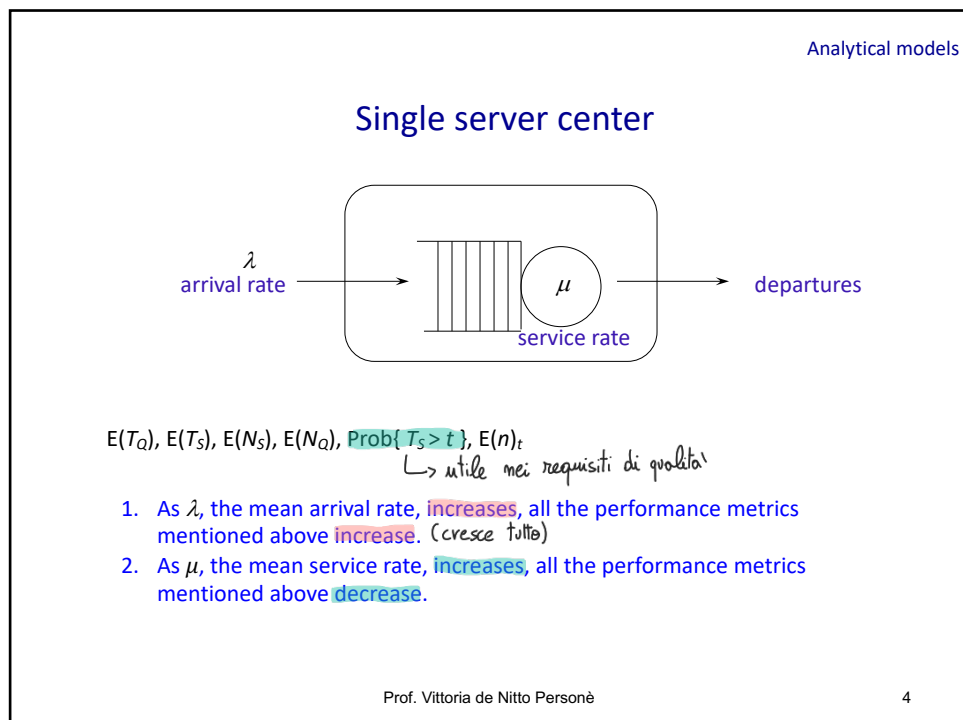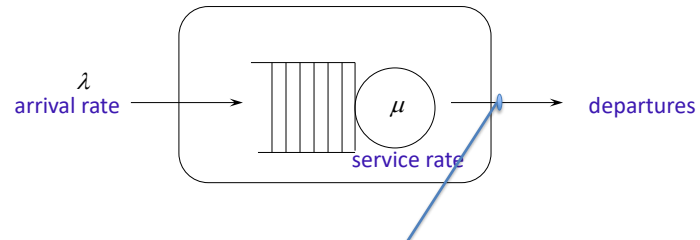# Single server center

$\lambda$
arrival rate → → $\mu$ → departures
service rate

Terminology
size ≅ service time $S$ (es: Job A chiede 5s)     $S=1/\mu$
time in the queue $T_Q$     waiting time
time in the system $T_S$     residence/response time (da quando arrivo a quando esco)
number in the system $N_S$
number in the queue $N_Q$
number in the service $U, \rho$

Prof. Vittoria de Nitto Personè                3

3

# Single server center

$\lambda$
arrival rate → → $\mu$ → departures
service rate

$E(T_Q)$, $E(T_S)$, $E(N_S)$, $E(N_Q)$, Prob{ $T_S > t$ }, $E(n)_t$
⤷ utile nei requisiti di qualità

1. As $\lambda$, the mean arrival rate, increases, all the performance metrics mentioned above increase. (cresce tutto)
2. As $\mu$, the mean service rate, increases, all the performance metrics mentioned above decrease.

Prof. Vittoria de Nitto Personè                4

4

2

# Single server center

$\lambda$
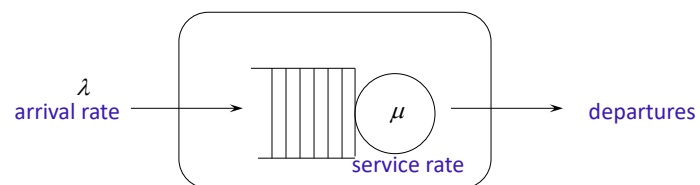arrival rate → | $\mu$ | → departures
service rate

$E(T_Q)$, $E(T_S)$, $E(N_S)$, $E(N_Q)$, Prob$\{ T_S > t \}$, $E(n)_t$

**Def.** throughput (produttività)

$t=1$, $E(n)_1$  n° of completions (departures) in the time unit

Non tempo totale !

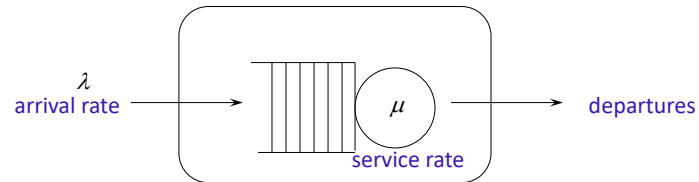Prof. Vittoria de Nitto Personè                    5

5

---

# Single server center

$\lambda$
arrival rate → | $\mu$ | → departures
service rate

**Def.** utilization

How can we "mathematically" define the utilization?

$$\rho = \lambda / \mu$$

Prof. Vittoria de Nitto Personè                    6

6

---

# Single server center

$\lambda$
arrival rate → → departures

$\mu$
service rate

E($T_Q$), E($T_S$), E($N_S$), E($N_Q$), Prob{ $T_S > t$ }, E($n$)$_t$

$$E\left(T_s\right) = E\left(T_Q\right) + E\left(S\right)$$

$T_{TOT}$    in coda    servizio

7

---

# Single server center

$\lambda$
arrival rate → → departures

$\mu$
service rate

E($T_Q$), E($T_S$), E($N_S$), E($N_Q$), Prob{ $T_S > t$ }, E($n$)$_t$

$$E\left(N_s\right) = E\left(N_Q\right) + E\left(number\ in\ service\right)$$

POP. TOTALE
servente singolo    in coda    $\rho \div$ media in corso di servizio

8

# Single server center



$\lambda$
arrival rate

$\mu$
service rate

departures

$E(T_Q)$, $E(T_S)$, $E(N_S)$, $E(N_Q)$, $\text{Prob}\{ T_S > t \}$, $E(n)_t$

$$E(N_s) = E(N_Q) + \rho$$

9

# Single server center

This server is faster

$\lambda = 1/6$  j/s

$\mu = 1/3$ j/s $> \lambda$ , il contrario non sarebbe bilanciato, accumulandosi nella coda

$\lambda = 1/6$  j/s

$\mu = 1/5$ j/s $> \lambda$

Which system has greater throughput?

10

09/03/21

---

Analytical models

## Single server center

**This server is faster**

$\lambda = 1/6$ j/s

$\mu = 1/3$ j/s

$\lambda = 1/6$ j/s

$\mu = 1/5$ j/s

$\mu > \lambda$

By assuming *job flow balance*, the throughput is the same !!
For both systems $X = \lambda = 1/6$ j/s

BUT the faster server shows the shorter queue and so shorter mean response time

In other words, improving the mean response time does not necessarily improve the throughput

*throughput* $X = \dfrac{Job\ completati\ C}{tempo\ \tau}$

$\rho = \dfrac{tempo\ Busy\ B}{tempo\ \tau\ osservazione}$

$\dfrac{C}{\tau} = \dfrac{C}{B} \cdot \dfrac{B}{\tau}$

$\dfrac{completati}{tempo\ occupato} \left[ \dfrac{job}{s} \right] = \mu$

$\rightarrow X = \mu \cdot \rho =$

$= \mu \cdot \dfrac{\lambda}{\mu} = \lambda$

Prof. Vittoria de Nitto Personè                                    11

11

---

Analytical models
*basic laws*

## Single server center

*random*

If the center is in stochastic equilibrium (stationary condition),

$$\lambda < \mu, \quad \rho = \lambda / \mu < 1$$

$E(n)_1 = X = \lambda$

throughput $= MIN(\lambda, \mu)$

Throughput is independent of the service rate $\mu$

If the center is NOT in stochastic equilibrium, (NO bilanciamento flusso)

$\lambda > \mu$, (nella coda c'è sempre qualcuno)

$E(n)_1 = X = \mu$

the center cannot work off the arrival rate, the queue grows unlimited

Prof. Vittoria de Nitto Personè                                    12

12

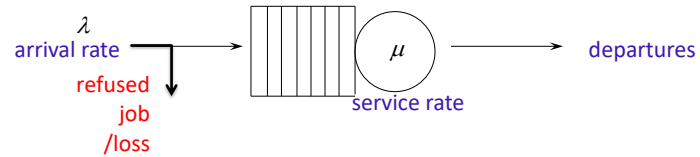6

# Single server center

cooda cresce

What's up if $\lambda > \mu$ ?

the center cannot work off the arrival rate, the queue grows unlimited

media in coda
nel tempo T

$$\mathrm{E}(N_Q \text{ in } T) \geq \lambda T - \mu T = T(\lambda - \mu) \;\; \xrightarrow{>0} \;\; \infty \;\; \text{as} \;\; T \to \infty$$

entrati in T    usati in T

13

---

# Single server center with finite buffer (coda finita)

$\lambda$
arrival rate

refused
job
/loss

$\mu$

service rate

departures

Each arrival when the queue is *full* will be lost

Which is the throughput?

14

## Slide 15

### Single server center with finite buffer

*coda finita SEMPRE STAZIONARIA, la coda non è infinita*

$\lambda$
arrival rate

refused job /loss

$\mu$
service rate

departures

Each arrival when the queue is *full* will be lost

Which is the throughput? (*ricave MENO di $\lambda$*)

$X = \lambda$     $\rho = \lambda / \mu$

No!
On the contrary

$X < \lambda$

Prof. Vittoria de Nitto Personè          15

15

## Slide 16

### Single server center with finite buffer

$\lambda' = X$ (*quelli effettivi*)

$\lambda$
arrival rate

refused job /loss

$\mu$
service rate

departures

Each arrival when the queue is *full* will be lost

Which is the throughput?

$X = \lambda$     $\rho = \lambda / \mu$

No!
On the contrary

$X < \lambda$

Prof. Vittoria de Nitto Personè          16

16

P.6 discrete

## Multi Server Queue $\left(\begin{array}{l}\text{ha } m \text{ serventi} \\ \text{in } /\!/ \text{ identici}\end{array}\right)$

$\mu$

1

$\lambda$
arrivals → queue

2
⋮

departures →

m

$\mu \; \forall$ servente

CONSERVATIVO

$E(S_i)$ $E(S)$ $N_S = \begin{cases} \{0, \dots, m\} & \text{se } N_a = \emptyset \quad \text{(coda vuota)} \\ \\ N_a + m & \text{se } N_a > \emptyset \end{cases}$

$\underbrace{}_{\text{(quanti ce ne sono?)}}$

"come è la media?"

Nel servente singolo $N_S = N_a + \{0, 1\}$

↳ se serve o meno qualcuno

media: $E(N_S) = E(N_a) + \rho$

Prof. Vittoria de Nitto Personè          17

17

---

definisco

• $\rho_i = \left(\dfrac{\lambda}{m}\right) \cdot \dfrac{1}{\mu}$

• $\rho_{GLOBALE} \doteq$ def. come sempre:

## Multi Server Queue

$\dfrac{\text{tasso medio IN} \;\nwarrow \lambda}{\text{tasso MAX OUT}} = \dfrac{\lambda}{m\mu}$

$\searrow m\mu$

$\mu$

1

$\lambda$
arrivals → queue

2
⋮

departures →

m

$\mu$

anche se uguali:

▸ $\rho_i$ dice server "i" quanto usato,

▸ con >1 server mi dice n° server in media su m totali
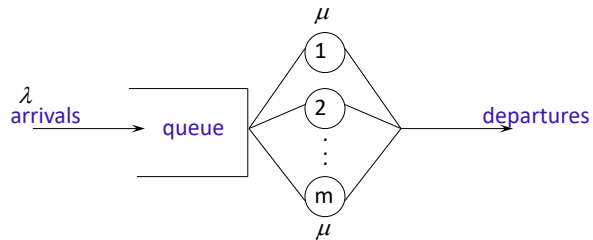
$$E(S_i) = 1/\mu \qquad E(S) = 1/m\mu = E(S_i)/m$$

se m=10, $\rho$=0,5 allora:

· server singolo: usato al 50%

· server multipli: di 10 server, mediamente occupati 0,5·10 = 5

NB: se popolazione $N_S \to \infty$, sempre 'm' saranno occupati, e calcolo $\rho$

Prof. Vittoria de Nitto Personè          18

18

## Multi Server Queue



$\mu$

1

$\lambda$
arrivals → queue

2
⋮
m

departures

$\mu$

$$E\left(N_s\right) = \begin{cases} E\left(N_Q\right) + \rho & if\ m = 1 \quad (singolo) \\ E\left(N_Q\right) + m\rho & f\ m > 1 \quad (multiple) \end{cases}$$
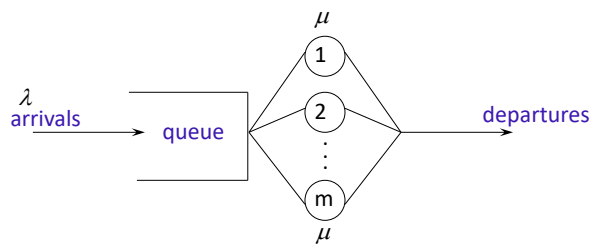
But how is the utilization defined for the multiserver case?

Prof. Vittoria de Nitto Personè                              19

19

---

## Multi Server Queue



$\mu$

1

$\lambda$
arrivals → queue

2
⋮
m

departures

$\mu$

$$\rho_i = \frac{\lambda_i}{\mu} = \frac{\lambda}{m\mu} \qquad \rho_{glob} = \frac{\lambda}{\mu_{glob}} = \frac{\lambda}{m\mu}$$

$\mu_1 = \mu_2 = \cdots$

Prof. Vittoria de Nitto Personè                              20

20

## Slide 21

# Multi Server Queue



$$\rho = \begin{cases} \dfrac{\lambda}{\mu} = \lambda E(S_i) & \text{if m = 1} \\[2ex] \dfrac{\lambda}{m\mu} = \dfrac{\lambda E(S_i)}{m} & \text{if m > 1} \end{cases}$$
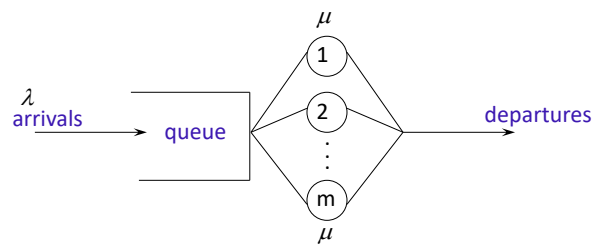
$$\frac{1}{\underset{[s]}{\underset{\downarrow}{\frac{\mu}{1}}}} = E(S_i) \quad [s]$$

Prof. Vittoria de Nitto Personè  21

21

## Slide 22

# Multi Server Queue



$$\rho_i = \rho_{glob} = \frac{\lambda}{m\mu}$$

Prof. Vittoria de Nitto Personè  22
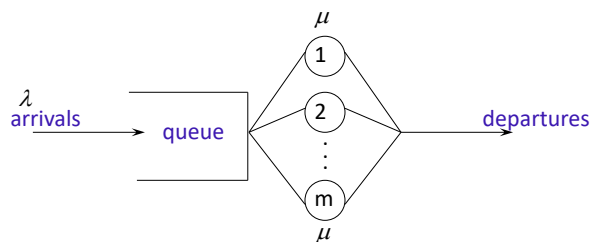
22

11

## Slide 23

# Multi Server Queue

μ
①
$\lambda$
arrivals → queue
②
⋮
departures
m
μ

tempo risposta:

atteso per sedermi

$E(T_s) = E(T_Q) + E(S_i)$

quando arriva a quando va via

$\frac{1}{\mu}, \forall i$ : qualsiasi servente (tutti uguali), da quando uno si siede, fino ad uscire, calcolo media, è $\frac{1}{\mu}$

• introduco $E(S) = \frac{1}{m\mu}$ , da quando uno entra e uno qualsiasi (diverso da chi è entrato) esce, "dopo quanto tempo se ne libera uno di quegli m, in medio"

• $E(T_q)$ = tempo attesa medio (arrivo - siede sul server)

Prof. Vittoria de Nitto Personè                23

23

## Slide 24

# Multi Server Queue

μ
①
$\lambda$
arrivals → queue
②
⋮
departures
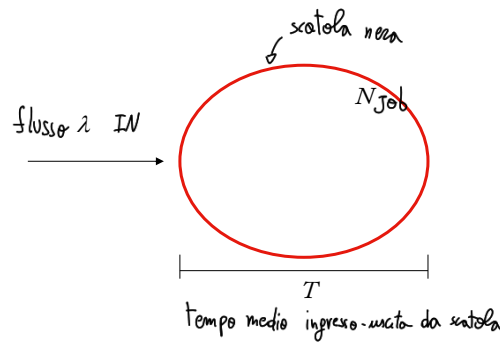m
μ

$$E(T_s) = E(T_Q) + E(S_i)$$

$$E(T_Q) = f(\lambda, \rho, E(S))$$

Prof. Vittoria de Nitto Personè                24

24

Little's law is very important for its broad applicability.
In general, we can see Little's law as applied at a black box:
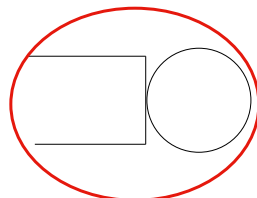it states relations between mean values

scatola nera

$N$ Job

flusso $\lambda$ IN

$T$
tempo medio ingresso-uscita da scatola

Little's Law (1961)

(a) queue discipline is FIFO,
(b) service node capacity is infinite,
(c) flow balance
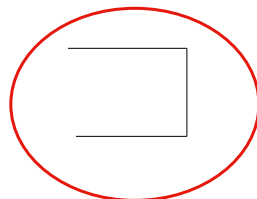
$N = \lambda T$   (sono medie)

If $\lambda$ is the mean arrival rate, $T$ is the mean residence time in the black box,
$N$ is the mean population in the black box, the theorem asserts that, if the system is "stable",
the mean population is given by the "mean arrival flow" multiplied the mean time the jobs
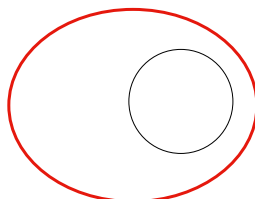spend in the black box

25

---

If the black box is the whole center, the theorem is applied to the center mean population:

$$E(N_S) = \lambda E(T_S)$$

If the black box is just the queue, the theorem is applied to the queue mean population:

$$E(N_Q) = \lambda E(T_Q)$$

If the black box is just the server, the theorem is applied to the server "mean population", in other words to the utilization:

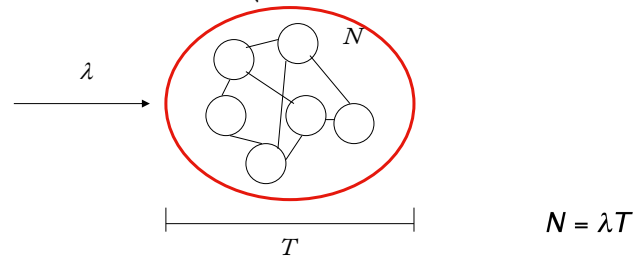$$\rho = \lambda E(S) \quad (\text{legge utilizzazione})$$

tempo che ci sto dentro

26

---

But if the black box is a network of centers, anyway interconnected,

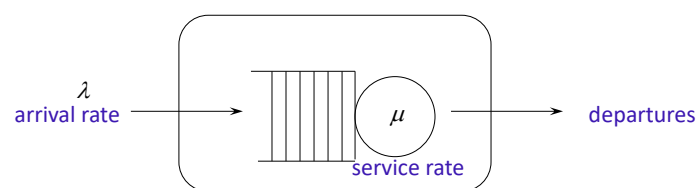↳ sia singoli che multi !

( Non mi interessa dentro come sia fatto)

$N$

$\lambda$

$T$

$$N = \lambda T$$

The theorem is applied to the entire network!!

27

---

# Single server center

$\lambda$
arrival rate

$\mu$

service rate

departures

$$E(T_s) = E(T_Q) + E(S)$$

Little's law · aggiunge

$$E(N_s) = \lambda E(T_s)$$

$$E(N_s) = E(N_Q) + \rho$$

$$E(N_Q) = \lambda E(T_Q)$$

$$E(T_s) = \frac{E(N_s)}{\lambda}$$

$$E(T_Q) = \frac{E(N_Q)}{\lambda}$$

28

28

14

## Multi Server Queue



$$E(T_s) = E(T_Q) + E(S_i)$$

↳ ottesa del "singolo"

Little's law

$$E(N_s) = \lambda E(T_s)$$

$$E(N_Q) = \lambda E(T_Q)$$

tempi ↗ Pop. media

$$E(T_s) = \frac{E(N_s)}{\lambda}$$

$$E(T_Q) = \frac{E(N_Q)}{\lambda}$$

29

---

# #2)

Consider a web server with a mean processing rate of 1.2 job/s.
If the server receives requests with a rate of 0.45 job/s and it has 0.225
enqueued jobs on average, determine: (serv. singolo)
(in media, nella coda)

a) the average utilization
b) the average response time.

picco

During rush hours the arrival rate grows of 20% and the average number of
enqueued jobs becomes 0.3681818.
Determine:
c) the performance metrics a) and b) (calcola di nuovo stessi indici)
d) which further increasing in arrival rate makes the server collapsing (coda ⟶ ∞)
e) the performance metrics a) and b) for the limiting case d).

30

15

# #1)

## SVOLGIMENTO

- "mean processing rate" $\doteq$ service rate $\mu = 1.2$ J/s
- "server receivs requests" $\doteq$ arrival rate $\lambda = 0.45$ J/s
- "enqueued Job average" $\doteq$ popolazione media in coda $E[N_a] = 0.225$

## Richieste

**a)** average utilization $\rho = \dfrac{\lambda}{\mu} = 0.375$

**b)** average response time (tempo TOT nel sistema) $E[T_s] = \dfrac{E[N_s]}{\lambda}$

dove $E[N_s] = E[N_q] + \rho = 0,225 + 0,375 = 0,6$ , allora $E[T_s] = \dfrac{0,6}{0,45} = 1.\overline{3}$

        persone in coda     persone in server

$\lambda$ incrementa del 20% $\rightarrow$ $\lambda' = \lambda \cdot 1,2 = 0,54$ J/s ; $\mu = 1.2$ J/s e $E[N_q'] = 0,3681818$

**c)** ricalcola le metriche

**c.a)** average utilization $\rho' = \dfrac{\lambda'}{\mu} = 0.45$

**c.b)** average response time (tempo TOT nel sistema) $E[T_s'] = \dfrac{E[N_s']}{\lambda'}$

dove $E[N_s'] = E[N_q'] + \rho' = 0,3681818 + 0,45 = 0,8181818$ , allora $E[T_s'] = \dfrac{0.81818}{0.54} = 1.\overline{51}$

        persone in coda     persone in server

Analogamente $E[T_a] = \dfrac{E[N_a]}{\lambda'} = 0,6818188$ s    e    $E[T_s] = E[T_a] + E[S] = 0,681818 s + \dfrac{1}{1.2} = 1,5\overline{1}$

                                                                    $\| \quad \dfrac{1}{\mu}$

**d)** aumento in % tale che server collassa e coda $\rightarrow \infty$

Vorrei $\rho \rightarrow 1$, cioè $\lambda' \rightarrow \mu$ , quindi $\lambda' \cdot \underset{\uparrow}{x} \overset{?}{=} \mu$ , trovo $x = \dfrac{1.2}{0.54} = 2,\overline{2}$ (aumento 120%)

                                        aumento %

**e)** in tale condizione, $\rho = \dfrac{\lambda}{\mu} = 1$    e $E[T_s] = \dfrac{E[N_s]}{\lambda'} = \infty$ (arrivi coda > partenze server)

# #2)

Let us consider a server that processes jobs with rate 0.8 jobs/s.
By assuming that the server receives jobs with a rate depending on the time slot as follows:

    8.00 a.m. – 12.00 a.m. average arrival rate 1.5 jobs/s
    12.00 a.m. – 2.00 p.m. average arrival rate 0.5 jobs/s
    2.00 p.m. – 7.00 p.m. average arrival rate 1.5 jobs/s
    7.00 p.m. – 9.00 p.m. average arrival rate 0.5 jobs/s
    9.00 p.m. – 8.00 a.m. average arrival rate 0.05 jobs/s

Determine:
    a) average arrival rate per day (24 hours)
    b) average utilization per day (˚ ˚)
    c) average throughput per day
    d) average throughput for each time slot (∀ fascia oraria)

Please, justify and comment the results by indicating the used laws.

Soluzione in lect5Dex1intro AM.pdf

Prof. Vittoria de Nitto Personè        31