



# 31/03/2022

## Performance Modeling of Computer Systems and Networks

Prof. Vittoria de Nitto Personè

### The multi-server queue

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

Copyright © Vittoria de Nitto Personè, 2021  
<https://creativecommons.org/licenses/by-nc-nd/4.0/> | CC BY-NC-ND 4.0

1

P. 37 slide      • t. servizio esp. ( $\mu$ ) (continuo)      . M = "exponentially continuous, memoryless"  
                   • arrivi Poisson ( $\lambda$ ) (discreto)       $\hookrightarrow$  Non invecchia      Analytical models  
                   Erlang, 1917      the multiserver queue

**M/M/m abstract scheduling**  
 multiserver

$E(N_Q)$  Erlang

$P_{i+1} = P_i$ , cambio solo il significato  
 che gli do.

probabilità marginale di avere  $n$  job nel centro.  $P(n) = \begin{cases} \frac{1}{n!} (m\rho)^n p(0) & \text{for } n=1, \dots, m \\ \frac{m^m}{m!} \rho^n p(0) & \text{for } n > m \text{ (si forma coda)} \end{cases}$

$p(0) = \left[ \sum_{i=0}^{m-1} \frac{(m\rho)^i}{i!} + \frac{(m\rho)^m}{m!(1-\rho)} \right]^{-1}$

Flusso in arrivo ripartito egualmente ( $\lambda_1 = \lambda_2 = \dots = \lambda_m = \lambda/m$ )

identici

Prof. Vittoria de Nitto Personè

2

Analytical models  
the multiserver queue

### The Erlang-C formula

Prob. tutti i server pieni

$$\begin{aligned}
 P_Q &\cong \Pr\{n \geq m\} = \sum_{n=m}^{\infty} p(n) \\
 &= \sum_{n=m}^{\infty} \frac{m^m}{m!} \rho^n p(0) = \frac{m^m}{m!} p(0) \sum_{n=m}^{\infty} \rho^n \\
 &= \frac{m^m}{m!} p(0) \sum_{n=0}^{\infty} \rho^{n+m} = \frac{m^m}{m!} p(0) \rho^m \sum_{n=0}^{\infty} \rho^n
 \end{aligned}$$

$\frac{1}{1-\rho}$

Prof. Vittoria de Nitto Personè

3

3

Analytical models  
the multiserver queue

### The Erlang-C formula

m serveri tutti occupati  $P_Q = \frac{(m\rho)^m}{m!(1-\rho)} p(0)$

$E(N_Q)$  Erlang =  $P_Q \frac{\rho}{1-\rho}$      $\xrightarrow{\text{tutti server pieni}}$   $E(N_S)$  =  $P_Q \frac{\rho}{1-\rho} + m\rho$

Little's law

$$E(T_Q) = \frac{E(N_Q)}{\lambda} \quad E(T_Q) = P_Q \frac{\rho}{\lambda(1-\rho)} = \frac{P_Q E(S)}{1-\rho}$$

Prof. Vittoria de Nitto Personè

4

4

Analytical models  
the multiserver queue

### The Erlang formula

*La erlang e' per l'esponenziale!*

**M/M/m**

$$E(T_Q)_{\text{Erlang}} = \frac{P_Q E(S)}{1 - \rho}$$

**M/M/1**

$$E(T_Q)_{KP} = \frac{\rho E(S)}{1 - \rho} = \frac{E(S_{rem})}{1 - \rho}$$

$E(S) = \frac{E(S_i)}{m} = \frac{1}{m\mu}$  *se ne libera una qualsiasi*

Prof. Vittoria de Nitto Personè

5

5

Analytical models  
the multiserver queue

### The Multi Server Queue

$c \equiv$  busy servers  $(0, \dots, m)$

$E(c) = \sum_{n=0}^{m-1} np(n) + \sum_{n=m}^{\infty} mp(n) = m\rho$  *probabilità di 'n' server occupati sono sempre 'm'*

$\rho = \sum_{n=0}^{m-1} \frac{n}{m} p(n) + \sum_{n=m}^{\infty} p(n) = \sum_{n=0}^{m-1} \frac{n}{m} p(n) + P_Q \rightarrow \rho \geq P_Q$  *"P\_Q (tutti pieni)"*

$$E(T_Q)_{\text{Erlang}} = \frac{P_Q E(S)}{1 - \rho} \leq E(T_Q)_{KP} = \frac{\rho E(S)}{1 - \rho}$$

Prof. Vittoria de Nitto Personè

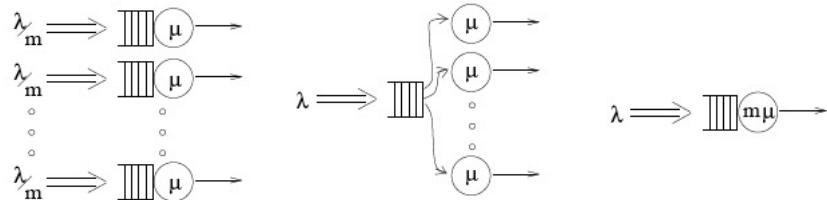
6

6

P. 289 perf , p. 39 slide

Analytical models  
server organizations

## Server Organizations

quale conviene? DIPENDE SEMPRE  
dal contesto!

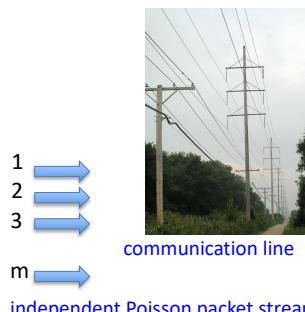
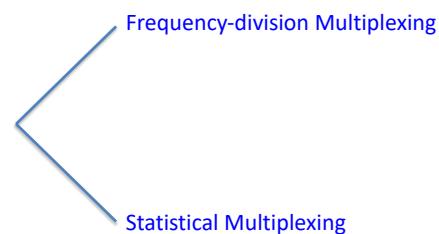
- service rate  $m\mu$
- $\rho = \frac{\lambda}{m\mu}$  per tutt.

Prof. Vittoria de Nitto Personè

7

Analytical models  
server organizations

## Communication systems

each with an arrival rate  $\lambda/m$   
packets per secondthe transmission  
time for each packet Exponential( $1/\mu$ )

Prof. Vittoria de Nitto Personè

8

8

Analytical models  
server organizations

## Communication systems

**Frequency-division Multiplexing** (*separo canale, do garanzie!*)

divide the transmission capacity into  $m$  equal portions

channel

subchannel

magnitude

1  
2  
3  
 $m$

separated streams

Prof. Vittoria de Nitto Personè

9

9

Analytical models  
server organizations

## Communication systems

**Frequency-division Multiplexing**

$\lambda_m \rightarrow \boxed{\mu}$

$\lambda_m \rightarrow \boxed{\mu}$

...

$\lambda_m \rightarrow \boxed{\mu}$

*1° modello visto*

Prof. Vittoria de Nitto Personè

10

10

Analytical models  
server organizations

## Communication systems

### Statistical Multiplexing

keep the transmission capacity as a whole

1  
2  
3  
 $m$   
merge into a single stream

**Multiplexer**

$n$  links  
any speed

1 link, any speed

Prof. Vittoria de Nitto Personè

11

11

Analytical models  
server organizations

## Communication systems

### Frequency-division Multiplexing

$\lambda_m \Rightarrow \boxed{\mu}$

$\lambda_m \Rightarrow \boxed{\mu}$

○  
○  
○

$\lambda_m \Rightarrow \boxed{\mu}$

**Statistical multiplexing**  
condivide la capacità comunicativa, c'è un "merge"  
 $\lambda \Rightarrow \boxed{m\mu}$

How do the two approaches compare with respect to mean response time?

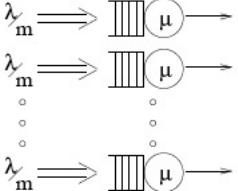
Prof. Vittoria de Nitto Personè

12

12

Analytical models  
server organizations

## Communication systems

<p><b>Frequency-division Multiplexing</b></p>  $E(T_S) = \frac{\rho E(S)}{1-\rho} + E(S) = \frac{E(S)}{1-\rho}$ $E(T_S) = \frac{1}{\mu \left(1 - \frac{\lambda}{\mu}\right)} = \frac{1}{\mu - \lambda}$	<p><b>Statistical multiplexing</b></p>  $\lambda \Rightarrow \boxed{M/M/1}$ $\lambda / \mu$
--	---

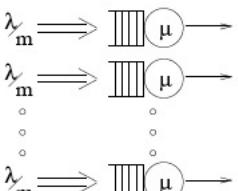
Prof. Vittoria de Nitto Personè

13

13

Analytical models  
server organizations

## Communication systems ( tutti e 3 i centri sono M/M/1 )

<p><b>Frequency-division Multiplexing</b></p>  $E(T_S)^{FDM} = \frac{1}{\mu - \frac{\lambda}{m}} = \frac{m}{m\mu - \lambda}$	<p><b>Statistical multiplexing</b></p>  $E(T_S)^{SM} = \frac{1}{m\mu - \lambda}$
---	--

*FDM shows a response time  $m$  times greater than for SM !*

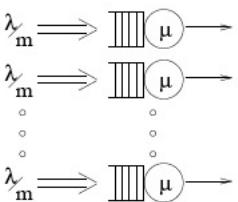
Prof. Vittoria de Nitto Personè

14

14

Analytical models  
server organizations

## Communication systems

<b>Frequency-division Multiplexing</b>  <p><math>\lambda_m \rightarrow \text{[parallel bars]} (\mu) \rightarrow</math></p> <p><math>\lambda_m \rightarrow \text{[parallel bars]} (\mu) \rightarrow</math></p> <p style="text-align: center;">⋮</p> <p><math>\lambda_m \rightarrow \text{[parallel bars]} (\mu) \rightarrow</math></p>	<b>Statistical multiplexing</b>  <p><math>\lambda \rightarrow \text{[parallel bars]} (m\mu) \rightarrow</math></p>
---	---

**1**      QoS guaranteed for each stream:  
           a specific service rate to each stream

**2**      If the original  $m$  streams were very regular (not Poisson), i.e., they were much less variable than Poisson, by merging them, we introduce lots of variability into the arrival stream.  
           This leads to problems if the application requires a low variability in delay, e.g., voice or video.

Prof. Vittoria de Nitto Personè

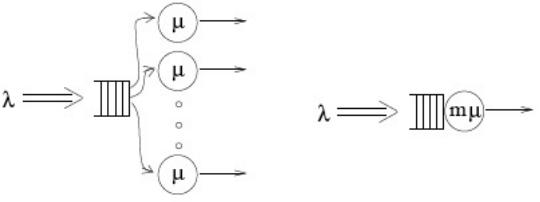
15

15

Analytical models  
server organizations

## Server Organizations

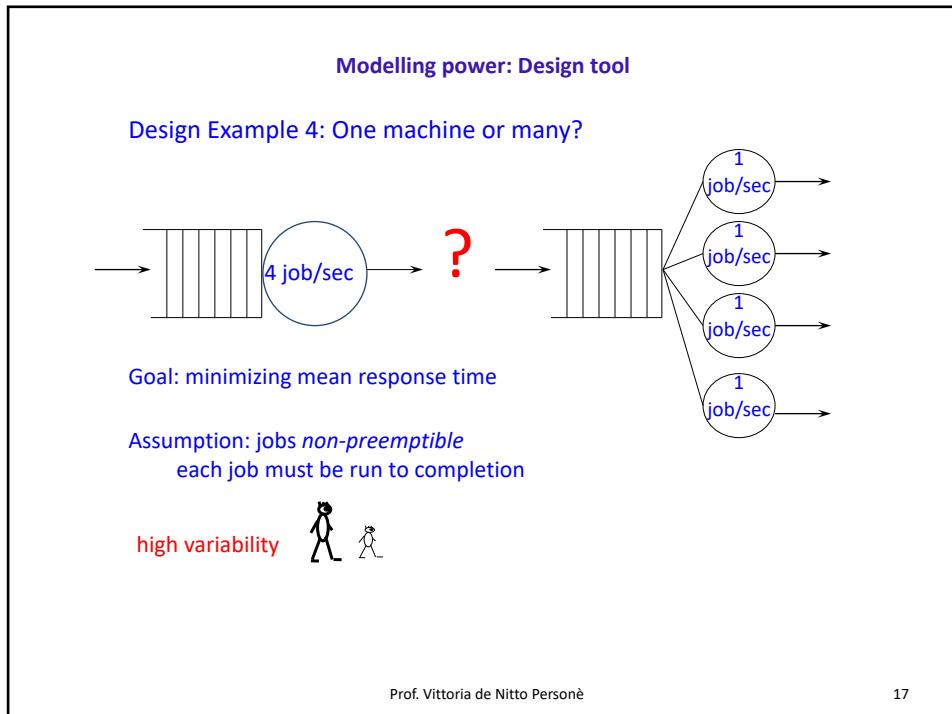
*unica da confrontare*

 <p><math>\lambda \rightarrow \text{[parallel bars]} \xrightarrow{\text{[parallel bars]}} \text{[parallel bars]} (\mu) \rightarrow</math></p> <p><math>\lambda \rightarrow \text{[parallel bars]} \xrightarrow{\text{[parallel bars]}} \text{[parallel bars]} (\mu) \rightarrow</math></p> <p style="text-align: center;">⋮</p> <p><math>\lambda \rightarrow \text{[parallel bars]} \xrightarrow{\text{[parallel bars]}} \text{[parallel bars]} (\mu) \rightarrow</math></p>	 <p><math>\lambda \rightarrow \text{[parallel bars]} (m\mu) \rightarrow</math></p>
--	--

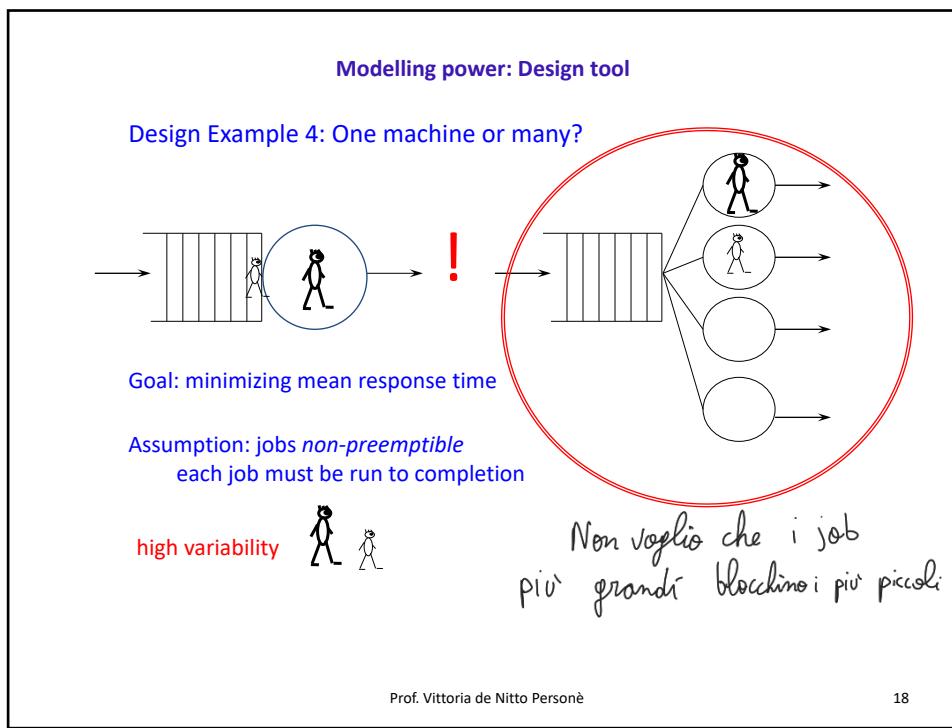
Prof. Vittoria de Nitto Personè

16

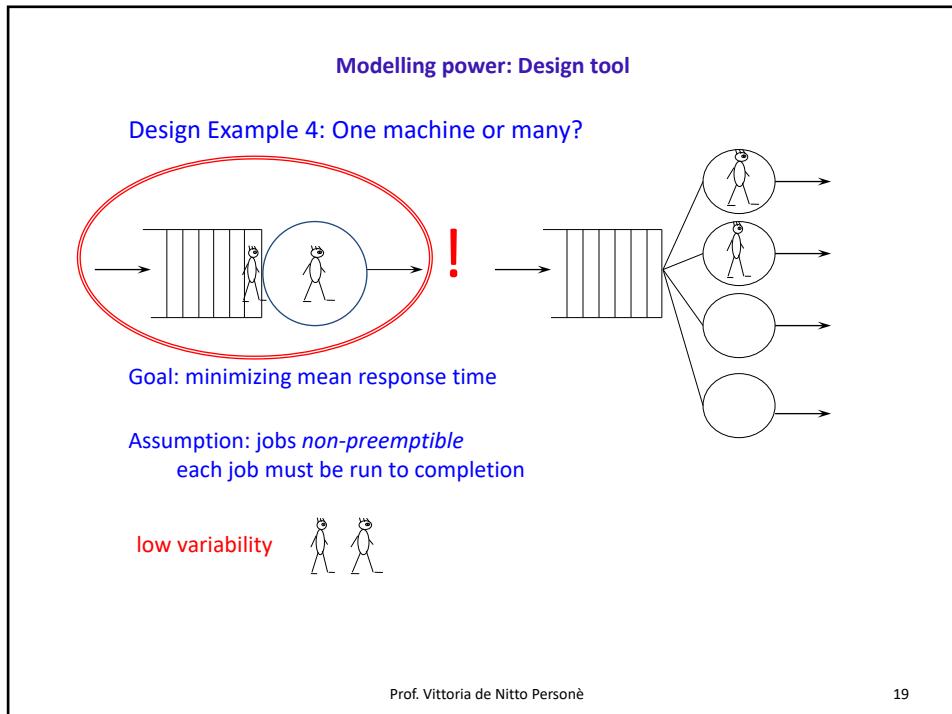
16



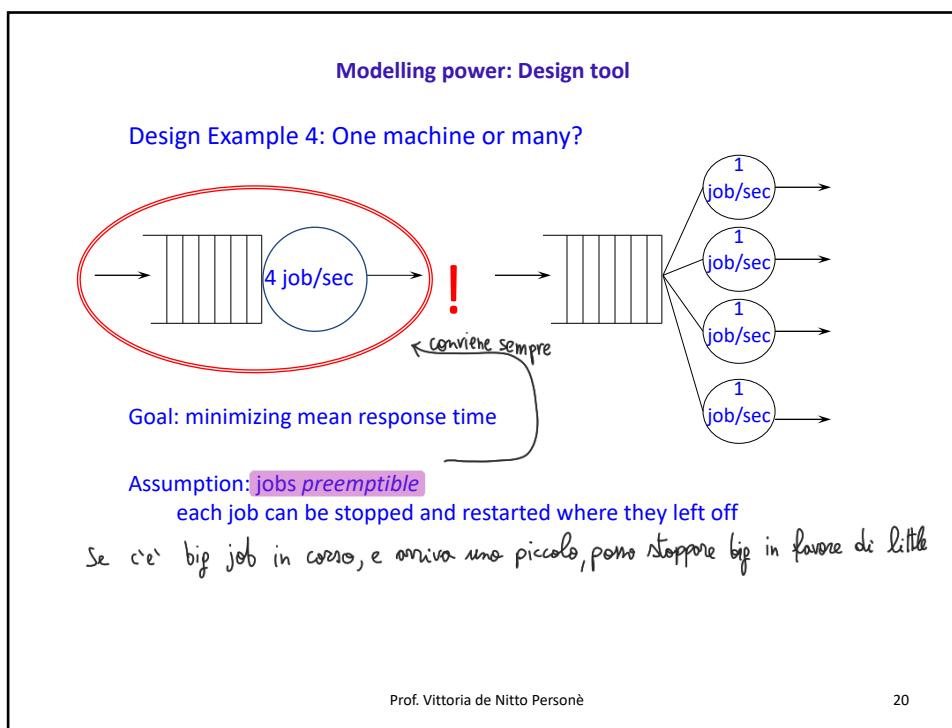
17



18



19



20

Analytical models  
server organizations

## Server Organizations

*multi server, ottendo meno*

*SINGLE WORKER, SERVIZIO più veloce*

*Prob. tutti pieni*

*tempo attesa media*

$E(T_Q)_{\text{Erlang}} = \frac{P_Q E(S)}{1 - \rho}$

$E(T_Q)_{\text{KP}} = \frac{\rho E(S)}{1 - \rho}$

$\rho \geq P_Q$

*servente singolo*

from the waiting time perspective the distributed capacity solution produces an improvement in the user perceived QoS

Tempo servizio rimanente

Prof. Vittoria de Nitto Personè

21

21

Analytical models  
server organizations

## Server Organizations

*concentrato*

What about the response time perspective??

$E(T_S)_{\text{Erlang}} = \frac{P_Q E(S)}{1 - \rho} + E(S_i)$

$E(T_S)_{\text{KP}} = \frac{\rho E(S)}{1 - \rho} + E(S)$

$E(S_i) = \frac{1}{\mu} = m \frac{1}{m\mu} = mE(S)$

Prof. Vittoria de Nitto Personè

22

22

Analytical models  
server organizations

## Server Organizations

What about the response time perspective??

(distribuita)  $E(T_S)_{Erlang} = \frac{P_Q E(S)}{1-\rho} + mE(S)$

Decreasing less than linear ^ v linear growth

$$E(T_S)_{KP} = \frac{\rho E(S)}{1-\rho} + E(S) \quad \text{chi influenza di più?}$$

Prof. Vittoria de Nitto Personè

23

23

Analytical models  
server organizations

## Server Organizations

Performance goal:

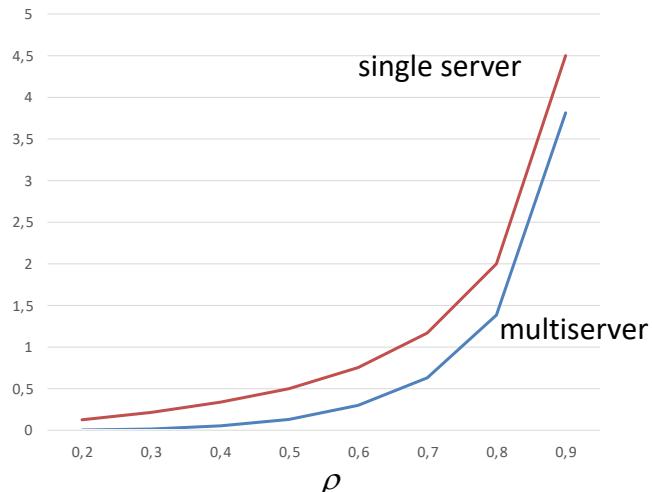
Waiting time perspective	Response time perspective	Distributed capacity
		$\rho \rightarrow 0$ distrib. capac. gives an $m$ times slower organization <small>(altro non c'è quasi mai, meglio servizio breve)</small>
		$\rho \rightarrow 1$ approximately the same response time

Prof. Vittoria de Nitto Personè (Nonno a  $\infty$  le attese nella coda) 24

24

## Waiting time perspective

$$E(S)=0,5 \text{ s}$$



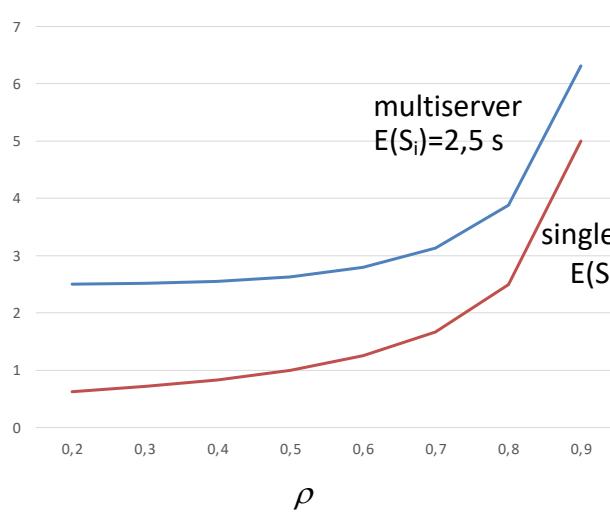
Prof. Vittoria de Nitto Personè

25

25

## Response time perspective

multiserver  
 $E(S_i)=2,5 \text{ s}$   
 single server  
 $E(S)=0,5 \text{ s}$



Prof. Vittoria de Nitto Personè

26

26

Analytical models  
server organizations

### Scaling factor

What about waiting and response time?

$$\left. \begin{array}{l} \rho = \frac{\lambda}{m\mu} \\ E(S_i) = \frac{1}{\mu} \quad E(S) = \frac{1}{m\mu} \end{array} \right| \quad \left. \begin{array}{l} \rho = \frac{a\lambda}{ma\mu} = \frac{\lambda}{m\mu} \\ E(S_i) = \frac{1}{a\mu} \quad E(S) = \frac{E(S_i)}{m} = \frac{1}{am\mu} \end{array} \right|$$

Prof. Vittoria de Nitto Personè      27

27

Analytical models  
server organizations

### Scaling factor

**Mean waiting time**  
*tempo servizio  $\mu t$*

$$E(T_Q)_{m,a}^{\text{ERLANG}} = \frac{P_Q E(S)_{m,a}}{1-\rho} = \frac{1 \cdot P_Q}{ma\mu(1-\rho)} = \frac{1}{a} \frac{P_Q E(S)_{m,1}}{(1-\rho)} = \frac{1}{a} E(T_Q)_{m,1}$$

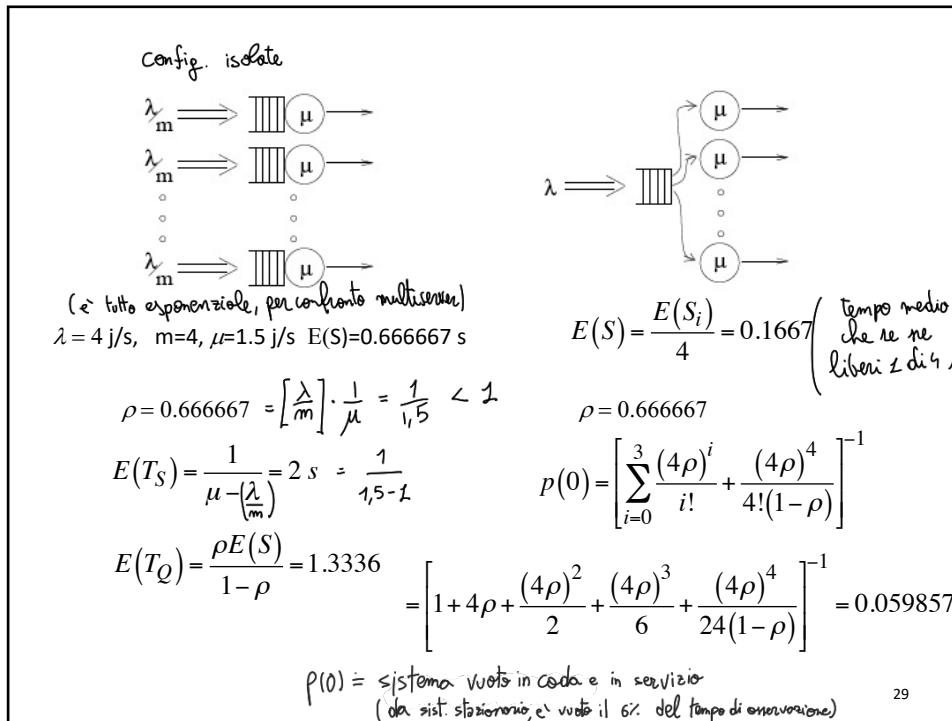
*sorventi fattori scalari*

$$E(S) = \frac{E(S_i)}{m} = \frac{1}{am\mu} \quad E(T_S)_{m,a} = E(T_Q)_{m,a} + \frac{1}{am\mu}$$

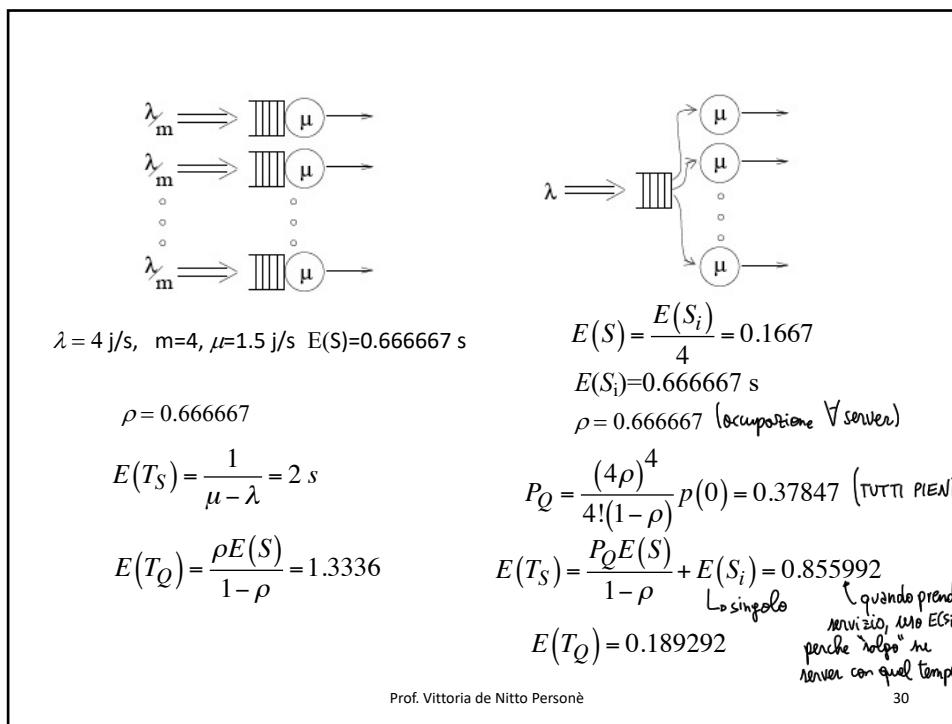
$$= \frac{1}{a} (E(T_S)_{m,1})^{28}$$

Prof. Vittoria de Nitto Personè

28



29



30

$$\lambda \implies \text{III}(\mu) \rightarrow$$

$$\lambda = 4 \text{ j/s}, \quad m\mu = 4 \times 1.5 = 6 \text{ j/s} \quad E(S) = 0.166667 \text{ s}$$

$\rho = 0.666667$  (come prima, ma me lo aspettavo)

$$E(T_S) = \frac{1}{m\mu - \lambda} = 0.5$$

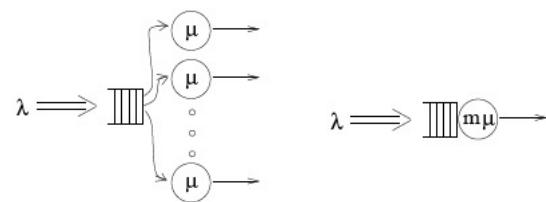
$$E(T_Q) = 0.3334 \left( \text{tempo attesa } \frac{\rho E(S)}{1-\rho} \right)$$

Prof. Vittoria de Nitto Personè

31

31

$$\begin{aligned} \lambda_m &\implies \text{III}(\mu) \rightarrow \\ \lambda_m &\implies \text{III}(\mu) \rightarrow \\ \dots &\dots \\ \lambda_m &\implies \text{III}(\mu) \rightarrow \end{aligned}$$



$$\rho = 0.666667$$

$$E(T_S) = \frac{1}{\mu - \lambda} = 2 \quad \geq \quad E(T_S) = 0.855992 \quad \geq \quad E(T_S) = \frac{1}{\mu - \lambda} = 0.5 \text{ s}$$

$$E(T_Q) = \frac{\rho E(S)}{1-\rho} = 1.3336 \quad \geq \quad E(T_Q) = 0.189292 \quad \leq \quad E(T_Q) = 0.3334$$

attesa minima multiserver

ESERCIZIO :  $\rho = 0.533334$   
PER CASA  $(20\% \text{ in meno})$

$\rho = 0.8$   
 $(20\% \text{ in più})$

Prof. Vittoria de Nitto Personè

32

32