

# Performance Modeling of Computer Systems and Networks

Prof. Vittoria de Nitto Personè

## Multiserver and Priority scheduling

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/>



1

### esempio:

Analytical models  
priority scheduling

#### Assumptions:

- Arrival rate 1 j/s random
- Average demand  $Z=4 \times 10^5$  ocat, expo, do not know size

(Non è che tutti chiedono uguale)

#### Possible configurations:

- 1 server of capacity  $C=10^6$  ocat/s
- Dual-core of  $C/2$  each one

#### QoS requirements:

- Average waiting  $T_Q < 0.15$  s
- For at least 35% of arrivals average response time  $T_S < 0.5$  s

Def.  $E(S) = Z/C = 0.4$  s (  $Z$  indipendente da  $C$ ,  $C$  è caratteristica fisica,  $Z$  no!)  
 $\hookrightarrow$  costante  $\hookrightarrow$  variabile

Capacità processamento macchina != quanto chiede job.  
alcuni elementi sono caratteristiche fisiche della macchina,  
non modellabile.

prof. Vittoria de Nitto Personè

2

2

"il job" = domanda media, quanto chiede?

QoS requirements:

- Average waiting  $T_Q < 0.15$  s

Analytical models  
priority scheduling

$\lambda = 1$  j/s,  $E(S) = 0.4$  s  $\Rightarrow \rho = 0.4$

- 1 server of capacity  $C=10^6$  operat/s  
(M/M/1 KP)  $E(T_Q) = 0.26$  s  $= \frac{\rho E(S)}{1-\rho}$   $E(T_Q)^{\text{Abstract-P}} = 0.2243$  s
- Dual-core of  $C/2$  each one

ogni core è dimezzato!

$$E(S_i) = \frac{Z}{C} = 2 \frac{Z}{C} = 2E(S) = 0.8$$

$$E(T_Q)_{\text{Erlang}} = \frac{\rho E(S)}{1-\rho} = 0.15238 \text{ s} > 0.15$$

NB:  $\frac{\lambda}{2} E(S^2) = \rho E(S)$  se EXP

$P_0 =$  tutti e due pieni;  $m=2$ ,  $\rho(0) = \sum_{i=0}^2 \left[ \frac{(2\rho)^i}{i!} + \frac{(2\rho)^2}{2!(1-\rho)} \right]^{-1} = \left[ 1 + 2\rho + \frac{(2\rho)^2}{2!(1-\rho)} \right]^{-1}$  Prob. centro vuoto

e calcolo  $P_0 = \frac{(m\rho)^m}{m!(1-\rho)}$   $\rho(0)$ , utile per  $E(T_{\text{erlang}})$

prof. Vittoria de Nitto Personè

3

in  $E(T_Q)$  erlang calcolo  $P_q$  e  $P_0$  con le formule citate, usando  $\rho = 0.4$  ed  $E(s) = 0.4$

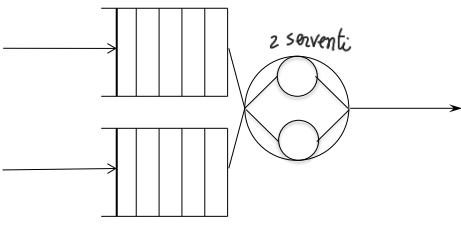
QoS requirements:

- Average waiting  $T_Q < 0.15$  s

Analytical models  
priority scheduling

$\lambda = 1$  j/s,  $E(S) = 0.4$  s  $\Rightarrow \rho = 0.4$

- Dual-core of  $C/2$  each one, classi priorità con multi server

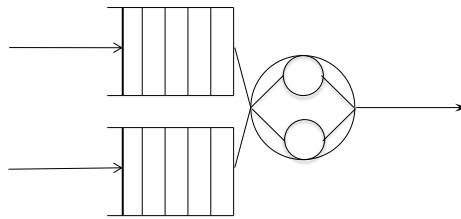


Con due code verrebbe molto più complesso, sconsigliato.

prof. Vittoria de Nitto Personè

4

## Multiserver with priority classes



5

Analytical models  
the multiserver queue

## The Erlang formula

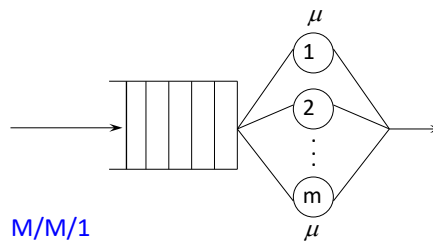
**M/M/m**

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

**M/M/1**

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

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



Prof. Vittoria de Nitto Personè

6

6

## Multiserver with priority classes

$$E(T_Q) = p_1 \frac{\rho_1 E(S)}{(1-\rho_1)} + p_2 \frac{\rho E(S)}{(1-\rho)(1-\rho_1)}$$

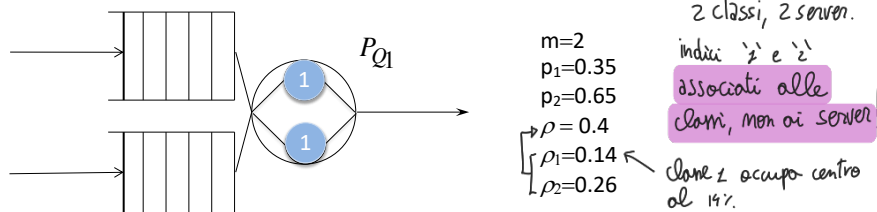


$$E(T_Q) = p_1 \frac{P_{Q1} E(S)}{(1-\rho_1)} + p_2 \frac{P_Q E(S)}{(1-\rho)(1-\rho_1)}$$

C'è prelazione, poichè il primo rapporto vede  $P_{Q1}$ , ed è in funzione di  $\rho_1$ , non  $\rho$  generico.

7

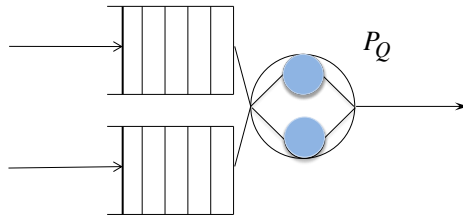
## Multiserver with priority classes



$P_{Q1} = \text{Erlang}(\rho_1) = 0.03438 =$  **tutti server** occupati da job di classe 1, usando erlang.

8

## Multiserver with priority classes



$$P_{Q1} = \text{Erlang}(\rho_1) = 0.03438 \quad P_Q = 0.22857 \quad \left( \begin{array}{l} \text{pieni indipendenti dalle classi:} \\ \text{solo clone 1, solo clone 2, misti...} \end{array} \right)$$

$$E(T_Q) = \rho_1 \frac{P_{Q1} E(S)}{(1-\rho_1)} + \rho_2 \frac{P_Q E(S)}{(1-\rho)(1-\rho_1)} = 0.12077 \quad \left( \begin{array}{l} \text{il più usato per } P_a \\ \text{una sempre } P_2 \end{array} \right)$$

(dual core + prelazione) QoS requirements:

- Average waiting  $T_Q < 0.15 \text{ s}$  !! "bound" rispettato
- GLOBALMENTE non sotto il requisito, non solo la clone 1.

← clone 1, con prelazione, vede solo il "suo"  $P_1$ , gli altri non li vede.

9

QoS requirements: (altro requisito)

- For at least 35% of arrivals average response time  $T_S < 0.5 \text{ s}$

Analytical models  
priority scheduling

$$\lambda = 1 \text{ j/s}, E(S) = 0.4 \text{ s} \quad \longrightarrow \quad \rho = 0.4$$

- 1 server of capacity  $C=10^6$  operat/s

$$E(T_Q) = 0.26 \text{ s} + 0.4 > 0.5, \text{ deve dividere le code!}$$

- Dual-core of  $C/2$  each one

$$E(S_i) = \frac{Z}{\frac{C}{2}} = 2 \frac{Z}{C} = 2E(S) = 0.8 \text{ s} > 0.5, \text{ senza contare } E(T_a)$$

10

## QoS requirements:

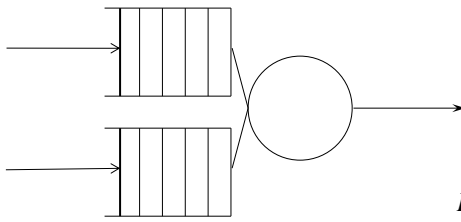
- For at least 35% of arrivals average response time  $T_S < 0.5$  s

Analytical models  
priority scheduling

$$\lambda = 1 \text{ j/s}, E(S) = 0.4 \text{ s} \quad \Rightarrow \quad \rho = 0.4$$

- 1 server of capacity  $C=10^6$  operat/s

## Abstract-P



$$\begin{aligned} p_1 &= 0.35 \\ p_2 &= 0.65 \\ \rho &= 0.4 \\ \rho_1 &= 0.14 \\ \rho_2 &= 0.26 \end{aligned}$$

$$E(T_{S1}) = 0.4651162$$

servente singolo +  
abstract priority

∀ obiettivo ∃ soluzione diversa.

prof. Vittoria de Nitto Personè

11