

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

Analytical models
priority scheduling

Assumptions:

- Arrival rate 1 j/s random
- Average demand $Z=4 \times 10^5$ oxat, expo, do not know size (astratto)
 Z = quanto job chiede, op/job

Possible configurations:

- 1 server of capacity $C=10^6$ oxat/s capacità server, non è v.a.
- Dual-core of $C/2$ each one dual core equivalente, ciascun proc ha capacità dimezzata.

QoS requirements:

- Average waiting $T_Q < 0.15$ s
- For at least 35% of arrivals average response time $T_S < 0.5$ s
la percentuale viene fornita dal testo

Def.

$E(S) = Z/C = 0.4$ s operazioni richiesta/operazioni server nell'unità di tempo

Z e C sono indipendenti, poichè C è una caratteristica fisica dell'hardware, costante; Z è una variabile, è quanto chiede un singolo job (varia da job a job), e mediamente è Z .

prof. Vittoria de Nitto Personè

2

2

QoS requirements:

- Average waiting $T_Q < 0.15$ 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

$$E[T_Q] = \frac{\lambda \cdot E(S)}{1 - \rho} \quad E(T_Q) = 0.26 \text{ s} \quad E(T_Q)^{\text{Abstract-P}} = 0.2243 \text{ s}$$

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

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

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

prof. Vittoria de Nitto Personè

3

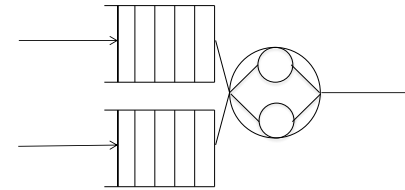
QoS requirements:

- Average waiting $T_Q < 0.15$ s

Analytical models
priority scheduling

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

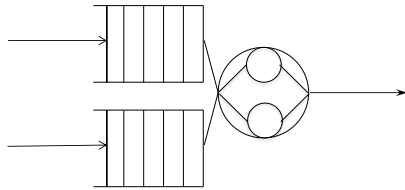
- Dual-core of $C/2$ each one



prof. Vittoria de Nitto Personè

4

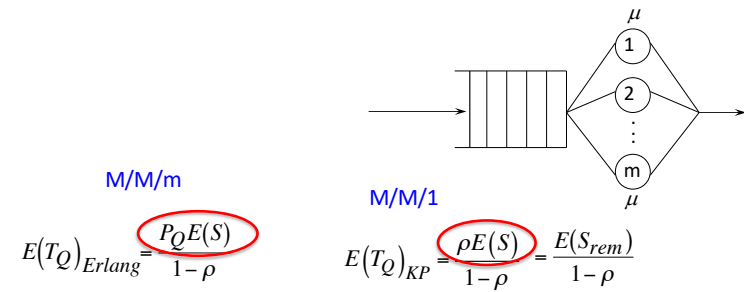
Multiserver with priority classes



5

Analytical models
the multiserver queue

The Erlang formula



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

Prof. Vittoria de Nitto Personè

6

6

Multiserver with priority classes

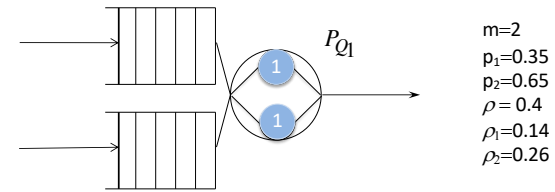
$$E(T_Q) = p_1 \frac{\rho_1 E(S)}{(1-\rho_1)} + p_2 \frac{\rho E(S)}{(1-\rho)(1-\rho_1)}$$



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

7

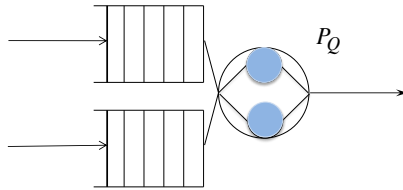
Multiserver with priority classes



$$P_{Q1} = \text{Erlang}(\rho_1) = 0.03438$$

8

Multiserver with priority classes



$$P_{Q1} = \text{Erlang}(\rho_1) = 0.03438 \quad P_Q = 0.22857$$

$$E(T_Q) = p_1 \frac{P_{Q1} E(S)}{(1 - \rho_1)} + p_2 \frac{P_Q E(S)}{(1 - \rho)(1 - \rho_1)} = 0.12077$$

QoS requirements:

- Average waiting $T_Q < 0.15$ s !!

9

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

$$E(T_Q) = 0.26 \text{ s}$$

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

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

prof. Vittoria de Nitto Personè

10

10

QoS requirements:

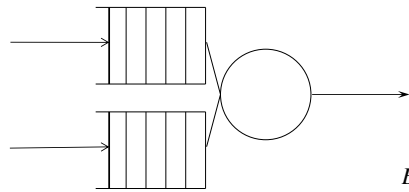
- For at least 35% of arrivals average response time $T_S < 0.5$ s

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



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

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