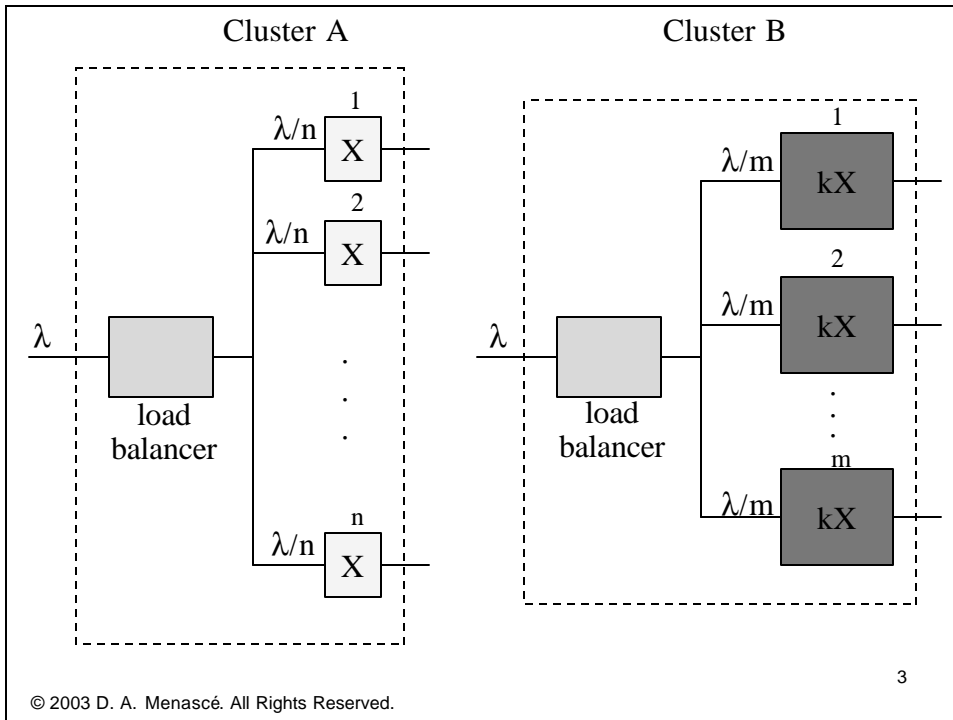


# Tradeoffs in Designing Web Clusters

Daniel A. Menascé  
Department of Computer Science  
George Mason University

## Typical Questions

- Should I use a large number of low-capacity inexpensive servers or a small number of high-capacity costly ones?
- How many servers of a given type are required to provide a certain performance level at a given cost?
- How many servers are needed to build a Web site with a given reliability?



## Comparison Criteria

- Equal Average Response Time
- Equal Cluster Capacity:  $nX = mkX$ .
  - $m = n / k$ .
- Equal Cost:  $n C(X) = m C(kX)$ .
  - $m = n C(x) / C(kX)$
- Equal Reliability.

## Reliability Considerations

$$R_A = 1 - (1 - r_A)^n$$

$$R_B = 1 - (1 - r_B)^m$$

$$R_A = R_B \quad \Rightarrow$$

$$m = \frac{n \log(1 - r_A)}{\log(1 - r_B)}$$

## Basic Performance Model (M/G/1)

$$T = S + \frac{U \times S(1 + C^2)}{2(1 - U)}$$

- T: average response of a Web request
- S: average request service time
- C: coefficient of variation of the service time
- U: server utilization, equal to  $I_w S$

The utilization has to be  $< 1$  ...

$$U_A = \frac{I}{n} \times \frac{1}{X} < 1 \quad \Rightarrow I < nX$$

$$U_B = \frac{I}{m} \times \frac{1}{kX} < 1 \quad \Rightarrow I < mkX$$

## Response Time Equations

$$T_A = \frac{1}{X} + \frac{\frac{I}{n} \times \left(\frac{1}{X}\right)^2 (1+C^2)}{2 \left(1 - \frac{I}{n} \times \frac{1}{X}\right)}$$

$$T_B = \frac{1}{kX} + \frac{\frac{I}{m} \times \left(\frac{1}{kX}\right)^2 (1+C^2)}{2 \left(1 - \frac{I}{m} \times \frac{1}{kX}\right)}$$

## Average Number of Requests

- From Little's Law:

$$N_A = IT_A$$

$$N_B = IT_B$$

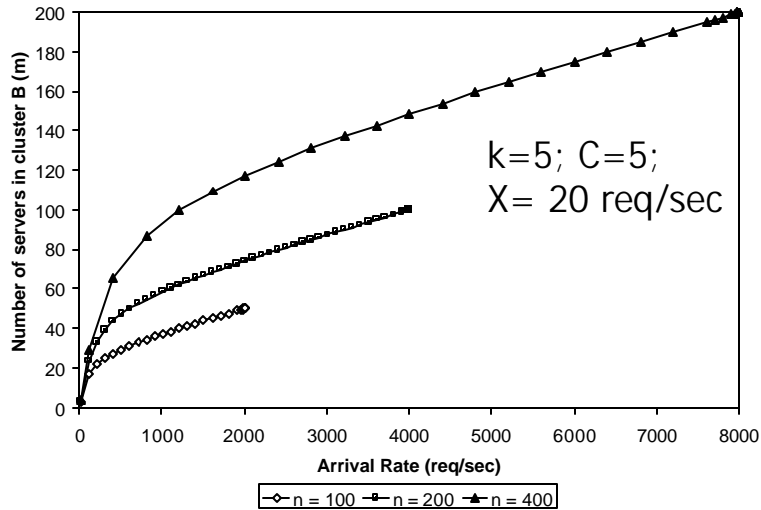
## Equal Response Time Case

$$T_A = T_B \quad \Rightarrow$$

$$m = \frac{1}{kX} \left[ I + \frac{1}{\frac{2(k-1)}{I(1+C^2)} + \frac{k}{nX-1}} \right]$$

$$\lim_{I \rightarrow nX} m = n/k$$

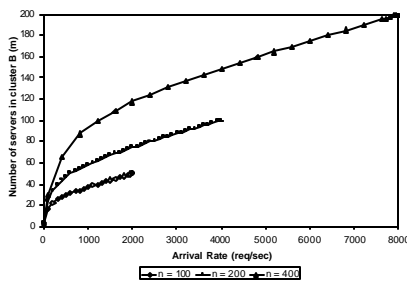
# Equal Response Time Case



© 2003 D. A. Menascé. All Rights Reserved.

11

# Equal Response Time Case



- When cluster A has 400 servers and  $\lambda = 4,810 \text{ req/sec}$ , cluster B needs 160 servers to obtain the same avg. response time of 1.03 sec.

- For a sufficiently large value of  $\lambda$ ,  $m$  increases linearly with  $\lambda$  at the rate of  $(k-1)/(k^2 X)$

© 2003 D. A. Menascé. All Rights Reserved.

12

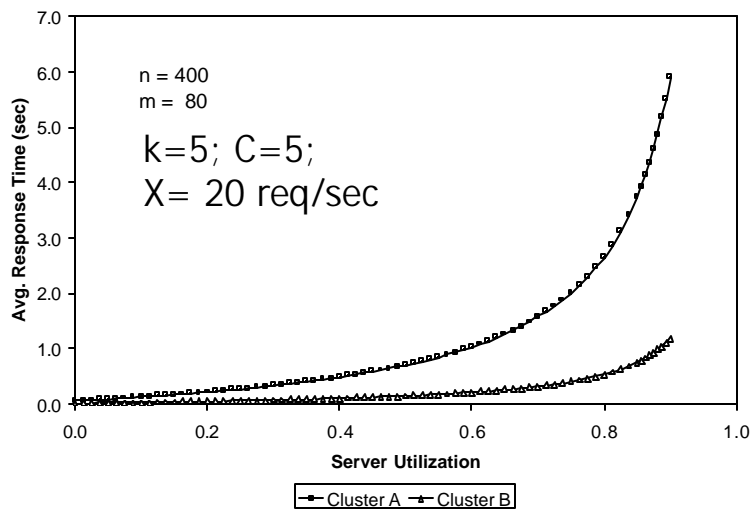
## Equal Capacity Case

- $m = n/k$
- Thus,  $T_B = T_A / k$        $N_B = N_A / k$
- Cluster B average response time is always  $k$  times less than that of cluster A.
- Cluster B can handle  $1/k$  of the requests that cluster A can handle.

© 2003 D. A. Menascé. All Rights Reserved.

13

## Equal Capacity Case



© 2003 D. A. Menascé. All Rights Reserved.

14

## Equal Cost Case

- If the cost is proportional to the capacity, then  $m$  has the same value as in the equal capacity case.
- Sublinear cost function: there is some economy of scale (i.e., the cost per unit of capacity decreases with the capacity).
- Superlinear cost function: the server cost per unit of capacity increases with the capacity.

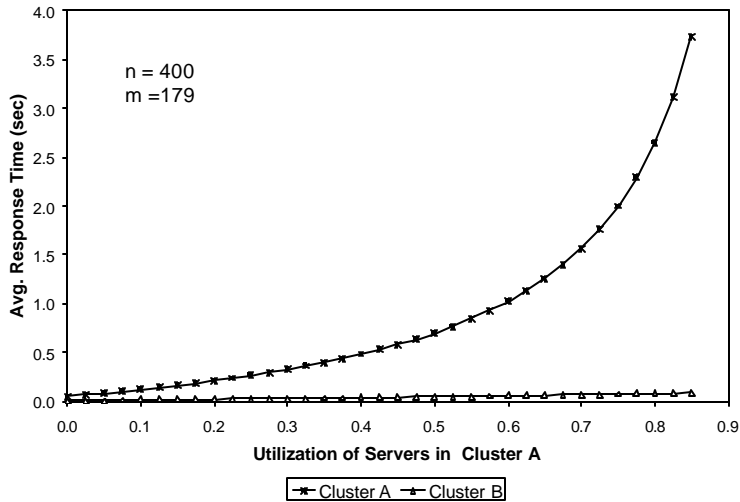
## Equal Cost Case: sub-linear cost function

$$C(x) = a\sqrt{x} \Rightarrow m = m/\sqrt{k}$$

$$U_A = U_B \sqrt{k}$$



## Equal Cost Case: sublinear cost function



© 2003 D. A. Menascé. All Rights Reserved.

17

## Equal Cost Case: superlinear cost function

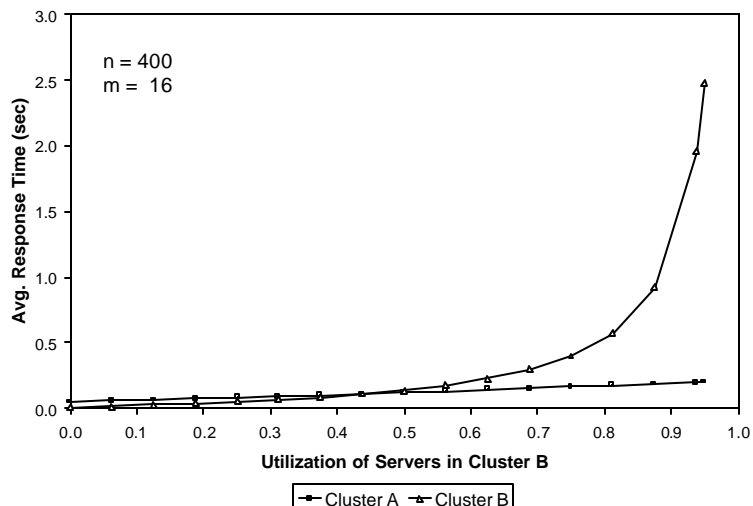
$$C(x) = ax^2 \Rightarrow m = m/k^2$$

$$U_B = kU_A$$

© 2003 D. A. Menascé. All Rights Reserved.

18

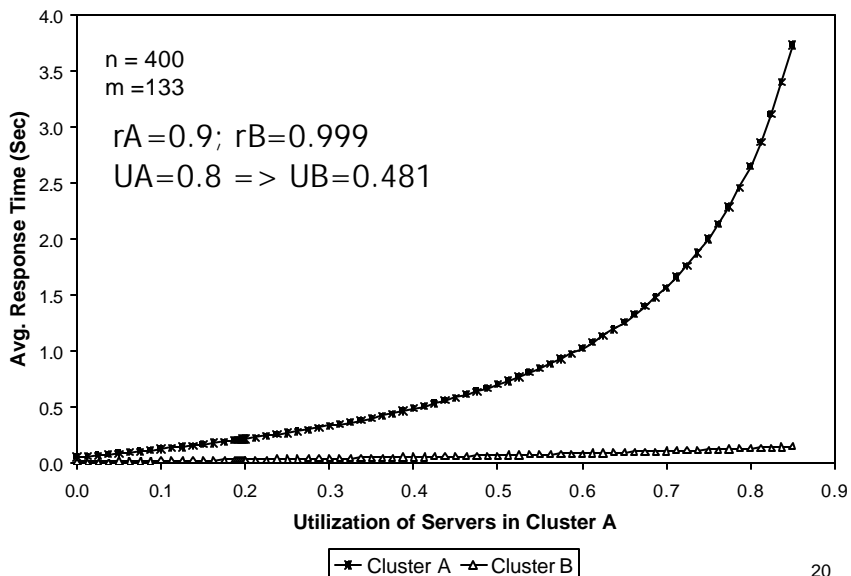
## Equal Cost Case: superlinear cost function



© 2003 D. A. Menascé. All Rights Reserved.

19

## Equal Reliability Case



© 2003 D. A. Menascé. All Rights Reserved.

20

Comparison:  $\lambda = 4,800$  req/sec;  
 $TA = 1.025$  sec;  $UA = 0.6$ ;  $n = 400$

$T_B$ (sec)	$m$	$U_B$	Case
1.025	54	89%	Equal response time.
0.205	80	60%	Equal total capacity
0.058	179	27%	Equal cost. $C(x) = vx$
8	16	> 100%	Equal cost. $C(x) = x^2$
0.083	133	36%	Equal cluster reliability.

21