

From Systems to Descriptive Models

Prof. Daniel A. Menascé
Department of Computer Science
George Mason University
www.cs.gmu.edu/faculty/menasce.html

1

© 2004 D. A. Menascé. All Rights Reserved.

Copyright Notice

- Most of the figures in this set of slides come from the book “Performance by Design: computer capacity planning by example,” by Menascé, Almeida, and Dowdy, Prentice Hall, 2004. It is strictly forbidden to copy, post on a Web site, or distribute electronically, in part or entirely, any of the slides in this file.

2

© 2004 D. A. Menascé. All Rights Reserved.

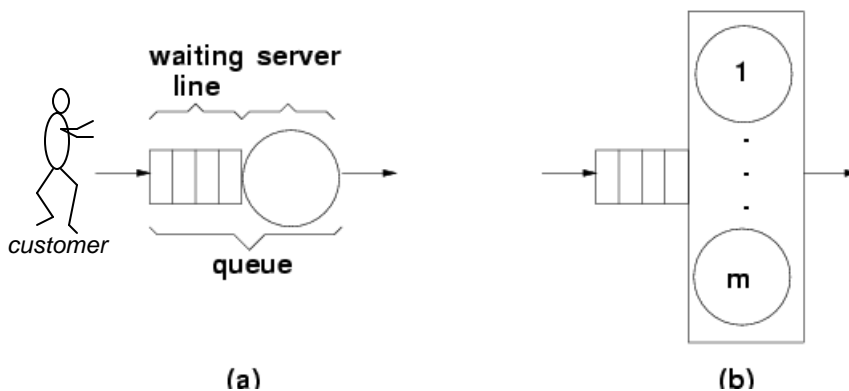
Modeling

- Abstraction of a real system.
- Should capture enough details to satisfy goals of the study.
- Types of models:
 - Simulation
 - Analytic
 - Hybrid

3

© 2004 D. A. Menascé. All Rights Reserved.

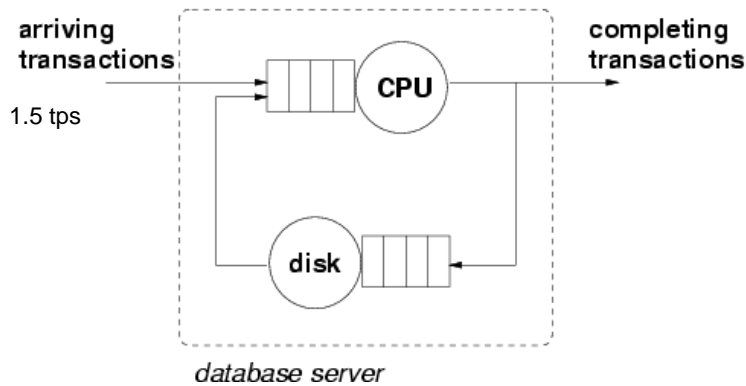
Graphical Notation for Queues



4

© 2004 D. A. Menascé. All Rights Reserved.

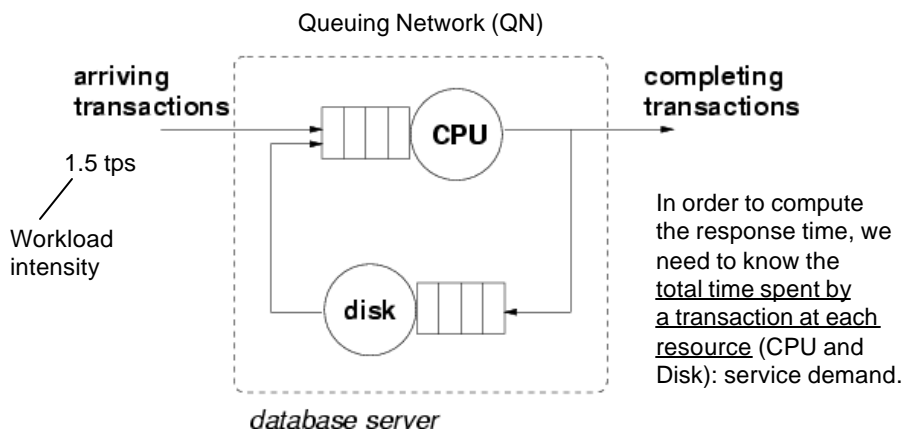
Simple Database Server Example: Open Workload



5

© 2004 D. A. Menascé. All Rights Reserved.

Simple Database Server Example: Open Workload



6

© 2004 D. A. Menascé. All Rights Reserved.

Multiclass DB Example

transaction group	percent of total	avg. CPU time (sec)	avg. no. I/Os
Trivial	45%	0.04	5.5
Medium	25%	0.18	28.9
Complex	30%	1.20	85

Each transaction group is assigned to a customer class in the Queuing Network.

7

© 2004 D. A. Menascé. All Rights Reserved.

When do we need multiple classes?

- Heterogeneous service demands.
- Different types of workloads.
- Different service level objectives.

8

© 2004 D. A. Menascé. All Rights Reserved.

Open Class

- Workload intensity specified by an arrival rate (usually independent of the system state).
- Unbounded number of customers in the system.
- Throughput is an input parameter, which is equal to the arrival rate in equilibrium.

9

© 2004 D. A. Menascé. All Rights Reserved.

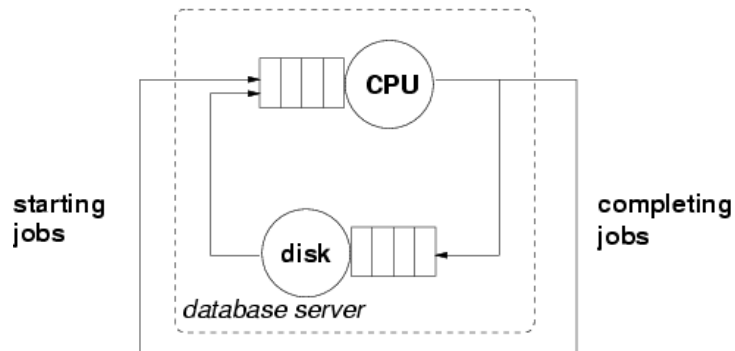
Closed Class

- Workload intensity specified by customer population (i.e., concurrency level)
- Bounded and known number of customers in the system.
- Throughput is an output parameter.

10

© 2004 D. A. Menascé. All Rights Reserved.

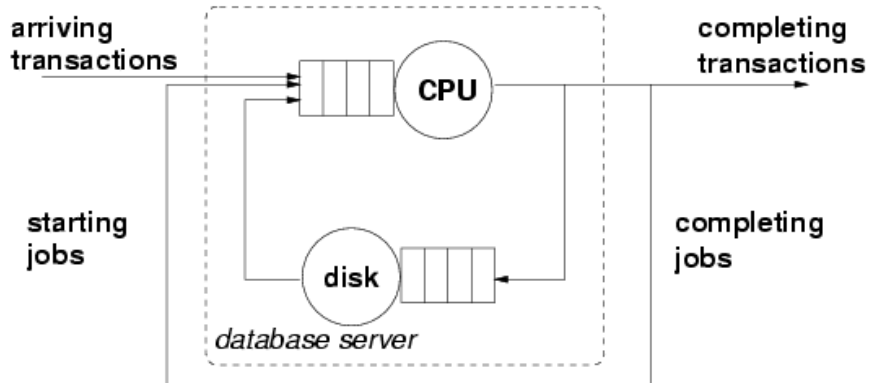
Closed Workload Example



11

© 2004 D. A. Menascé. All Rights Reserved.

Mixed Workload Example



12

© 2004 D. A. Menascé. All Rights Reserved.

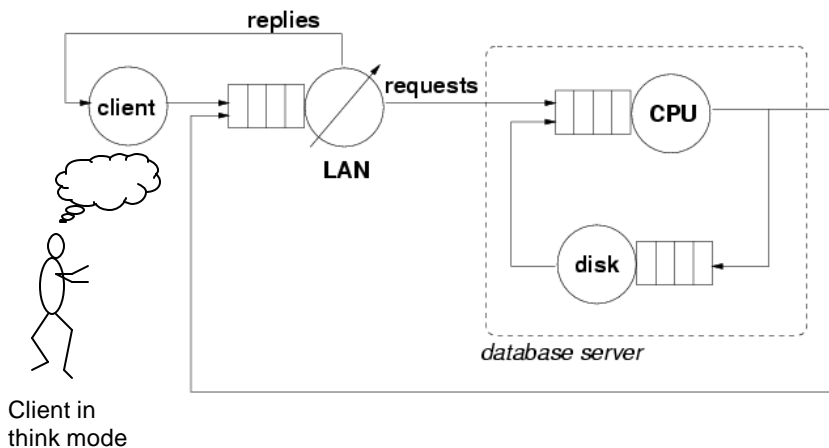
Service Level Agreements (SLA)

- Agreement between users and providers of computing services on the levels of various performance metrics.
 - 99.99% availability during 8:00am-11:00pm period and 99.9% at other times
 - Less than 4 sec page download time for requests over non-secure connections less than 6 sec for requests over secure connections
 - Minimum throughput of 2,000 page downloads/sec.

13

© 2004 D. A. Menascé. All Rights Reserved.

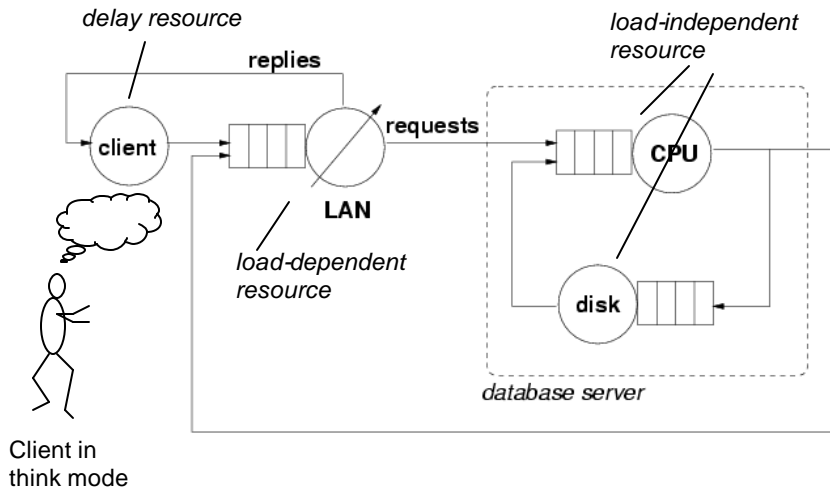
Different Types of Resources



14

© 2004 D. A. Menascé. All Rights Reserved.

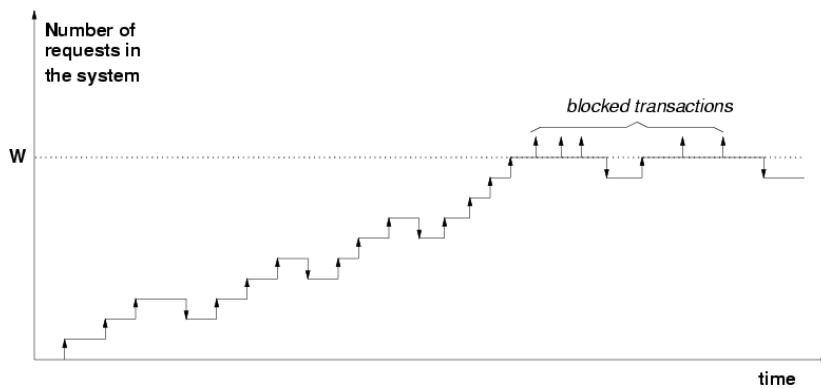
Different Types of Resources



15

© 2004 D. A. Menascé. All Rights Reserved.

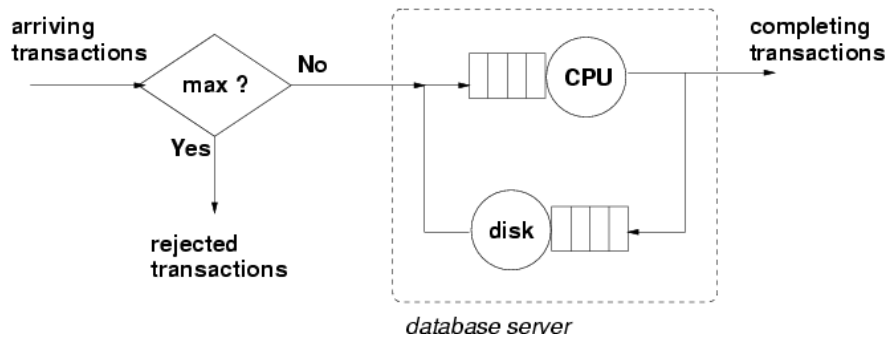
Blocking (Admission Control)



16

© 2004 D. A. Menascé. All Rights Reserved.

System with Admission Control

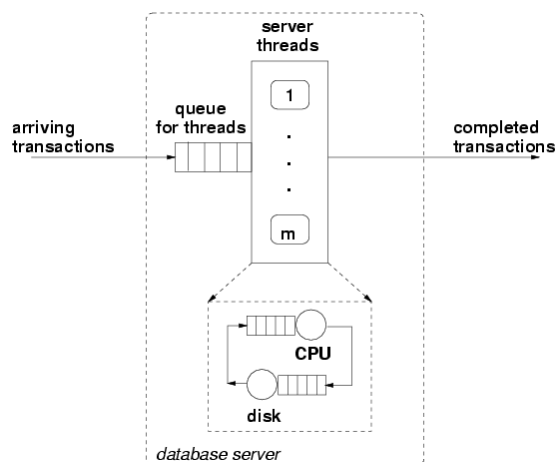


$$\text{Throughput} = \text{Arrival Rate} \times (1 - \text{Probability of Rejection})$$

17

© 2004 D. A. Menascé. All Rights Reserved.

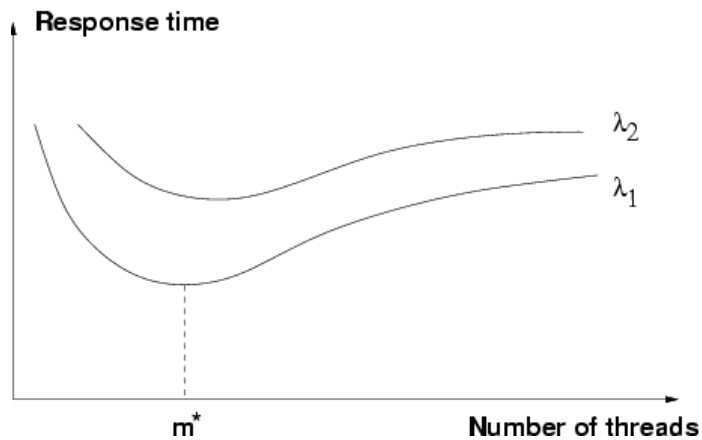
Database Server with Software Contention



18

© 2004 D. A. Menascé. All Rights Reserved.

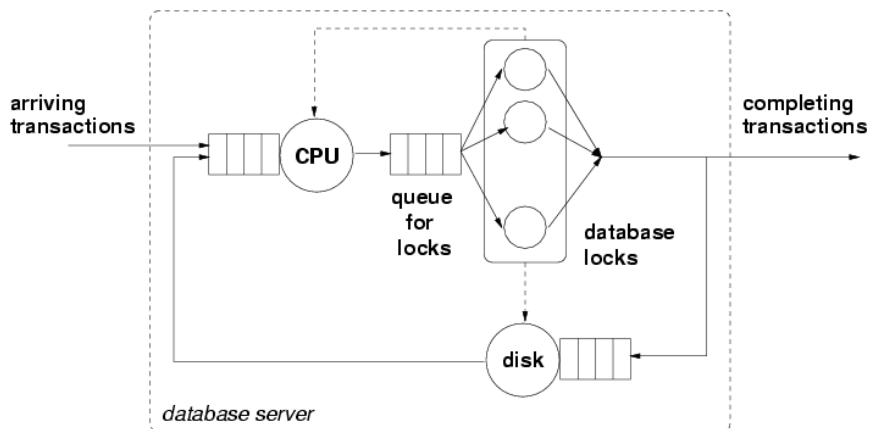
Effect of Software Contention



19

© 2004 D. A. Menascé. All Rights Reserved.

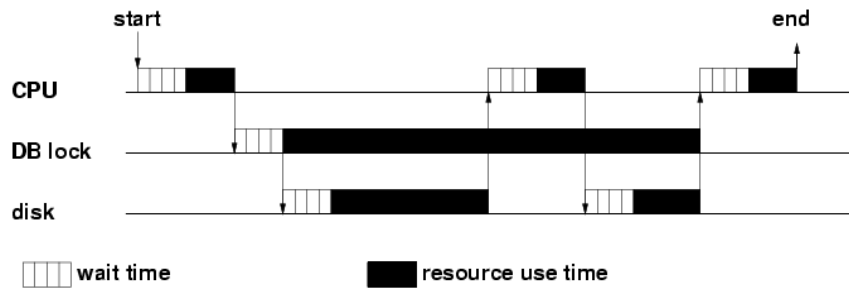
A DB Server with Lock Contention



20

© 2004 D. A. Menascé. All Rights Reserved.

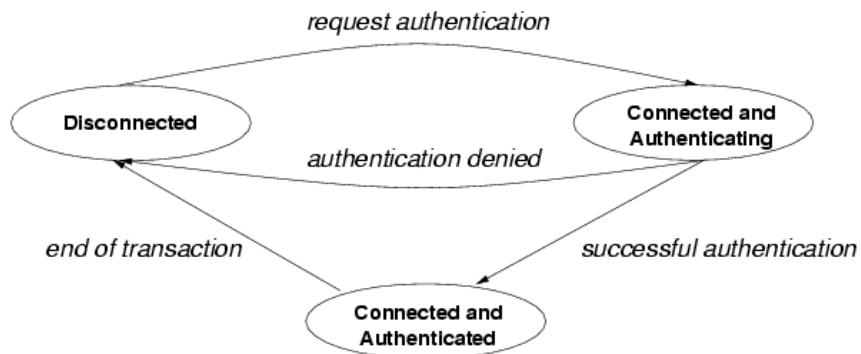
Simultaneous Resource Possession Time Axes



21

© 2004 D. A. Menascé, All Rights Reserved.

Class Switching



Authentication can be significantly more CPU-intensive than DB access.

22

© 2004 D. A. Menascé, All Rights Reserved.

Queuing Disciplines

- First Come First Served (FCFS)
- Priority Queuing (FCFS breaks the tie):
 - Non-preemptive
 - Preemptive resume
 - Preemptive repeat.
- Round-Robin (RR)
- Processor Sharing (PS)
- Last Come First Served-Preemptive Resume (LCFS-PR)
- Shortest Job First (SJF)

23

© 2004 D. A. Menascé. All Rights Reserved.

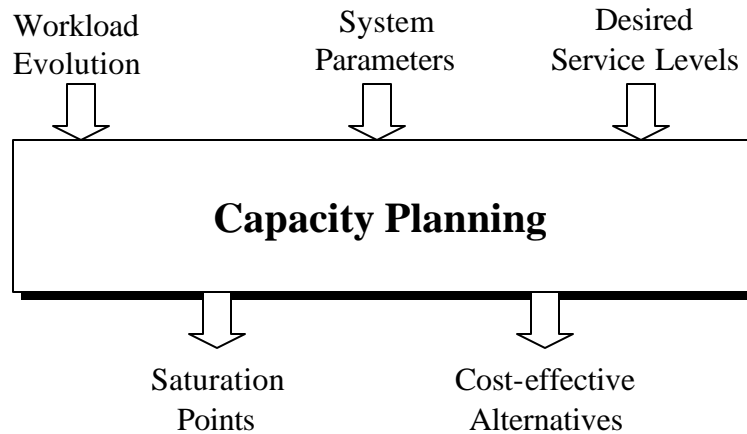
QN Models

- K queues ($i = 1, \dots, K$)
- R classes ($r = 1, \dots, R$)
- Input parameters:
 - Workload intensity
$$\vec{I} = (I_1, \dots, I_r, \dots, I_R) \quad \text{for open classes}$$
 - $$\vec{N} = (N_1, \dots, N_r, \dots, N_R) \quad \text{for closed classes}$$
 - Service demands (D_{ir})

24

© 2004 D. A. Menascé. All Rights Reserved.

Capacity Planning Input and Output Variables

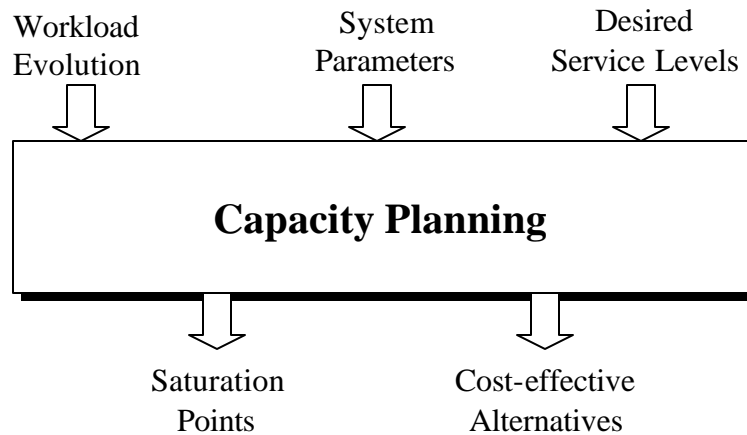


25

© 2004 D. A. Menascé. All Rights Reserved.

Capacity Planning Input and Output Variables

- intensity*
- nature*

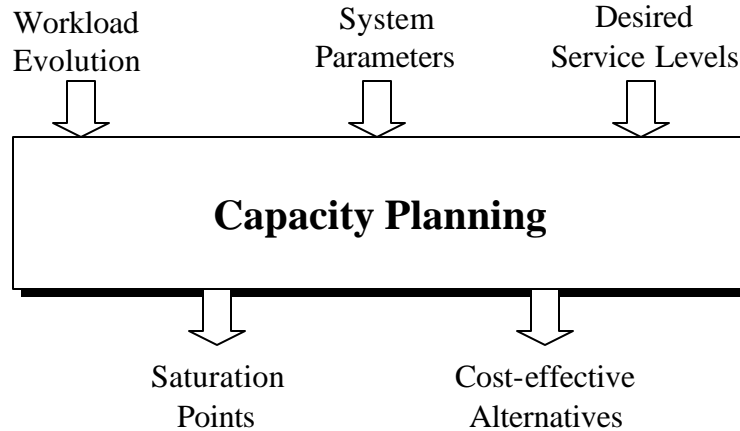


26

© 2004 D. A. Menascé. All Rights Reserved.

Capacity Planning Input and Output Variables

- *intensity*
- *nature*
- *processors, disks, networks*
- *max. no. connections*

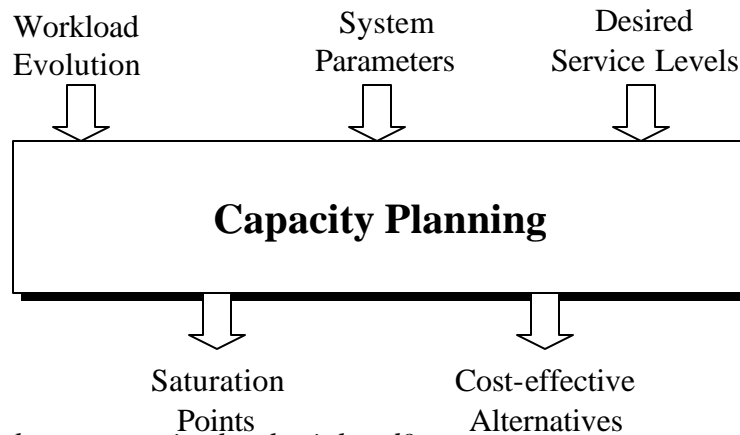


© 2004 D. A. Menascé. All Rights Reserved.

27

Capacity Planning Input and Output Variables

- *intensity*
- *nature*
- *processors, disks, networks*
- *max. no. connections*
- *RT ≤ 2 sec*
- *tput > 10 tps*

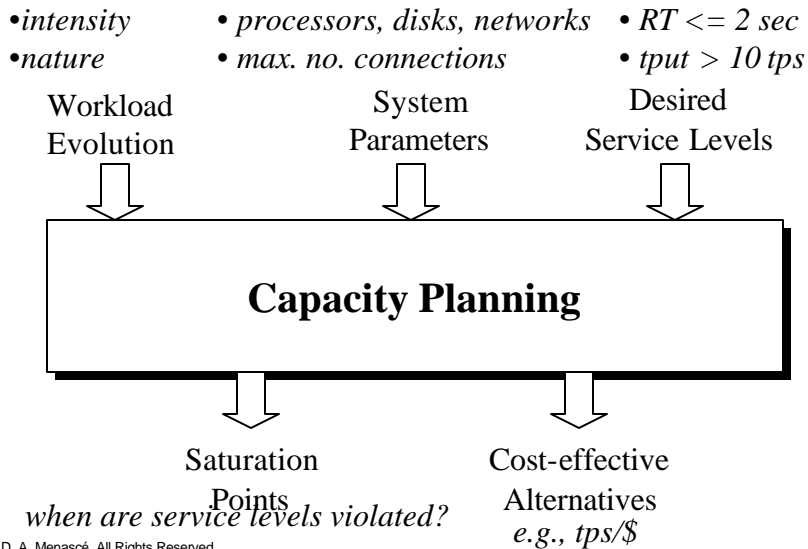


when are service levels violated?

© 2004 D. A. Menascé. All Rights Reserved.

28

Capacity Planning Input and Output Variables



29

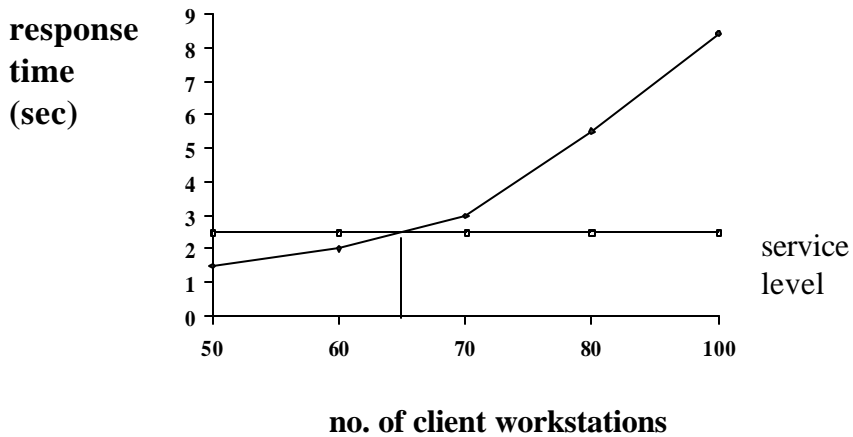
Capacity Planning Definition

Capacity Planning is the process of *predicting* when the *service levels* will be violated as a function of the *workload evolution*, as well as the determination of the most cost-effective way of delaying system saturation.

30

© 2004 D. A. Menascé. All Rights Reserved.

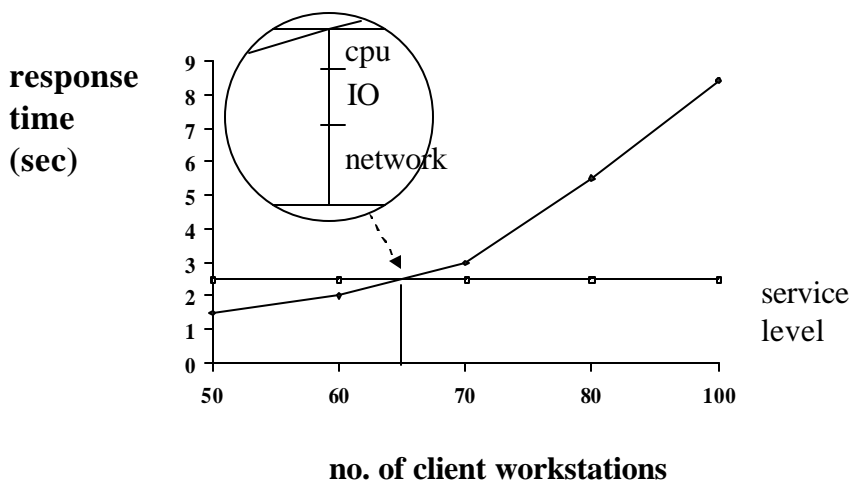
Capacity Planning Concept



31

© 2004 D. A. Menascé. All Rights Reserved.

Capacity Planning Concept



32

© 2004 D. A. Menascé. All Rights Reserved.

Typical Capacity Planning Questions

- Situation: migrating from a mainframe based to a C/S system.
- Questions:
 - how many clients will the new system support with acceptable response time?
 - How many servers and how should they be configured to handle the load?
 - Should I use a two-tier or a three-tier architecture?

33

© 2004 D. A. Menascé. All Rights Reserved.