

CS 672

System Level Performance Models of Computer Systems

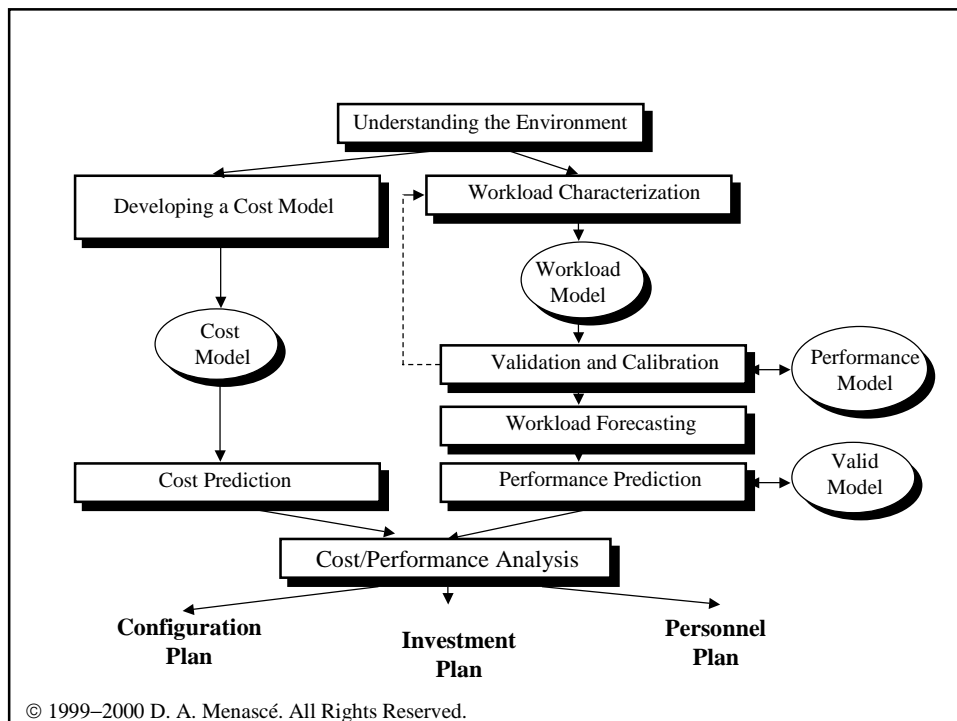
Dr. Daniel A. Menascé

<http://www.cs.gmu.edu/faculty/menasce.html>

Department of Computer Science

George Mason University

© 1999–2000 D. A. Menascé. All Rights Reserved.



© 1999–2000 D. A. Menascé. All Rights Reserved.

Part V: Learning Objectives

Characterize system-level models

Present State Transition Diagram (STD) technique

Show general solution to STDs

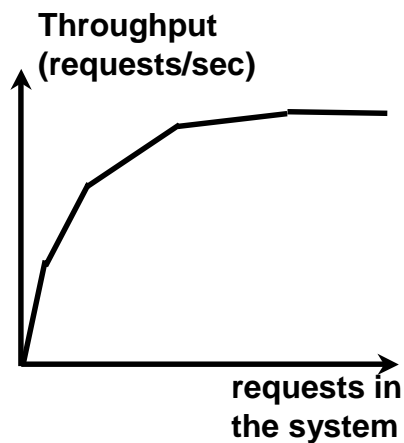
Show how to obtain performance metrics from the solution of STDs

3

© 1999–2000 D. A. Menascé. All Rights Reserved.

System-level Models

- System is seen as a black box.
- Only its input-output characteristics are considered.
- Inputs: arrivals of requests
- Output: throughput.



4

© 1999–2000 D. A. Menascé. All Rights Reserved.

System-level Example

- A Web server receives 10 requests/sec.
- The maximum number of requests in the server is 3.
- Requests that arrive and find three requests being processed are rejected.

5

© 1999–2000 D. A. Menascé. All Rights Reserved.

System-level Example

- The measured throughput as a function of the number of requests is:

Number of requests	Throughput (req/sec)
0	0
1	12
2	15
3	16

6

© 1999–2000 D. A. Menascé. All Rights Reserved.

System-level Example: a few questions

- Q1: What is the probability that an incoming request is rejected?
- Q2: What is the average number of requests in execution?
- Q3: What is the average throughput of the Web server?
- Q4: What is the average time spent by an HTTP request in the Web server?

7

© 1999–2000 D. A. Menascé. All Rights Reserved.

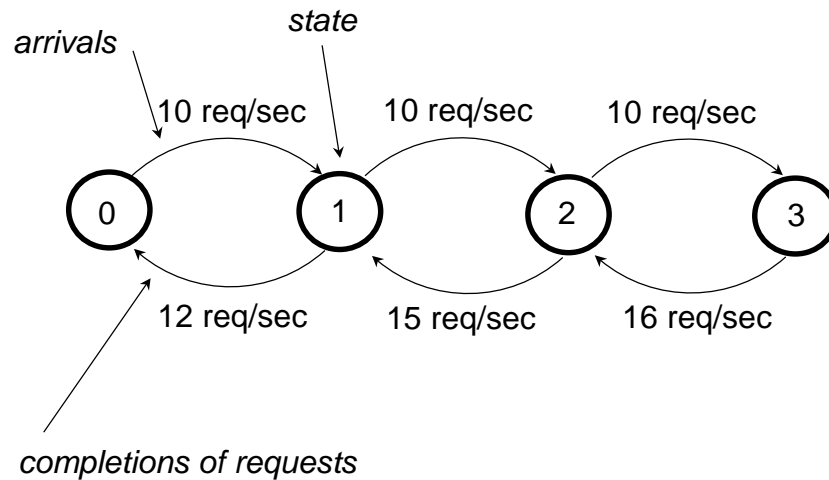
System-level Example

- Characterize the Web server by its state, i.e., the number k of requests in the Web server.
- Assumptions made:
 - homogeneous workload: all requests are equivalent
 - memoryless: how the system arrived at system k does not matter.
 - operational equilibrium: no. requests at beginning of interval = no. request at the end.

8

© 1999–2000 D. A. Menascé. All Rights Reserved.

System-level Example



© 1999–2000 D. A. Menascé. All Rights Reserved.

9

System-level Example

- Assume we are able to find the values of:
 - P_k = probability that there are k requests in the Web server.
- Question: can we answer all the questions posed before as a function of the P_k 's?

© 1999–2000 D. A. Menascé. All Rights Reserved.

10

System-level Example: a few questions

- Q1: What is the probability that an incoming request is rejected?
- A:

11

© 1999–2000 D. A. Menascé. All Rights Reserved.

System-level Example: a few questions

- Q1: What is the probability that an incoming request is rejected?
- A: It is the probability that an arriving HTTP request finds 3 requests already being processed. The answer is then P_3 .

12

© 1999–2000 D. A. Menascé. All Rights Reserved.

System-level Example: a few questions

- Q2: What is the average number of requests in execution?
- A: using the definition of average:

13

© 1999–2000 D. A. Menascé. All Rights Reserved.

System-level Example: a few questions

- Q2: What is the average number of requests in execution?
- A: using the definition of average:

$$n_{\text{req}} = 0 \times P_0 + 1 \times P_1 + 2 \times P_2 + 3 \times P_3$$

14

© 1999–2000 D. A. Menascé. All Rights Reserved.

System-level Example: a few questions

- Q3: What is the average throughput of the Web server?
- A: again, using the definition of average:

15

© 1999–2000 D. A. Menascé. All Rights Reserved.

System-level Example: a few questions

- Q3: What is the average throughput of the Web server?
- A: again, using the definition of average:

$$X = 0 \times P_0 + 12 \times P_1 + 15 \times P_2 + 16 \times P_3$$

throughput value at each state



16

© 1999–2000 D. A. Menascé. All Rights Reserved.

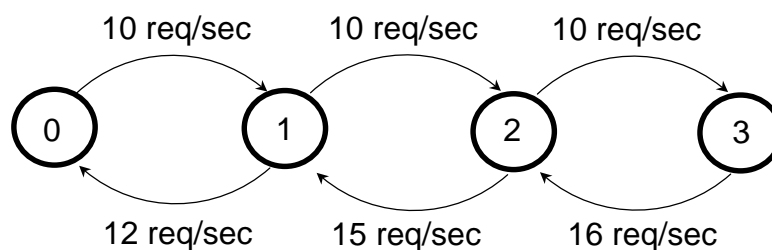
System-level Example: a few questions

- Q4: What is the average time spent by an HTTP request in the Web server?
- A: It will be a function of the average number of requests, n_{req} , and the average throughput X . More on this later...

17

© 1999–2000 D. A. Menascé. All Rights Reserved.

System-level Example: computing the P_k 's

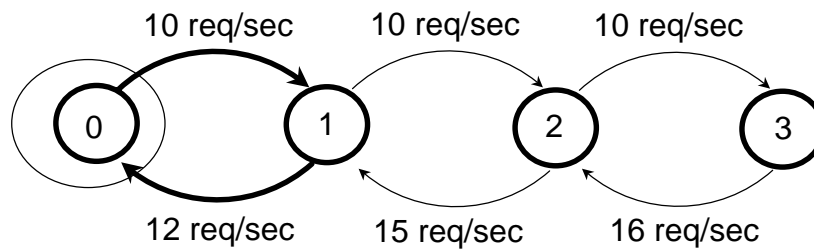


- use the flow in = flow out principle: the flow into a set of states is equal to the flow out of this set of states in equilibrium.

18

© 1999–2000 D. A. Menascé. All Rights Reserved.

System-level Example: computing the P_k 's



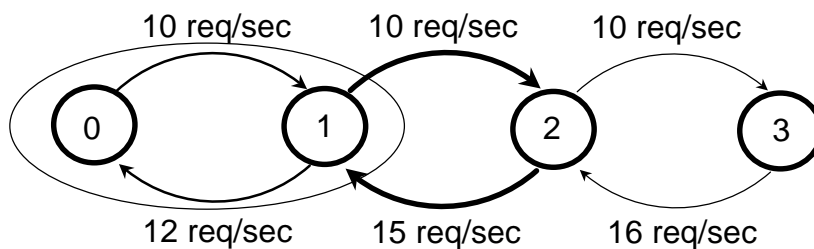
flow in = flow out

$$12 \times P_1 = 10 \times P_0$$

19

© 1999–2000 D. A. Menascé. All Rights Reserved.

System-level Example: computing the P_k 's



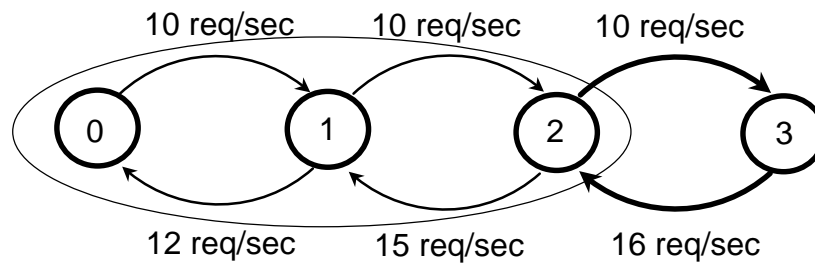
flow in = flow out

$$15 \times P_2 = 10 \times P_1$$

20

© 1999–2000 D. A. Menascé. All Rights Reserved.

System-level Example: computing the P_k 's



flow in = flow out

$$16 \times P_3 = 10 \times P_2$$

21

© 1999–2000 D. A. Menascé. All Rights Reserved.

System-level Example: computing the P_k 's

- Putting it all together:

$$12 \times P_1 = 10 \times P_0 \Rightarrow P_1 = 10/12 P_0$$

$$15 \times P_2 = 10 \times P_1 \Rightarrow P_2 = 10/15 P_1 \\ = \frac{10 \times 10}{15 \times 12} P_0$$

$$16 \times P_3 = 10 \times P_2 \Rightarrow P_3 = 10/16 P_2 \\ = \frac{10 \times 10 \times 10}{16 \times 15 \times 12} P_0$$

22

© 1999–2000 D. A. Menascé. All Rights Reserved.

System-level Example: computing the P_k 's

- Putting it all together:

$$P_1 = 10/12 P_0; \quad P_2 = \frac{10 \times 10}{15 \times 12} P_0; \text{ and}$$

$$P_3 = \frac{10 \times 10 \times 10}{16 \times 15 \times 12} P_0$$

- But, the Web server has to be in one of the four states at any time. So,
 $P_0 + P_1 + P_2 + P_3 = 1.$

23

© 1999–2000 D. A. Menascé. All Rights Reserved.

System-level Example: computing the P_k 's

- Solving for P_0 and then for the other P_k 's we get:

k	P_k
0	0.365
1	0.305
2	0.203
3	0.127

24

© 1999–2000 D. A. Menascé. All Rights Reserved.

System-level Example: answering the questions

- Q1: What is the probability that an incoming request is rejected?
- A: It is the probability that an arriving HTTP request finds 3 requests already being processed. The answer is then

$$P_3 = 0.127 = 12.7\%.$$

25

© 1999–2000 D. A. Menascé. All Rights Reserved.

System-level Example: answering the questions

- Q2: What is the average number of requests in execution?
- A: using the definition of average:

$$\begin{aligned}n_{\text{req}} &= 0 \times 0.365 + 1 \times 0.305 + \\ &\quad 2 \times 0.203 + 3 \times 0.127 \\ &= 1.091 \text{ requests}\end{aligned}$$

26

© 1999–2000 D. A. Menascé. All Rights Reserved.

System-level Example: answering the questions

- Q3: What is the average throughput of the Web server?
- A: again, using the definition of average:

$$\begin{aligned} X &= 0 \times 0.365 + 12 \times 0.305 + \\ &\quad 15 \times 0.203 + 16 \times 0.127 \\ &= 8.731 \text{ requests/sec.} \end{aligned}$$

27

© 1999–2000 D. A. Menascé. All Rights Reserved.

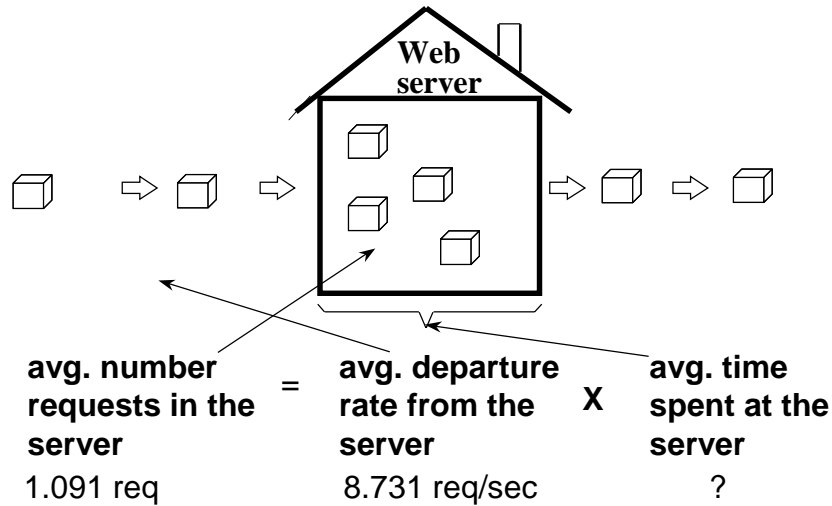
System-level Example: answering the questions

- Q4: What is the average time spent by an HTTP request in the Web server?
- A: It is a function of the average number of requests, n_{req} , and the average throughput X . We need Little's Law to answer this question.

28

© 1999–2000 D. A. Menascé. All Rights Reserved.

Little's Law



29

© 1999–2000 D. A. Menascé. All Rights Reserved.

System-level Example: answering the questions

- Q4: What is the average time spent by an HTTP request in the Web server?
- A: From Little's Law,

$$R = n_{\text{req}} / X = 1.091 / 8.731 = 0.125 \text{ sec.}$$

30

© 1999–2000 D. A. Menascé. All Rights Reserved.

Practice Drill

Using Models for Decision Making

- What happens if the maximum number of allowed TCP connections changes from 3 to 10?
- What if the load on the server doubles?
- What is the impact of a threefold increase in the server's capacity?

© 1999–2000 D. A. Menascé. All Rights Reserved.

Types of System-level Models

- **Population Size:**
 - infinite
 - finite
- **Service Rate:**
 - fixed
 - variable
- **Maximum Queue Size:**
 - unlimited
 - limited

© 1999–2000 D. A. Menascé. All Rights Reserved.

32

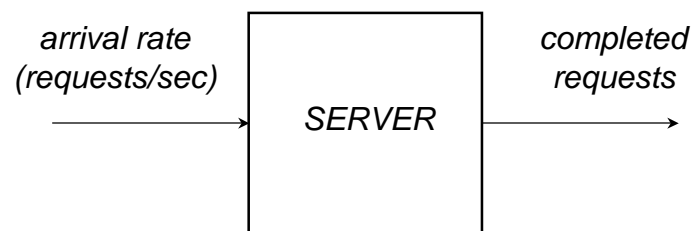
Types of System-level Models (population size)

- Infinite Population: the number of clients is very large. The rate at which requests arrive to the system does not depend on the number of requests in the system.
 - e.g., requests arriving from the Internet to a public Web server.

33

© 1999–2000 D. A. Menascé. All Rights Reserved.

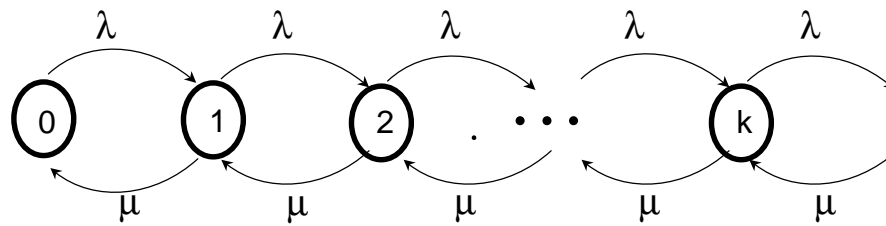
Types of System-level Models (infinite population)



34

© 1999–2000 D. A. Menascé. All Rights Reserved.

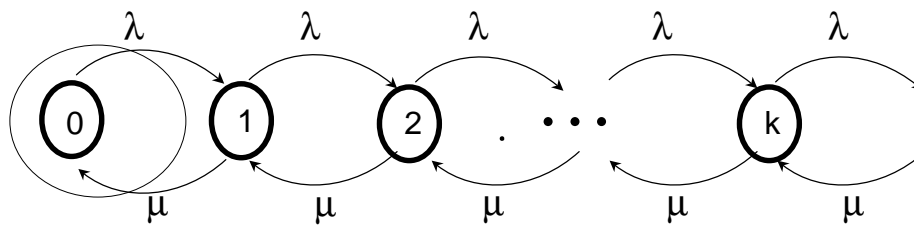
Infinite Population/Infinite Queue



flow in = flow out

© 1999–2000 D. A. Menascé. All Rights Reserved.

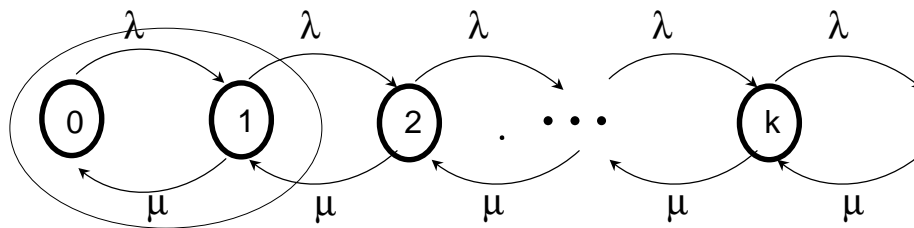
Infinite Population/Infinite Queue



flow out = flow in
 $\lambda P_0 = \mu P_1$

© 1999–2000 D. A. Menascé. All Rights Reserved.

Infinite Population/Infinite Queue

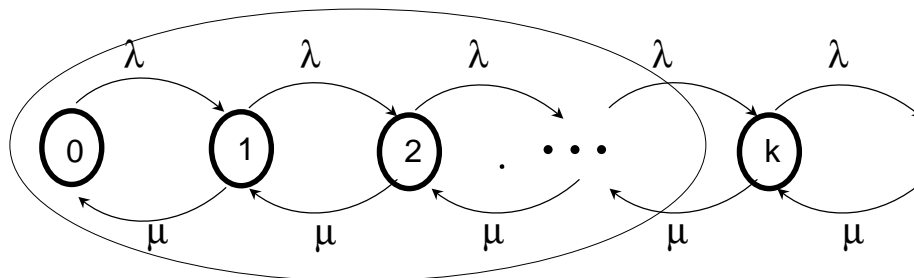


flow out = flow in

$$\begin{aligned}\lambda P_0 &= \mu P_1 \\ \lambda P_1 &= \mu P_2\end{aligned}$$

© 1999–2000 D. A. Menascé. All Rights Reserved.

Infinite Population/Infinite Queue



flow out = flow in

$$\begin{aligned}\lambda P_0 &= \mu P_1 \\ \lambda P_1 &= \mu P_2\end{aligned}$$

$$\dots$$

$$\lambda P_{k-1} = \mu P_k$$

© 1999–2000 D. A. Menascé. All Rights Reserved.

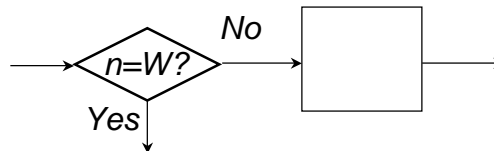
Infinite Population/Infinite Queue Example

- A DB server receives 30 req/sec. Each request takes 0.02 sec on the average. Find:
 - Fraction of requests in the DB server?
 - Average response time.
 - Average response time for a server twice as fast.
 - Average response time for a server twice as fast for twice the arrival rate.

© 1999–2000 D. A. Menascé. All Rights Reserved.

Types of System-level Models (maximum queue size)

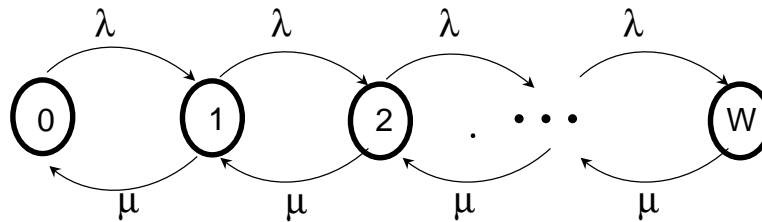
- Unlimited Queue Size: all arriving requests are queued for service. No requests are rejected!
- Limited Queue Size: requests that find more than W requests waiting for service are rejected.



© 1999–2000 D. A. Menascé. All Rights Reserved.

40

Infinite Population/Finite Queue



- arriving requests that find the server in state W are lost.

© 1999–2000 D. A. Menascé. All Rights Reserved.

Infinite Population/Finite Queue Example

- A DB server receives 30 req/sec. Each request takes 0.02 sec on the average. At most 4 request can be queued. Find:
 - Fraction of requests in the DB server?
 - Average response time.
 - Average response time for a server twice as fast.
 - Average response time for a server twice as fast for twice the arrival rate.

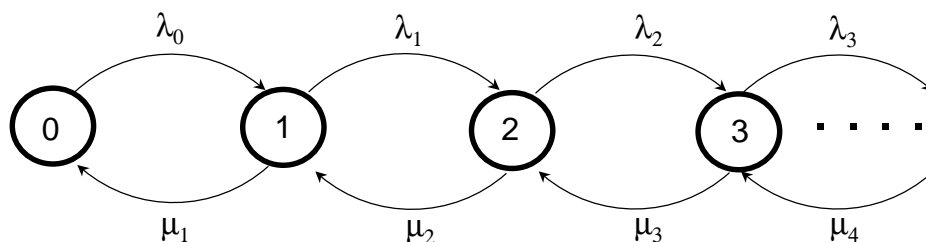
© 1999–2000 D. A. Menascé. All Rights Reserved.

Infinite Population/Finite Queue Example (cont'd)

- What is the maximum value for the maximum number of request queued so that less than 1% of the requests are rejected?

© 1999–2000 D. A. Menascé. All Rights Reserved.

Generalized System-level Models



*Generalized System-level Models can be solved
using the **flow in = flow out** principle!*

44

© 1999–2000 D. A. Menascé. All Rights Reserved.

Generalized System-level Models

$$p_k = \prod_{i=0}^{k-1} \frac{\lambda_i}{\mu_{i+1}}$$

$$p_0 = \left[\sum_{k=0}^{\infty} \prod_{i=0}^{k-1} \frac{\lambda_i}{\mu_{i+1}} \right]^{-1}$$

45

© 1999–2000 D. A. Menascé. All Rights Reserved.

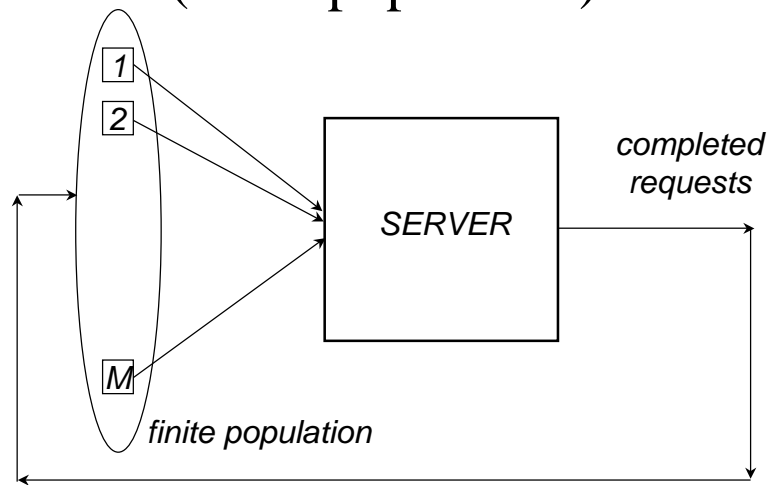
Types of System-level Models (population size)

- Finite Population: the number of clients is limited. The rate at which requests arrive to the system depends on how many have already arrived.
 - e.g., requests arriving to an intranet Web server from a known number of clients within the organization.

46

© 1999–2000 D. A. Menascé. All Rights Reserved.

Types of System-level Models (finite population)

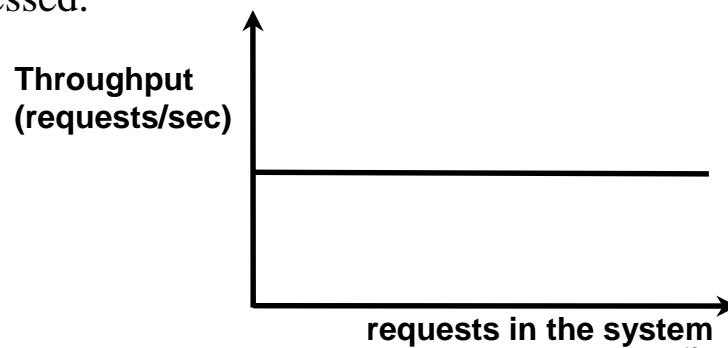


© 1999–2000 D. A. Menascé. All Rights Reserved.

47

Types of System-level Models (service rate)

- Fixed Service Rate: the throughput does not vary with the number of requests being processed.

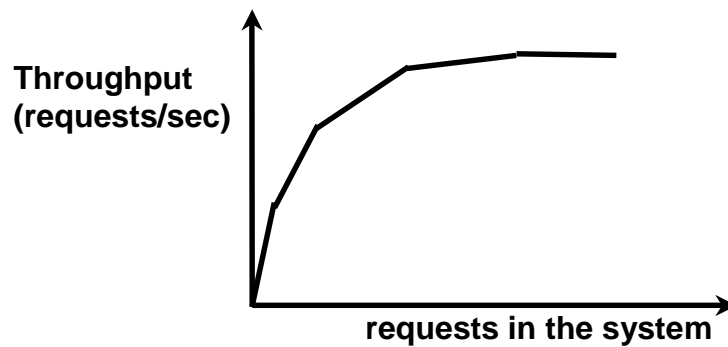


© 1999–2000 D. A. Menascé. All Rights Reserved.

48

Types of System-level Models (service rate)

- Variable Service Rate: the throughput depends on the number of requests being processed.



© 1999–2000 D. A. Menascé. All Rights Reserved.

49

Types of System-level Models

Population	Service Rate	Queue Size
infinite	fixed	unlimited
infinite	fixed	limited
infinite	variable	unlimited
infinite	variable	finite
finite	fixed	
finite	variable	

© 1999–2000 D. A. Menascé. All Rights Reserved.

50

System-level Models

Example

A Web server receives 30 requests/sec. Its throughput function is given below. The server queue is limited to five requests. What is the server utilization, avg. throughput, avg. no. requests, avg. response time, and fraction of lost requests?

No. of requests	Throughput (req/sec)
1	18
2	35
3 or more	50

51

© 1999–2000 D. A. Menascé. All Rights Reserved.

System-level Models

Example (cont'd)

Using the Generalized System-level model equations we get that

Average Number of Requests	1.850
Server Utilization	82.700%
Server Average Throughput	28.4 req/sec
Fraction of Lost Requests	0.05343

From Little's Law,

$$\begin{aligned}\text{avg. response time} &= \text{avg. no requests} / \\ &\quad \text{avg. throughput} = \\ &\quad 1.85 / 28.4 = 0.065 \text{ sec.}\end{aligned}$$

52

© 1999–2000 D. A. Menascé. All Rights Reserved.

System-level Models Example (cont'd)

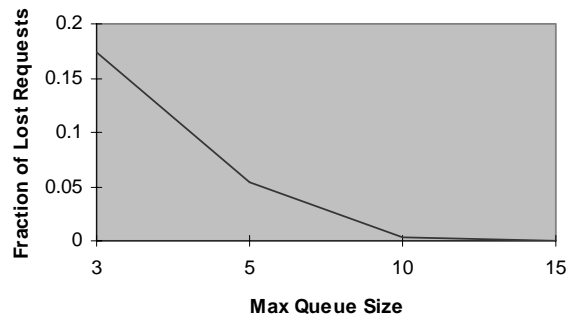
Max. Queue Size	3	5	10	15
Average Number of Requests	1.433	1.850	2.212	2.264
Server Utilization	79.81%	82.7%	83.9%	83.96%
Avg. Server Throughput (req/sec)	24.8	28.4	29.9	30
Average Response Time (sec)	0.058	0.065	0.074	0.075
Fraction of Lost Requests	0.173077	0.05343	0.003869	0.000299

53

© 1999–2000 D. A. Menascé. All Rights Reserved.

System-level Models Example (cont'd)

Fraction of Lost Requests vs. Max Queue Size



54

© 1999–2000 D. A. Menascé. All Rights Reserved.

Summary

System-level models view a server as a black box. Only its arrival process and throughput functions are relevant.

State Transition Diagrams (STDs) can be used to find the probability that k requests are in the server. Use the *flow in = flow out* principle.

Little's Law can be used to compute the response time from the average number of requests and from the throughput.

55

© 1999–2000 D. A. Menascé. All Rights Reserved.