

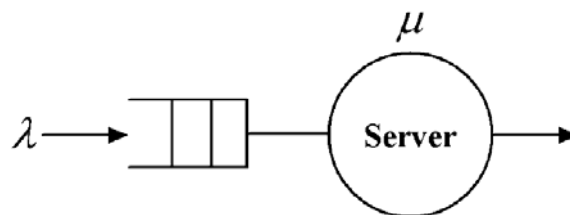
Analytic Modeling

Stochastic Model

1

Infinite Population Model

- Single server case



2

Notation

L = length of observation period

n = number of customers arrived and served in L

$\lambda = n/L$ = arrival rate

x_j = service time of j^{th} customer

$\bar{x} = \frac{1}{n} \sum_{j=1}^n x_j$ = mean service time

3

Performance Bounds

Model parameters

- mean interarrival time, $1/\lambda$
- mean service time, \bar{x}

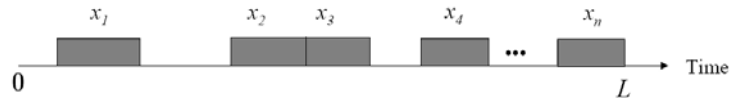
Performance measures

- utilization factor, U
- mean response time, \bar{T}
- throughput, R

4

Utilization Factor –Single Server Case

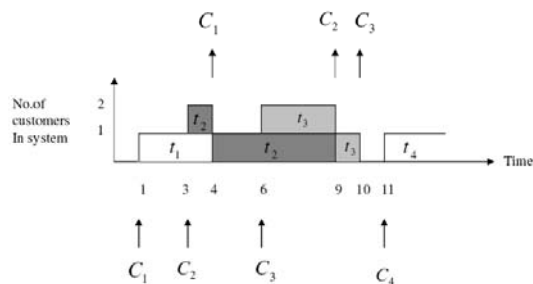
Consider the following case:



- total busy time = $\sum_{j=1}^n x_j$
- $U = \% \text{ of time server is busy} = \frac{1}{L} \sum_{j=1}^n x_j = \lambda \bar{x}$
- note: $U \leq 1$
- results can be extended to multiple servers

5

Mean Response Time (Little's Formula)



- Sum of response time = $\sum_{j=1}^n t_j = \text{total area under curve} = L\bar{N}$

- Thus, $\frac{n}{L} \frac{1}{n} \sum_{j=1}^n t_j = \bar{N}$ or

- $\lambda \bar{T} = \bar{N}$ - Little's Formula

6

Throughput

- $R = n/L = \lambda$ where n is the number of customers served in time L
- Since $\lambda = \frac{U}{\bar{x}}$, then $R = \frac{U}{\bar{x}}$

7

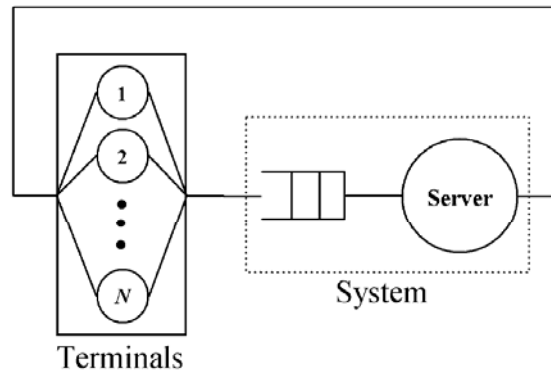
Some Properties

- Utilization factor, throughput, and Little's formula are not affected by scheduling discipline
- Little's formula can also be applied to an arbitrary queuing system (including queuing networks)
- Note: the fundamental results are based on the system returning to the idle state at time L , an indication that the system is stable

8

Finite Population Model

Single server case



9

Performance Bounds

Model parameters

- mean think time, \bar{h}
- mean service time, \bar{x}

Performance measures

- utilization factor, U
- mean response time, \bar{T}

10

Utilization Factor (Upper Bound)

$$\text{When } N = 1, U = \frac{\bar{x}}{\bar{x} + \bar{h}}$$

$$\text{Let } N^* = \frac{(\bar{x} + \bar{h})}{\bar{x}}$$

$$U \leq \begin{cases} N \frac{\bar{x}}{\bar{x} + \bar{h}} & 1 \leq N \leq N^* \\ 1 & N \geq N^* \end{cases}$$

11

Mean Response Time (Lower Bound)

- Arrival rate to system (queue and server) = arrival rate to user terminals = λ
- Apply Little's formula to system
$$\lambda \bar{T} = \bar{N}(s) = \text{mean number in system}$$
- Apply Little's formula to user terminals
$$\lambda \bar{h} = \bar{N}(t) = \text{mean number in thinking state}$$
- Since $\bar{N}(s) + \bar{N}(t) = N$, $\lambda(\bar{T} + \bar{h}) = N$ or $\bar{T} = \frac{N}{\lambda} - \bar{h}$
- Since $U = \lambda \bar{x} \leq 1$, $\bar{T} \geq N\bar{x} - \bar{h}$ (lower bound)

12

Summary

- Little's formula
 - mathematical relationship between λ , \bar{N} , and \bar{T} , two of them are performance measures
 - derivation is based on infinite population model, but result is application to other models, e.g., finite population model, queueing networks
- Infinite population model
 - exact analytic result for utilization factor is available and is a function of input parameters

13

Summary

- Finite population model
 - exact analytic results are generally not available
 - performance bounds are available -upper bound for utilization and lower bound for mean response time
 - performance bounds are functions of input parameters
- In general, exact analytic results are available if certain assumptions are made on input parameters such as interarrival time, service time and think time
 - for the case of exponential distribution, results are available for infinite population model, finite population model and queueing networks
 - for non-exponential distributions, selected cases where analytic results are available

14