1

ECE560: Computer Systems Performance Evaluation



Lecture 14 -Embedded Markov-Chain Queueing Systems & Priority Queueing Systems

Instructor: Dr. Liudong Xing

## Administration Issues (3/27/2024)

- Homework #5 assigned
  - Due April 1, Monday
- Project final report
  - Due: April 19, Friday
  - Refer to Report Guidelines



### Agenda

- Embedded Markov Chain Queueing Systems
  - M/G/1,
  - M/D/1,
  - GI/M/1
- Priority Queueing Systems

### Embedded Markov Chain Queueing Systems

- More general arrival process or service times are allowed
  - <u>M/G/1</u>,
  - M/D/1,
  - GI/M/1
- Solution (see extra notes posted for detailed derivations):
  - Constructing an embedded Markov chain
  - And applying z-transform and Laplace-Stieltjes transform methods

### M/G/1 Queueing Systems

- Assume
  - Poisson arrival process with rate  $\lambda$
  - General service time distribution with
    - different customers have independent service times
    - E[s] and E[s<sup>2</sup>] exist ( in order to calculate L,W)

5

### Steady-State State Probability

• Steady-state probability distribution:

$$\begin{cases} \pi_0 = 1 - \rho \\ \pi_i = \pi_0 k_i + \sum_{j=1}^{i+1} \pi_j k_{i-j+1} \end{cases}$$

where  $k_n = \Pr[n \text{ customers arrive during one service interval}]$ 

$$k_n = \int_0^\infty e^{-\lambda t} \frac{(\lambda t)^n}{n!} dW_s(t), \quad n = 0, 1, 2, ...$$

 $W_s(t)$ : cdf of service time s

### Other state probabilities

- π<sub>i</sub>: steady state probability that a departing customer leaves *i* customers behind π<sub>i</sub> = Pr[X=i]
- *p<sub>i</sub>*: probability that there are *i* customers in the system at arbitrary times - *p<sub>i</sub>* = *Pr[N=i]*
- $r_i$ : probability that an arriving customer finds *i* customers already in the system
- It can be shown (by Klienrock) that for M/G/1 systems:

 $\pi_i = p_i = r_i$ 

7















### Renewal Processes (Chapter 4.5)

- A Poisson process can be characterized as a counting process for which the inter-arrival times (times between successive events) are *i.i.d*, exponential *r.v.*s
- A renewal process is a generalization of the Poisson process
  - A counting process for which the interarrival times (times between successive events) are *i.i.d* r.v.s

GI/M/1







## $p_n$ for GI/M/1

- Distinction between
  - $\pi_n$ : the steady-state probability that an arriving customer finds *n* customers in the system
    - From "an arriving customer" point of view
  - p<sub>n</sub>: the steady-state probability that there are *n* customers in the system
    - From "a random observer" point of view









### Hands-On Problem

Considering a computer subsystem that can be modeled as the GI/M/1 queueing system. Specifically, the service time is exponentially distributed with a constant rate of 60 jobs per second. The Laplace-Stieltjes transform of the job inter-arrival time  $\tau$  is assumed to be  $A^*[\theta] = \frac{\theta + \lambda}{3\lambda}$  with  $\lambda = 30$ .

- What is the probability that an arriving job finds the system is busy?
- What is the probability that an arriving job finds 3 customers in the system?
- What is the average number of jobs in the system queue?
- What is the average time a job spends in the system?

GI/M/1

Dr.	Xing	©	
-----	------	---	--

# Agenda • Embedded Markov Chain Queueing Systems - M/G/1, - GI/M/1 • Priority Queueing Systems

### Kendall Notation (review) Standard notation for queueing systems: A/B/c/K/m/Z• A: arrival process or inter-arrival time distribution - 'M' = Poisson arrival process - 'D' = Deterministic (constant) arrival rate - 'G' = General arrival process • **B**: service process or service time dist. - 'M' = Exponential service time dist. - 'D' = Deterministic (constant) service time - 'G' = General service time • c: number of servers • K: the capacity of the system (queue+server(s)) (default: $\infty$ ) • m: total job/customer population (default: ∞) • Z: scheduling discipline (default: FIFO) Dr. Xing © GI/M/1 26







### **Control Policies**

• Control policies to resolve the situation wherein a customer of class *i* arrives to find a customer of class *j* in service (i < j)

### • Non-preemptive

- The newly arrived customer always waits until the customer in service completes service before gaining access to the service facility  $\rightarrow$  a head-of-the-line (HOL) system
- GI/M/1 30

### Control Policies (Cont'd)

#### • Preemptive

- service of the customer of class *j* is interrupted, and the newly arrived customer of higher priority begins service
- The interrupted customer returns to the head of the *j*th class
- The interrupted customer resumes the service at the point of interruption *preemptive-resume*, or
- The interrupted customer repeats the entire service from the beginning -preemptive-repeat

Dr. Xing ©

GI/M/1

31

### M/G/1 Priority Q Systems

 Consider an M/G/1 queueing system with any queue discipline that chooses customers by an algorithm that does not consider customer service times or any measure of them. Then the performance measures L, W, Lq, Wq will be the same as for the FCFS queue discipline (shown by Kleinrock)

$$L = \rho + \frac{\rho^2 (1 + C_s^2)}{2(1 - \rho)}, \text{ where } C_s = \frac{\sqrt{Var[s]}}{E[s]}$$

$$W = L/\lambda = W_s + \frac{\rho W_s (1 + C_s^2)}{2(1 - \rho)}$$

$$L_q = L - \rho = \frac{\rho^2 (1 + C_s^2)}{2(1 - \rho)}$$

$$W_q = L_q/\lambda = \frac{\rho^2 (1 + C_s^2)}{2(1 - \rho)\lambda} = \frac{\rho W_s (1 + C_s^2)}{2(1 - \rho)}$$
Dr. Xing © GI/M/1 32

### M/M/c Priority Q Systems • Also called <u>multi-server priority systems</u> • Consider M/M/c non-preemptive systems with *n* priority classes - The customers arrival to class *i* is Poisson with rate $\lambda_i$ . Then the overall arrival pattern is Poisson with mean (superposition property) $\lambda = \lambda_1 + \lambda_2 + \dots + \lambda_n$ - Each customer has the same exponential service time requirement with a mean of $1/\mu$ - Cobham showed that $\rho = \frac{\lambda}{\mu c} = \frac{\lambda W_s}{c}$ $E[q_1] = \frac{C[c,\alpha]W_s}{c(1-\lambda_1 W_s/c)}$ Ref.(M/M/c with FCFS in L#13): $W_q = E[q] = \frac{C[c, \alpha]W_s}{c(1-\rho)}$ $C[c,\alpha]W_s$ $E[q_i] = -$ , j = 2, ..., nDr. Xing © GI/M/1 33

