

Fundamentals of Computer Networking: Introduction to Queuing Theory

Guevara Noubir
CSG150, lecture 3

Reading: Textbook chapter 3.

F2003, CSG150

Fundamentals of Computer Networking

Lecture 3, 1

Outline

- Delay Models
- Little's Theorem
- The $M/M/1$ queuing system
- The $M/G/1$ queuing system

F2003, CSG150

Fundamentals of Computer Networking

Lecture 3, 2

Delay Models

- *Delay* (or *latency*) of data packet is an important measure of the performance of a network

Delay = *PropagationDelay* + *TransmissionDelay* + *QueuingDelay*

PropagationDelay = *Distance/SpeedOfLight* (independent of message size)

TransmissionDelay = *MessageSize/Bandwidth* (Bandwidth = data-rate here)

QueuingDelay = delay due to time spent waiting in queues (**most important delay**)

- The queuing delay depends on several parameters:
 - Arrival process
 - Service discipline
 - Processing delay
 - Others: bandwidth of the link, buffer size

F2003, CSG150

Fundamentals of Computer Networking

Lecture 3, 3

Queuing Theory Framework

- Queuing system:
 - Servers (one or several): e.g., router, computer processor
 - Customers: e.g., users, packets, program tasks
 - Queues: customers wait in queues before getting services
- Assumption: queues are unbounded.

F2003, CSG150

Fundamentals of Computer Networking

Lecture 3, 4

Little Theorem (1961)

- Measurement quantities of interest:
 - T : average delay incurred by a customer
 - N : average number of customers in the system
- Little's Theorem:
 - $N = \lambda T$ where λ is rate of the arrival process
- Little's Theorem provides a general and fundamental relation between N , T , and λ . It is independent of the nature of the arrival process or of the service time distribution.

F2003, CSG150

Fundamentals of Computer Networking

Lecture 3, 5

Proof of Little's Theorem

- Notation:
 - $\alpha(t)$: number of users that arrived before time t
 - $\beta(t)$: number of users that departed before time t
 - T_i time spent by user i within the system
 - $N(t)$ number of users in the system at time t
- Arrival rate: $\lambda(t) = \frac{\alpha(t)}{t}$
 - λ is the limit
- Average time within the system: $T(t) = \frac{1}{\alpha(t)} \sum_{i=0}^{\alpha(t)} T_i$
- Average number of users at time t : $N_t = \frac{1}{t} \int_0^t N(\tau) d\tau$

F2003, CSG150

Fundamentals of Computer Networking

Lecture 3, 6

Proof (Cont'd)

- The usage of the system can be bounded:

$$\sum_{i=0}^{\beta(t)} T_i \leq \int_0^t N(\tau) d\tau \leq \sum_{i=0}^{\alpha(t)} T_i$$

$$\frac{\beta(t)}{t} \frac{\sum_{i=0}^{\beta(t)} T_i}{\beta(t)} \leq \frac{1}{t} \int_0^t N(\tau) d\tau \leq \frac{\alpha(t)}{t} \frac{\sum_{i=0}^{\alpha(t)} T_i}{\alpha(t)}$$

- Taking the limit when $t \rightarrow +\infty$

$$\lambda T \leq N \leq \lambda T$$

F2003, CSG150

Fundamentals of Computer Networking

Lecture 3, 7

Applications: Flow Control

- Sliding window flow control (e.g., Go-Back-N or Selective Repeat) with window size: W
- The number of packets in the system is always less than W : $\lambda T = N \leq W$
- Conclusion:
 - for a given window size, if T increases, then the arrival rate has to be decreased
 - for a given arrival rate, if T increases, then the window size has to be increased
 - for a given T , if the arrival rate increases, then the window size has to be increased

F2003, CSG150

Fundamentals of Computer Networking

Lecture 3, 8

Application: Throughput in a Time Sharing System

- Assumptions:
 - N terminals, one processor, one queue
 - Terminals are always occupied
 - System activity: users log-on, reflection (R on average), submit task to the processor, tasks are queued, tasks execution takes on average P units of time
- The delay for a user task is on average T s.t.:
$$R+P \leq T \leq R + NP$$

Using Little's Theorem: $\frac{N}{R+NP} \leq \lambda \leq \frac{N}{R+P}$

F2003, CSG150

Fundamentals of Computer Networking

Lecture 3, 9

Time Sharing (Cont'd)

- The processor is also a queuing system where $N \leq 1$
- In the steady state mode: the arrival rate in the system is the same as for the processor
- Using Little's Theorem a second time: $\lambda P \leq 1$

- Combining these two bounds we get:

$$\lambda \leq \min\left\{\frac{1}{P}, \frac{N}{R+P}\right\} = \frac{1}{P} \min\left\{1, \frac{N}{1+R/P}\right\}$$

- The smallest term indicates the bottleneck

F2003, CSG150

Fundamentals of Computer Networking

Lecture 3, 10

The $M/M/1$ queuing system

- Little's Theorem is a general tool that allows us to calculate the steady-state average delay of a queuing system
- Notation:
 - M : memoryless, G : general, D : deterministic, 1: number of servers in the system
- $M/M/1$:
 - Arrival rate is Poisson distributed
 - Service time is exponentially distributed
 - These two processes are independent

F2003, CSG150

Fundamentals of Computer Networking

Lecture 3, 11

Poisson Process

- A Poisson process with arrival rate λ :

- The probability distribution function (pdf):

$$\Pr(n \text{ arrivals in interval } [t, t + \tau]) = \frac{e^{-\lambda\tau} (\lambda\tau)^n}{n!}$$

- The arrival distribution of two disjoint intervals is independent

- Properties:

- expected number of arrivals in a length- τ interval is: $\lambda\tau$.

F2003, CSG150

Fundamentals of Computer Networking

Lecture 3, 12

Poisson Process (Cont'd)

- Probabilities for small intervals:
 - $\Pr(0 \text{ arrival}) = e^{-\lambda\delta}/0! = 1 - \lambda\delta + \alpha(\delta)$
 - $\Pr(1 \text{ arrival}) = \lambda\delta e^{-\lambda\delta}/1! = \lambda\delta + \alpha(\delta)$
 - $\Pr(2 \text{ arrivals}) = (\lambda\delta)^2 e^{-\lambda\delta}/2! = \alpha(\delta)$If δ tends to 0, then we have $\Pr(0 \text{ arrivals}) = 1 - \lambda\delta$ and $\Pr(1 \text{ arrival}) = \lambda\delta$
- Inter-arrival times:
 - Let t_n be the arrival time of the n^{th} customer and $\tau_n = t_{n+1} - t_n$
 - Then: $\Pr(\tau_n > s) = e^{-\lambda s}$ (\Rightarrow exponential distribution)

F2003, CSG150

Fundamentals of Computer Networking

Lecture 3, 13

Other Properties

- Poisson processes are used to model the traffic of a large number of similar and independent users
- If n independently and identically distributed packet arrival processes (rate λ/n) occur at the head of a link then the aggregated process can be shown to be well approximated by a Poisson process of rate λ . n is considered to be a large value.
- The aggregation of k independent Poisson processes of rates $\lambda_1, \lambda_2, \dots, \lambda_k$ yields a Poisson process of rate: $\lambda_1 + \lambda_2 + \dots + \lambda_k$

F2003, CSG150

Fundamentals of Computer Networking

Lecture 3, 14

Exponential Service Time

- Let s_n denote the service time for the n^{th} customer. The service time distribution is exponential with parameter μ if:
 $\Pr[s_n \leq s] = 1 - e^{-\mu s}$
- The expected service time for a job is: $1/\mu$
- The exponential service time is memoryless in the sense that:
 $\Pr(s_n > r + t \mid s_n > t) = \Pr(s_n > r)$
- Poisson processes are closely related to exponential distributions: inter-arrival times of a Poisson process with rate λ have an exponential distribution with parameter λ .

F2003, CSG150

Fundamentals of Computer Networking

Lecture 3, 15

Analysis of the M/M/1 Queuing System

- The state of the system is captured by the number of customers in the system at time t
- We consider a discrete version of the process evolution:
 - Time: $0, \delta, 2\delta, 3\delta, \dots, k\delta, \dots$
 - N_k : number of customers at time $k\delta$
- Properties:
 - $\Pr[N_{k+1} = l \mid N_k = l] = \sum_{i \geq 0} \Pr[i \text{ arrivals and } i \text{ departures in } \delta \text{ interval}]$
 - $\Pr[N_{k+1} = 0 \mid N_k = 0] \approx 1 - \lambda\delta + (\lambda\delta)(\mu\delta) \approx 1 - \lambda\delta$

F2003, CSG150

Fundamentals of Computer Networking

Lecture 3, 16

M/M/1 Analysis (Cont'd)

- Let: $P_{i,j} = \Pr[N_{k+1}=j \mid N_k = i]$
- $P_{0,0} = 1 - \lambda\delta + \alpha(\delta) \approx 1 - \lambda\delta$
- $P_{i,i} = 1 - \lambda\delta - \mu\delta + \alpha(\delta) \approx 1 - \lambda\delta - \mu\delta$ (for $i \geq 1$)
- $P_{i,i+1} = \lambda\delta + \alpha(\delta) \approx \lambda\delta$ (for $i \geq 0$)
- $P_{i,i-1} = \mu\delta + \alpha(\delta) \approx \mu\delta$ (for $i \geq 1$)
- $P_{i,j} = \alpha(\delta) \approx \alpha(\delta)$ (for $j \neq i, i+1, i-1$)
- The state transitions represent a Markov chain

F2003, CSG150

Fundamentals of Computer Networking

Lecture 3, 17

Stationary Distribution of a System

- After a long period of time the system reaches a *steady state*
- Let: $p_i = \lim_{k \rightarrow +\infty} \Pr[N_k = i]$
- From the Markov chain diagram we have:
 - $p_i = p_{i-1}(\lambda\delta) + p_i(1 - \lambda\delta - \mu\delta) + p_{i+1}(\mu\delta)$
 - Hence: $(p_i - p_{i+1}) = \rho(p_{i-1} - p_i)$, where $\rho = \lambda/\mu$
 - Let $\Delta_i = p_i - p_{i+1}$, then $\Delta_i = \rho\Delta_{i-1}$ (for $i > 0$)
 - We also have: $\Delta_0 = (1 - \rho)p_0$

F2003, CSG150

Fundamentals of Computer Networking

Lecture 3, 18

Stationary Distribution of a System

$$p_i = p_0 - \sum_{j=0}^{i-1} \Delta_j$$

$$p_i = p_0 - (1 - \rho)p_0 \frac{1 - \rho^i}{1 - \rho}$$

$$p_i = \rho^i p_0$$

- Since: $\sum_{i \geq 0} p_i = 1$, then $p_0 = 1 - \rho$, and $p_i = \rho^i (1 - \rho)$

F2003, CSG150

Fundamentals of Computer Networking

Lecture 3, 19

Steady State Averages

- Steady state average number of customers:

$$\sum_{i=0}^{\infty} i p_i = \sum_{i=0}^{\infty} i \rho^i (1 - \rho) = \frac{\rho}{1 - \rho}$$

- Average delay T (using Little's Th.):

$$T = \frac{1}{\mu - \lambda}$$

- Average waiting time W :
(delay-service time)

$$W = \frac{\lambda}{\mu(\mu - \lambda)}$$

F2003, CSG150

Fundamentals of Computer Networking

Lecture 3, 20

Applications

- Scaling up the arrival rate and service rate
 - If we increase the arrival and service rates by the same factor then average number of customers in the system stays the same, while the average delay goes down
- Multiplexing several connections on one link
 - Benefit of statistical multiplexing

F2003, CSG150

Fundamentals of Computer Networking

Lecture 3, 21

App 1: Blocking Probability

- Consider a queuing system with:
 - K servers
 - $N \geq K$ in system customers (in service + waiting)
 - Departing customers are immediately replaced by new customers
 - \bar{X} is the average customer service time
- Average customer time in the system T ?
 - $T = N/\lambda$ and $\bar{K} = \lambda \bar{X}$
 - Thus: $T = N \bar{X} / K$

F2003, CSG150

Fundamentals of Computer Networking

Lecture 3, 22

App 1: Blocking Probability (Cont'd)

- Assume that customers are blocked (and lost) if the system is full:
 - β is the proportion of customers that are blocked
 - The system may go through moments where less than K servers are active
 - Then:

$$\bar{K} = (1 - \beta)\lambda\bar{X}$$
$$\beta = 1 - \frac{\bar{K}}{\lambda\bar{X}} \geq 1 - \frac{K}{\lambda\bar{X}}$$

F2003, CSG150

Fundamentals of Computer Networking

Lecture 3, 23

App 2: Terminal Concentrator

- Consider a terminal concentrator:
 - 4 input lines, each line of 64 kbps
 - 1 output line of 128 kbps
 - Mean packet size is 12800 bits
 - Each of the four input lines delivers Poisson traffic with $\lambda_i = 2$ pkts/s
- Mean delay of a packet within the concentrator:
 - $\lambda = 8$ pkts/s, $\mu = 10$ pkts/s, $T = 1/(\mu - \lambda) = 500$ ms
- Average number of packets within the concentrator:
 - $N = \rho/(1 - \rho) = 4$

F2003, CSG150

Fundamentals of Computer Networking

Lecture 3, 24

App 2: Terminal Concentrator (Cont'd)

- Remarks:
 - The output line is capable of handling the generated traffic (128Kbps > 12800 * 8), but a substantial input queue builds up.
 - The reason is the randomness of the arrivals
- Usefulness of modeling and analysis:
 - Delay estimation
 - Buffers dimensioning

F2003, CSG150

Fundamentals of Computer Networking

Lecture 3, 25

App 3: Statistical Multiplexing vs. Dedicated Channels

- Let a system consist of:
 - Two computers connected using a 64Kbps line
 - 8 parallel sessions
 - Each session generates Poisson traffic with $\lambda_i = 2$ pkts/s
 - Packets length is exponentially distributed with mean 2000 bits.
- Two possible strategies:
 - Give each session a dedicated portion of the channel (e.g. TDM or FDM)
 - Have all the packets compete for the shared channel

F2003, CSG150

Fundamentals of Computer Networking

Lecture 3, 26

App 3: Statistical Multiplexing vs. Dedicated Channels (Cont'd)

- Dedicated channels (8*8Kbps):
 - $\lambda = 2$ pkts/s, $\mu = 4$ pkts/s, $T = 1/(\mu - \lambda) = 500$ ms
- Statistical multiplexing:
 - $\lambda = 16$ pkts/s, $\mu = 32$ pkts/s, $T = 1/(\mu - \lambda) = 62.5$ ms
- Explanation: because of the randomness of the arrival rate, some of the dedicated channel may be unused (because the corresponding session is idle) while packets are queued for other sessions

F2003, CSG150

Fundamentals of Computer Networking

Lecture 3, 27

The *M/G/1* System

- *M/G/1* system:
 - Arrival rate is Poisson
 - Service time has a general distribution
- It is not possible to derive a closed-form stationary distribution (as in *M/M/1*) but we can derive other results
- Assume that:
 - Customers are served on a FCFS basis
 - X_i (service time of i^{th} arrival) identically distributed, mutually independent, and independent of the inter-arrival times

F2003, CSG150

Fundamentals of Computer Networking

Lecture 3, 28

P-K Formula

- Average service time: $\bar{X} = E\{X\} = \frac{1}{\mu}$
- Second moment of service time: $\overline{X^2} = E\{X^2\}$
- *Pollaczek-Khinchin* (P-K) formula:

$$W = \frac{\lambda \overline{X^2}}{2(1-\rho)}$$
- Then: $T = \bar{X} + \frac{\lambda \overline{X^2}}{2(1-\rho)}$
- Using Little's Theorem: $N_Q = \frac{\lambda^2 \overline{X^2}}{2(1-\rho)}$; $N = \rho + \frac{\lambda^2 \overline{X^2}}{2(1-\rho)}$

F2003, CSG150

Fundamentals of Computer Networking

Lecture 3, 29

Verification of P-K Formula for Exponentially Distributed Service Time

- When service times are exponentially distributed as in the *M/M/1* system:

$$\bar{X} = 1/\mu; \overline{X^2} = 2/\mu^2$$

$$W = \frac{\rho}{\mu(1-\rho)}; T = \frac{1}{\mu - \lambda}$$

- When the service time is identical for all customers:

$$M/D/1: \bar{X} = 1/\mu; \overline{X^2} = 1/\mu^2$$

$$W = \frac{\rho}{2\mu(1-\rho)}$$

M/D/1 provides lower bounds for W , T , N_Q and N

F2003, CSG150

Fundamentals of Computer Networking

Lecture 3, 30

Proof of the P-K Formula

- We use the concept of *mean residual service time*
- Notation:
 - W_i : waiting time of customer i
 - R_i : residual time to completion of the current customer at instant when i arrives ($R_i=0$, if no customer is being serviced)
 - Q_i : number of customers waiting in queue when i arrives
- Since customers are serviced in order, we have:

$$W_i = R_i + \sum_{j=i-Q_i}^{i-1} X_j$$

$$W = R + N_Q \bar{X}$$

F2003, CSG150

Fundamentals of Computer Networking

Lecture 3, 31

Proof of the P-K Formula (Cont'd)

- From Little's Theorem: $N_Q = W\lambda$, then: $W = R/(1-\rho)$

$$R = \frac{1}{t} \int_0^t R(\tau) d\tau$$

$$R = \frac{1}{t} \sum_{i=1}^{\beta(t)} \frac{1}{2} X_i^2$$

$$R = \frac{\lambda \overline{X^2}}{2}$$

- Thus the P-K formula

F2003, CSG150

Fundamentals of Computer Networking

Lecture 3, 32

Why Poisson Assumption?

- Where did we use Poisson Process arrivals assumption?
 - At the moment when a packet arrives the queue is typical
 - $\lim P\{N(t) = n \mid \text{an arrival occurred just after } t\} = \lim \{N(t) = n\}$
 - Section 3.3.2
 - If arrival not Poisson:
 - Inter-arrival: uniformly distributed between 2 and 4 seconds
 - Customer service time is: 1 second
 - => An arriving customer finds the queue empty
 - => but an external customer sees a average queue length of 1/3

F2003, CSG150

Fundamentals of Computer Networking

Lecture 3, 33

Unstable $M/G/1$ Systems

- For several probability distribution functions the second moment is finite (proportional to the square of the mean): e.g., exponential, constant, uniform. However, it is not general to all distributions.
- Let X be the random variable representing the service time for a customer s.t.:
 - $\Pr[X=1] = 2/3$; $\Pr[X=2^i]=1/4^i$ (for $i>0$)
 - The mean of X is finite, but the second moment is infinite
 - In this kind of systems we may have an accumulations of arrivals that exceeds the service capability

F2003, CSG150

Fundamentals of Computer Networking

Lecture 3, 34

Applications of $M/G/1$: GBN ARQ

- Simplified analysis of Go-Back-n ARQ:
 - No-modulus, all acknowledgements are received
 - If the lowest number in the window is not acked by the end of the window the sender assumes that the error occurred and starts retransmitting
 - Errors are independent from one to another
 - All frames take a unit of time to be transmitted
- The service time distribution is:
 - $\Pr[X=1+ni] = p^i(1-p)$ ($i \geq 0$)

F2003, CSG150

Fundamentals of Computer Networking

Lecture 3, 35

Applications of $M/G/1$: GBN ARQ

- If the packets are generated at the sender by a Poisson process, then we have an $M/G/1$ system:

$$\bar{X} = 1 + \frac{np}{1-p} \quad W = \frac{\lambda \bar{X}^2}{2(1-\lambda \bar{X})}$$

$$\bar{X}^2 = 1 + \frac{2np}{1-p} + \frac{n^2(p+p^2)}{(1-p)^2} \quad T = \bar{X} + W$$

Formulas good to know:

$$\sum_{k=0}^{\infty} p^k = \frac{1}{1-p}, \quad \sum_{k=0}^{\infty} kp^k = \frac{p}{(1-p)^2}, \quad \sum_{k=0}^{\infty} k^2 p^k = \frac{p+p^2}{(1-p)^3}$$

F2003, CSG150

Fundamentals of Computer Networking

Lecture 3, 36

M/G/1 with Priorities

- System with priority:
 - Customers are divided into classes: 1 ... k
 - Customers in class i are given priority over customers of class j (for any $j > i$)
 - Non-preemptive:
 - Customers are served in their order of arrival
- Notation:
 - Arrival process for class i : Poisson with rate λ_i
 - Service time of customers of class i : X_i
 - W_i average waiting time for a customer in class i
 - R average residual time
 - Q_i average number of customers of class i waiting in queue

F2003, CSG150

Fundamentals of Computer Networking

Lecture 3, 37

M/G/1 with Priorities (Cont'd)

Class 1	Class 2
$W_1 = R + Q_1 \bar{X}_1$	$W_2 = R + Q_1 \bar{X}_1 + Q_2 \bar{X}_2 + \lambda_1 W_2 \bar{X}_1$
$Q_1 = W_1 \lambda_1$	$Q_i = W_i \lambda_i$
$\Rightarrow W_1 = \frac{R}{1 - \rho_1}$	$\Rightarrow W_2 = \frac{R}{(1 - \rho_1)(1 - \rho_1 - \rho_2)}$
Class i	
$W_i = \frac{R}{(1 - \rho_1 - \dots - \rho_{i-1})(1 - \rho_1 - \dots - \rho_k)}$	

F2003, CSG150

Fundamentals of Computer Networking

Lecture 3, 38

M/G/1 with Priorities (Cont'd)

- As for the P-K formula: $R = \frac{1}{2} \sum_{i=1}^k \lambda_i \bar{X}_i^2$
- Thus the average waiting time for a customer in class i :

$$W_i = \begin{cases} \frac{\sum_{i=1}^k \lambda_i \bar{X}_i^2}{2(1 - \rho_1)} & \text{if } i = 1 \\ \frac{\sum_{i=1}^k \lambda_i \bar{X}_i^2}{2(1 - \rho_1 - \dots - \rho_{i-1})(1 - \rho_1 - \dots - \rho_i)} & \text{if } i > 1 \end{cases}$$

F2003, CSG150

Fundamentals of Computer Networking

Lecture 3, 39

M/M/m Markov System

- Steady state probabilities:

$$p_n = \begin{cases} p_0 \frac{(m\rho)^n}{n!} & n \leq m \\ p_0 \frac{m^m \rho^n}{m!} & n > m \end{cases} \quad \rho = \frac{\lambda}{m\mu} < 1$$

$$p_0 = \left[\sum_{n=0}^{m-1} \frac{(m\rho)^n}{n!} + \frac{(m\rho)^m}{m!(1-\rho)} \right]^{-1}$$

Erlang C Formula: probability of having to wait for service

$$P\{\text{Queuing}\} = P_Q = \frac{p_0 (m\rho)^m}{m!(1-\rho)} \quad W = \frac{\rho P_Q}{\lambda(1-\rho)}$$

F2003, CSG150

Fundamentals of Computer Networking

Lecture 3, 40

M/M/m/m Markov System

- Steady state probabilities:

$$p_n = p_0 \left(\frac{\lambda}{\mu}\right)^n \frac{1}{n!}; p_0 = \left[\sum_{n=0}^m \left(\frac{\lambda}{\mu}\right)^n \right]^{-1}$$

Blocking probability: Erlang B formula

$$p_m = \frac{(\lambda / \mu)^m / m!}{\sum_{n=0}^m (\lambda / \mu)^n / n!}$$