

Schwarz (Sections 2.1-3), Molloy, Kleinrock.

- Read any of the queuing theory references, e.g.
 - The M/M/1 Queue
 - Poisson Arrival Model
 - Basic Single Queue Model

INSIDE A ROUTER

MODELING AND ANALYSIS, PART III: NETWORK LAYER PERFORMANCE

- Part III: Network Analysis.

* Part II: Inside a Router.

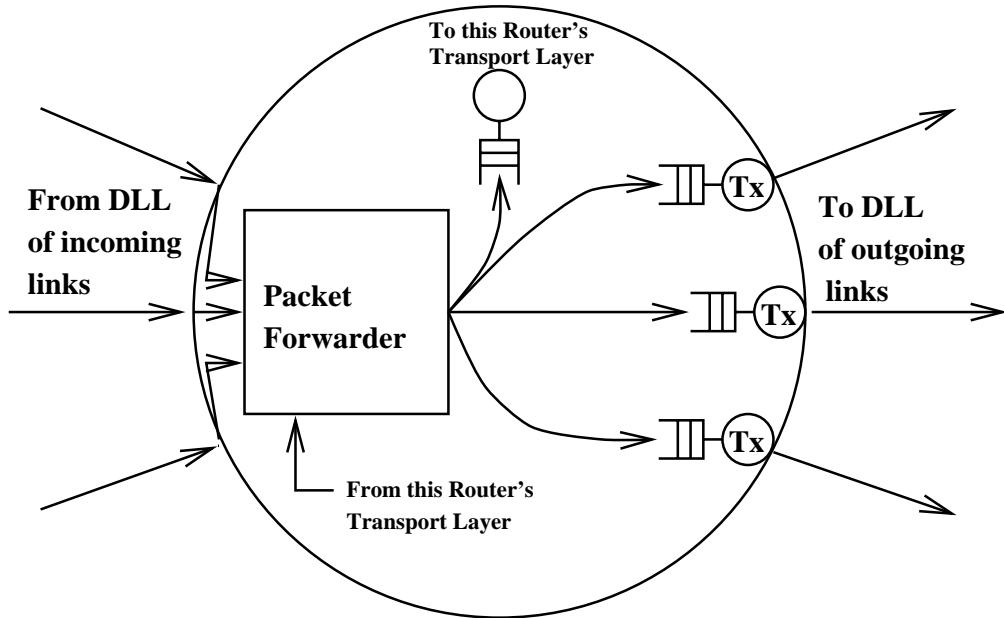
- Part I: Essentials of Probability.

• Three parts.

MODELING AND ANALYSIS

NETWORK LAYER PERFORMANCE

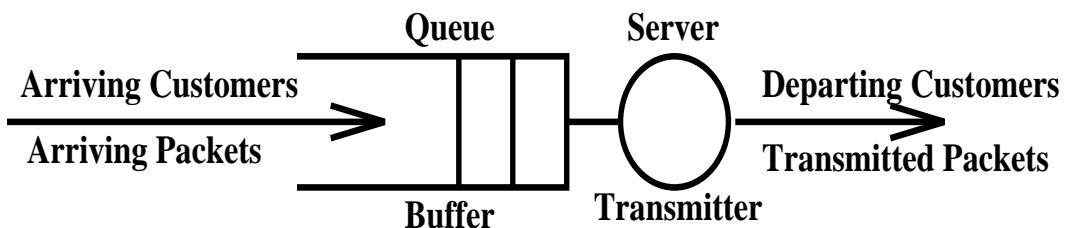
Queueing in the Network Layer at a Router



CCN, ECSE-4670: Performance II: Inside a Router, August 11, 1998, © K.S. Vastola, RPI 3

Basic Single Queue Model

- Classical queueing theory can be applied to an output link in a router.



- For example, a 56 kbps transmission line can “serve” 1000-bit packets at a rate of

$$\frac{56,000 \text{ bits/sec}}{1000 \text{ bits/packet}} = 56 \text{ packets/sec}$$

CCN, ECSE-4670: Performance II: Inside a Router, August 11, 1998, © K.S. Vastola, RPI 4

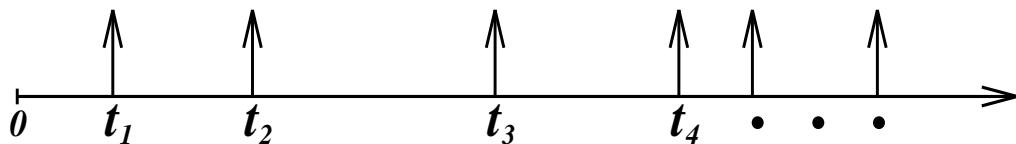
Applications of Queueing Analysis Outside of Networking

- Checkout line in a supermarket.
- Waiting for a teller in a bank.
- Batch jobs waiting to be processed by the CPU.
- “That’s the way the whole thing started,
Silly but it’s true,
Thinking of a sweet romance
Beginning in a queue.”
—G. Gouldman, “Bus Stop,” The Hollies

CCN, ECSE-4670: Performance II: Inside a Router, August 11, 1998, © K.S. Vastola, RPI 5

The Poisson Arrival Model

- A Poisson process is a sequence of events “randomly spaced in time.”



- Examples
 - Customers arriving to a bank.
 - Packets arriving to a buffer.
- The rate λ of a Poisson process is the average number of events per unit time (over a long time).

CCN, ECSE-4670: Performance II: Inside a Router, August 11, 1998, © K.S. Vastola, RPI 6

τ_1 , the time until the first arrival,
has an exponential distribution!

$$F_{\tau_1}(t) = P(\tau_1 \leq t) = 1 - e^{-\lambda t} \quad \text{and} \quad f_{\tau_1}(t) = \lambda e^{-\lambda t}$$

- So

$$P(\tau_1 < t) = P^0(t) = e^{-\lambda t}$$

- Let τ_1 = the time until the next arrival.
- Pick an arbitrary starting point in time (call it 0).

Interarrival Times of a Poisson Process

- For 2 disjoint (non-overlapping) intervals, (s_1, s_2) and (s_3, s_4) , (i.e. $s_1 < s_2 \leq s_3 < s_4$), the number of arrivals in (s_1, s_2) is independent of the number of arrivals in (s_3, s_4) .

$$P^n(t) = \frac{n!}{(s_1-s_0)^n} e^{-\lambda(s_1-s_0)}$$

- For a length of time t , the probability of n arrivals in t units of time is

Properties of a Poisson Process

Interarrival Times of a Poisson Process (cont.)

- Let τ_2 = the length of time between the first and second arrival.
- We can show that

$$P(\tau_2 > t \mid \tau_1 = s) = P(\tau_2 > t) = e^{-\lambda t} \quad \text{for any } s, t > 0$$

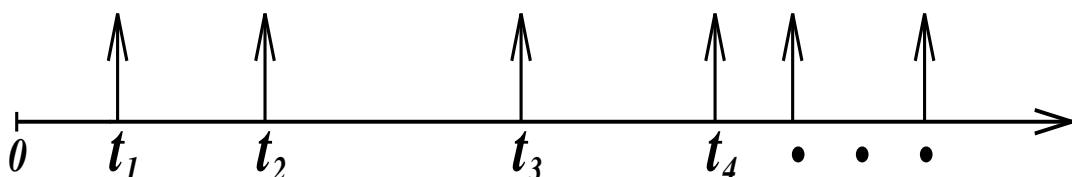
i.e. τ_2 is exponential and independent of τ_1 !

- Similarly define τ_3 as the time between the second and third arrival; τ_4 as the time between the third and fourth arrival; . . .
- The random variables $\tau_1, \tau_2, \tau_3, \dots, \tau_n, \dots$ are called the interarrival times of the Poisson process.

CCN, ECSE-4670: Performance II: Inside a Router, August 11, 1998, © K.S. Vastola, RPI 9

Interarrival Times of a Poisson Process (cont.)

- The interarrival time random variables, $\tau_1, \tau_2, \tau_3, \dots$
 - Are (pair-wise) independent.
 - Each has an exponential distribution with mean $1/\lambda$.



CCN, ECSE-4670: Performance II: Inside a Router, August 11, 1998, © K.S. Vastola, RPI 10

- k = number of buffer slots (omitted when $k = \infty$).
 - m = number of servers.
 - G = General (or arbitrary).
 - M = exponential, D = deterministic.
 - γ is a symbol representing the service distribution
 - D = Deterministic (constant τ).
 - M = Poisson (exponential interarrival times τ).
 - X is a symbol representing the interarrival process
- notation: $X/\gamma/m/k$, where
“ $M/M/1$ ” is a special case of more general (Kendall)

Queueing Notation

- The $M/M/1$ queue is the most basic and important queueing model.
- An infinite length buffer.
- One (1) server.
- the “service rate”.
- Exponential service times (with mean $1/\mu$, so μ is
- Poisson arrivals (with rate λ).
- An $M/M/1$ queue has

The $M/M/1$ Queue

Aside: The D/D/1 Queue

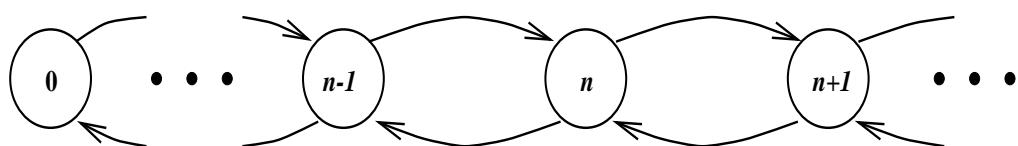
- The D/D/1 queue has
 - Deterministic arrivals (periodic with period = $1/\lambda$).
 - Deterministic service times (each service takes exactly $1/\mu$).
 - As well as 1 server and an infinite length buffer.
- If $\lambda < \mu$ then there is no waiting in a D/D/1 queue.

**Randomness is a major cause
of delay in a network node!**

CCN, ECSE-4670: Performance II: Inside a Router, August 11, 1998, © K.S. Vastola, RPI 13

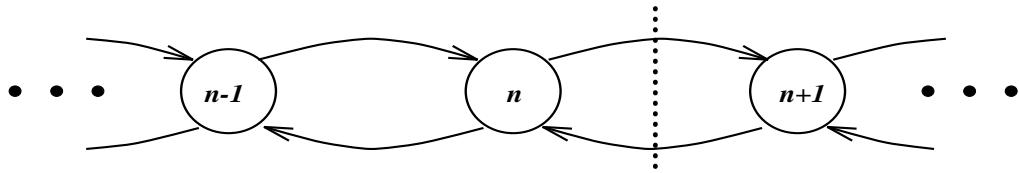
State Analysis of an M/M/1 Queue

- Let n be the state of the system = the number of packets in the system (including the server).
- Let p_n be the steady state probability of finding n customers waiting in the system (including the server).
- How to find p_n ? The state diagram:



CCN, ECSE-4670: Performance II: Inside a Router, August 11, 1998, © K.S. Vastola, RPI 14

State Analysis of an M/M/1 Queue (cont.)



- If this system is stable (i.e. $p_n \neq 0$ for each n), then in steady state it will drift back and forth across the dotted line. So,
- the number of transitions from left to right
= the number of transitions from right to left.
- Thus we obtain the balance equations

$$p_n \lambda = p_{n+1} \mu \quad \text{for each } n \geq 0$$

CCN, ECSE-4670: Performance II: Inside a Router, August 11, 1998, © K.S. Vastola, RPI 15

State Analysis of an M/M/1 Queue (cont.)

- Lets solve the balance equations: $p_n \lambda = p_{n+1} \mu$
- For $n = 0$ we get
- If we let $\rho = \lambda/\mu$, this becomes

$$p_1 = \rho p_0$$

- Similarly

$$p_2 = \rho p_1 = \rho^2 p_0$$

- And in general

$$p_n =$$

CCN, ECSE-4670: Performance II: Inside a Router, August 11, 1998, © K.S. Vastola, RPI 16

- Finally note that $p_n = (1-p)p_0$, $n = 0, 1, 2, 3, \dots$ is a geometric distribution.
- So p is sometimes called the "server utilization" = probability that the server is working = probability that the queueing system is NOT empty
- Also $p = 1 - p_0$ makes intuitive sense.
- Note that requiring $p < 1$ for stability (i.e. $\lambda < \mu$)

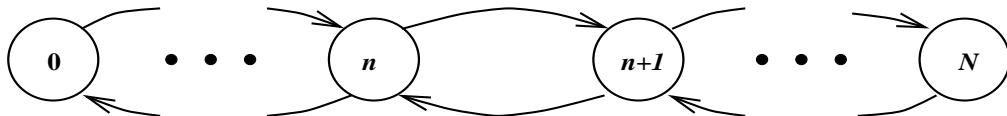
State Analysis of an M/M/1 Queue (cont.)

- $p_0 = 1 - p$ and $p_n = (1-p)p_0$ for $n = 1, 2, 3, \dots$
- So we must have
 - We obtain
 - We need to solve for p_0 , so we need one more equation. Use
 - We have $p_n = p_n p_0$ for $n = 1, 2, 3, \dots$
 - We need to solve for p_0 , so we need one more equation. Use

State Analysis of an M/M/1 Queue (cont.)

The Finite Buffer Case: M/M/1/N

- Infinite buffer assumption is unrealistic in practice.
- N = total number of buffer slots (including server).
- New state diagram:



- Get the same balance equations $p_n\lambda = p_{n+1}\mu$ but now only for $n = 0, 1, 2, \dots, N - 1$ with $N < \infty$. So

$$p_n = \rho p_{n-1} = \rho^n p_0 \quad \text{for } n = 0, 1, 2, \dots, N$$

as before, but we get a different p_0 .

The Finite Buffer Case: M/M/1/N (cont.)

- From $p_n = \rho^n p_0$ for $n = 0, 1, 2, \dots, N < \infty$ and $\sum_{n=0}^N p_n = 1$ we get

$$p_0 = 1 - \sum_{n=1}^N \rho^n p_0$$

- So

$$p_0 = \frac{1}{1 + \sum_{n=1}^N \rho^n} = \frac{1}{1 + \frac{\rho(1-\rho^N)}{(1-\rho)}} = \cdots = \frac{1 - \rho}{1 - \rho^{N+1}}$$

- Note that this holds for any $\rho \geq 0$. No need to assume $\rho < 1$. We always have stability in finite buffer case.

packets.

- Thus, if we desire a blocking probability less than 10^{-6} , we need a buffer capable of holding 19 or more while $p_N < 10^{-6}$ for $N \geq 19$.
- Example: For $p = 0.5$, $p_N > 10^{-6}$ for $N \leq 18$,
- We can use P_B to choose the correct buffer size.
- P_B is very important!

Blocking Probability and Buffer Size (cont.)

- P_B is called the blocking probability.
- Since arrivals are independent of buffer state, we have at an arbitrary point in time,
- Note that p_N is the probability that the buffer is full away due to a full buffer.
- $p_N = P_B =$ probability an arriving packet is turned away due to a full buffer.
- So in the finite buffer case,

$$p_n = \frac{1 - p_{N+1}}{(1 - p)^{n+1}} \quad \text{for } n = 0, 1, 2, \dots, N$$

Blocking Probability and the Right Size Buffer

- So the average rate = $\gamma = \mu(1 - p_0) + 0p_0$
- When the server is idle, the output rate = 0
- When the server is busy, the output rate = μ
- $P(\text{server is busy}) = 1 - p_0$

Look at the output side.

Alternate way to compute throughput of M/M/1/N:

Throughput in the Finite Buffer Case (cont.)

- The throughput γ of any queuing system is the rate at which customers successfully leave the system.
- For the M/M/1 infinite buffer case, $\gamma = \lambda$ if the system is stable. (Everything that arrives and is not blocked must eventually depart.)
- For the M/M/1/N finite buffer case, $\gamma = \lambda(1 - P_B)$. (Everything that arrives and is not blocked must eventually depart.)

Throughput in the Finite Buffer Case

The infinite buffer model is a very good approximation of a finite buffer system.
Even for moderate buffer sizes!

- For $p = 0.8, N = 32$, the difference is only 0.06%.
- For $p = 0.8$ and $N = 16$ packets, these probabilities differ by less than 2.3%.
- For a finite buffer, $p_n = (1-p)p_n / (1-p_{N+1})$
- For an infinite buffer, $p_n = (1-p)p_n$

by the Infinite Buffer Model

Approximation of a Finite Buffer System

- Isn't that neat?
 - Solving for P_B we get
 - Equating our two formulas for γ we get
- $$\mu(1 - p_0) = \alpha(1 - P_B)$$
- Aside: Derivation of $p_N = P_B$ Using Throughput

- Little's Formula holds for very general queuing systems (not just M/M/1). Even whole networks!
- where λ is the "arrival rate for customers eventually served" (which we had called γ).

$$\lambda E(T) = E(n)$$

- Little's Formula says
- $E(T) =$ the average delay for a customer.
- Let $T =$ time spent by a customer in a queuing system (waiting and being served).

Little's Formula and Queuing Delay

$$\begin{aligned} &= (1 - \rho) \frac{(1 - \rho)^2}{\rho} = \\ &= (1 - \rho) \sum_{n=0}^{\infty} n \rho^n = (1 - \rho) \sum_{n=0}^{\infty} n \end{aligned}$$

- So the average number in the system is
- $n =$ the number in the system (including the server).
- Let's look again at the M/M/1 queuing system.

How Long is That Line?

- Sometimes we consider the waiting time W , i.e. the time spent waiting in the queue (not in service). So, which is unitless.

$$\mu E(n) = \frac{\rho}{\lambda} = \frac{\rho(1-p)}{1-p}$$

more convenient to consider

- $E(T)$ is measured in units of time. Sometimes it is

$$E(T) = \frac{\lambda}{\mu} = \frac{\lambda(1-p)}{1-p}$$

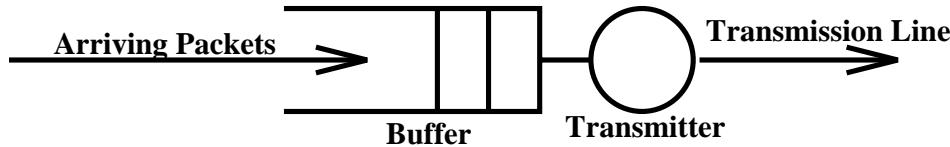
- Let's apply Little to the M/M/1 queue

Little's Formula and Queueing Delay (cont.)

- In steady state, the average number of customers left found on arrival, i.e. $\lambda E(T) = E(n)$ behind on departure should equal the average number
- When it leaves the system, it has been in the system during its time in the system.
- When it arrives to the queueing system, it should find for $E(T)$. Thus, $\lambda E(T)$ customers should have arrived
- Little's formula is either deep or obvious. Intuition:
- Pick a "typical customer".
- Little's Formula is either deep or obvious. Intuition:

Little's Formula and Queueing Delay (cont.)

Single Link Example



- Poisson packet arrivals with rate $\lambda = 2000 \text{ p/s}$.
- Fixed link capacity $C = 1.544 \text{ Mb/s}$ (**T1 Carrier rate**).
- We approximate the packet length distribution by an exponential with mean $L = 515 \text{ b/p}$.
- Thus the service time is exponential with mean

$$\frac{1}{\mu} = \frac{L}{C} = \frac{515 \text{ b/p}}{1.544 \text{ Mb/s}} \approx 0.33 \text{ ms/p}$$

i.e. packets are served at a rate of $\mu = 3000 \text{ p/s}$.

Single Link Example (cont.)

- Using our formulas for an M/M/1 queue

$$\rho = \frac{\lambda}{\mu} = 0.67$$

So,

$$E(n) = \frac{\rho}{1 - \rho} = 2.0 \text{ packets}$$

and

$$E(T) = \frac{E(n)}{\lambda} = 1.0 \text{ ms}$$

$$P_B = \frac{\sum_{i=0}^k p_i / i!}{\sum_{i=0}^k p_i / i!}$$

Erlang Loss) Formula

- Blocking probability is given by the Erlang B (or Erlang C) formula.
- Any customer (a call) which doesn't get a circuit is blocked (gets a busy signal).
- Models a trunk line with k circuits available.
- Important model in circuit switched networks.
- (except one in each server).
- $M/M/k$ for $k \geq 1$. One or more servers, no buffers

Other Queueing Models (cont.)

- Has worse performance at lower loads than $M/M/1$ plexed with $k = 24$.
- (e.g. a T1 carrier is typically time division multiplexed with same total capacity.
- Good model of a link which is made up of multiple channels, either physically or through multiplexing which are useful in networking.
- $M/M/k$ for $k > 1$. Multiple servers.

Other Queueing Models

- M/D/1. Deterministic service times (packet length).
 - Under heavy load ($\rho \approx 1$), M/D/1 has half the delay of an M/M/1.
 - This is one motivation for fixed-packet-length sys-tems like ATM.
 - Special case of M/G/1 with $\sigma^2 = 0$

$$E(n) = \left(\frac{\rho}{\rho - 1} \right) \left(1 - \frac{\rho}{2} \right) \quad \rho < 1$$

Other Queueing Models (cont.)

- M/G/1. Arbitrary service (packet length) distribution.
 - Can still compute the mean number in the system via the Pollaczek-Khinchine (P-K) formula
 - Can still compute the mean delay in the system.
 - Where σ^2 is the variance of the service time distribution. Again, variability (randomness) causes delay.
 - Can apply Little's Formula to get the mean delay.

Other Queueing Models (cont.)

- Can also model and analyze other queuing systems
 - With priority.
 - With general arrival processes.
 - With "vacations."
 - Many others.
- See Schwartz (Ch. 2), Kleinrock (Vol. I & II) or take ECSE-6820/DSES-6820, Queuing (sic) Systems & Applications.
- Queuing theory is also used in analysis of Operating Systems, e.g. in CSCI-6140.

Other Queuing Models (cont.)