Packet Switching: queueing delay, loss



queuing and loss:

- If arrival rate (in bits) to link exceeds transmission rate of link for a period of time:
 - packets will queue, wait to be transmitted on link
 - packets can be dropped (lost) if memory (buffer) fills up

How do loss and delay occur?

packets queue in router buffers

- packet arrival rate to link (temporarily) exceeds output link capacity
- packets queue, wait for turn



Four sources of packet delay



d_{proc}: nodal processing

- check bit errors
- determine output link
- typically < msec

d_{queue}: queueing delay

- time waiting at output link for transmission
- depends on congestion level of router

Four sources of packet delay





* Check out the Java applet for an interactive animation on trans vs. prop delay

Caravan analogy



- cars "propagate" at 100 km/hr
- toll booth takes 12 sec to service car (bit transmission time)
- car~bit; caravan ~ packet
- Q: How long until caravan is lined up before 2nd toll booth?

- time to "push" entire caravan through toll booth onto highway = 12*10 = 120 sec
- time for last car to propagate from 1st to 2nd toll both: 100km/(100km/hr)= 1 hr
- A:62 minutes

Caravan analogy (more)



- suppose cars now "propagate" at 1000 km/hr
- and suppose toll booth now takes one min to service a car
- Q: Will cars arrive to 2nd booth before all cars serviced at first booth?
 - <u>A: Yes</u> after 7 min, 1st car arrives at second booth; three cars still at 1st booth.

Queueing delay (revisited)



^{*} Check out the Java applet for an interactive animation on queuing and loss

Introduction 1-7

La/R -> 1

Queueing Theory Basics

- Each 'node' or 'station' or router called a queue
- Each packet called a 'job'
- A queue has a servicing/processing station and a buffer or queue where jobs wait for service
- The behaviour of a queue is determined by the queueing policy (e.g. FIFO) and the service time (e.g. proportional to packet length or fixed)
- The performance (throughput, delay etc) depends on the queue parameters and the arrival process of jobs

Analysis

- Analysis of a single queue is difficult
- Analysis of networks of queues is even more difficult.
- The best-known results are derived with striong assumptions on all parameters.
- The standard naming scheme of queues is of the form X/Y/k/b where X = arrival process, Y = service time process, k= number of service stations, b = length of buffer
- ♦ We will only look at M/M/1/∞ queues (M=markovian)
- ✤ For networks of M/M/1/∞ queues, it is enough to analyze single queues. Network performance can be very easily obtained from individual queue performance.

$M/M/1/\infty$ queues

- The first M: Poisson arrival process. Probability of N(t) packets arriving in any interval of time t is P(N(t)=k) = (λt)^k exp(-λt)/k!, k = 0,1,2,....
- The second M: Exponential interarrival times Probability of job k arriving t units after job k-1 is P(x=t) = μexp(-μt) if t>0 and 0 otherwise. It follows that E[x] = 1/μ, variance[x] = 1/μ²
- * Under these assumptions, utilization = Prob(queue is non-empty) = ρ where $\rho = \lambda/\mu$
- * So when λ approaches μ (cannot exceed μ), utilization goes towards 100%

$M/M/1/\infty$ queues - contd.

- However, expected number of jobs in the queue is = $\rho /(1 - \rho)$ where $\rho = \lambda/\mu$
- So when λ approaches
 μ the number of jobs in
 the queue approaches
 infinity!!
- As a result delay goes up.
- Therefore most systems cannot be driven at capacity.



Little's Law

One of the very few general laws:

The average number of customers in a (stable) queueing system L is equal to the long-term average effective arrival rate, λ , multiplied by the average time a customer spends in the system, W; or L = λ W.

Applies to single queues or networks

So average delay seen by a packet (from previous slide) = $\rho / [\lambda (1 - \rho)]$

"Real" Internet delays and routes

- what do "real" Internet delay & loss look like?
- * traceroute program: provides delay measurement from source to router along end-end Internet path towards destination. For all *i*:
 - sends three packets that will reach router *i* on path towards destination
 - router *i* will return packets to sender
 - sender times interval between transmission and reply.



"Real" Internet delays, routes

traceroute: gaia.cs.umass.edu to www.eurecom.fr



* Do some traceroutes from exotic countries at www.traceroute.org

Packet loss

- queue (aka buffer) preceding link in buffer has finite capacity
- packet arriving to full queue dropped (aka lost)
- lost packet may be retransmitted by previous node, by source end system, or not at all



* Check out the Java applet for an interactive animation on queuing and loss

Throughput

- *throughput:* rate (bits/time unit) at which bits transferred between sender/receiver
 - instantaneous: rate at given point in time
 - average: rate over longer period of time

