

#### **College of Information Studies**

University of Maryland Hornbake Library Building College Park, MD 20742-4345

# TCP

#### Session 22 INST 346 Technologies, Infrastructure and Architecture

# Improving Support for Learning

- Less lecture, more discussion, slow down
- More examples and lab-style homework
- In-class activities
- Discuss homework and quizzes in class
- Solicit topics in advance for exam review
- More readings (!)
- Jump around less on the slides
- More extensive exam study guide
- More quiz questions (!)

## Goals

- TCP
  - Connection Setup
  - Reliable Transfer
  - Timeout Setting
  - Flow Control
  - Disconnection

• Maybe: BGP

## TCP: Overview RFCs: 793,1122,1323, 2018, 2581

- point-to-point:
  - one sender, one receiver
- reliable, in-order byte steam:
  - no "message boundaries"
- pipelined:
  - TCP congestion and flow control set window size

- full duplex data:
  - bi-directional data flow in same connection
  - MSS: maximum segment size
- connection-oriented:
  - handshaking (exchange of control msgs) inits sender, receiver state before data exchange
- flow controlled:
  - sender will not overwhelm receiver

#### **TCP** segment structure



#### **Connection Management**

before exchanging data, sender/receiver "handshake":

- agree to establish connection (each knowing the other willing to establish connection)
- agree on connection parameters





Socket connectionSocket =
 welcomeSocket.accept();

#### TCP 3-way handshake



## TCP sender events:

#### data rcvd from app:

- create segment with seq #
- seq # is byte-stream number of first data byte in segment
- start timer if not already running
  - think of timer as for oldest unacked segment
  - expiration interval: TimeOutInterval

#### timeout:

- retransmit segment that caused timeout
- restart timer ack rcvd:
- if ack acknowledges previously unacked segments
  - update what is known to be ACKed
  - start timer if there are still unacked segments

# TCP seq. numbers, ACKs

#### sequence numbers:

 byte stream "number" of first byte in segment's data

#### acknowledgements:

- seq # of next byte expected from other side
- cumulative ACK
- Q: how receiver handles out-of-order segments
  - A: TCP spec doesn't say,
    - up to implementor

#### outgoing segment from sender



### TCP seq. numbers, ACKs



simple telnet scenario

#### **TCP:** retransmission scenarios



lost ACK scenario

#### **TCP:** retransmission scenarios



#### TCP ACK generation [RFC 1122, RFC 2581]

event at receiver	TCP receiver action
arrival of in-order segment with expected seq #. All data up to expected seq # already ACKed	delayed ACK. Wait up to 500ms for next segment. If no next segment, send ACK
arrival of in-order segment with expected seq #. One other segment has ACK pending	immediately send single cumulative ACK, ACKing both in-order segments
arrival of out-of-order segment higher-than-expect seq. # . Gap detected	immediately send <i>duplicate ACK</i> , indicating seq. # of next expected byte
arrival of segment that partially or completely fills gap	immediate send ACK, provided that segment starts at lower end of gap

# TCP fast retransmit

- time-out period often relatively long:
  - long delay before resending lost packet
- detect lost segments via duplicate ACKs.
  - sender often sends many segments backto-back
  - if segment is lost, there will likely be many duplicate ACKs.
- TCP fast retransmit if sender receives 3 ACKs for same data ("triple duplicate ACKs"), resend unacked segment with smallest seq #
  - likely that unacked segment lost, so don't wait for timeout

## TCP fast retransmit



## TCP round trip time, timeout

- Q: how to set TCP timeout value?
- Ionger than RTT
  - but RTT varies
- too short: premature timeout, unnecessary retransmissions
- too long: slow reaction to segment loss

- <u>Q:</u> how to estimate RTT?
- SampleRTT: measured time from segment transmission until ACK receipt
  - ignore retransmissions
- SampleRTT will vary, want estimated RTT "smoother"
  - average several recent measurements, not just current SampleRTT

## TCP round trip time, timeout

EstimatedRTT =  $(1 - \alpha)$ \*EstimatedRTT +  $\alpha$ \*SampleRTT

- exponential weighted moving average
- influence of past sample decreases exponentially fast
- typical value:  $\alpha = 0.125$



## TCP round trip time, timeout

- timeout interval: EstimatedRTT plus "safety margin"
  - large variation in EstimatedRTT -> larger safety margin
- estimate SampleRTT deviation from EstimatedRTT:

```
DevRTT = (1-\beta)*DevRTT +
\beta*|SampleRTT-EstimatedRTT|
(typically, \beta = 0.25)
```

```
TimeoutInterval = EstimatedRTT + 4*DevRTT
```

\* Check out the online interactive exercises for more examples: http://gaia.cs.umass.edu/kurose\_ross/interactive/

## **TCP flow control**



# TCP flow control

- receiver "advertises" free buffer space by including rwnd value in TCP header of receiver-to-sender segments
  - RcvBuffer size set via socket options (typical default is 4096 bytes)
  - many operating systems autoadjust RcvBuffer
- sender limits amount of unacked ("in-flight") data to receiver's rwnd value
- guarantees receive buffer will not overflow



#### TCP: closing a connection

- client, server each close their side of connection
  - send TCP segment with FIN bit = I
- respond to received FIN with ACK
  - on receiving FIN, ACK can be combined with own FIN
- simultaneous FIN exchanges can be handled

### TCP: closing a connection



#### Inter-AS routing is different

#### policy:

- intra-AS: single admin, so single consistent policy
- inter-AS: each admin wants control over how its traffic routed and who routes through its AS

#### performance:

- intra-AS: can focus on performance
- inter-AS: policy may dominate over performance

#### Inter-AS tasks

- suppose router in AS1 receives datagram destined outside of AS1:
  - router should forward packet to gateway router, but which one?

#### AS1 must:

- learn which dests are reachable through AS2, which through AS3
- 2. propagate this reachability info to all routers in AS1



#### Internet inter-AS routing: BGP

- BGP (Border Gateway Protocol): the de facto inter-domain routing protocol
  - "glue that holds the Internet together"
- BGP provides each AS a means to:
  - eBGP: obtain subnet reachability information from neighboring ASes
  - **iBGP:** propagate reachability information to all ASinternal routers.
  - determine "good" routes to other networks based on reachability information and policy
- allows subnet to advertise its existence to rest of Internet: "1 am here"

## eBGP, iBGP connections





gateway routers run both eBGP and iBGP protools

## **BGP** basics

- BGP session: two BGP routers ("peers") exchange BGP messages over semi-permanent TCP connection:
  - advertising paths to different destination network prefixes (BGP is a "path vector" protocol)
- when AS3 gateway router 3a advertises path AS3,X to AS2 gateway router 2c:
  - AS3 promises to AS2 it will forward datagrams towards X



## Path attributes and BGP routes

- advertised prefix includes BGP attributes
  - prefix + attributes = "route"
- two important attributes:
  - AS-PATH: list of ASes through which prefix advertisement has passed
  - NEXT-HOP: indicates specific internal-AS router to nexthop AS
- Policy-based routing:
  - gateway receiving route advertisement uses import policy to accept/decline path (e.g., never route through AS Y).
  - AS policy also determines whether to *advertise* path to other other neighboring ASes

## BGP path advertisement



- AS2 router 2c receives path advertisement AS3,X (via eBGP) from AS3 router 3a
- Based on AS2 policy, AS2 router 2c accepts path AS3,X, propagates (via iBGP) to all AS2 routers
- Based on AS2 policy, AS2 router 2a advertises (via eBGP) path AS2, AS3, X to AS1 router 1c

## BGP path advertisement



gateway router may learn about multiple paths to destination:

- AS1 gateway router 1C learns path AS2,AS3,X from 2a
- AS1 gateway router 1C learns path AS3,X from 3a
- Based on policy, AS1 gateway router 1C chooses path AS3, X, and advertises path within AS1 via iBGP

#### BGP: achieving policy via advertisements



Suppose an ISP only wants to route traffic to/from its customer networks (does not want to carry transit traffic between other ISPs)

- A advertises path Aw to B and to C
- B chooses not to advertise BAw to C:
  - B gets no "revenue" for routing CBAw, since none of C, A, w are B's customers
  - C does not learn about CBAw path
- C will route CAw (not using B) to get to w

#### BGP: achieving policy via advertisements



Suppose an ISP only wants to route traffic to/from its customer networks (does not want to carry transit traffic between other ISPs)

- A,B,C are provider networks
- X,W,Y are customer (of provider networks)
- X is dual-homed: attached to two networks
- policy to enforce: X does not want to route from B to C via X
  - .. so X will not advertise to B a route to C

## **BGP** route selection

- router may learn about more than one route to destination AS, selects route based on:
  - I. local preference value attribute (policy decision)
  - 2. shortest AS-PATH
  - 3. closest NEXT-HOP router (hot potato routing)
  - 4. additional criteria

## Hot Potato Routing



- 2d learns (via iBGP) it can route to X via 2a or 2c
- hot potato routing: choose local gateway that has least intradomain cost (e.g., 2d chooses 2a, even though more AS hops to X): don't worry about inter-domain cost!