# End-to-End Protocols: UDP and TCP

Hui Chen, Ph.D.

Dept. of Engineering & Computer Science

Virginia State University

Petersburg, VA 23806

# Acknowledgements

- Some pictures used in this presentation were obtained from the Internet

- The instructor used the following references
  - Larry L. Peterson and Bruce S. Davie, Computer Networks: A Systems Approach, 5th Edition, Elsevier, 2011
  - Andrew S. Tanenbaum, Computer Networks, 5th Edition, Prentice-Hall, 2010
  - James F. Kurose and Keith W. Ross, Computer Networking: A Top-Down Approach, 5th Ed., Addison Wesley, 2009
  - Larry L. Peterson's (http://www.cs.princeton.edu/~llp/) Computer Networks class web site

# Acknowledgements

- ☐ Animations in the PDF version of the slides is produced using
  - ■ PPspliT
  - ■ http://www.dia.uniroma3.it/~rimondin/downloads.php

# Outline

❑ User Datagram Protocol

❑ Transmission Control Protocol
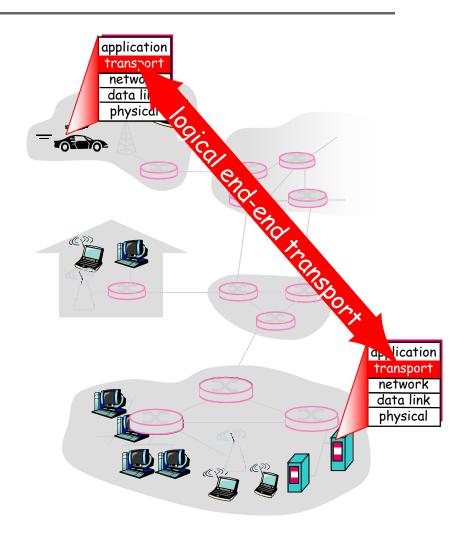
# Network Applications

Network

- Users make use of networks via network applications at hosts
- A hosts can run many network applications simultaneously
- Each application is one or more running programs (processes)
- Q: How processes share the underlying network layers?

# Transport Layer Services and Protocols

- provide *logical communication* between application processes running on different hosts
- transport protocols run in end systems
  - **send side**
    - breaks app messages into segments, passes to network layer
  - **receive side:**
    - reassembles segments into messages, passes to app layer
- more than one transport protocol available to applications
  - **Internet: TCP and UDP**



application
transport
network
data link
physical

logical end-end transport

application
transport
network
data link
physical

# Transport vs. Network Layer (1)

- *network layer:* logical communication between hosts

- *transport layer:* logical communication between processes
  - relies on, enhances, network layer services

<span style="color:red">Household analogy:</span>

*12 kids sending letters among themselves via their parents*

- processes = kids

- application messages = letters in envelopes

- hosts = houses

- transport protocol = Ann and Bill (parents)

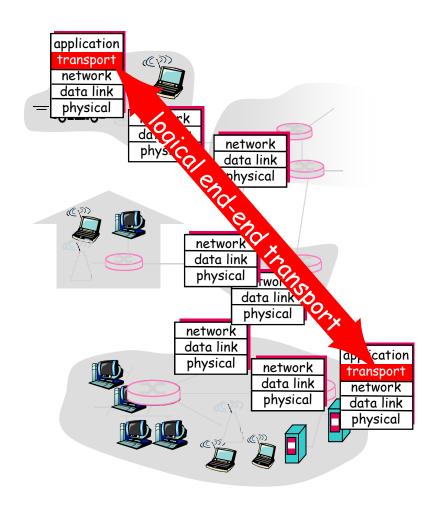- network-layer protocol = postal service

# Transport vs. Network Layer (2)

- Network layer: Underlying best-effort network
  - drop messages
  - re-orders messages
  - delivers duplicate copies of a given message
  - limits messages to some finite size
  - delivers messages after an arbitrarily long delay

- Transport Layer: Common end-to-end services
  - guarantee message delivery
  - deliver messages in the same order they are sent
  - deliver at most one copy of each message
  - support arbitrarily large messages
  - support synchronization
  - allow the receiver to flow control the sender
  - support multiple application processes on each host

# Internet Transport-Layer Protocols

- ☐ Reliable, in-order delivery (TCP)
  - ■ congestion control
  - ■ flow control
  - ■ connection setup
- ☐ Unreliable, unordered delivery: UDP
  - ■ no-frills extension of "best-effort" IP
- ☐ Services not available:
  - ■ delay guarantees
  - ■ bandwidth guarantees

# Multiplexing/Demultiplexing

Host-to-host delivery ⬅➡ process-to-process delivery

# Multiplexing/Demultiplexing
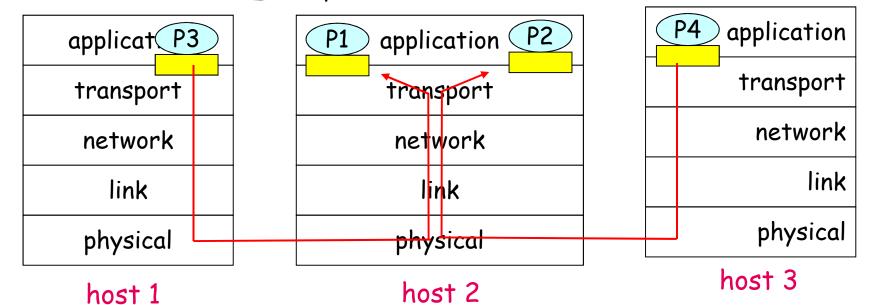
Host-to-host delivery ←→ process-to-process delivery

**Demultiplexing at rcv host:**
delivering received segments
to correct socket

**Multiplexing at send host:**
gathering data from multiple
sockets, enveloping data with
header (later used for
demultiplexing)

▭ = socket ⬭ = process



host 1    host 2    host 3

# Simple Demultiplexer (1)

- ❑ Need to know to or from which process the data is sent or come
  - ▪ Identify processes on hosts
- ❑ How to identify processes on hosts?
  - ▪ Introduce concept of "port"
  - ▪ *Q: why not to use process id?*

# Processes ID: Windows Example



CSCI 445 – Fall 2016

# Processes ID: Linux Example

# Simple Demultiplexer  (2)

- How to identify processes on hosts?
    - *Q: why not to use process id?*
    - Introduce concept of "port"
        - Endpoints identified by ports
        - servers have well-known ports
        - see /etc/services on Unix/Linux
        - see C:\WINDOWS\system32\drivers\etc\services on MS Windows



Process 8

Process 3

Host 1

Process 3

Process 8

Host 2

# Simple Demultiplexer: UDP

□ Adds multiplexing to Internet Protocol

▪ Endpoints identified by ports (UDP ports)

▪ Demultiplex via ports on hosts

▪ Nothing more is added

    □ Unreliable and unordered datagram service

    □ No flow control

▪ User Datagram Protocol (UDP)

    □ A process is identified by <host, port>

    □ Connectionless model

□ Header format

▪ Optional checksum

    □ psuedo header + UDP header + data

    □ pseudo header = protocol number + source IP address and destination IP address + UDP length field  → From IP header

From UDP header

| 0 | 16 | 31 |
|---|---|---|
| SrcPort | | DstPort |
| Length | | Checksum |
| Data | | |

# In-Class Exercise L15-1



|     | 0 | 16 | 31 |

| SrcPort | DstPort |
| Length | Checksum |
| Data |

□ Q1: How many UDP ports are there?

□ Q2: How big are UDP headers?

□ Q3: How much data does a UDP datagram can carry?

□ Turn your work in before you leave!

# Transmission Control Protocol (TCP)

- Connection-oriented
- Byte-stream
  - applications writes bytes
  - TCP sends segments
  - applications reads bytes
- Full duplex
- Flow control: keep sender from overrunning receiver
- Congestion control: keep sender from overrunning network

# Data Link Versus Transport

- Potentially connects many different hosts
  - need explicit connection establishment and termination
- Potentially different RTT
  - need adaptive timeout mechanism
- Potentially long delay in network
  - need to be prepared for arrival of very old packets

- Potentially different capacity at destination
- need to accommodate different node capacity
- Potentially different network capacity
- need to be prepared for network congestion

# Segment Format (1)

# Segment Format (2)

- Each connection identified with 4-tuple:
  - (SrcPort, SrcIPAddr, DsrPort, DstIPAddr)
- Sliding window + flow control
  - acknowledgment, SequenceNum, AdvertisedWinow
- Flags
  - SYN, FIN, RESET, PUSH, URG, ACK
- Checksum
  - pseudo header + TCP header + data



Data (SequenceNum)

Sender

Receiver

Acknowledgment + AdvertisedWindow

# Sequence and Acknowledgement Numbers (1)

❑ Host A sends a file of 500,000 bytes over a TCP connection with Maximum Segment Size (MSS) as 1,000 bytes to host B

- ■ How many segments? 500,000/1,000 = 500

- ■ Sequence number assignments

   - ❑ Sequence number of 1st segment? 0

   - ❑ Sequence number of 2nd segment? 1,000

   - ❑ Sequence number of 3rd segment?  2,000

   - ❑ ……

# Sequence and Acknowledgement Numbers (2)

- Scenario 1
  - Host B received all bytes numbered 0 to 1,999 from host A
  - What would host B put in the acknowledgement number field of the segment it sends to A?
    - 2,000: the sequence number of the next byte host B is expecting
- Scenario 2
  - Host B received two segments containing bytes from 0-999, and 2,000-2,999, respectively?
  - What would host B put in the acknowledgement number field of the segment it sends to A?
    - 1000: TCP only acknowledges bytes up to the first missing byte in the stream, and it is the next byte host B is expecting
- Scenario 3
  - Host B received 1st segment containing bytes from 0-999. Somehow, next it received 3rd segment containing bytes from 2,000-2,999.
  - What does host B in this case that the segments arrive out of order?
    - TCP does not specify how to deal with this situation. Hence, it is up to the implementation.
      - Option 1: Host B immediately discards out-of-order segment → simple receiver design
      - Option 2: Host B keeps the out-of-order segment and waits for missing bytes to fill in the gaps → more efficient on bandwidth utilization → taken in practice

# TCP is Connection-Oriented

☐ Keep track of states of receiver and sender

- Connection Establishment

- Connection Termination

- TCP finite state machine and state transition

# Connection Establishment



Active participant (client)　　　　　　Passive participant (server)

SYN, SequenceNum = $x$

SYN + ACK, SequenceNum = $y$, Acknowledgment = $x + 1$

ACK, Acknowledgment = $y + 1$

# Connection Termination

# State Transition Diagram



CSCI 445 – Fall 2016

# Connection Establishment and State Transition

# Connection Establishment and State Transition

client

# Connection Establishment and State Transition

client



CLOSED

Passive open | Close

Active open/SYN

Close

LISTEN

SYN/SYN + ACK        Send/SYN

SYN_RCVD          SYN/SYN + ACK          SYN_SENT

ACK          SYN + ACK/ACK

Close/FIN          ESTABLISHED          **client**

# Connection Establishment and State Transition

client          server



CLOSED

Active open/SYN

Passive open | Close

Close

LISTEN

SYN/SYN + ACK          Send/SYN
SYN_RCVD          SYN/SYN + ACK          SYN_SENT

ACK          SYN + ACK/ACK

Close/FIN          ESTABLISHED          **client**

# Connection Establishment and State Transition



client                    server

## client

CLOSED

Passive open | Close

Active open/SYN

Close

LISTEN

SYN/SYN + ACK | Send/SYN

SYN/SYN + ACK

SYN_RCVD                    SYN_SENT

ACK | SYN + ACK/ACK

Close/FIN

ESTABLISHED

## server

CLOSED

Passive open | Close

Active open/SYN

Close

LISTEN

SYN/SYN + ACK | Send/SYN

SYN/SYN + ACK

SYN_RCVD                    SYN_SENT

ACK | SYN + ACK/ACK

Close/FIN

ESTABLISHED

# Connection Establishment and State Transition

client    server

closed



**client**            **server**

# Connection Establishment and State Transition

client     server
closed         closed



**client**



**server**

# Connection Establishment and State Transition

client        server

closed          closed

-------------- Action: passive open



**client**          **server**

# Connection Establishment and State Transition

client    server

closed

------------ Action: passive open
listen



**client**

**server**

# Connection Establishment and State Transition

client        server

closed

Action: active open

Action: passive open
listen



client

server

# Connection Establishment and State Transition



client    server

closed

Action: active open

Action: passive open

listen

SYN, SequenceNum = x

**client**

**server**

# Connection Establishment and State Transition

client    server

Action: active open

SYN_SENT

*SYN, SequenceNum = x*

listen

Action: passive open



**client**

**server**

# Connection Establishment and State Transition

client    server

Action: active open

SYN_SENT

SYN, SequenceNum = x

Action: passive open

listen



**client**

**server**

# Connection Establishment and State Transition

client    server

Action: active open

SYN_SENT

SYN, SequenceNum = x

SYN_RECV



**client**    **server**

# Connection Establishment and State Transition



Action: active open

SYN_SENT

SYN, SequenceNum = x

SYN_RECV

SYN+ACK, SequenceNum = y
Acknowledgement = x+1

**client**

**server**

# Connection Establishment and State Transition



Action: active open

client          server

SYN_SENT

SYN, SequenceNum = x

SYN_RECV

SYN+ACK, SequenceNum = y
Acknowledgement = x+1

# Connection Establishment and State Transition



Action: active open

client                                    server

SYN, SequenceNum = x

SYN_RECV

SYN+ACK, SequenceNum = y
Acknowledgement = x+1

Established

**client**

**server**

# Connection Establishment and State Transition

client      server

Action: active open

*SYN, SequenceNum = x*

SYN_RECV

*SYN+ACK, SequenceNum = y, Acknowledgement = x+1*

Established

*ACK, Acknowledgement = y+1*

**client**

**server**

# Connection Establishment and State Transition



Action: active open

client                                server

SYN, SequenceNum = x

SYN_RECV

SYN+ACK, SequenceNum = y
Acknowledgement = x+1

Established    ACK, Acknowledgement = y+1

# Connection Establishment and State Transition



client          server

Action: active open

SYN, SequenceNum = x

SYN+ACK, SequenceNum = y
Acknowledgement = x+1

Established

ACK, Acknowledgement = y+1

Established

**client**

**server**

# Connection Termination and State Transition (1)

client        server

timed wait

**Client closes first**

closed

| | Client | | | | Server | | |
|---|---|---|---|---|---|---|---|
| Close/FIN | ESTABLISHED | | | Close/FIN | ESTABLISHED | | |
| | Close/FIN    FIN/ACK | | | | Close/FIN    FIN/ACK | | |
| FIN_WAIT_1 | | | CLOSE_WAIT | FIN_WAIT_1 | | | CLOSE_WAIT |
| ACK | FIN/ACK | | Close/FIN | ACK | FIN/ACK | | Close/FIN |
| FIN_WAIT_2 | ACK + FIN/ACK | CLOSING | LAST_ACK | FIN_WAIT_2 | ACK + FIN/ACK | CLOSING | LAST_ACK |
| FIN/ACK | | ACK  Timeout after two segment lifetimes | ACK | FIN/ACK | | ACK  Timeout after two segment lifetimes | ACK |
| | TIME_WAIT | → | CLOSED | | TIME_WAIT | → | CLOSED |

# Connection Termination and State Transition (1)

client          server

timed wait

Client closes
first          closed

# Connection Termination and State Transition (1)

client          server

timed wait

**Client closes first**

closed

# Connection Termination and State Transition (1)

client      server

close

timed wait

**Client closes first**    closed

# Connection Termination and State Transition (1)

client    server

close ——— FIN ——→

timed wait

Client closes first    closed



*Close*/FIN  ESTABLISHED  FIN/ACK

*Close*/FIN

FIN_WAIT_1 ← — — — *Close*/FIN    CLOSE_WAIT

ACK    FIN/ACK    *Close*/FIN

FIN_WAIT_2    CLOSING    LAST_ACK

ACK    *Timeout after two segment lifetimes*    ACK

FIN/ACK    TIME_WAIT    CLOSED

---

*Close*/FIN  ESTABLISHED  FIN/ACK

*Close*/FIN

FIN_WAIT_1    CLOSE_WAIT

ACK    FIN/ACK    *Close*/FIN

FIN_WAIT_2    CLOSING    LAST_ACK

ACK    *Timeout after two segment lifetimes*    ACK

FIN/ACK    TIME_WAIT    CLOSED

# Connection Termination and State Transition (1)

client     server

close ————————————→
                FIN

timed wait

**Client closes first**     closed

| | | | | | |
|---|---|---|---|---|---|

Client side state diagram:
- *Close*/FIN → ESTABLISHED
- ESTABLISHED → *Close*/FIN → FIN_WAIT_1
- ESTABLISHED → FIN/ACK → CLOSE_WAIT
- FIN_WAIT_1 → ACK → FIN_WAIT_2
- FIN_WAIT_1 → FIN/ACK → CLOSING
- FIN_WAIT_1 → ACK + FIN/ACK → TIME_WAIT
- CLOSE_WAIT → *Close*/FIN → LAST_ACK
- FIN_WAIT_2 → FIN/ACK → TIME_WAIT
- CLOSING → ACK → TIME_WAIT
- TIME_WAIT → *Timeout after two segment lifetimes* → CLOSED
- LAST_ACK → ACK → CLOSED

Server side state diagram:
- *Close*/FIN → ESTABLISHED
- ESTABLISHED → *Close*/FIN → FIN_WAIT_1
- ESTABLISHED → FIN/ACK → CLOSE_WAIT
- FIN_WAIT_1 → ACK → FIN_WAIT_2
- FIN_WAIT_1 → FIN/ACK → CLOSING
- FIN_WAIT_1 → ACK + FIN/ACK → TIME_WAIT
- CLOSE_WAIT → *Close*/FIN → LAST_ACK
- FIN_WAIT_2 → FIN/ACK → TIME_WAIT
- CLOSING → ACK → TIME_WAIT
- TIME_WAIT → *Timeout after two segment lifetimes* → CLOSED
- LAST_ACK → ACK → CLOSED

# Connection Termination and State Transition (1)

client          server

close ─────────── FIN ──────────►

timed wait

Client closes
first

closed

| | | | | | |
|---|---|---|---|---|---|
| Close/FIN | | | | Close/FIN | ESTABLISHED |
| ↓ | ESTABLISHED | | | ↓ | |
| | Close/FIN | FIN/ACK | | | Close/FIN  FIN/ACK |
| FIN_WAIT_1 | | CLOSE_WAIT | | FIN_WAIT_1 | CLOSE_WAIT |
| ACK | FIN/ACK | Close/FIN | | ACK  FIN/ACK | Close/FIN |
| FIN_WAIT_2 | ACK+FIN/ACK  CLOSING | LAST_ACK | | FIN_WAIT_2  ACK+FIN/ACK  CLOSING | LAST_ACK |
| FIN/ACK | ACK  Timeout after two segment lifetimes | ACK | | FIN/ACK  ACK  Timeout after two segment lifetimes | ACK |
| | TIME_WAIT | CLOSED | | TIME_WAIT | CLOSED |

# Connection Termination and State Transition (1)



client                    server

close                        FIN

                             ACK

timed wait

Client closes
first

closed

ESTABLISHED

*Close*/FIN

*Close*/FIN          FIN/ACK

FIN_WAIT_1                    CLOSE_WAIT

FIN/ACK                    *Close*/FIN

ACK              ACK + FIN/ACK

FIN_WAIT_2       CLOSING          LAST_ACK

                 ACK   *Timeout after two       ACK
                       segment lifetimes*

FIN/ACK

TIME_WAIT                     CLOSED

ESTABLISHED

*Close*/FIN

*Close*/FIN          FIN/ACK

FIN_WAIT_1                    CLOSE_WAIT

FIN/ACK                    *Close*/FIN

ACK              ACK + FIN/ACK

FIN_WAIT_2       CLOSING          LAST_ACK

                 ACK   *Timeout after two       ACK
                       segment lifetimes*

FIN/ACK

TIME_WAIT                     CLOSED

# Connection Termination and State Transition (1)

client            server

close

FIN

ACK

timed wait

Client closes first

closed

**Client state diagram (left):**

Close/FIN → ESTABLISHED

ESTABLISHED

Close/FIN → FIN_WAIT_1

FIN/ACK → CLOSE_WAIT

FIN_WAIT_1
- ACK → FIN_WAIT_2
- FIN/ACK → CLOSING
- ACK + FIN/ACK → TIME_WAIT

FIN_WAIT_2
- FIN/ACK → TIME_WAIT

CLOSE_WAIT
- Close/FIN → LAST_ACK

CLOSING
- ACK → TIME_WAIT

LAST_ACK
- ACK → CLOSED

TIME_WAIT → Timeout after two segment lifetimes → CLOSED

**Server state diagram (right):**

Close/FIN → ESTABLISHED

ESTABLISHED

Close/FIN → FIN_WAIT_1

FIN/ACK → CLOSE_WAIT

FIN_WAIT_1
- ACK → FIN_WAIT_2
- FIN/ACK → CLOSING
- ACK + FIN/ACK → TIME_WAIT

FIN_WAIT_2
- FIN/ACK → TIME_WAIT

CLOSE_WAIT
- Close/FIN → LAST_ACK

CLOSING
- ACK → TIME_WAIT

LAST_ACK
- ACK → CLOSED

TIME_WAIT → Timeout after two segment lifetimes → CLOSED

# Connection Termination and State Transition (1)

client        server

close  —FIN→

←ACK—

timed wait

**Client closes first**

closed

# Connection Termination and State Transition (1)

client          server

close ──FIN──▶

◀──ACK──

timed wait

closed

**Client closes first**



ESTABLISHED
Close/FIN
Close/FIN → FIN_WAIT_1
FIN/ACK → CLOSE_WAIT

FIN_WAIT_1
ACK
FIN/ACK
ACK + FIN/ACK
FIN_WAIT_2
FIN/ACK → TIME_WAIT

CLOSING
ACK → TIME_WAIT

CLOSE_WAIT
Close/FIN → LAST_ACK
ACK → CLOSED

TIME_WAIT
Timeout after two segment lifetimes → CLOSED

ESTABLISHED
Close/FIN
Close/FIN → FIN_WAIT_1
FIN/ACK → CLOSE_WAIT

FIN_WAIT_1
ACK
FIN/ACK
ACK + FIN/ACK
FIN_WAIT_2
FIN/ACK → TIME_WAIT

CLOSING
ACK → TIME_WAIT

CLOSE_WAIT
Close/FIN → LAST_ACK
ACK → CLOSED

TIME_WAIT
Timeout after two segment lifetimes → CLOSED

# Connection Termination and State Transition (1)

client              server

close ──────── FIN ────────→

←──────── ACK ────────

close

←──────── FIN ────────

timed wait

closed

Client closes first

# Connection Termination and State Transition (1)

client    server

close

FIN

ACK

close

FIN

timed wait

Client closes first

closed

# Connection Termination and State Transition (1)



client     server

close

FIN

ACK

close

FIN

timed wait

closed

Client closes first

# Connection Termination and State Transition (1)

client         server

close

FIN

ACK

close

FIN

timed wait

ACK

**Client closes first**

closed

## Client state diagram

| ESTABLISHED |

Close/FIN → FIN_WAIT_1

Close/FIN (FIN_WAIT_1)

FIN/ACK → CLOSE_WAIT

FIN_WAIT_1 → ACK → FIN_WAIT_2

ACK + FIN/ACK

CLOSE_WAIT → Close/FIN → LAST_ACK

CLOSING

FIN_WAIT_2 → FIN/ACK → TIME_WAIT

CLOSING → ACK → TIME_WAIT

TIME_WAIT → Timeout after two segment lifetimes → CLOSED

LAST_ACK → ACK → CLOSED

## Server state diagram

| ESTABLISHED |

Close/FIN → FIN_WAIT_1

FIN/ACK → CLOSE_WAIT

FIN_WAIT_1 → ACK → FIN_WAIT_2

ACK + FIN/ACK → FIN_WAIT_1

FIN/ACK

CLOSE_WAIT → Close/FIN → LAST_ACK

CLOSING

FIN_WAIT_2 → FIN/ACK → TIME_WAIT

CLOSING → ACK → TIME_WAIT

TIME_WAIT → Timeout after two segment lifetimes → CLOSED

LAST_ACK → ACK → CLOSED

# Connection Termination and State Transition (1)



Client closes first

# Connection Termination and State Transition (1)

# Connection Termination and State Transition (1)

client     server

close

FIN

ACK

FIN     close

timed wait

ACK

**Client closes first**

closed

# Connection Termination and State Transition (1)

client          server

close ——— FIN ———→

←——— ACK ———

close

←——— FIN ———

timed wait

——— ACK ———→

**Client closes first**

closed

# Connection Termination and State Transition (1)



client    server

close
FIN
ACK
FIN    close
timed wait
ACK
closed

Client closes first

# Connection Termination and State Transition (2)

- ❏ This side closes first
  - ▪ ESTABLISHED → FIN_WAIT_1 → FIN_WAIT_2 → TIME_WAIT
- ❏ Other side closes first
  - ▪ ESTABLISHED → CLOSE_WAIT → LAST_ACK → CLOSED
- ❏ Both sides close at the same time
  - ▪ ESTABLISHED → FIN_WAIT_1 → CLOSING → TIME_WAIT → CLOSED

# TCP Sliding Window: Why Different?

- Potentially connects many different hosts
    - need explicit connection establishment and termination
- Potentially different RTT
    - need adaptive timeout mechanism
- Potentially long delay in network
    - need to be prepared for arrival of very old packets

- Potentially different capacity at destination
    - need to accommodate different node capacity
- Potentially different network capacity
    - need to be prepared for network congestion

# TCP Sliding Window: Reliable and Ordered Delivery

TCP uses cumulative acknowledgements to acknowledge receiving of all the bytes up to the first missing byte



- ❑ Sending side
  - LastByteAcked ≤ LastByteSent
  - LastByteSent ≤ LastByteWritten
  - buffer bytes between LastByteAcked and LastByteWritten

Receiving side
LastByteRead < NextByteExpected
NextByteExpected ≤ LastByteRcvd +1
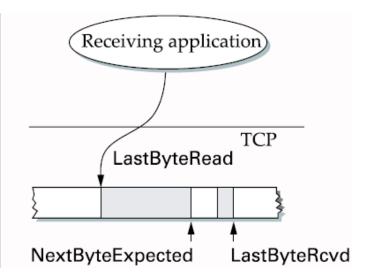buffer bytes betweenNextByteRead and LastByteRcvd

# TCP Flow Control (1)

- receive side of TCP connection has a receive buffer

- app process may be slow at reading from buffer

- speed-matching service: matching the send rate to the receiving app's drain rate

flow control

sender won't overflow receiver's buffer by transmitting too much, too fast

# TCP Flow Control (2)

- Send buffer size: MaxSendBuffer
- Receive buffer size: MaxRcvBuffer
- Receiving side
  - LastByteRcvd - LastByteRead $\leq$ MaxRcvBuffer
  - AdvertisedWindow = MaxRcvBuffer – ((NextByteExpected -1) - LastByteRead)) → maximum possible free space remaining in the buffer
- Sending side
  - LastByteSent - LastByteAcked $\leq$ AdvertisedWindow
    - LastByteSent – LastByteAcked: unacknowledged bytes sender has put in TCP
    - Otherwise, the sender may overrun the receiver
  - EffectiveWindow = AdvertisedWindow - (LastByteSent -LastByteAcked) → how much data it can sent
  - LastByteWritten - LastByteAcked $\leq$ MaxSendBuffer
  - If the sender tries to write y bytes to TCP
    - block sender if (LastByteWritten - LastByteAcked) + y > MaxSenderBuffer
- Always send ACK in response to arriving data segment
- Persist when AdvertisedWindow = 0

# Flow Control and Buffering (3)

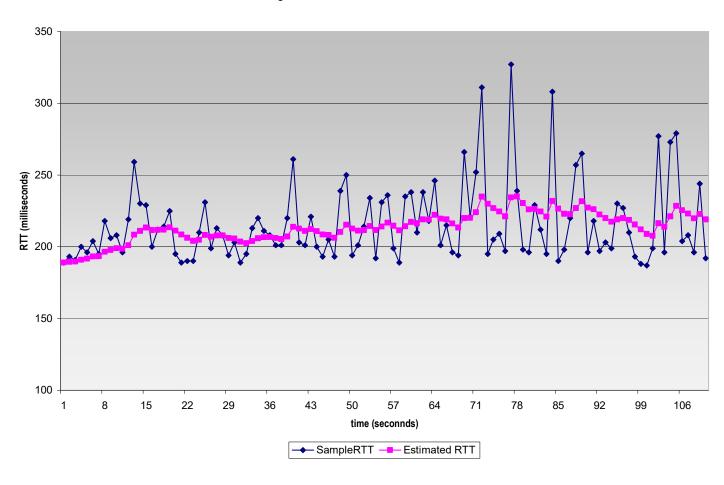| | A | Message | B | Comments |
|---|---|---|---|---|
| 1 | → | < request 8 buffers> | → | A wants 8 buffers |
| 2 | ← | <ack = 15, buf = 4> | ← | B grants messages 0-3 only |
| 3 | → | <seq = 0, data = m0> | → | A has 3 buffers left now |
| 4 | → | <seq = 1, data = m1> | → | A has 2 buffers left now |
| 5 | → | <seq = 2, data = m2> | • • • | Message lost but A thinks it has 1 left |
| 6 | ← | <ack = 1, buf = 3> | ← | B acknowledges 0 and 1, permits 2-4 |
| 7 | → | <seq = 3, data = m3> | → | A has 1 buffer left |
| 8 | → | <seq = 4, data = m4> | → | A has 0 buffers left, and must stop |
| 9 | → | <seq = 2, data = m2> | → | A times out and retransmits |
| 10 | ← | <ack = 4, buf = 0> | ← | Everything acknowledged, but A still blocked |
| 11 | ← | <ack = 4, buf = 1> | ← | A may now send 5 |
| 12 | ← | <ack = 4, buf = 2> | ← | B found a new buffer somewhere |
| 13 | → | <seq = 5, data = m5> | → | A has 1 buffer left |
| 14 | → | <seq = 6, data = m6> | → | A is now blocked again |
| 15 | ← | <ack = 6, buf = 0> | ← | A is still blocked |
| 16 | • • • | <ack = 6, buf = 4> | ← | Potential deadlock |

Dynamic buffer allocation.  The arrows show the direction of transmission.  An ellipsis (…) indicates a lost TCP segment

# Adaptive Retransmission: Original Algorithm

◻ Measure SampleRTT for each segment/ACK pair

◻ Compute weighted average of RTT

- EstimatedRTT = α x EstimatedRTT + β x SampleRTT

- where α **+** β **= 1**

  ◻ α between 0.8 and 0.9

  ◻ β between 0.1 and 0.2

- Set timeout based on EstimatedRTT

  ◻ TimeOut = 2 x EstimatedRTT
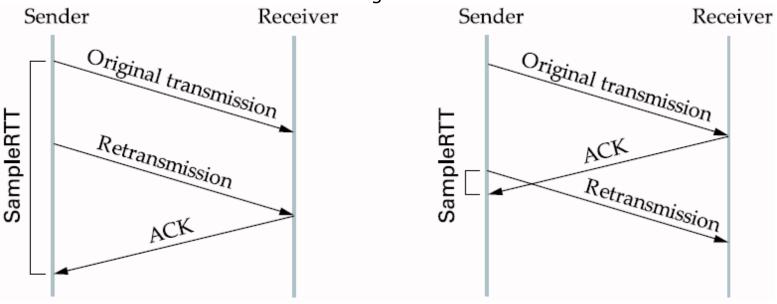
# Example RTT estimation:

**RTT: gaia.cs.umass.edu to fantasia.eurecom.fr**

# Adaptive Retransmission: Karn/Partridge Algorithm

Problem with original algorithm
ACK does not really acknowledge a transmission, it acknowledges the receipt of data → can not distinguish an ACK is for which transmission/retransmission of a segment



- ☐ Do not sample RTT when retransmitting
- ☐ **Double timeout after each retransmission**
  - ■ Congestion is the most likely cause of lost segments → TCP should not react too aggressively to a timeout

# Jacobson/ Karels Algorithm

- Previous approaches did not take the variance of the sample RTT into account
  - If no variance, Estimated RTT is good enough, $2 \times$ Estimated RTT is too pessimistic
  - If variance large, timeout value should not be too dependent on Estimated RTT
- New Calculations for average RTT
  - Difference = SampleRTT – EstimtaedRTT
  - EstimatedRTT = EstimatedRTT + ($\delta$ x Difference)
  - Deviation = Deviation + $\delta$( |Difference| - Deviation)
    - where $\delta$ is a factor between 0 and 1
  - Consider variance when setting timeout value
    - TimeOut = $\mu$ x EstimatedRTT + $\phi$ x Deviation
    - where $\mu$ = 1 and $\phi$ = 4
- Notes
  - algorithm only as good as granularity of clock (500ms on Unix)
  - accurate timeout mechanism important to congestion control

# TCP: Sequence Number Wrap Around

| Bandwidth | Time until Wraparound |
|---|---|
| T1 (1.5 Mbps) | 6.4 hours |
| Ethernet (10 Mbps) | 57 minutes |
| T3 (45 Mbps) | 13 minutes |
| Fast Ethernet (100 Mbps) | 6 minutes |
| OC-3 (155 Mbps) | 4 minutes |
| OC-12 (622 Mbps) | 55 seconds |
| OC-48 (2.5 Gbps) | 14 seconds |

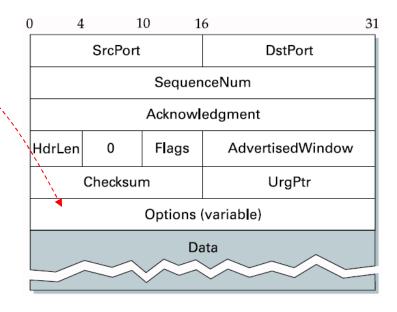**Time until 32-bit sequence number space wraps around**

# TCP: Can Keep Pipe Full?

| Bandwidth | Delay × Bandwidth Product |
|---|---|
| T1 (1.5 Mbps) | 18 KB |
| Ethernet (10 Mbps) | 122 KB |
| T3 (45 Mbps) | 549 KB |
| Fast Ethernet (100 Mbps) | 1.2 MB |
| OC-3 (155 Mbps) | 1.8 MB |
| OC-12 (622 Mbps) | 7.4 MB |
| OC-48 (2.5 Gbps) | 29.6 MB |

**Required window size for 100-ms RTT.**

# Solution: TCP Extensions

- Implemented as header options
- Store timestamp in outgoing segments → measure RTT
- Extend sequence space with 32-bit timestamp → protected against sequence number wrap-around
- Shift (scale) advertised window → keep the pipe full
- Selective acknowledgement (SAC) → acknowledge any additional (out-of-order) blocks of received data

| 0 | 4 | 10 | 16 | 31 |
|---|---|---|---|---|
| SrcPort | | | DstPort | |
| SequenceNum | | | | |
| Acknowledgment | | | | |
| HdrLen | 0 | Flags | AdvertisedWindow | |
| Checksum | | | UrgPtr | |
| Options (variable) | | | | |
| Data | | | | |

TCP Extensions for High Performance
http://tools.ietf.org/html/rfc1323

# Summary

- User Datagram Protocol
  - Multiplexer/Demultiplexer for IP
- Transmission Control Protocol
  - Reliable Byte Stream
    - Connection-oriented
      - Connection establishment
      - Connection termination
    - Automatics Repeated-Request: ACKs and NACKs
    - Flow-control
    - Timeout value estimation
    - Extensions
- *Congestion control (future discussions)*