


Copyright © William C. Cheng



## CS551

# Basic TCP Mechanisms

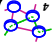
**Bill Cheng**

<http://merlot.usc.edu/cs551-f12>

---

Computer Communications - CSCI 551

Copyright © William C. Cheng



## What Does TCP Provide?

- ↳ Does TCP Provide...
- ↳ Connection establishment?
- ↳ Connectionless communication?
- ↳ Congestion control?
- ↳ Differentiated services?
- ↳ Duplicate packet detection?
- ↳ Flow control?
- ↳ Loss recovery?
- ↳ Message or record boundaries?
- ↳ Ordered data delivery?
- ↳ Out-of-order data delivery?
- ↳ Quality of service guarantees?
- ↳ Urgent data indication?

---


## Where And Why Is TCP Used?

- ↳ Where: anywhere reliable communication is needed
- ↳ everywhere (60-90% of traffic is TCP)
- ↳ file transfer/http, http, p2p, DNS (zone transfers), BGP, SMTP, remote login/telnet
- ↳ Why?
- ↳ connection oriented (reliability, ordering, ...)
- ↳ has the right features (congestion control, flow control, ...)
- ↳ widely deployed ⇒ interoperability, understood, exhaustively studied, pretty good implementations
- ↳ doing your own protocol is a lot of work

---

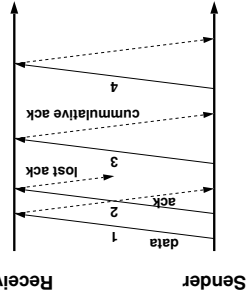
Computer Communications - CSCI 551

Copyright © William C. Cheng



## TCP Reliability Mechanism

- ↳ Summary or mechanisms
- ↳ window-based flow control
- ↳ sequence numbers, 3-way handshake, reliability (ACK, retransmission policies)
- ↳ congestion control
- ↳ RTT estimation

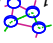


Sender Receiver

---

Computer Communications - CSCI 551

Copyright © William C. Cheng



## CS551

# Basic TCP Mechanisms


**Bill Cheng**

<http://merlot.usc.edu/cs551-f12>

---

Computer Communications - CSCI 551

Copyright © William C. Cheng



## What Does TCP Provide?

- ↳ Does TCP Provide...
- ↳ Connection establishment? **Y**
- ↳ Connectionless communication? **N**
- ↳ Congestion control? **Y**
- ↳ Differentiated services? **Y (sort of)**
- ↳ Duplicate packet detection? **Y**
- ↳ Flow control? **Y**
- ↳ Loss recovery? **Y**
- ↳ Message or record boundaries? **N**
- ↳ Ordered data delivery? **Y**
- ↳ Out-of-order data delivery? **N**
- ↳ Quality of service guarantees? **N**
- ↳ Urgent data indication? **Y**

---


## TCP Summary

- ↳ Communication abstraction:
- ↳ Reliable
- ↳ Ordered
- ↳ Point-to-point
- ↳ Byte-stream
- ↳ Protocol implemented entirely at the ends
- ↳ Assumes unreliable, non-sequenced delivery
- ↳ Fate sharing
- ↳ Operations
- ↳ OPEN/LISTEN, CONNECT, SEND, RECEIVE, ABORT

---

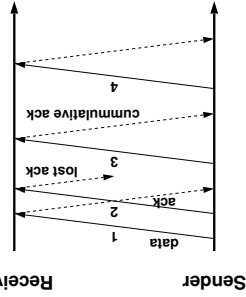
Computer Communications - CSCI 551

Copyright © William C. Cheng



## TCP Reliability Mechanism

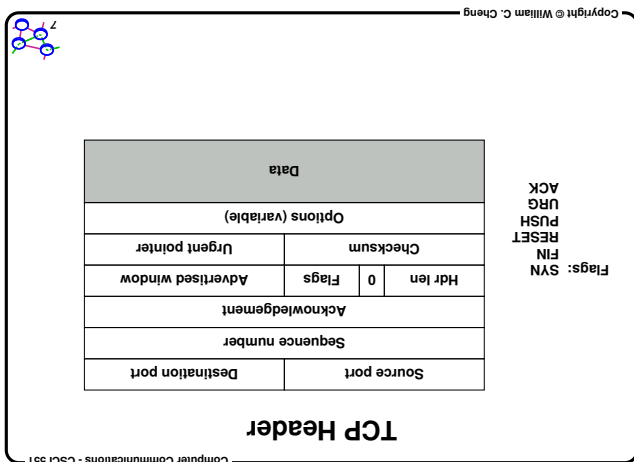
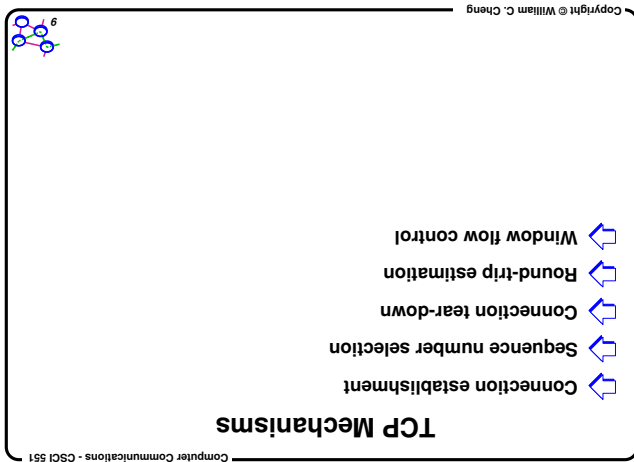
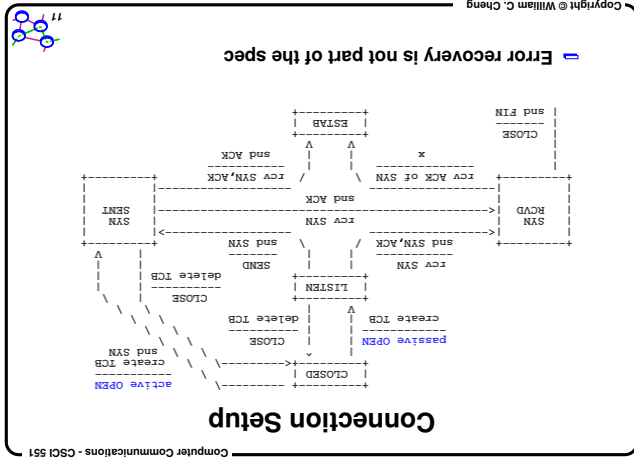
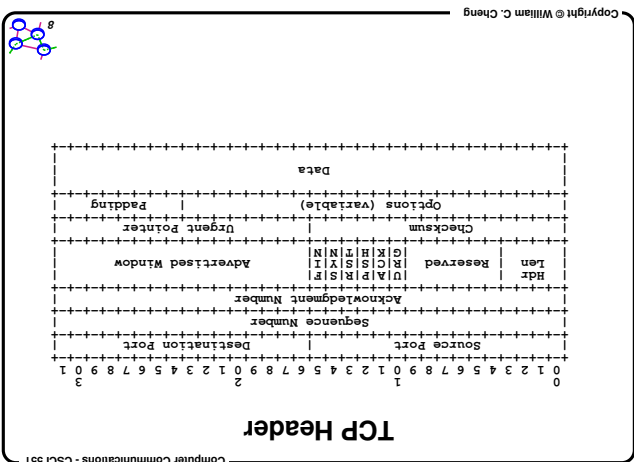
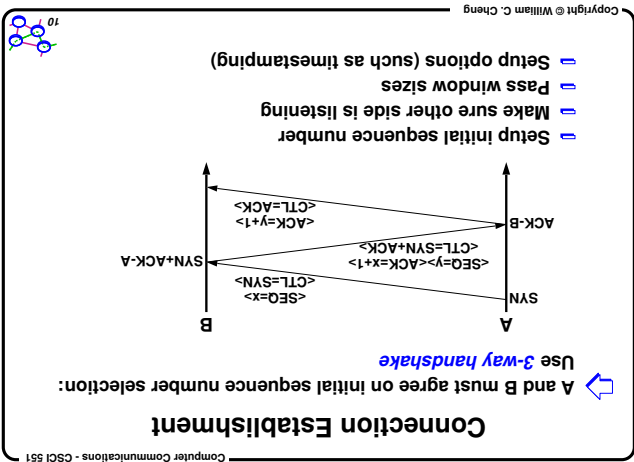
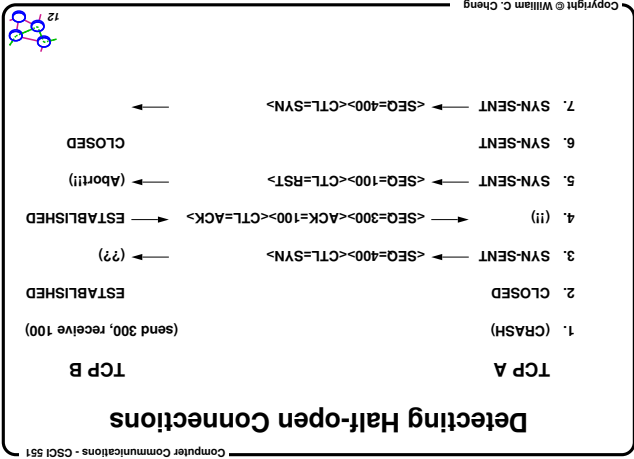
- ↳ Summary or mechanisms
- ↳ window-based flow control
- ↳ sequence numbers, 3-way handshake, reliability (ACK, retransmission policies)
- ↳ congestion control
- ↳ RTT estimation



Sender Receiver

---

Computer Communications - CSCI 551





Copyright © William C. Cheng

20

### Window Flow Control: Sender

Application writes here

Advised window (from receiver)

send buffer (from receiver)

Effective window

Last byte sent

Last byte ACKed

Sent but not acked

Not yet sent

Sequence numbers

Computer Communications - CSC1 551

Copyright © William C. Cheng

22

### Window Advancement Issues

- Advancing a full window
  - Slow receiver: receive buffer fills up (all data ACKed), last advertised window size is 0, sender is not allowed to send, how does it know the window has opened up?
  - solution: sender sends one data byte (probe) when advertised window size is 0
- Silly window syndrome
  - Fast sender, slow receiver: receiver with small buffer advertise it, sender quickly fills it with small amount of data
  - solution: delay ACK at receiver, Nagle's algorithm at sender
  - Nagle's algorithm -- send 1st partial packet, but not more until it is ACKed or have a full packet

(basically, stop & wait for small packets)

Computer Communications - CSC1 551

Copyright © William C. Cheng

24

### Timestamp Extension

- Used to improve timeout mechanism by more accurate measurement of RTT
- When sending a packet, insert current timestamp into option
- Receiver echoes timestamp in ACK

Computer Communications - CSC1 551

Copyright © William C. Cheng

19

### Flow Control

- Problem: Fast sender can overrun receiver
- Packet loss, unnecessary retransmissions
- TCP's solution:
  - Window sizes are passed in every packet to avoid sender overrunning the receiver
  - Sender transmits at pre-negotiated rate
  - Sender limited to a window's worth of unacknowledged data
- Flow control different from congestion control

Computer Communications - CSC1 551

Copyright © William C. Cheng

21

### Window Flow Control: Receiver

Application reads here

Receive buffer (possible window)

Missing Data

Receiver (Advised window)

Last byte received

Next byte expected

Acked but not delivered to user

Received but not ACKed yet

Sequence numbers

Computer Communications - CSC1 551

Copyright © William C. Cheng

23

### TCP Extensions

- Needed for high-bandwidth delay connections
- Accurate round-trip time estimation
- Window size limitations
- Impact of loss
- Implemented using TCP options
- Timestamp
- Protection from sequence number wraparound
- Large windows

Computer Communications - CSC1 551

Copyright © William C. Cheng

## Large Windows

- Recall that the AdvertisedWindow field is 16-bit long
- Delay × Bandwidth vs. link speed, assuming 100ms RTT
  - 1.5Mbps (T1): 18KB
  - 10Mbps (Ethernet): 122KB
  - 45Mbps (T3): 549KB
  - 100Mbps (FDDI): 1.2MB
  - 155Mbps (STS-3): 1.8MB
  - 622Mbps (STS-12): 7.4MB
  - 1.2Gbps (STS-24): 14.8MB
- Apply scaling factor to advertised window
  - Specifies how many bits window must be shifted to the left
- Scaling factor exchanged during connection setup

Computer Communications - CSCI 551

Copyright © William C. Cheng

## Protection From Wraparound

- Wraparound time vs. link speed:
  - 1.5Mbps (T1): 6.4 hours
  - 10Mbps (Ethernet): 57 minutes
  - 45Mbps (T3): 13 minutes
  - 100Mbps (FDDI): 6 minutes
  - 155Mbps (STS-3): 4 minutes
  - 622Mbps (STS-12): 55 seconds
  - 1.2Gbps (STS-24): 28 seconds
- Recall that MSL is 2 minutes in TCP
- Use timestamp (32-bit) to distinguish sequence number wraparound

Computer Communications - CSCI 551