# Lecture 7

The Transport Layer (2)

# Recap: Reliability via "Stop-and-Wait"

- **Sender key ideas:**

  - Sender sends one packet of application data, then **waits** for **acknowledgment** that receiver got it (ACK)

  - Upon receiving an ACK for that packet, sender can move on to send next data packet

    - Which packet an ACK corresponds to is determined based on **sequence number**

  - If the ACK isn't received within a certain **timeout**, sender retransmits

- **Receiver key ideas:**

  - Upon receiving a data packet, the receiver sends an ACK **for that packet** (i.e. carrying same **sequence number**)

  - If this is a **new** data packet (i.e. next expected sequence number), delivers it "up" to application

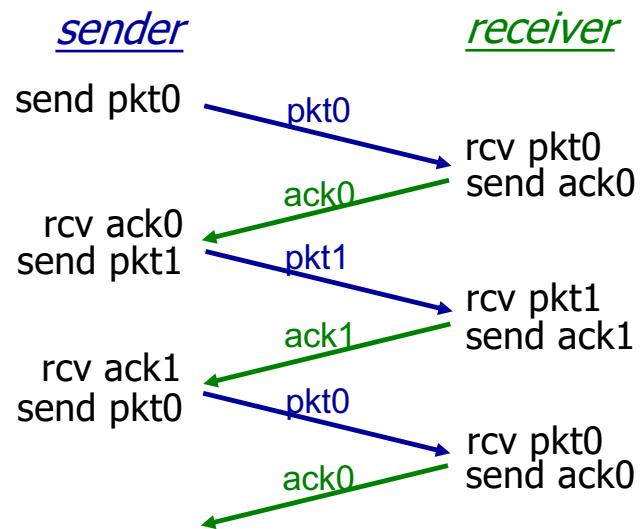# Complete Reliable Data Transfer Protocol Using the Stop-and-Wait Approach

**Sender:**
- Set seq=0
- Wait for data to send **(1)**
- Create and send packet(seq, data)
- Start timer
- Wait for ACK or Timeout **(2)**
  - If ACK && not corrupt && ACK_seq == seq:
    - Stop timer
    - seq = (seq+1) % 2
    - Wait for more data (1)
  - Else if Timeout:
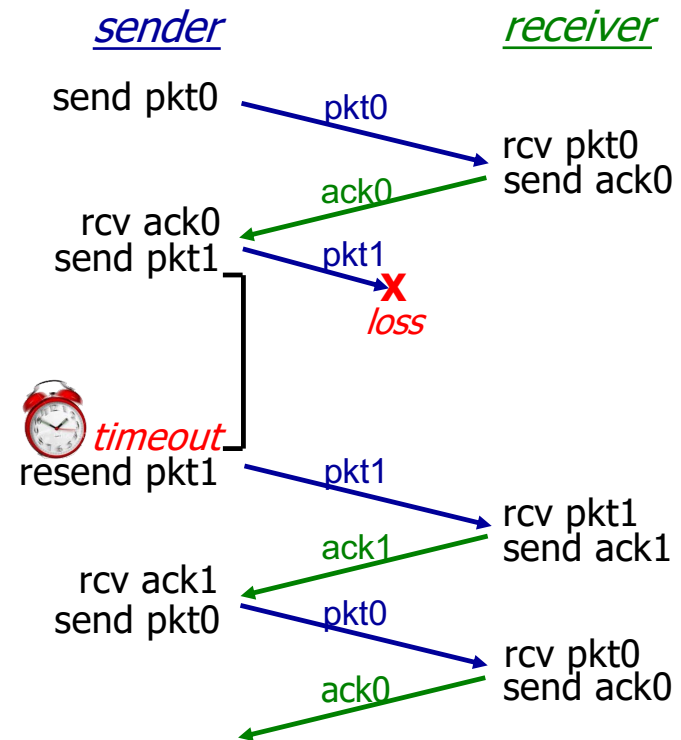    - Resend packet
    - Restart timer
    - Wait for ACK (2)

**Receiver:**
- Set seq=0
- Wait for packet from sender **(1)**
  - If not corrupt:
    - Send ACK(recvd_seq)
    - If recvd_seq == seq:
      - Extract and deliver to application
      - seq = (seq+1) % 2
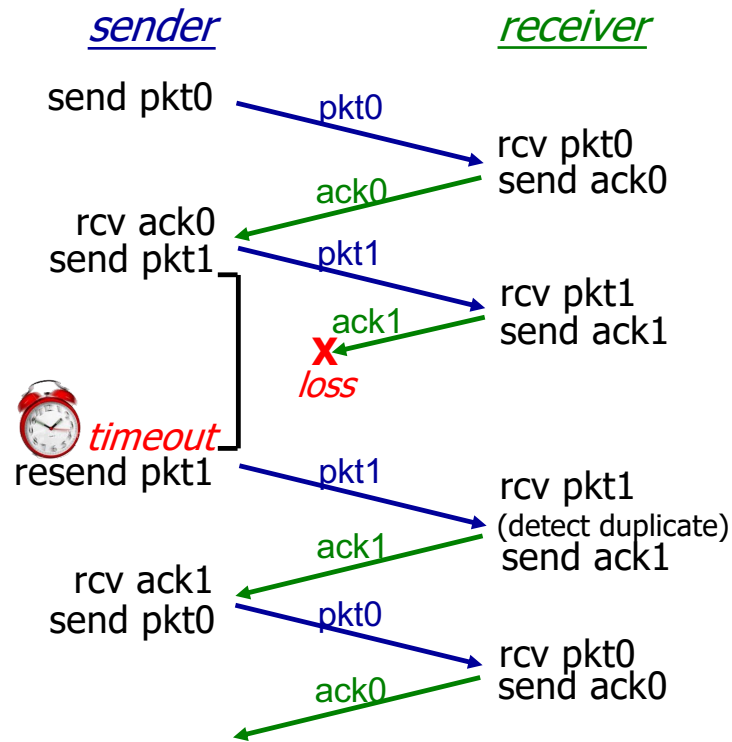- Wait for new packet from sender (1)
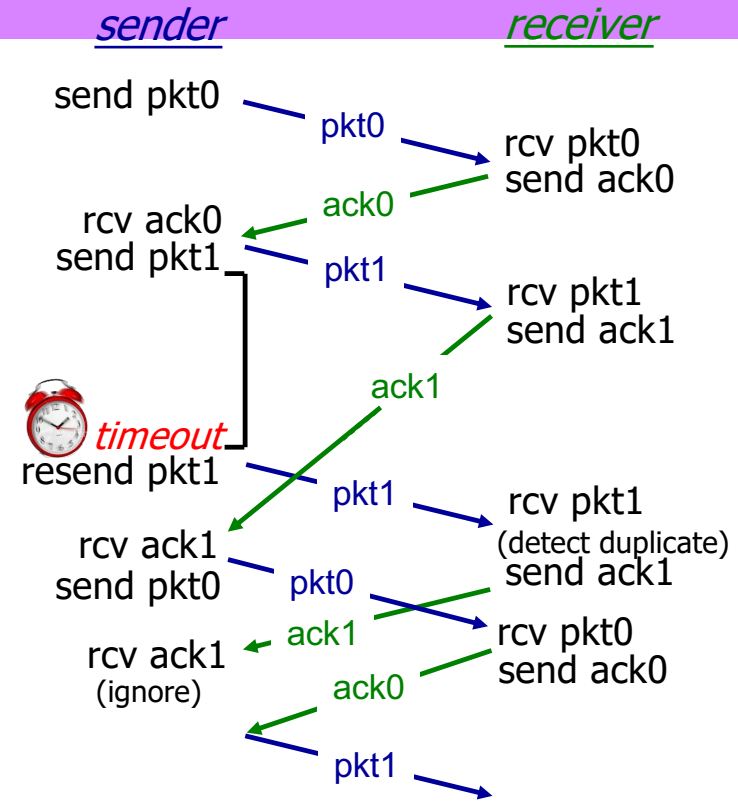
3

# Stop-and-Wait: Example Executions

sender                                                receiver

send pkt0 ──── pkt0 ────►
                                                rcv pkt0
                                                send ack0
rcv ack0   ◄──── ack0 ────
send pkt1  ──── pkt1 ────►
                                                rcv pkt1
                                                send ack1
rcv ack1   ◄──── ack1 ────
send pkt0  ──── pkt0 ────►
                                                rcv pkt0
                                                send ack0
           ◄──── ack0 ────

(a) no loss

sender                                                receiver

send pkt0 ──── pkt0 ────►
                                                rcv pkt0
                                                send ack0
rcv ack0   ◄──── ack0 ────
send pkt1  ──── pkt1 ──►✗
                                                loss

🕐 *timeout*
resend pkt1 ──── pkt1 ────►
                                                rcv pkt1
                                                send ack1
rcv ack1   ◄──── ack1 ────
send pkt0  ──── pkt0 ────►
                                                rcv pkt0
                                                send ack0
           ◄──── ack0 ────

(b) packet loss

# Stop-and-Wait: Example Executions

sender                          receiver

send pkt0 — pkt0 →
                                rcv pkt0
                                send ack0
rcv ack0 ← ack0
send pkt1 — pkt1 →
                                rcv pkt1
                                send ack1
                        ack1
                        X
                        loss
timeout
resend pkt1 — pkt1 →
                                rcv pkt1
                                (detect duplicate)
                                send ack1
rcv ack1 ← ack1
send pkt0 — pkt0 →
                                rcv pkt0
                                send ack0
            ← ack0

(c) ACK loss

sender                          receiver

send pkt0 — pkt0 →
                                rcv pkt0
                                send ack0
rcv ack0 ← ack0
send pkt1 — pkt1 →
                                rcv pkt1
                                send ack1
                        ack1
timeout
resend pkt1 — pkt1 →
                                rcv pkt1
                                (detect duplicate)
rcv ack1                        send ack1
send pkt0 — pkt0 →
rcv ack1 ← ack1
(ignore)                        rcv pkt0
                        ack0    send ack0
                        — pkt1 →

(d) premature timeout/ delayed ACK

# Complete Reliable Data Transfer Protocol Using the Stop-and-Wait Approach

- **Why don't we use duplicate ACKs as evidence of loss?**

- In the preliminary versions of this protocol that we discussed last week, we said the receiver could send a NAK or duplicate ACK after receiving a corrupted message to prompt the sender to retransmit. The final version relies on timeouts only. **Why**?

- i.e. Why not retransmit immediately in the case where ACK seq doesn't match our expected seq?

# Stop-and-Wait RDT Operation: Premature Timeouts

sender                              receiver

send pkt0
                    pkt0
                              rcv pkt0
                              send ack0
rcv ack0            ack0
send pkt1
                    pkt1
                              rcv pkt1
                              send ack1
                    ack1

timeout
resend pkt1
                    pkt1
                              rcv pkt1
                              (detect duplicate)
rcv ack1                      send ack1
send pkt0           pkt0
                    ack1
rcv ack1                      rcv pkt0
(ignore)                      send ack0
                    ack0
                    pkt1

**What would happen if we didn't ignore this ack?**

- Sender retransmits pkt0
- Receiver gets pkt0, assumes ack was lost, retransmits ack0
- Sender gets ack0 after having sent pkt1, so retransmits pkt1
- Receiver gets pkt1 retransmission, retransmits ack1
- ...

**If** the pkt was really lost, retransmitting on duplicate ack lets us respond faster.
**But**, if not, generates a lot of extra traffic (without additional protocol changes)
**And,** sender cannot tell which case it's in

7

# Reliable Data Transfer: Mechanism Summary

- **Checksums**: detect errors
- **ACKs/NAKs**: provide sender with feedback about what has been received
- **Retransmissions**: sender resends lost/corrupted packets
- **Sequence numbers**: identify packets, allow de-duplication
- **Timeouts**: detect (probable) loss, decide when to retransmit

# Analyzing Stop-and-Wait

- Our protocol is finally correct…but is it good?

# Analyzing Stop-and-Wait

- Our protocol is finally correct…but is it good?

sender            receiver

first packet bit transmitted, $t = 0$

first packet bit arrives

last packet bit arrives, send ACK

RTT

ACK arrives, send next packet, $t = RTT + L / R$
(assumes 0 transmission delay for ACK)

# Analyzing Stop-and-Wait

- Our protocol is finally correct…but is it good?

RTT ~60ms
Bandwidth ~1Gbps
Segment size ~1500bytes

data packet →

(a) a stop-and-wait protocol in operation

**Transmission delay**:
$(1500*8)/(10^9) = 12$ microsec

Assuming negligible transmission delay for ACK and no queuing/processing, we wait 60.012 ms to send next segment

**How long would it take to transfer a 1Gb file at this rate?**

# Improving Performance

- **How can we improve our protocol to get better performance?**
  - Think back to HTTP discussion…what strategies did we use there?

# Improving Performance

- **How can we improve our protocol to get better performance?**
  - Pipelining – instead of stopping and waiting after **every** packet, we can send **multiple** packets to "fill the pipe" before waiting for acknowledgment

(a) a stop-and-wait protocol in operation          (b) a pipelined protocol in operation

# Pipelined Reliable Data Transfer

- Instead of only 1 unacknowledged packet, sender can have **up to N unacknowledged packets** at a time
  - N packets "in-flight"
  - Requires more than 1-bit sequence number

- Sender needs to **buffer** sent but not-yet-acknowledged packets so they can be retransmitted as needed

- Sender buffer typically operates as a "**sliding window**"

base    N

Sent and ACK'd

Sent but not ACK'd

Not yet sent

*sequence number* →

next_seq

# Sliding Window

- Sender buffer typically operates as a "**sliding window**"
  - base is the sequence number of the start of window (last ACK'd packet + 1)
  - next_seq is the sequence number that will be used for the next packet to send

# Sliding Window

- Sender buffer typically operates as a "**sliding window**"
  - base is the sequence number of the start of window (last ACK'd packet + 1)
  - next_seq is the sequence number that will be used for the next packet to send



base    N

sequence number →

next_seq

usable sequence numbers

Sent and ACK'd

Sent but not ACK'd

Not yet sent

Cannot be used (yet)

# Sliding Window

- Sender buffer typically operates as a "**sliding window**"
  - once the window is full (we've used up all the sequence numbers between base and base+N), no new packets can be sent…

base      N

sequence number →

next_seq

Sent and ACK'd

Sent but not ACK'd

Not yet sent

# Sliding Window

- Sender buffer typically operates as a "**sliding window**"
  - …until a new ACK lets us "slide" the window forward

base         N

*sequence number* →

next_seq

■ Sent and ACK'd

■ Sent but not ACK'd

□ Not yet sent

# A Pipelined Reliable Data Transfer Protocol: "Go-Back-N"

- Go-Back-N protocol key ideas:
  - Sender can have **up to N unacknowledged packets** "in-flight" at a time
  - Receiver sends **cumulative acknowledgments**
    - Cumulative acknowledgement for sequence number X means I have received everything **up to and including** sequence number X
      - sometimes called ARU for "all received up-to"
    - Acknowledgments let the sender "slide its window forward" and send new packets until there are again N unacknowledged packets
  - When a **timeout** occurs, the sender will **"Go-Back-N"** and resend **ALL not-yet-acknowledged packets** (of which there can be up to N)

# Go-Back-N: receiver

- ACK-only: always send ACK for correctly-received packet so far, with highest *in-order* seq #
  - may generate duplicate ACKs
  - need only remember `rcv_base`

- on receipt of out-of-order packet:
  - can discard (don't buffer) or buffer: an implementation decision
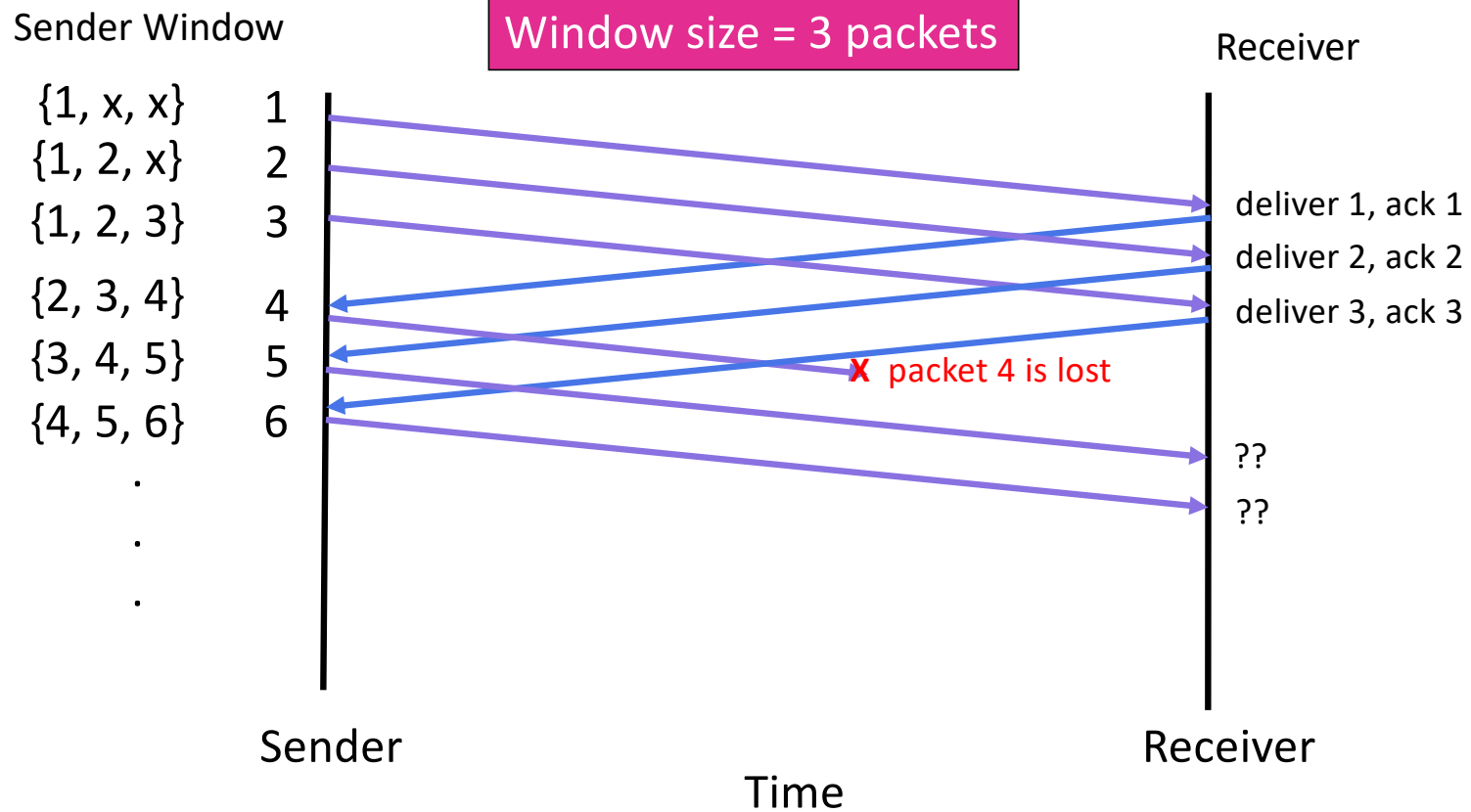  - re-ACK pkt with highest in-order seq #
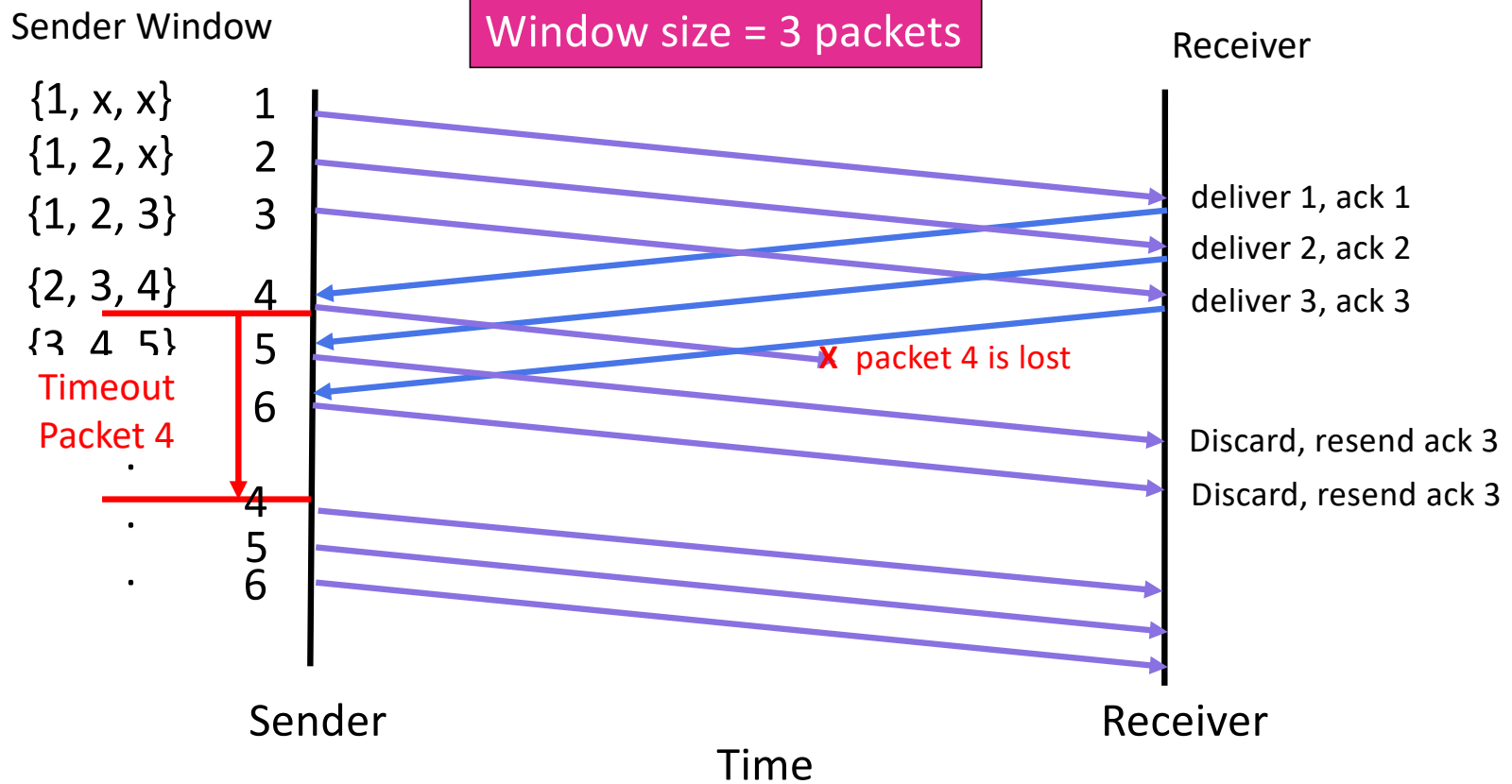
Receiver view of sequence number space:

... |||||||||||||| ...

rcv_base

| received and ACKed

| Out-of-order: received but not ACKed

| Not received

# Go-Back-N Normal Case Operation

Sender Window

Window size = 3 packets

Receiver

| Sender Window | | |
|---|---|---|
| {1, x, x} | 1 | |
| {1, 2, x} | 2 | |
| {1, 2, 3} | 3 | |
| {2, 3, 4} | 4 | |
| {3, 4, 5} | 5 | |
| {4, 5, 6} | 6 | |

deliver 1, ack 1
deliver 2, ack 2
deliver 3, ack 3

deliver 4, ack 4
deliver 5, ack 5
deliver 6, ack 6

Sender

Receiver

Time

# Go-Back-N Operation with Loss

Sender Window

Window size = 3 packets

Receiver

| Sender Window | | |
|---|---|---|
| {1, x, x} | 1 | |
| {1, 2, x} | 2 | |
| {1, 2, 3} | 3 | |
| {2, 3, 4} | 4 | |
| {3, 4, 5} | 5 | |
| {4, 5, 6} | 6 | |

deliver 1, ack 1
deliver 2, ack 2
deliver 3, ack 3

X packet 4 is lost

??
??

Sender

Receiver

Time

# Go-Back-N Operation with Loss



Sender Window

Window size = 3 packets

Receiver

{1, x, x}  1
{1, 2, x}  2
{1, 2, 3}  3

{2, 3, 4}  4
{3, 4, 5}  5

Timeout
Packet 4

4
5
6

deliver 1, ack 1
deliver 2, ack 2
deliver 3, ack 3

X packet 4 is lost

6

Discard, resend ack 3
Discard, resend ack 3

Sender

Receiver

Time

23

# Go-Back-N: Details

- Sender
  - Sends up to N unacknowledged packets
  - Maintains timer for **oldest unacknowledged packet**
    - Set timer when sending first packet in empty window
    - Reset timer when receiving ack that moves up window (or stop if no more un-ack'd packets)
    - (and reset timer after resending window)
  - Resends ALL un-ack'd packets on timeout
- Receiver
  - Sends cumulative ACK upon receiving next expected packet
  - Resends last cumulative ACK upon receiving out-of-order packet
  - Discards packets received out of order (some variants allow these to be buffered – why?)

# Go-Back-N: Extended FSM (Sender)

rdt_send(data)
_____

if (nextseqnum < base+N) {
   sndpkt[nextseqnum] = make_pkt(nextseqnum,data,chksum)
   udt_send(sndpkt[nextseqnum])
   if (base == nextseqnum)
     start_timer
   nextseqnum++
   }
else
 refuse_data(data)

if there is space in window, send new data
(starting timer if not already set)

$\Lambda$
_____
base=1
nextseqnum=1

Wait

timeout
_____
start_timer
udt_send(sndpkt[base])
udt_send(sndpkt[base+1])
…
udt_send(sndpkt[nextseqnum-1])

reset timer and resend all packets in window

rdt_rcv(rcvpkt)
  && corrupt(rcvpkt)
_____
$\Lambda$

rdt_rcv(rcvpkt) &&
  notcorrupt(rcvpkt)
_____
base = getacknum(rcvpkt)+1
If (base == nextseqnum)
  stop_timer
 else
  start_timer

slide window forward and reset timer (if there are still unacknowledged packets in window)

25

# Go-Back-N: Extended FSM (Receiver)

any other event
————————————
udt_send(sndpkt)

resend cumulative ack

$\Lambda$
————————————
expectedseqnum=1
sndpkt =
  make_pkt(0,ACK,chksum)

Wait

rdt_rcv(rcvpkt)
  && notcorrupt(rcvpkt)
  && hasseqnum(rcvpkt,expectedseqnum)
————————————
extract(rcvpkt,data)
deliver_data(data)
sndpkt = make_pkt(expectedseqnum,ACK,chksum)
udt_send(sndpkt)
expectedseqnum++

received next expected
packet, so deliver to
application, and send ack

# Why limit it to *N* packets?

- Flow Control (later)

- Congestion Control (later)

- Fixed size for sequence numbers
  - How many possible sequence numbers if we have a *k* bit field?

# Remarks

- Receiver: Data delivered to application layer one segment at a time
  - If data from segment $n$ is sent upward, it means
    - Segment $n$ was received correctly, and all previous segments were received correctly and delivered to the application layer
- Receiver discards out of order segments… what a waste?
  - If segment $n$ is lost, no point in keeping segment $n + 1$
  - Remember the sender sends all packets in the window!
- Event based programming
  - Call from upper layers, timer interrupts, call from lower layers

# Go-Back-N in action

sender window (N=4)      sender      receiver

0 1 2 3 4 5 6 7 8    send pkt0

0 1 2 3 4 5 6 7 8    send pkt1

0 1 2 3 4 5 6 7 8    send pkt2      receive pkt0, send ack0

0 1 2 3 4 5 6 7 8    send pkt3    **X** loss    receive pkt1, send ack1

(wait)      receive pkt3, discard,
     (re)send ack1

0 1 2 3 4 5 6 7 8    rcv ack0, send pkt4

0 1 2 3 4 5 6 7 8    rcv ack1, send pkt5    receive pkt4, discard,
     (re)send ack1

ignore duplicate ACK    receive pkt5, discard,
     (re)send ack1

pkt 2 timeout

0 1 2 3 4 5 6 7 8    send pkt2

0 1 2 3 4 5 6 7 8    send pkt3

0 1 2 3 4 5 6 7 8    send pkt4    rcv pkt2, deliver, send ack2

0 1 2 3 4 5 6 7 8    send pkt5    rcv pkt3, deliver, send ack3

     rcv pkt4, deliver, send ack4

     rcv pkt5, deliver, send ack5

# Another Approach: "**Selective Repeat**"

- Sender can send up to N unacknowledged packets
  - *N* consecutive seq #s
  - limits seq #s of sent, unACKed packets
- Receiver **individually** acknowledges each correctly received packet and **buffers** out-of-order packets
- Sender maintains **separate** timer for each un-ACK'd packet to retransmit individually (i.e., repeat selectively)

# Another Approach: "Selective Repeat"

- Sender can send up to N unacknowledged packets

- Receiver **individually** acknowledges each correctly received packet and **buffers** out-of-order packets

- Sender maintains **separate** timer for each un-ACK'd packet to retransmit individually (i.e. repeat selectively)



base    N

sequence number →

next_seq

■ Sent and ACK'd

■ Sent but not ACK'd

□ Not yet sent

# Another Approach: "Selective Repeat"

- Sender buffer:

sendbase       N

*sequence number* →

next_seq

🟩 Sent and ACK'd

🟨 Sent but not ACK'd

⬜ Not yet sent

- Receiver buffer:

recvbase      N

*sequence number* →

🟩 Received and ACK'd

⬛ Not received (but expected)

⬜ Not received

# Another Approach: "Selective Repeat"

- Sender buffer:

sendbase       N

sequence number →

next_seq

- Receiver buffer:

recvbase     N

sequence number →

| | |
|---|---|
| 🟩 | Sent and ACK'd |
| 🟨 | Sent but not ACK'd |
| ⬜ | Not yet sent, not usable seq (outside window) |
| ⬜ | Not yet sent, usable seq (in window) |
| 🟩 | Received and ACK'd |
| ⬜ | Not received (but expected) |
| ⬜ | Not received, outside window |
| ⬜ | Not received, able to accept (in window) |

# Selective Repeat Operation with Loss

Sender Window

Window size = 3 packets

Receiver

| | |
|---|---|
| {1, x, x} | 1 |
| {1, 2, x} | 2 |
| {1, 2, 3} | 3 |
| {2, 3, 4} | 4 |
| {3, 4, 5} | 5 |
| {4, 5, 6} | 6 |

deliver 1, ack 1

deliver 2, ack 2

deliver 3, ack 3

**X** packet 4 is lost

??

??

Sender

Receiver

Time

# Selective Repeat Operation with Loss

Sender Window

Window size = 3 packets

Receiver

{1, x, x}   1
{1, 2, x}   2
{1, 2, 3}   3
                  deliver 1, ack 1
{2, 3, 4}   4
                  deliver 2, ack 2
{3, 4, 5}   5
                  deliver 3, ack 3
Timeout
Packet 4    6
                  **X** packet 4 is lost
            4
                  buffer 5, ack 5
{4, 5, 6}
                  buffer 6, ack 6
{4, 5, 6}
                  deliver 4,5,6 ack 4
{7, x, x}   7

Sender          Time          Receiver

# Selective Repeat Details

## sender

**data from above:**

- if next available seq # in window, send packet & start timer

**timeout($n$):**

- resend packet $n$, restart timer

**ACK($n$) in [sendbase,sendbase+N-1]:**

- mark packet $n$ as received
- if n smallest unACKed packet, advance window base to next unACKed seq #

## receiver

**packet $n$ in [rcvbase, rcvbase+N-1]**

- send ACK($n$)
- out-of-order: buffer
- in-order: deliver (also deliver buffered, in-order packets), advance window to next not-yet-received packet

**packet $n$ in [rcvbase-N,rcvbase-1]**

- ACK($n$)

**otherwise:**

- ignore

# Selective Repeat Details

## sender

data from above:

- if next available seq # in window

tir...

AC...

- if n smallest unACKed packet, advance window base to next unACKed seq #

**Why?**
- Why do we need to acknowledge old packets?
- Why don't we need to acknowledge more than N before our window start?

## receiver

packet *n* in [rcvbase, rcvbase+N-1]

- send ACK(*n*)
- out-of-order: buffer
- in-order: deliver (also deliver buffered, in-order packets), advance window to next not-yet-received packet

packet *n* in [rcvbase-N,rcvbase-1]

- ACK(*n*)

otherwise:

- ignore

# Asymmetric Knowledge...

- Why do we need to acknowledge old packets?
  - Same basic reason as in our previous protocols...ACKs can get lost, so sender may not have moved up window yet (and still needs to recv ACK to do so)

- Why don't we need to acknowledge more than N before our window start?
  - Sender window can't be **too** far behind receiver window (because it can't send more than N beyond its own window start)
  - e.g. if I receive seq 9 and window size is 4, sender can't still be waiting for anything before 6

# Subtle Points

- In reality, sequence numbers are **finite**
  - e.g. 16-bit sequence number -> $2^{16}$ = 65,536 distinct numbers;
  - 32-bit sequence number -> $2^{32}$ = 4,294,967,296 distinct numbers
- Consider sequence numbers from 0-5, window size = 4:

0 1 2 3 4 5 0 1 2 3 □□□□□□□□□  sender

0 1 2 3 4 5 0 1 2 3 4 5 0 1 □□□□□  receiver

# Subtle Points

- In reality, sequence numbers are **finite**
  - e.g. 16-bit sequence number -> 2^16 = 65,536 distinct numbers;
  - 32-bit sequence number -> 2^32 = 4,294,967,296 distinct numbers

- Consider sequence numbers from 0-5, window size = 4:



0 1 2 3 4 5 0 1 2 3      sender

0 1 2 3 4 5 0 1 2 3 4 5 0 1      receiver

receiver has received 0-3, so it moves up its window…

…but the sender doesn't know that yet! (e.g. because acks were lost)

So, sender will retransmit 0, 1, 2, 3
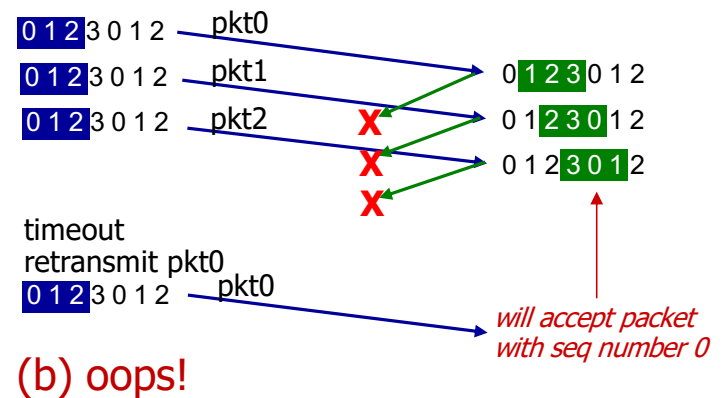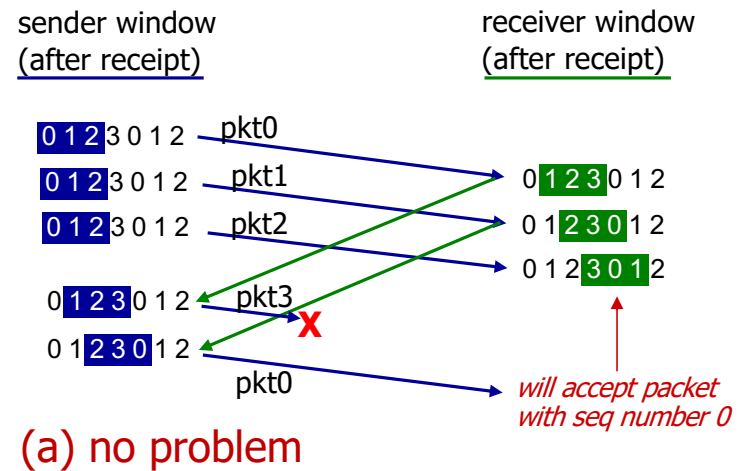
**What will the receiver do?**

# Subtle Points

- In reality, sequence numbers are **finite**
  - e.g. 16-bit sequence number -> $2^{16}$ = 65,536 distinct numbers;
  - 32-bit sequence number -> $2^{32}$ = 4,294,967,296 distinct numbers

- Consider sequence numbers from 0-5, window size = 4:

0 1 2 3 4 5 0 1 2 3         sender

0 1 2 3 4 5 0 1 2 3 4 5 0 1         receiver

receiver has received 0-3, so it moves up its window…

…but the sender doesn't know that yet! (e.g. because acks were lost)

So, sender will retransmit 0, 1, 2, 3

**Receiver can't differentiate retransmission from new message reusing its sequence number!**

# Subtle Points

- In reality, sequence numbers are **finite**
  - e.g. 16-bit sequence number -> 2^16 = 65,536 distinct numbers;
  - 32-bit sequence number -> 2^32 = 4,294,967,296 distinct numbers

- Consider sequence numbers from 0-5, window size = 4:

**3**

`0 1 2 3 4 5 0 1 2 3                    ` sender

`0 1 2 3 4 5 0 1 2 3 4 5 0 1            ` receiver

receiver has received 0-3, so it moves up its window…

…but the sender doesn't know that yet! (e.g. because acks were lost)

So, sender will retransmit 0, 1, 2, 3

**Solution: need to shrink window OR expand sequence range**

# Selective repeat: a dilemma!

example:

- seq #s: 0, 1, 2, 3 (base 4 counting)
- window size=3

sender window
(after receipt)

receiver window
(after receipt)

0 1 2 3 0 1 2  pkt0
0 1 2 3 0 1 2  pkt1
0 1 2 3 0 1 2  pkt2
0 1 2 3 0 1 2  pkt3   X
0 1 2 3 0 1 2
                pkt0

0 1 2 3 0 1 2
0 1 2 3 0 1 2
0 1 2 3 0 1 2
0 1 2 3 0 1 2

*will accept packet
with seq number 0*

## (a) no problem

0 1 2 3 0 1 2  pkt0
0 1 2 3 0 1 2  pkt1
0 1 2 3 0 1 2  pkt2   X
                      X
                      X
timeout
retransmit pkt0
0 1 2 3 0 1 2  pkt0

0 1 2 3 0 1 2
0 1 2 3 0 1 2
0 1 2 3 0 1 2

*will accept packet
with seq number 0*

## (b) oops!

# Selective repeat: a dilemma!

example:

- seq #s: 0, 1, 2, 3 (base 4 counting)
- window size=3

Q: what relationship is needed between sequence # size and window size to avoid problem in scenario (b)?

0 1 2 3 0 1 2    pkt0
0 1 2 3 0 1 2    pkt1
0 1 2 3 0 1 2    pkt2

0 1 2 3 0 1 2    pkt3

0 1 2 3 0 1 2

0 1 2 3 0 1 2

0 1 2 3 0 1 2

- *receiver can't see sender side*
- *receiver behavior identical in both cases!*
- *something's (very) wrong!*

*will accept packet with seq number 0*

0 1 2 3 0 1 2    pkt2

timeout
retransmit pkt0
0 1 2 3 0 1 2    pkt0

0 1 2 3 0 1 2

0 1 2 3 0 1 2

0 1 2 3 0 1 2

(b) oops!

*will accept packet with seq number 0*

# Subtle Points

- In reality, sequence numbers are **finite**
  - e.g. 16-bit sequence number -> 2^16 = 65,536 distinct numbers;
  - 32-bit sequence number -> 2^32 = 4,294,967,296 distinct numbers

- Consider sequence numbers from 0-5, window size = 3:

```
0 1 2 3 4 5 0 1 2 3                    sender

0 1 2 3 4 5 0 1 2 3 4 5 0 1            receiver
```

What if this message needed to be retransmitted multiple times, and some copy is **delayed** in the network such that it arrives now?

So far, we have assumed messages arrive in the **order** they were sent. **What happens if this is not true?**

# Summary

- We can build a reliable transport service on top of an unreliable network

- **Pipelining** is used to improve performance and make reliable data transfer practical

- We need to be careful in reasoning about sender and receiver behavior because they may have different views of the current state (asymmetric knowledge)

# Reliable Data Transfer - Recap

- Stop-and-wait
  - Send 1 packet, wait until it is acknowledged to send the next
- Go-Back-N
  - Sender can have up to N unacknowledged packets at any time
  - Receiver discards out of order packets and sends cumulative ACKs
  - Upon timeout for oldest unacknowledged packet, sender resends ALL unacknowledged packets
- Selective Repeat
  - Sender can have up to N unacknowledged packets at any time
  - Receiver buffers out of order packets (that fall within window) and sends selective ACKs
  - Upon timeout for a specific packet, sender resends that packet

# TCP Service Abstraction

- **TCP is a connection-oriented protocol that delivers a reliable, in-order, byte stream**

- **Connection-oriented**: two processes coordinate ("handshake") before beginning to send data to each other

- **Reliable**: TCP resends lost packets
  - Until it gives up and shuts down connection

- **In-order**: TCP only hands consecutive chunks of data to application

- **Byte stream**: TCP assumes there is an incoming stream of data, and attempts to deliver it to application

# TCP Mechanisms

- Builds on most of what we've seen so far:
  - Checksums
  - Sequence numbers (byte offsets)
  - Sender and receiver maintain a sliding window
  - Receiver sends cumulative acknowledgements (like **Go-Back-N**)
    - Sender maintains a single retransmission timer
  - Receivers can buffer out-of-sequence packets (like **Selective Repeat**)
    - technically, this is optional…not defined by TCP spec
- And we'll see a few more soon: fast retransmit, timeout estimation algorithms

# TCP Segment Structure

Used for
Multiplexing /
Demultiplexing
(similar to UDP,
but remember that
TCP
demultiplexing
actually uses
source port/IP)

32 bits

| Source port | Destination port |
|---|---|
| Sequence number | |
| Acknowledgment | |

| HdrLen | 0 | Flags | Advertised window |
|---|---|---|---|

| Checksum | Urgent pointer |
|---|---|
| Options (variable) | |
| Data | |

# TCP Segment Structure

Used for error detection (same as UDP)

<--- 32 bits --->

| Source port | Destination port |
|---|---|
| Sequence number | |
| Acknowledgment | |

| HdrLen | 0 | Flags | Advertised window |
|---|---|---|---|

| Checksum | Urgent pointer |
|---|---|

| Options (variable) |
|---|

| Data |
|---|

# TCP Segment Structure

As discussed last time, used for maintaining **ordered** delivery and differentiating new vs duplicate data

32 bits

| | |
|---|---|
| Source port | Destination port |
| Sequence number | |
| Acknowledgment | |
| HdrLen | 0 | Flags | Advertised window |
| Checksum | Urgent pointer |
| Options (variable) | |
| Data | |

# TCP Segment Structure

As discussed last time, used for maintaining **ordered** delivery and differentiating new vs duplicate data

**32 bits**

| Source port | Destination port |
|---|---|
| Sequence number | |
| Acknowledgment | |
| HdrLen | 0 | Flags | Advertised window |
| Checksum | Urgent pointer |
| Options (variable) | |
| Data | |

**But**, for TCP, the sequence number is a **byte offset**, **not** a packet id

# TCP Byte Stream Service…

Application @ Host A



Application @ Host B

# ...provided using TCP segments

Host A



TCP Data

*Segment* sent when:
1. Segment full (Max Segment Size),
2. Not full, but times out

TCP Data

Host B

# TCP Segments

- IP packet
  - No bigger than Maximum Transmission Unit (MTU)
  - E.g., up to 1500 bytes with Ethernet
- TCP segment
  - IP payload contains segment composed of TCP header and data
  - TCP header $\geq$ 20 bytes long
- TCP payload/data
  - No more than Maximum Segment Size (MSS) bytes
  - E.g., up to 1460 consecutive bytes from the stream
  - MSS = MTU – (IP header) – (TCP header)

| IP Data / TCP Segment | | |
|---|---|---|
| TCP Payload/Data | TCP Hdr | IP Hdr |

# TCP Sequence Numbers

ISN (Initial Sequence Number)

k bytes

Host A

Sequence number
= 1st byte in segment =
ISN + k

# TCP Sequence Numbers

ISN (Initial Sequence Number)

k

Host A

Sequence number
= 1st byte in segment =
ISN + k

TCP Data

TCP HDR

ACK sequence number
= next expected byte
= seqno + length(data)

TCP Data

TCP HDR

Host B

# TCP Segment Structure

32 bits

| Source port | Destination port |
|---|---|
| Sequence number | |
| Acknowledgment | |

| HdrLen | 0 | Flags | Advertised window |
|---|---|---|---|

| Checksum | Urgent pointer |
|---|---|

| Options (variable) |
|---|

| Data |
|---|

Starting byte offset of data carried by this segment

# TCP Segment Structure

Sequence of **next expected** byte (i.e. sequence number just **after** the last byte received **in order** so far)

← 32 bits →

| Source port | Destination port |
| --- | --- |
| Sequence number | |
| Acknowledgment | |

| HdrLen | 0 | Flags | Advertised window |
| --- | --- | --- | --- |

| Checksum | Urgent pointer |
| --- | --- |

| Options (variable) |
| --- |

| Data |
| --- |

# TCP ACK Behavior

- Sender sends segment
  - Data starts with sequence number X
  - Segment contains B bytes [X, X+1, X+2, …, X+B-1]

- Upon receipt of segment, receiver sends an ACK
  - If all data prior to X already received:
    - ACK acknowledges X+B (because that is next expected byte)
  - If highest in-order byte received is Y s.t. (Y+1) < X
    - ACK acknowledges Y+1
    - Even if this has been ACKed before

Why?

# TCP ACK Behavior: Normal (no loss) Case

- Sender: seqno=X, length=B

- Receiver: ACK=X+B

- Sender: seqno=X+B, length=B

- Receiver: ACK=X+2B

- Sender: seqno=X+2B, length=B


- Sequence number of next packet is same as last ACK field

# TCP ACK Behavior: Loss Case

- **Sender sends packets with 100B and sequence numbers**:
  - 100, 200, 300, 400, 500, 600, 700, 800, 900, …
- **Assume the fifth packet (seqno 500) is lost, but no others**
- **Stream of ACKs will be**:
  - 200, 300, 400, 500, 500 (seqno:600), 500 (seqno:700), 500 (seqno:800), 500 (seqno:900),…
  - Acknowledgements are **cumulative**

# Reliable Data Transfer with TCP

- To recover a lost packet, sender needs to **retransmit**

- What triggers retransmission?
  - Our usual mechanism: **Timeout**
  - Sender maintains a **single timer** for **oldest unacknowledged segment** (like Go-Back-N)
  - On timeout, retransmit **only** that oldest unacknowledged segment (closer to Selective Repeat)

# Reliable Data Transfer with TCP: Setting Timeouts

- **How long should the timeout be?**
  - Too long: slow reaction to loss
  - Too short: waste bandwidth retransmitting packets that were not really lost

- **Goal:** timeout should be *close* to RTT
  - Definitely can't be shorter (many unnecessary retransmissions)
  - But, much longer will make reactions slow

# RTT Estimation

- We want timeout close to RTT, but…how do we know what the RTT is?
- We can measure it!
  - `SampleRTT`: measured time from segment transmission until ACK receipt
  - But measurements will vary over time…

# RTT Estimation

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

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



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

sampleRTT
EstimatedRTT

RTT (milliseconds) vs time (seconds)

# RTT Estimation

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

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



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

◆ sampleRTT
■ EstimatedRTT

RTT (milliseconds) vs time (seconds)

Example:
- Our current RTT estimate is 100ms
- New SampleRTT is 200ms
- What is updated EstimatedRTT?

.875*100 + .125*200 = 112.5ms

# TCP Timeout Setting

- timeout interval: **EstimatedRTT** plus "safety margin"
  - large variation in **EstimatedRTT:** want a larger safety margin

`TimeoutInterval = EstimatedRTT + 4*DevRTT`

estimated RTT          "safety margin"

- **DevRTT**: EWMA of **SampleRTT** deviation from **EstimatedRTT**:

`DevRTT = (1-β)*DevRTT + β*|SampleRTT-EstimatedRTT|`

(typically, $\beta$ = 0.25)

# TCP Timeouts: Details

- **Retransmitted** segments are **ignored** for SampleRTT calculations
  - Eliminates ambiguity of whether ACK was sent in response to original or retransmission

- What value should the timeout start at? (before any SampleRTT is measured)
  - Recommendation: 1 second (https://tools.ietf.org/html/rfc6298)

- Timeout is **doubled** when it expires
  - More on congestion control soon...

# TCP Retransmission Scenarios



lost ACK scenario

premature timeout

# TCP Retransmission Scenarios

Host A                    Host B

Seq=92, 8 bytes of data

Seq=100, 20 bytes of data

ACK=100

X

ACK=120

Seq=120,  15 bytes of data

cumulative ACK
**covers** for earlier lost
ACK

# Reliable Data Transfer with TCP

- To recover a lost packet, sender needs to **retransmit**

- What triggers retransmission?
  - Our usual mechanism: **Timeout**
    - Sender maintains a **single timer** for **oldest unacknowledged segment** (like Go-Back-N)
    - On timeout, retransmit **only** that oldest unacknowledged segment (closer to Selective Repeat)
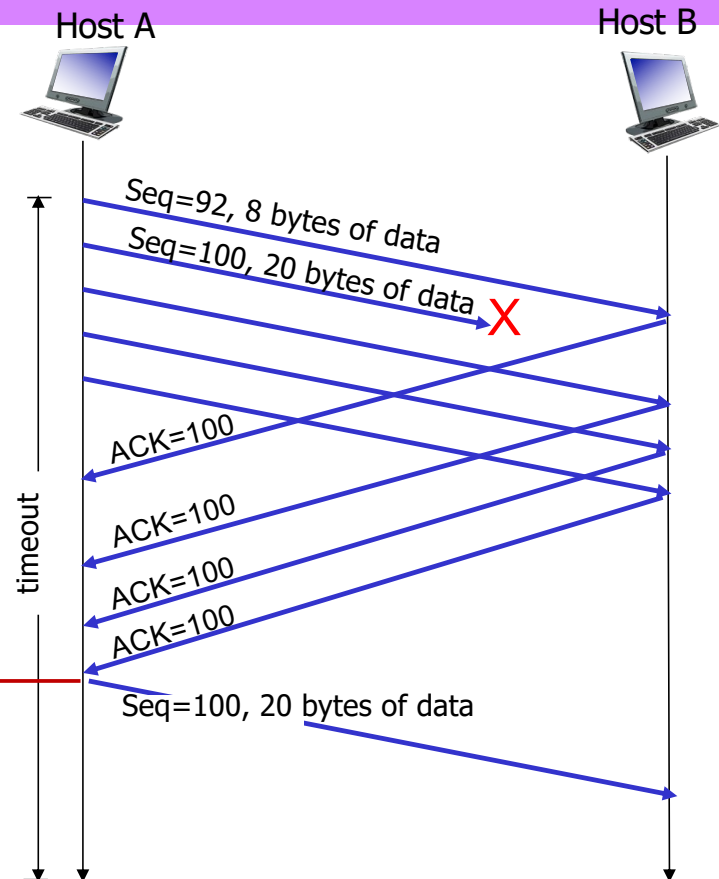
# Reliable Data Transfer with TCP

- To recover a lost packet, sender needs to **retransmit**

- What triggers retransmission?
  - Our usual mechanism: **Timeout**
    - Sender maintains a **single timer** for **oldest unacknowledged segment** (like Go-Back-N)
    - On timeout, retransmit **only** that oldest unacknowledged segment (closer to Selective Repeat)
  - An optimization: **Fast Retransmit**

# TCP Fast Retransmit

- Recall our loss scenario:
  - Sender sends packets with 100B and seqnos.:
    - 100, 200, 300, 400, 500, 600, 700, 800, 900, …
  - Assume the fifth packet (seqno 500) is lost, but no others
  - Stream of ACKs will be:
    - 200, 300, 400, 500 (seqno:600), 500 (seqno:700), 500 (seqno:800), 500 (seqno:900),…
- **Duplicate ACKs are a sign of isolated loss**
  - 500 clearly hasn't been delivered…but other segments are getting through

# TCP Fast Retransmit

- TCP will retransmit a segment **upon receiving 3 duplicate ACKs** for that segment
  - without waiting for timeout

💡 Receipt of three duplicate ACKs indicates 3 segments received after a missing segment – lost segment is likely. So retransmit!
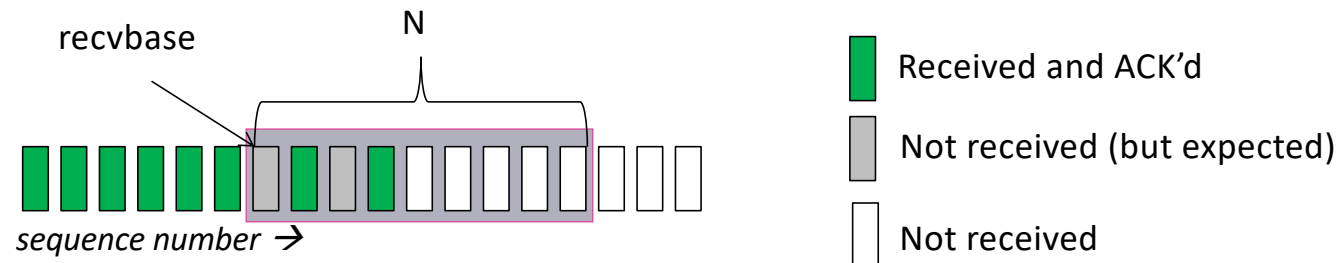
Host A

Host B

Seq=92, 8 bytes of data

Seq=100, 20 bytes of data

X

ACK=100

ACK=100

ACK=100

ACK=100

Seq=100, 20 bytes of data

timeout

# TCP Flow Control

- **Flow Control** key idea: sender should not transmit faster than the receiver can process!

# Flow Control

- Recall our discussion from last week on reliable data transfer
  - Receiver maintains a **buffer** that stores received but not-yet-delivered data



- Buffers have finite storage space, so what would happen if sender sends more data than the receiver can fit in its buffer?
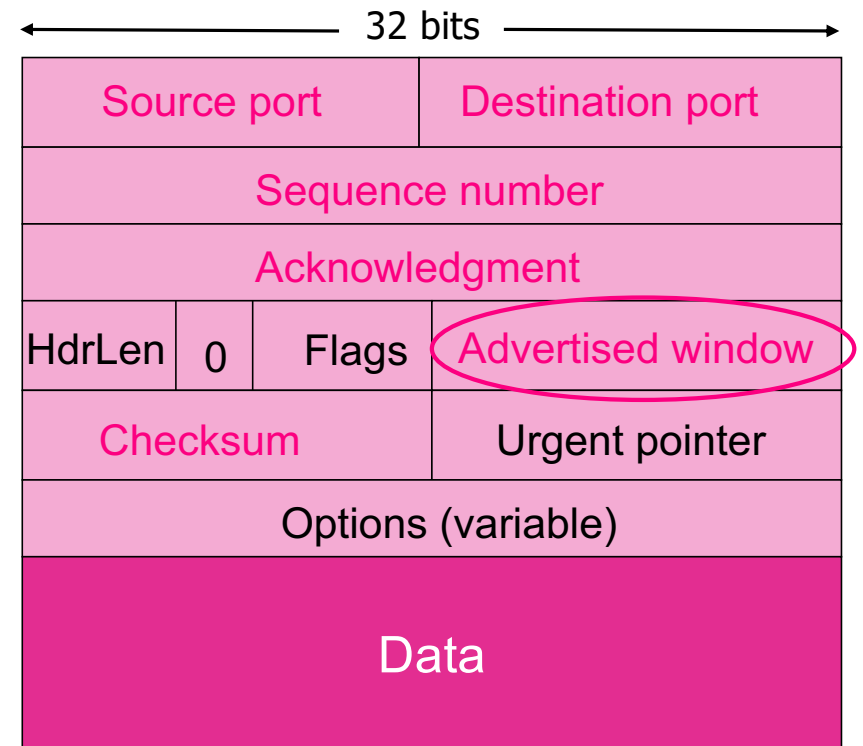
# Flow Control

- We saw a **limited form** of flow control in the Go-Back-N and Selective Repeat protocols we looked at:
  - Receiver maintains a **buffer** that **can hold N packets**
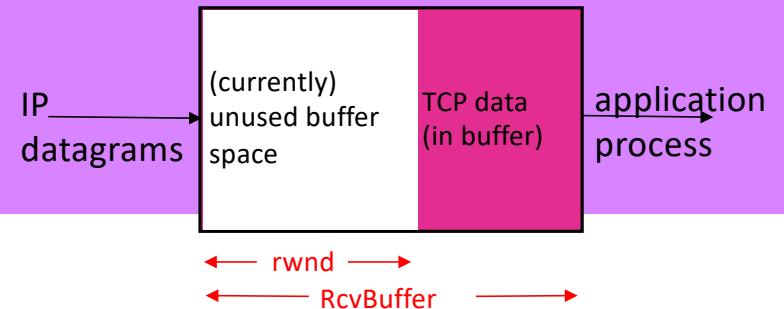  - Sender is allowed to have **up to N** unacknowledged packets out at any time



sequence number →

sequence number →

Sent and ACK'd

Sent but not ACK'd

Not yet sent

Received and ACK'd

Not received (but expected)

Not received

79

# TCP Flow Control

- More flexible, but same basic idea: **sender only sends as much data as it knows receiver can accept**

- Receiver uses an "Advertised Window" (RWND) to tell the sender how many bytes it can accept
  - Receiver indicates value of RWND in ACKs
  - Sender ensures that the total number of bytes in flight <= RWND

32 bits

| Source port | Destination port |
|---|---|
| Sequence number | |
| Acknowledgment | |
| HdrLen / 0 / Flags | Advertised window |
| Checksum | Urgent pointer |
| Options (variable) | |
| Data | |

# TCP Flow Control



IP datagrams → | (currently) unused buffer space | TCP data (in buffer) | → application process

← rwnd →
← RcvBuffer →

- **Advertised window limits rate**: Sender can send no faster than RWND/RTT bytes/sec

- Receiver only advertises more space when application has consumed old arriving data

- What happens when RWND=0?
  - Sender keeps probing by sending segments with one byte of data
  - Receiver can ACK, but won't increase RWND until application reads some data and buffer space opens up

```
rwnd = RcvBuffer – [LastByteRcvd – LastByteRead]
```

# Recall: TCP Service Abstraction

- **TCP is a connection-oriented protocol that delivers a reliable, in-order, byte stream**

- **Connection-oriented**: two processes coordinate ("handshake") before beginning to send data to each other

✓ - **Reliable**: TCP resends lost packets
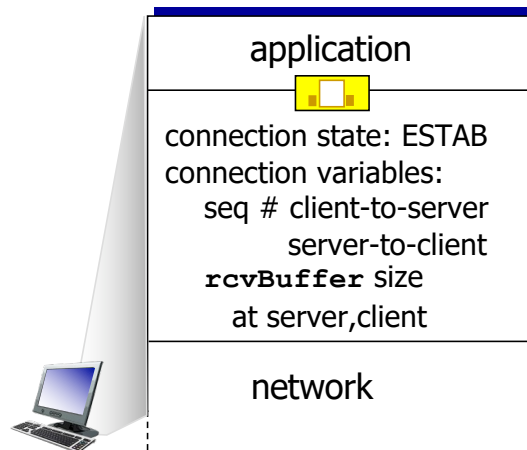  - Until it gives up and shuts down connection

✓ - **In-order**: TCP only hands consecutive chunks of data to application
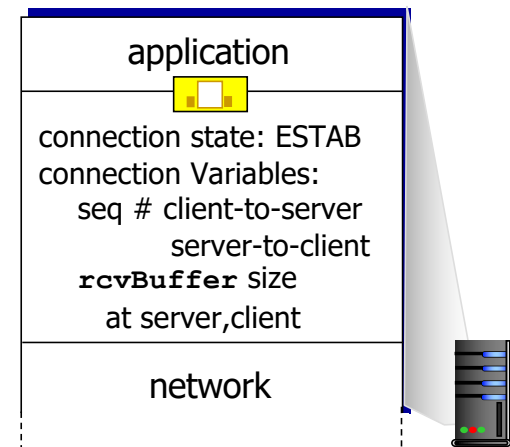
✓ - **Byte stream**: TCP assumes there is an incoming stream of data, and attempts to deliver it to application

# TCP Connection Management

- Before exchanging data, sender/receiver "handshake":
  - agree to establish connection (each knowing the other willing to establish connection)
  - agree on connection parameters (e.g., starting seq #s, initial buffer sizes)

application

connection state: ESTAB
connection variables:
   seq # client-to-server
     server-to-client
**rcvBuffer** size
   at server,client

network

application

connection state: ESTAB
connection Variables:
   seq # client-to-server
     server-to-client
**rcvBuffer** size
   at server,client

network

```
clientSocket.connect("hostname,
    "port number");
```
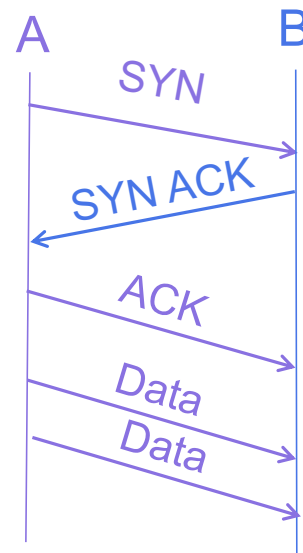
```
connectionSocket =
    welcomeSocket.accept();
```

# Connection Parameters: Initial Sequence Number

- Sequence number for the **very first byte**

- **Why not just use ISN = 0?**
  - Practical issue
    - IP addresses and port #s uniquely identify a connection
    - Eventually, though, these port #s do get used again; small chance an old packet is still in flight
    - Also, others might try to spoof your connection
  - Why does using ISN help?

- Hosts exchange ISNs when establishing connection

# Establishing a TCP Connection

- **Three-way handshake** to establish connection
  - Host A sends a SYN (open; "synchronize sequence numbers") to host B
  - Host B returns a SYN acknowledgment (SYN ACK)
  - Host A sends an ACK to acknowledge the SYN ACK

A       B

SYN

SYN ACK

ACK

Data

Data

# Establishing a TCP Connection

**Flags:**
SYN
ACK
FIN
RST
PSH
URG
CWR
ECE

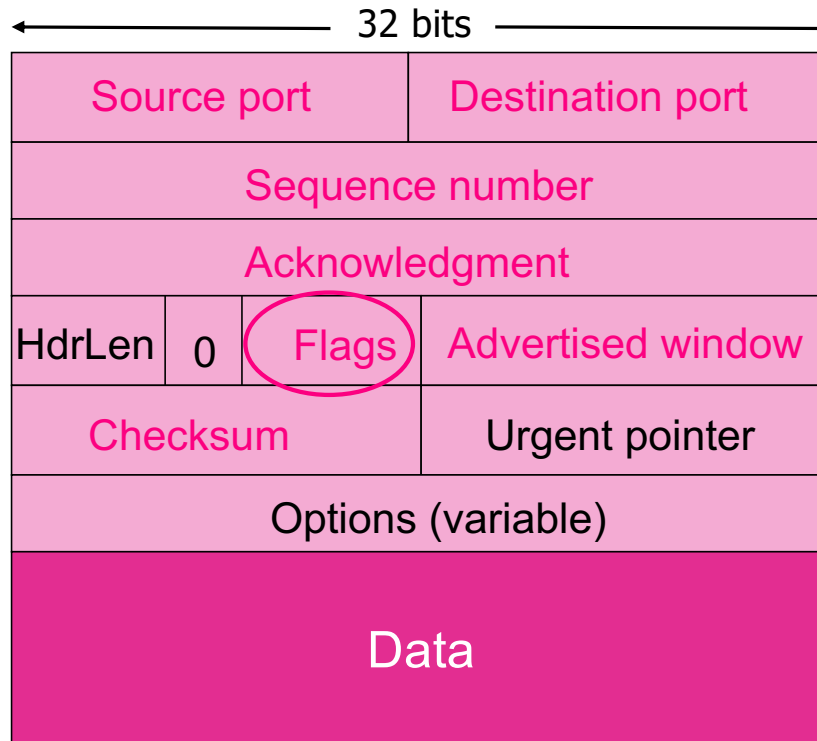| 32 bits | | |
|---|---|---|
| Source port | | Destination port |
| Sequence number | | |
| Acknowledgment | | |
| HdrLen | 0 | Flags | Advertised window |
| Checksum | | Urgent pointer |
| Options (variable) | | |
| Data | | |

# Establishing a TCP Connection

**Flags:**
SYN – connection setup
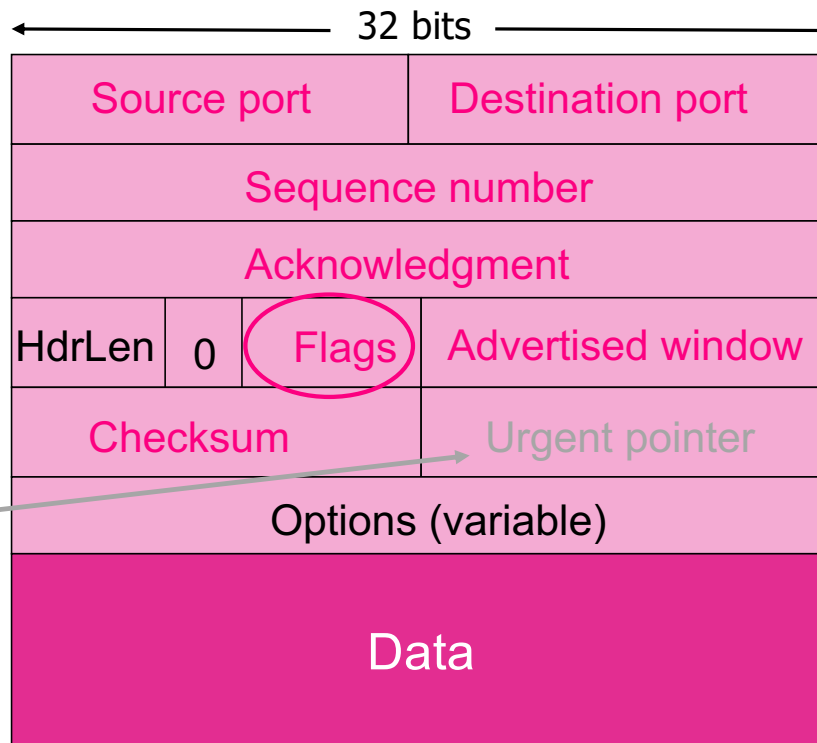ACK – contains valid ack
FIN – connection teardown
RST – connection teardown
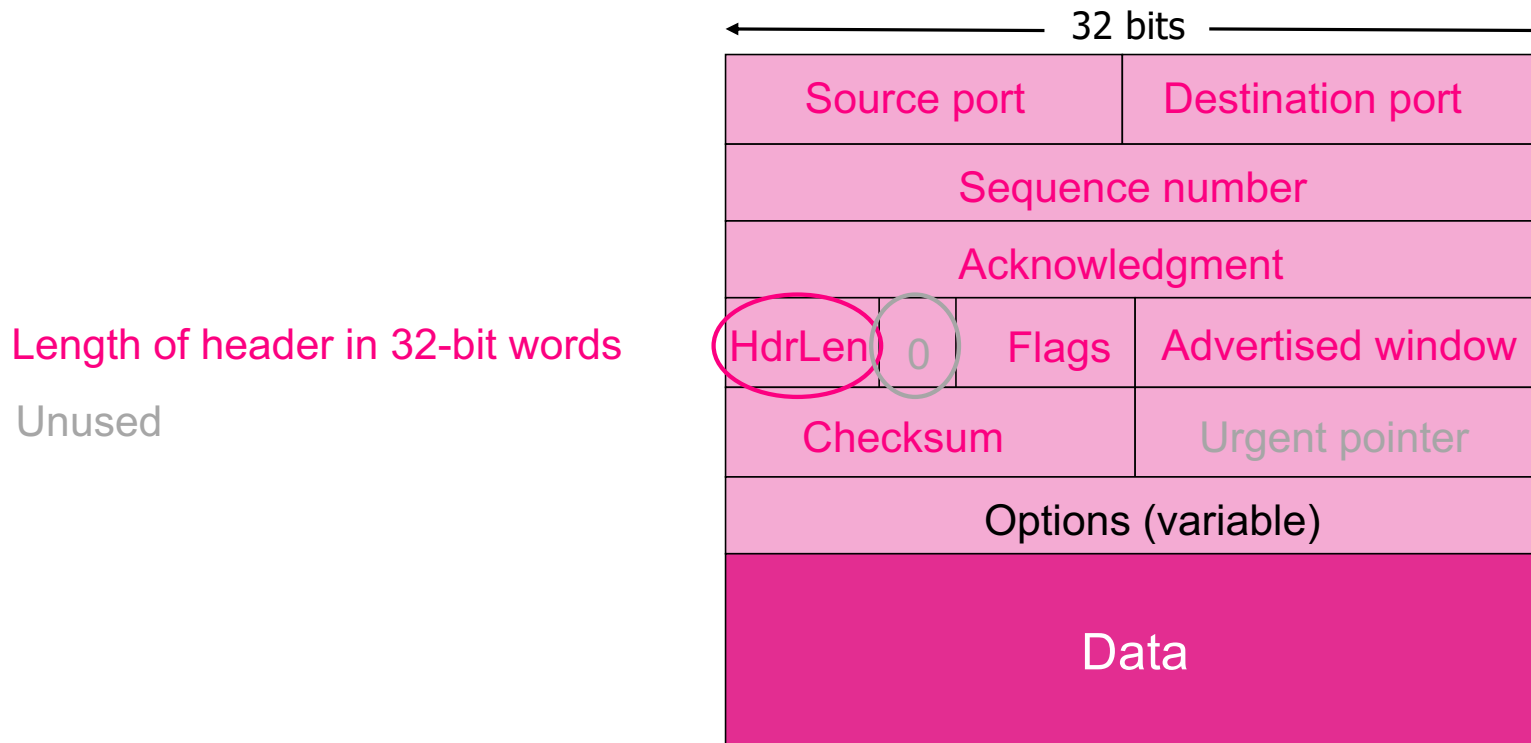PSH – pass data up immediately
URG – data marked urgent
CWR – congestion control
ECE – congestion control

# Wrapping up TCP Segment Structure

Length of header in 32-bit words

Unused

# Step 1: Host A's SYN

A tells B to open a connection

| A's port | B's port |
|----------|----------|
| A's Initial Sequence Number ||
| N/A ||

| 5 | 0 | SYN | Advertised window |
|---|---|-----|-------------------|

| Checksum | Urgent pointer |
|----------|----------------|

# Step 2: Host B's SYN-ACK

B tells A it accepts and is ready to accept next packet

| B's port | A's port |
|---|---|
| B's Initial Sequence Number | |
| ACK=A's ISN+1 | |

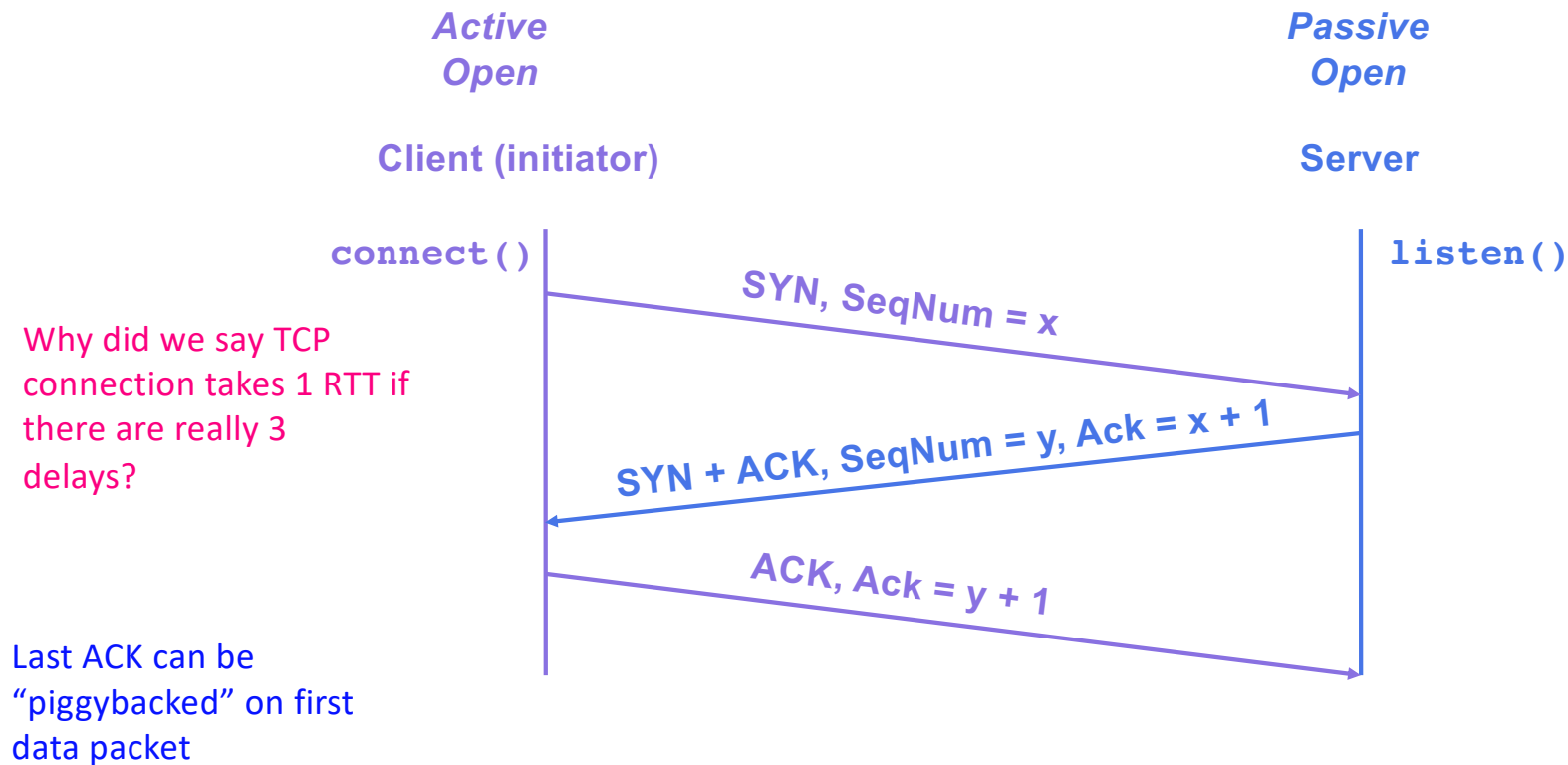| 5 | 0 | SYN|ACK | Advertised window |
|---|---|---|---|
| Checksum | | | Urgent pointer |

B's Initial Sequence Number = A's Initial Sequence Number?

# Step 3: Host A's ACK of B's SYN-ACK

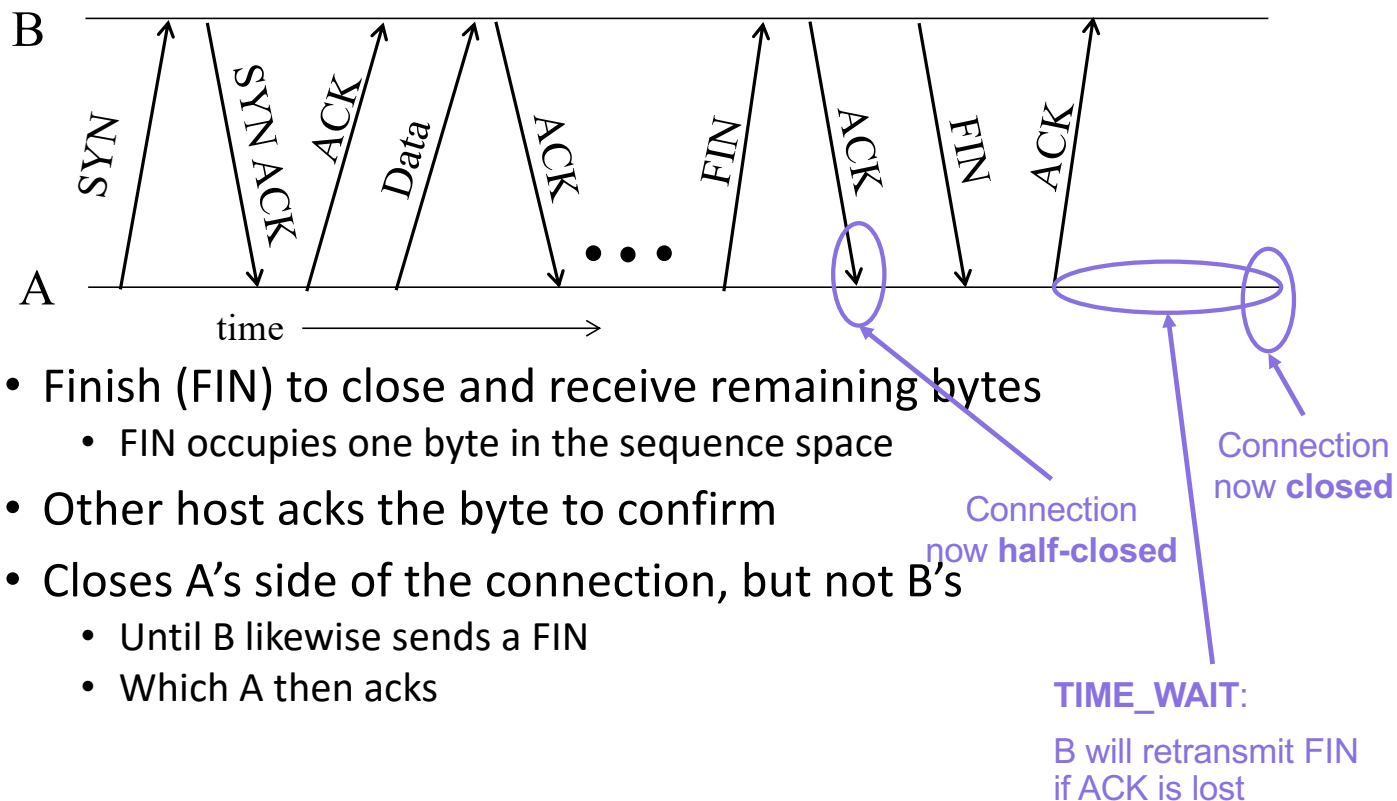A tells B to open a connection

| A's port | B's port |
|---|---|
| A's Initial Sequence Number | |
| ACK=B's ISN+1 | |

| 5 | 0 | ACK | Advertised window |
|---|---|---|---|
| Checksum | | | Urgent pointer |

# TCP 3-way Handshake

**Active Open**

**Passive Open**

**Client (initiator)**

**Server**

`connect()`

`listen()`

SYN, SeqNum = x

Why did we say TCP connection takes 1 RTT if there are really 3 delays?

SYN + ACK, SeqNum = y, Ack = x + 1

ACK, Ack = y + 1

Last ACK can be "piggybacked" on first data packet

# Closing a TCP Connection: Normal Termination Example



- Finish (FIN) to close and receive remaining bytes
  - FIN occupies one byte in the sequence space
- Other host acks the byte to confirm
- Closes A's side of the connection, but not B's
  - Until B likewise sends a FIN
  - Which A then acks

Connection now **half-closed**

Connection now **closed**

**TIME_WAIT**:

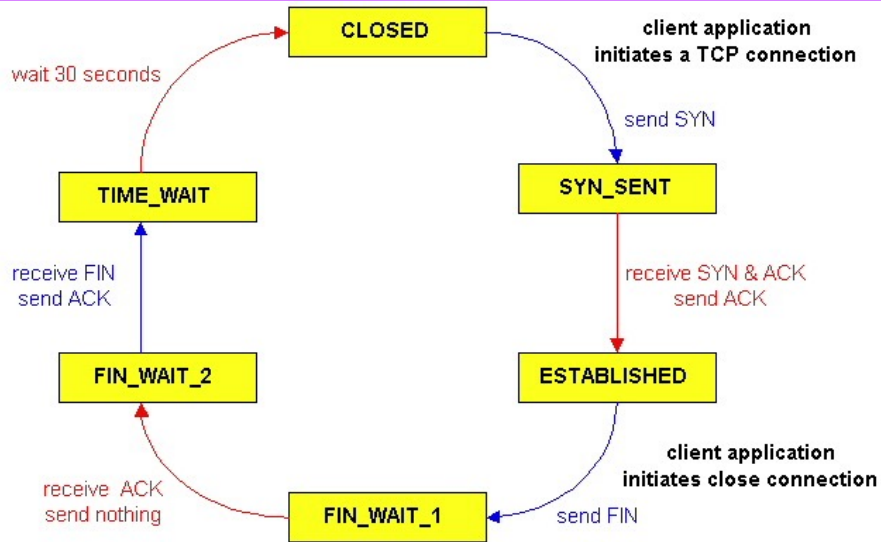B will retransmit FIN if ACK is lost
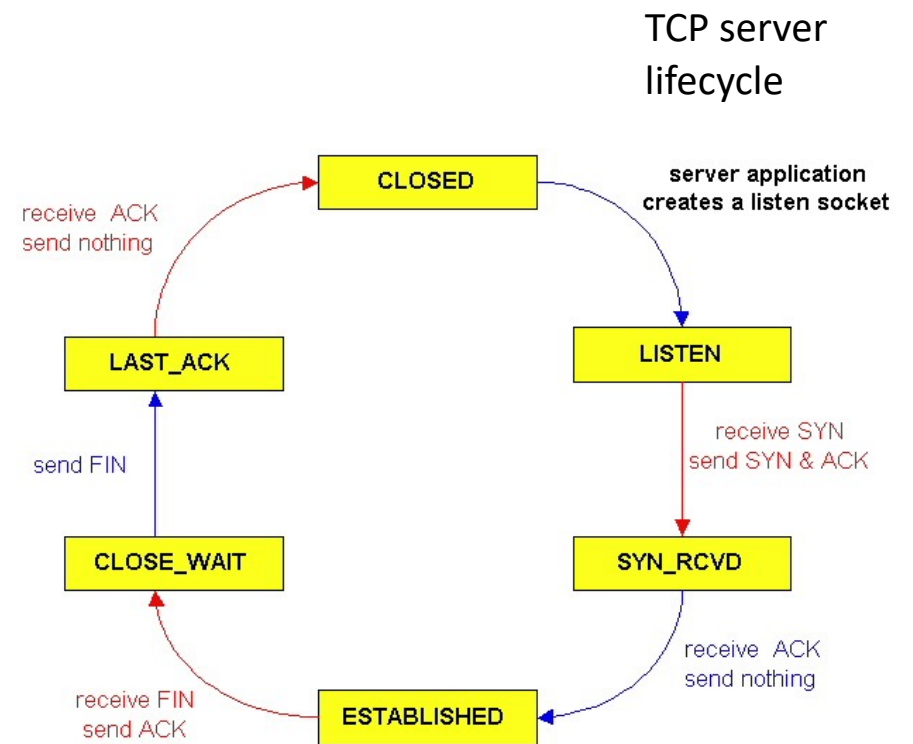
# Closing a TCP Connection: Abrupt Termination



- A sends a RESET (RST) to B
  - E.g., because application process on A crashed
- That's it
  - B does not ack the RST
  - Thus, RST is not delivered reliably, and any data in flight is lost
  - But: if B sends anything more, will elicit another RST

# TCP State Diagrams (Normal Case)



TCP client lifecycle

TCP server lifecycle

# Summary

- TCP: connection-oriented protocol delivering reliable, in-order byte stream
  - Reliability builds on mechanisms we've seen: checksums, sequence numbers (adapted to byte stream abstraction), cumulative acknowledgments, sliding window
    - New mechanisms: dynamic timeout estimation, fast retransmit
  - Flow control prevents buffer overruns
  - Handshaking procedures used to set up and cleanly teardown connections