# CSB051 – Computer Networks
## 電腦網路

# Chapter 3
# Transport Layer

吳俊興

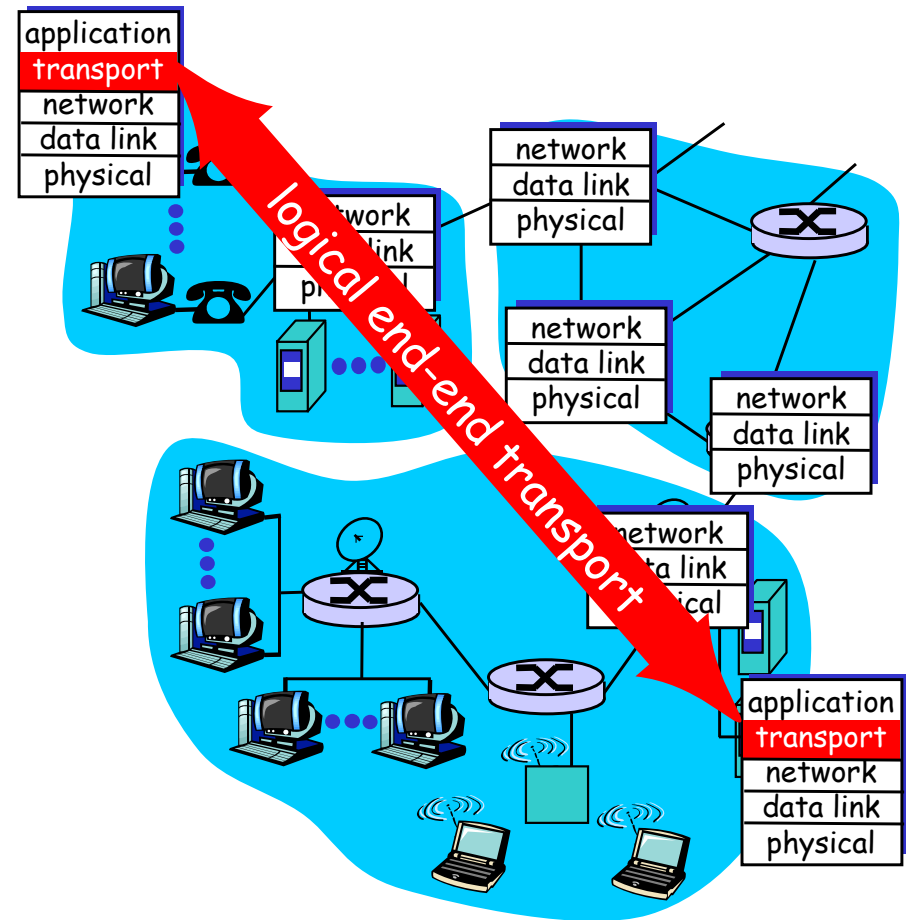國立高雄大學 資訊工程學系

# Chapter 3 outline

# 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

application
transport
network
data link
physical

network
data link
physical

network
data link
physical

network
data link
physical

network
data link
physical

network
data link
physical

application
transport
network
data link
physical

logical end-end transport

# Chapter 3 outline
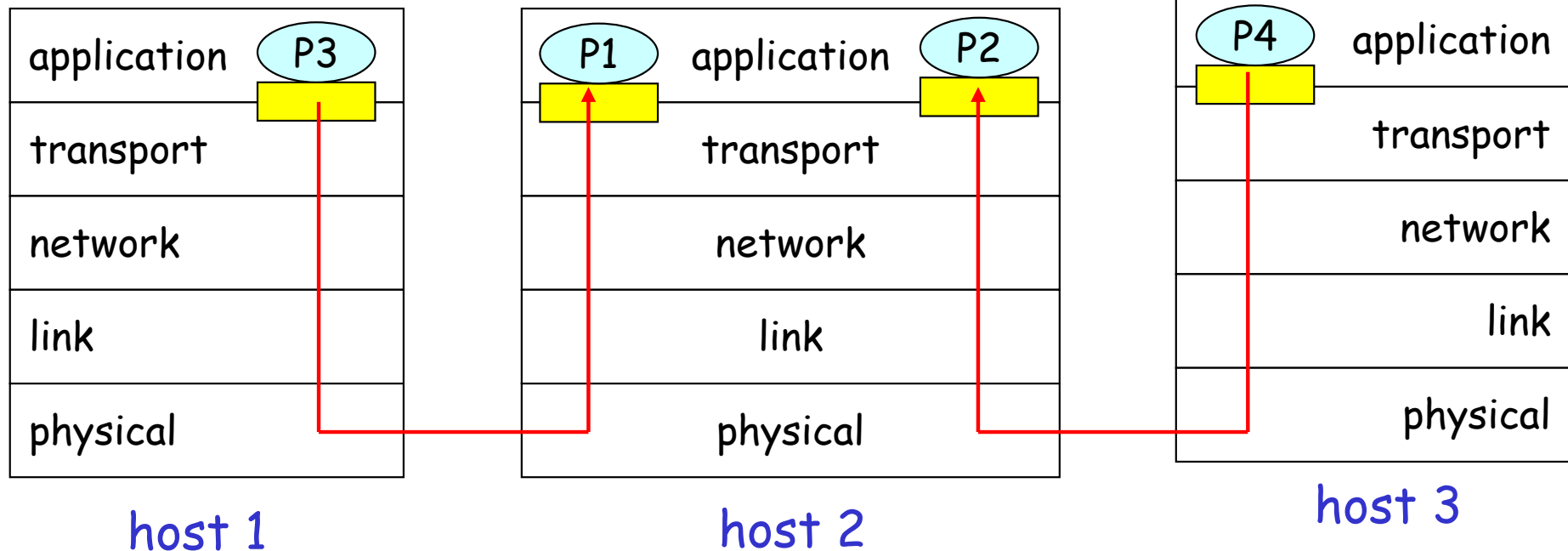
# Multiplexing/demultiplexing

**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

# How demultiplexing works

□ host receives IP datagrams
  ○ each datagram has source IP address, destination IP address
  ○ each datagram carries 1 transport-layer segment
  ○ each segment has source, destination port number (recall: well-known port numbers for specific applications)

□ host uses IP addresses & port numbers to direct segment to appropriate socket

←——————— 32 bits ———————→

| source port # | dest port # |
|:---:|:---:|
| other header fields | |
| application data (message) | |

TCP/UDP segment format

# Connectionless demultiplexing

□ Create sockets with port numbers:

```
DatagramSocket mySocket1 = new
    DatagramSocket(99111);

DatagramSocket mySocket2 = new
    DatagramSocket(99222);
```

□ UDP socket identified by two-tuple:
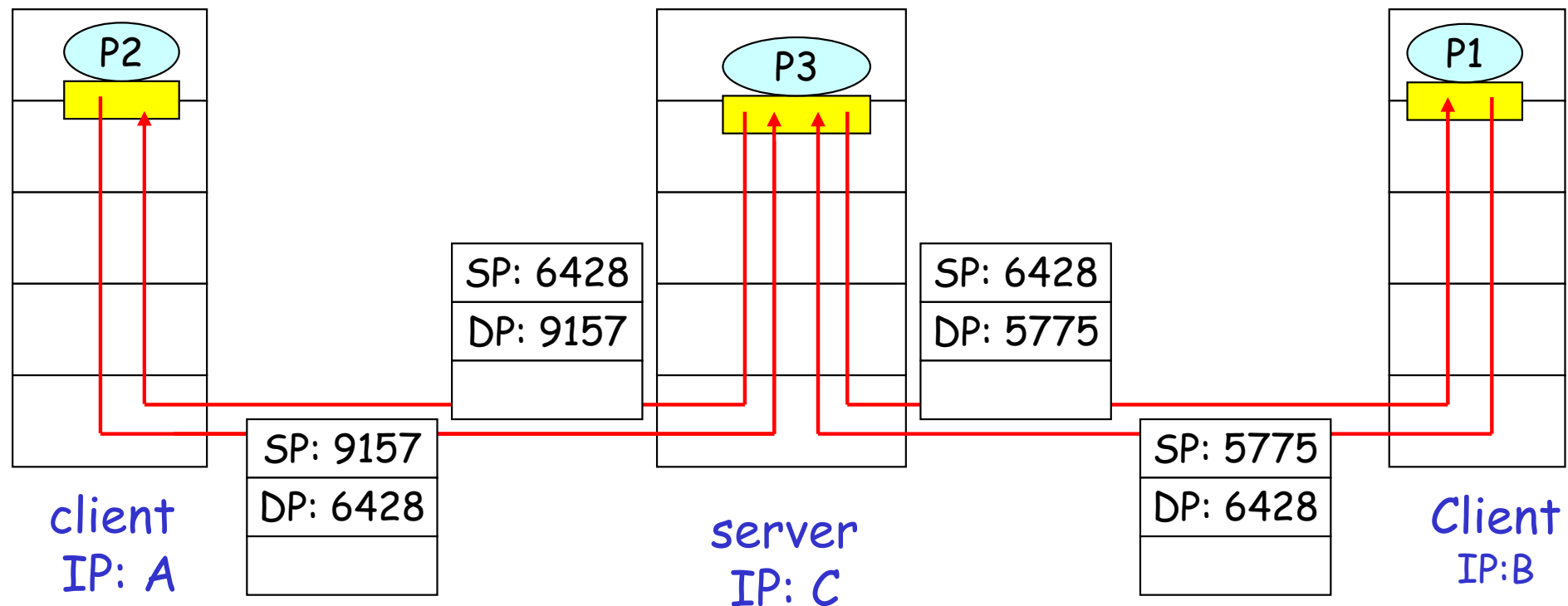
(dest IP address, dest port number)

□ When host receives UDP segment:
- checks destination port number in segment
- directs UDP segment to socket with that port number

□ IP datagrams with different source IP addresses and/or source port numbers directed to same socket

# Connectionless demux (cont)

`DatagramSocket serverSocket = new DatagramSocket(6428);`



SP provides "return address"

# Connection-oriented demux

- TCP socket identified by 4-tuple:
  - source IP address
  - source port number
  - dest IP address
  - dest port number
- recv host uses all four values to direct segment to appropriate socket

- Server host may support many simultaneous TCP sockets:
  - each socket identified by its own 4-tuple
- Web servers have different sockets for each connecting client
  - non-persistent HTTP will have different socket for each request

# Connection-oriented demux (cont)



P1

P4  P5  P6

P2  P3

SP: 5775
DP: 80
S-IP: B
D-IP:C

client
IP: A

SP: 9157
DP: 80
S-IP: A
D-IP:C

server
IP: C

SP: 9157
DP: 80
S-IP: B
D-IP:C

Client
IP:B

# Chapter 3 outline

# UDP: User Datagram Protocol [RFC 768]

- □ "no frills," "bare bones" Internet transport protocol
- □ "best effort" service, UDP segments may be:
  - ○ lost
  - ○ delivered out of order to app
- □ *connectionless:*
  - ○ no handshaking between UDP sender, receiver
  - ○ each UDP segment handled independently of others

## Why is there a UDP?

- □ finer control over what and when data to be sent
  - ○ no congestion control: UDP can blast away as fast as desired
- □ no connection establishment (which can add delay)
- □ simple: no connection state at sender, receiver
- □ small segment header
  - ○ TCP: 20bytes, UDP: 8bytes

# UDP: more

- □ often used for streaming multimedia apps
  - ○ loss tolerant
  - ○ rate sensitive
- □ **other UDP uses**
  - ○ DNS
  - ○ SNMP
- □ reliable transfer over UDP: add reliability at application layer
  - ○ application-specific error recovery!

Length, in bytes of UDP segment, including header

```
←————— 32 bits —————→

| source port # | dest port # |
| length        | checksum    |
|                             |
|       Application           |
|          data               |
|        (message)            |
|                             |
```

UDP segment format

# UDP checksum

**Goal:** detect "errors" (e.g., flipped bits) in transmitted segment

## Sender:

- □ treat segment contents as sequence of 16-bit integers
- □ checksum: addition (1's complement sum) of segment contents
- □ sender puts checksum value into UDP checksum field

## Receiver:

- □ compute checksum of received segment
- □ check if computed checksum equals checksum field value:
  - ○ NO - error detected
  - ○ YES - no error detected. *But maybe errors nonetheless?* More later ....

# Internet Checksum Example

□ Note

   ○ When adding numbers, a carryout from the most significant bit needs to be added to the result

□ Example: add two 16-bit integers

```
               1 1 1 0 0 1 1 0 0 1 1 0 0 1 1 0
               1 1 0 1 0 1 0 1 0 1 0 1 0 1 0 1
              ─────────────────────────────────
wraparound    ①1 0 1 1 1 0 1 1 1 0 1 1 1 0 1 1
              ─────────────────────────────────
      sum      1 0 1 1 1 0 1 1 1 0 1 1 1 1 0 0
 checksum      0 1 0 0 0 1 0 0 0 1 0 0 0 0 1 1
```

# Chapter 3 outline

# Principles of Reliable data transfer

- important in app., transport, link layers
- top-10 list of important networking topics!



(a) provided service

(b) service implementation

- characteristics of unreliable channel will determine complexity of reliable data transfer protocol (rdt)
    (udt=unreliable data transfer)

# 3.4.1 Reliable data transfer: getting started

**rdt_send():** called from above, (e.g., by app.). Passed data to deliver to receiver upper layer

**deliver_data():** called by **rdt** to deliver data to upper

**send side**

rdt_send()    data

reliable data transfer protocol (sending side)

udt_send()    packet

**receive side**

data    deliver_data()

reliable data transfer protocol (receiving side)

packet    rdt_rcv()

unreliable channel

**udt_send():** called by rdt, to transfer packet over unreliable channel to receiver

**rdt_rcv():** called when packet arrives on rcv-side of channel

# Reliable data transfer: getting started

We'll:

☐ incrementally develop sender, receiver sides of reliable data transfer protocol (rdt)

☐ consider only unidirectional data transfer
  ○ but control info will flow on both directions!

☐ use finite state machines (FSM) to specify sender, receiver

event causing state transition
actions taken on state transition

state: when in this "state" next state uniquely determined by next event

state 1

event
actions

state 2

# Developing reliable data transfer protocols

□ rdt1.0: over reliable channel
   ○ no bit errors, no packet losses

□ rdt2.x: channel with bit errors (stop-and-wait)
   ○ rdt2.0: ACK/NAK not corrupted
   ○ rdt2.1: garbled ACK/NAK – retransmit and duplicate
   ○ rdt2.2: NAK-free

□ rdt3.0: channel with errors and loss
   ○ time-based retransmission
   ○ alternating-bit protocol (stop-and-wait)

□ Pipelined protocols: send multiple packets without waiting for ACKs
   ○ Go-Back-N: sliding-window protocol
   ○ Selective Repeat

# Rdt1.0: reliable transfer over a reliable channel

- underlying channel perfectly reliable
  - no bit errors
  - no loss of packets
- separate FSMs for sender, receiver:
  - sender sends data into underlying channel
  - receiver read data from underlying channel

Wait for call from above

rdt_send(data)
_____

packet = make_pkt(data)
udt_send(packet)

**sender**

Wait for call from below

rdt_rcv(packet)
_____

extract (packet,data)
deliver_data(data)

**receiver**

# Rdt2.0: channel with bit errors

☐ **underlying channel may flip bits in packet**
  ○ checksum to detect bit errors

☐ *the* **question: how to recover from errors:**
  ○ *acknowledgements (ACKs):* receiver explicitly tells sender that pkt received OK
  ○ *negative acknowledgements (NAKs):* receiver explicitly tells sender that pkt had errors
  ○ sender retransmits pkt on receipt of NAK

☐ **new mechanisms in `rdt2.0` (beyond `rdt1.0`):**
  ○ error detection
  ○ receiver feedback: control msgs (ACK,NAK) rcvr->sender

# rdt2.0: FSM specification

**receiver**

rdt_send(data)
_____
snkpkt = make_pkt(data, checksum)
udt_send(sndpkt)

rdt_rcv(rcvpkt) &&
isNAK(rcvpkt)
_____
udt_send(sndpkt)

**Wait for call from above**

**Wait for ACK or NAK**

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

**sender**

rdt_rcv(rcvpkt) &&
corrupt(rcvpkt)
_____
udt_send(NAK)

**Wait for call from below**

rdt_rcv(rcvpkt) &&
notcorrupt(rcvpkt)
_____
extract(rcvpkt,data)
deliver_data(data)
udt_send(ACK)

# rdt2.0: operation with no errors

rdt_send(data)
_____
snkpkt = make_pkt(data, checksum)
udt_send(sndpkt)

rdt_rcv(rcvpkt) &&
isNAK(rcvpkt)
_____
udt_send(sndpkt)

Wait for call from above

Wait for ACK or NAK

rdt_rcv(rcvpkt) && isACK(rcvpkt)
_____
Λ

rdt_rcv(rcvpkt) &&
corrupt(rcvpkt)
_____
udt_send(NAK)

Wait for call from below

rdt_rcv(rcvpkt) &&
notcorrupt(rcvpkt)
_____
extract(rcvpkt,data)
deliver_data(data)
udt_send(ACK)

# rdt2.0: error scenario

rdt_send(data)
_____
snkpkt = make_pkt(data, checksum)
udt_send(sndpkt)

Wait for call from above

Wait for ACK or NAK

rdt_rcv(rcvpkt) && isNAK(rcvpkt)
_____
udt_send(sndpkt)

rdt_rcv(rcvpkt) && corrupt(rcvpkt)
_____
udt_send(NAK)

rdt_rcv(rcvpkt) && isACK(rcvpkt)
_____
Λ

Wait for call from below

rdt_rcv(rcvpkt) && notcorrupt(rcvpkt)
_____
extract(rcvpkt,data)
deliver_data(data)
udt_send(ACK)

# rdt2.0 has a fatal flaw!

**What happens if ACK/NAK corrupted?**

□ sender doesn't know what happened at receiver!

□ can't just retransmit: possible duplicate

**Handling duplicates:**

□ sender adds *sequence number* to each pkt

□ sender retransmits current pkt if ACK/NAK garbled

□ receiver discards (doesn't deliver up) duplicate pkt

**stop and wait**
Sender sends one packet, then waits for receiver response

# rdt2.1: sender, handles garbled ACK/NAKs

rdt_send(data)

---

sndpkt0 = make_pkt(0, data, checksum)
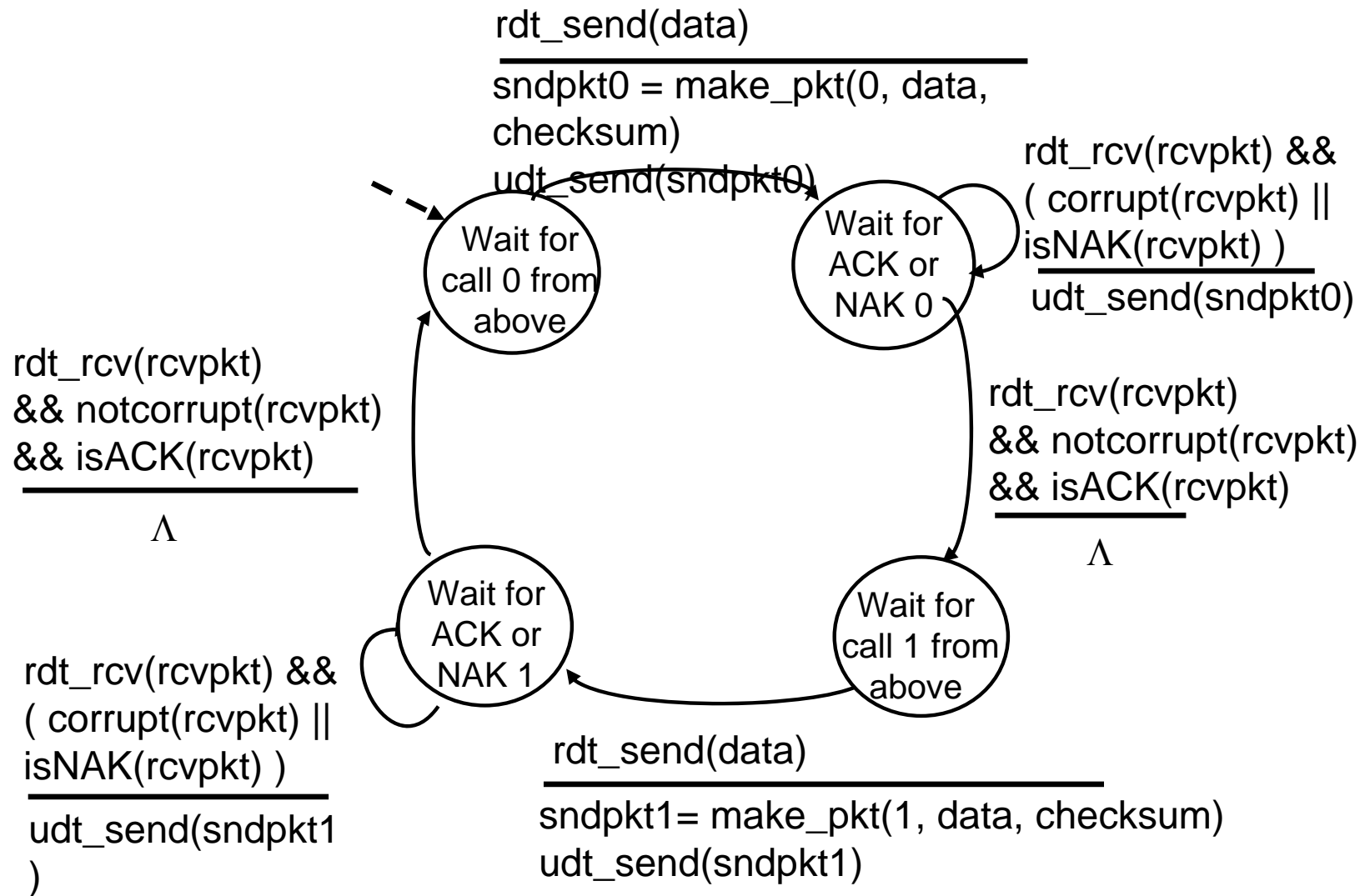udt_send(sndpkt0)

**Wait for call 0 from above**

**Wait for ACK or NAK 0**

rdt_rcv(rcvpkt) &&
( corrupt(rcvpkt) ||
isNAK(rcvpkt) )

---

udt_send(sndpkt0)

rdt_rcv(rcvpkt)
&& notcorrupt(rcvpkt)
&& isACK(rcvpkt)

---

$\Lambda$

rdt_rcv(rcvpkt)
&& notcorrupt(rcvpkt)
&& isACK(rcvpkt)

---

$\Lambda$

**Wait for ACK or NAK 1**

**Wait for call 1 from above**

rdt_rcv(rcvpkt) &&
( corrupt(rcvpkt) ||
isNAK(rcvpkt) )

---

udt_send(sndpkt1)

rdt_send(data)

---

sndpkt1= make_pkt(1, data, checksum)
udt_send(sndpkt1)

# rdt2.1: receiver, handles garbled ACK/NAKs

rdt_rcv(rcvpkt) && notcorrupt(rcvpkt)
 && has_seq0(rcvpkt)
_____
extract(rcvpkt,data)
deliver_data(data)
sndpkt = make_pkt(ACK, chksum)
udt_send(sndpkt)

rdt_rcv(rcvpkt) && (corrupt(rcvpkt)
_____
sndpkt = make_pkt(NAK, chksum)
udt_send(sndpkt)

rdt_rcv(rcvpkt) &&
  not corrupt(rcvpkt) &&
  has_seq1(rcvpkt)
_____
sndpkt = make_pkt(ACK, chksum)
udt_send(sndpkt)

rdt_rcv(rcvpkt) && (corrupt(rcvpkt)
_____
sndpkt = make_pkt(NAK, chksum)
udt_send(sndpkt)

rdt_rcv(rcvpkt) &&
  not corrupt(rcvpkt) &&
  has_seq0(rcvpkt)
_____
sndpkt = make_pkt(ACK, chksum)
udt_send(sndpkt)

**Wait for 0 from below**

**Wait for 1 from below**

rdt_rcv(rcvpkt) && notcorrupt(rcvpkt)
 && has_seq1(rcvpkt)
_____
extract(rcvpkt,data)
deliver_data(data)
sndpkt = make_pkt(ACK, chksum)
udt_send(sndpkt)

# rdt2.1: discussion

- seq # added to pkt
- two seq. #'s (0,1) will suffice.  Why?
- must check if received ACK/NAK corrupted
- twice as many states
  - state must "remember" whether "current" pkt has 0 or 1 seq. #

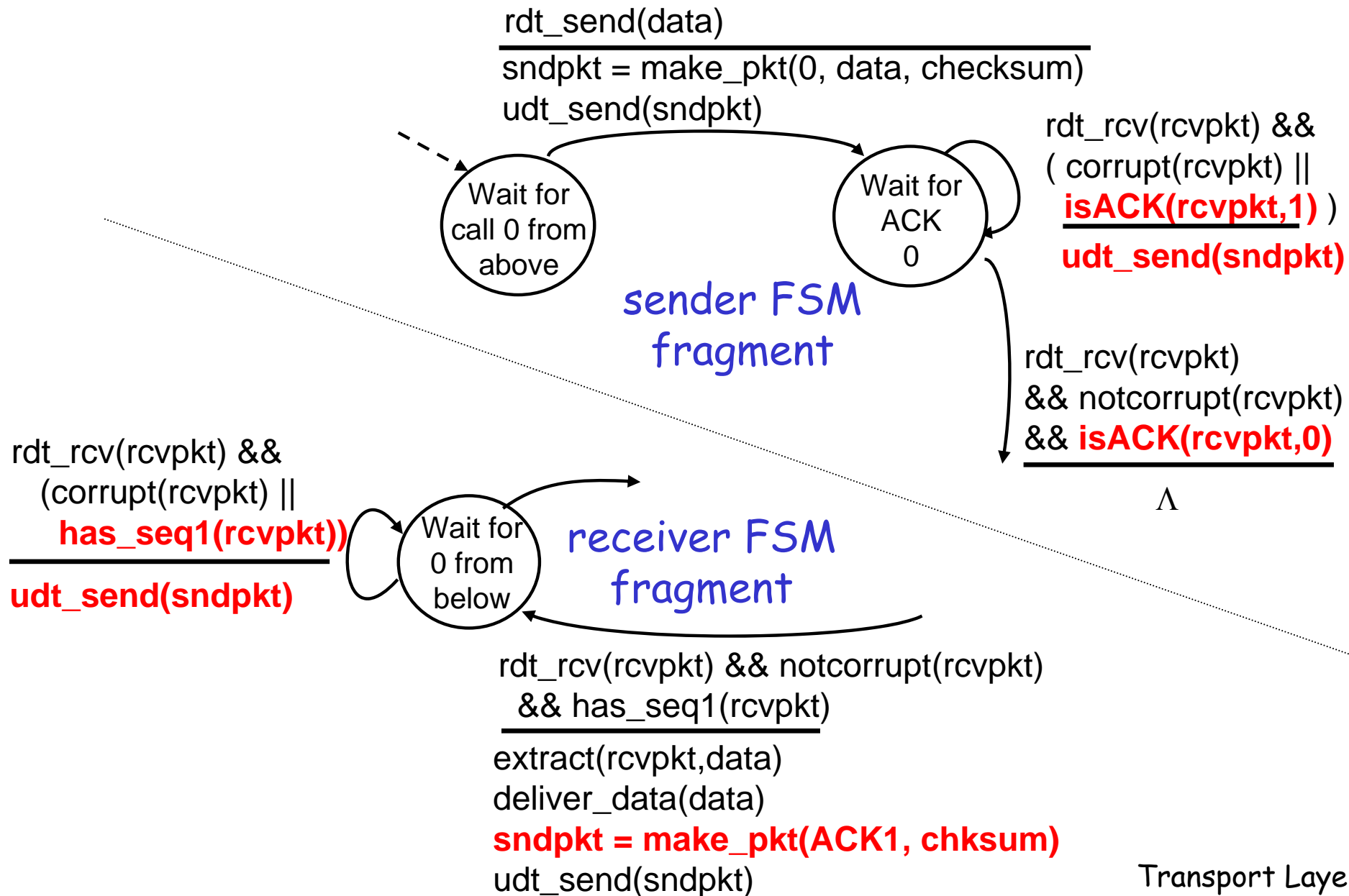- must check if received packet is duplicate
  - state indicates whether 0 or 1 is expected pkt seq #
- note: receiver can *not* know if its last ACK/NAK received OK at sender

# rdt2.2: a NAK-free protocol

- □ same functionality as rdt2.1, using ACKs only
- □ instead of NAK, receiver sends ACK for last pkt received OK
  - ○ receiver must *explicitly* include seq # of pkt being ACKed
- □ duplicate ACK at sender results in same action as NAK: *retransmit current pkt*

# rdt2.2: sender, receiver fragments

rdt_send(data)
_____
sndpkt = make_pkt(0, data, checksum)
udt_send(sndpkt)

rdt_rcv(rcvpkt) &&
( corrupt(rcvpkt) ||
  **isACK(rcvpkt,1)** )
**udt_send(sndpkt)**

Wait for
call 0 from
above

Wait for
ACK
0

**sender FSM**
**fragment**

rdt_rcv(rcvpkt)
&& notcorrupt(rcvpkt)
&& **isACK(rcvpkt,0)**
_____
$\Lambda$

rdt_rcv(rcvpkt) &&
  (corrupt(rcvpkt) ||
   **has_seq1(rcvpkt))**
_____
**udt_send(sndpkt)**

Wait for
0 from
below

**receiver FSM**
**fragment**

rdt_rcv(rcvpkt) && notcorrupt(rcvpkt)
  && has_seq1(rcvpkt)
_____
extract(rcvpkt,data)
deliver_data(data)
**sndpkt = make_pkt(ACK1, chksum)**
udt_send(sndpkt)

# rdt3.0: channels with errors *and* loss

**New assumption:**
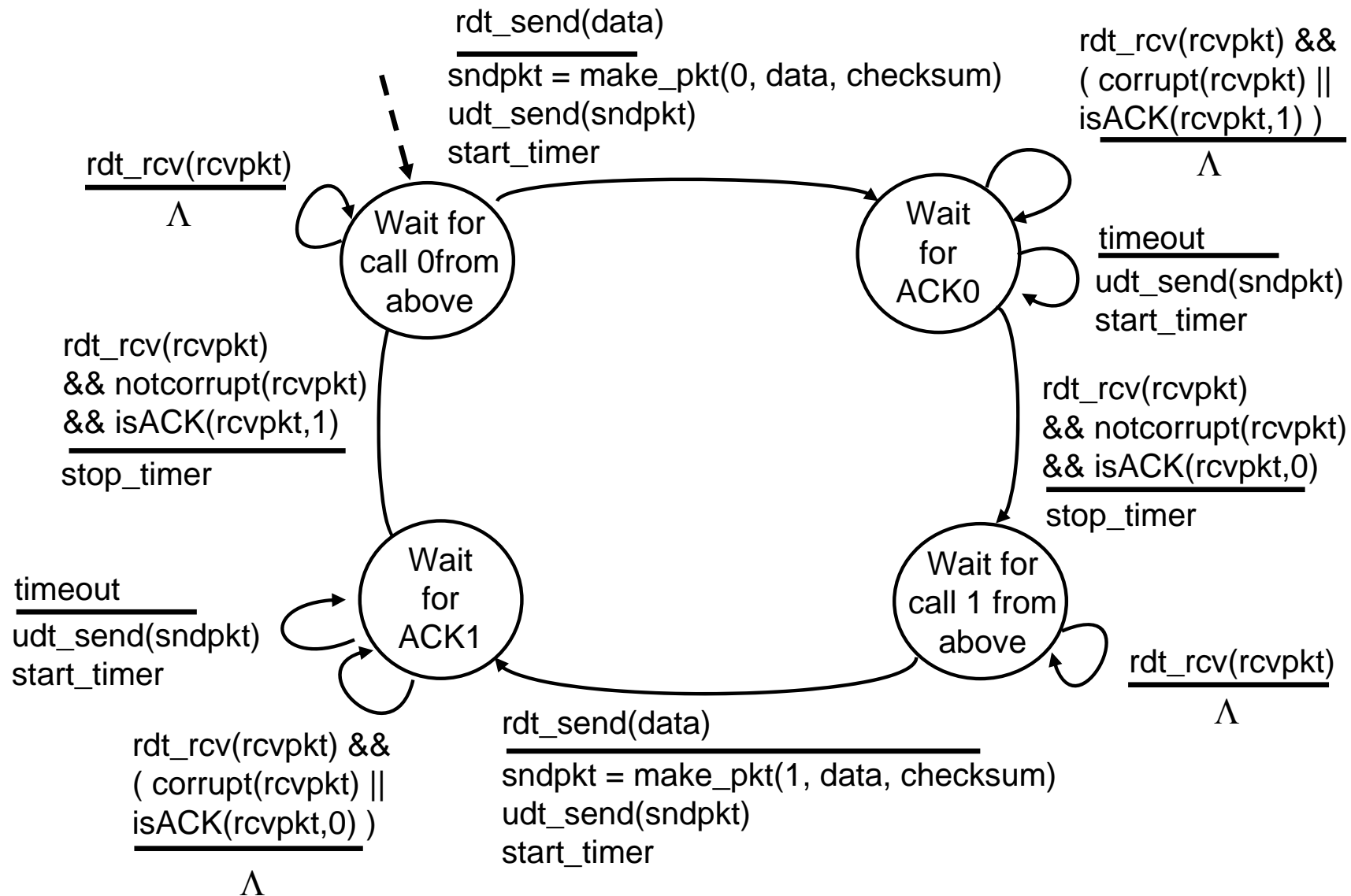underlying channel can also lose packets (data or ACKs)

- checksum, seq. #, ACKs, retransmissions will be of help, but not enough

**Approach:** sender waits "reasonable" amount of time for ACK
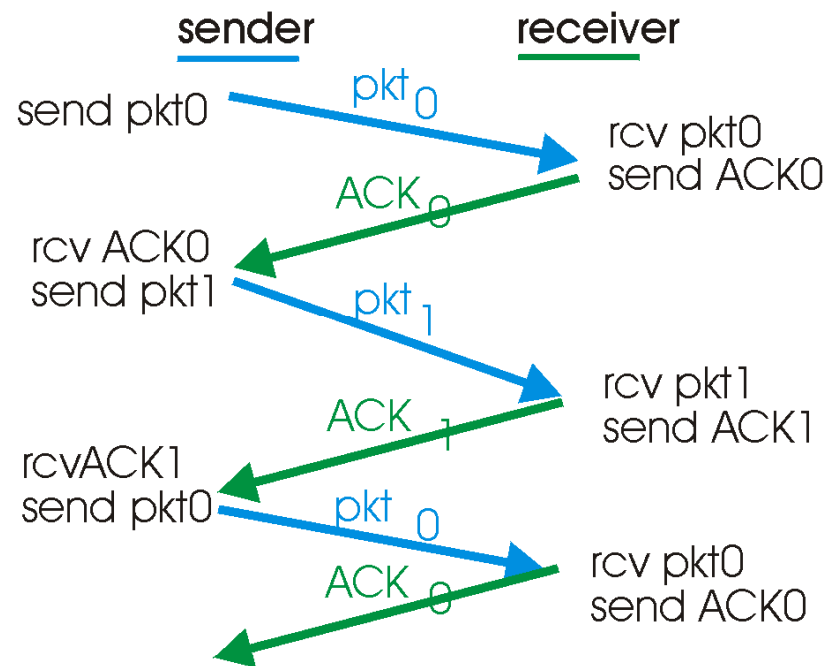
- □ retransmits if no ACK received in this time
- □ if pkt (or ACK) just delayed (not lost):
  - retransmission will be duplicate, but use of seq. #'s already handles this
  - receiver must specify seq # of pkt being ACKed
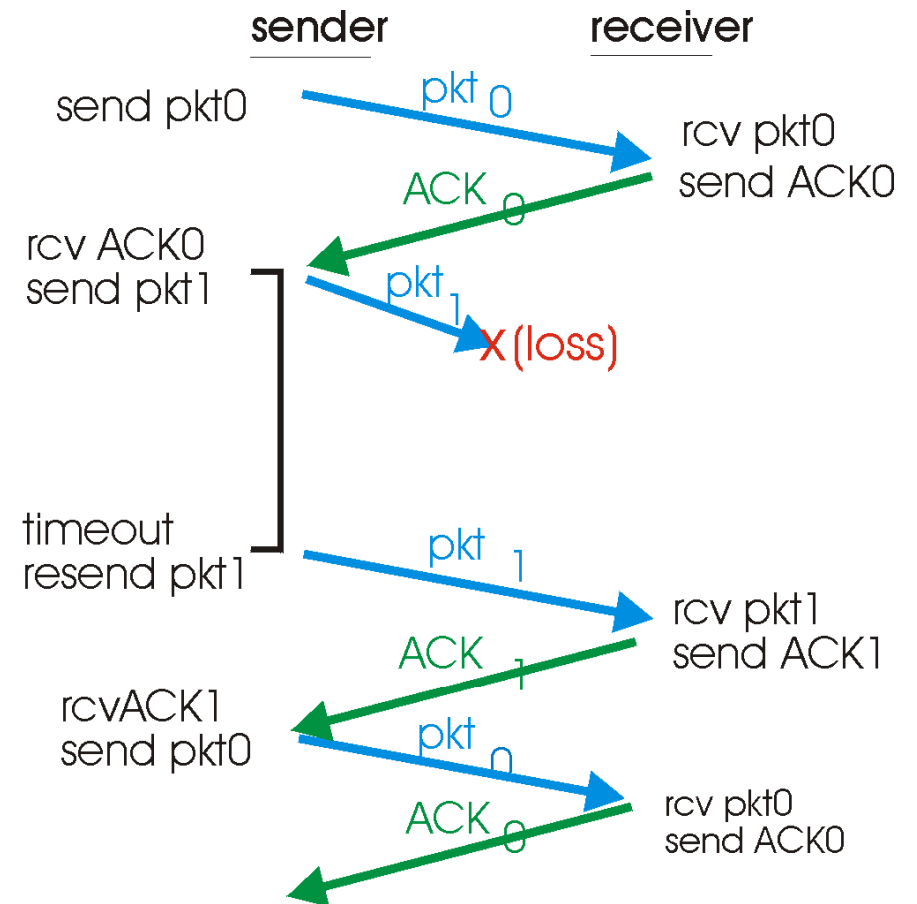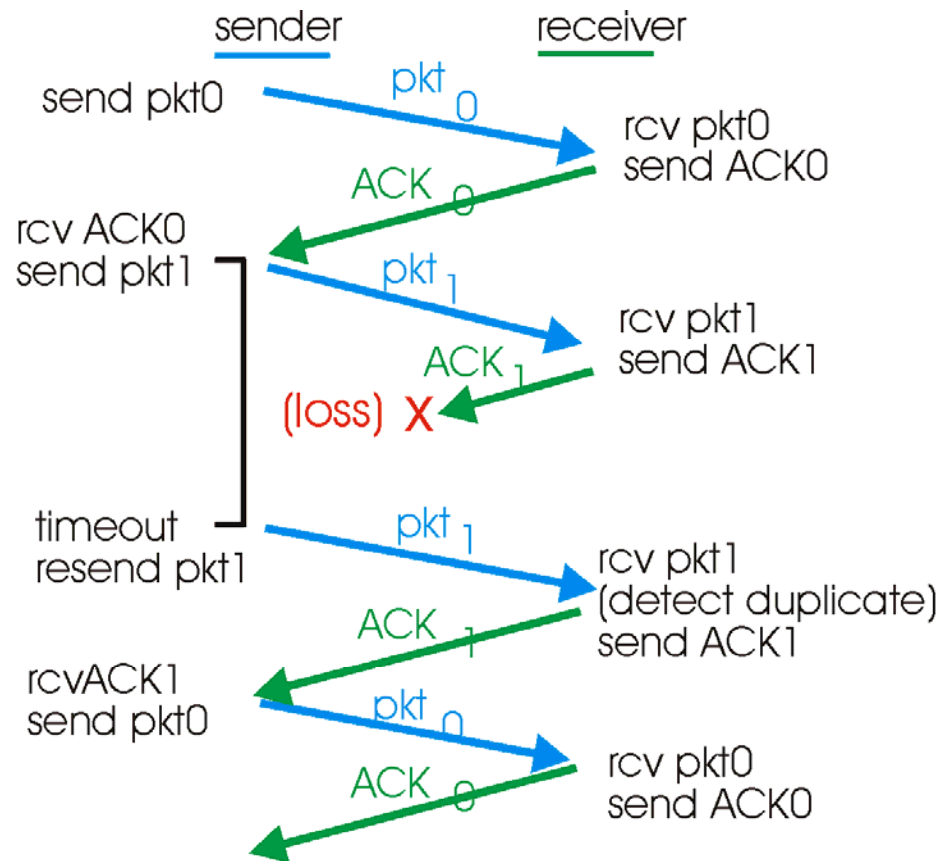- □ requires countdown timer

# rdt3.0 sender

rdt_send(data)
_____
sndpkt = make_pkt(0, data, checksum)
udt_send(sndpkt)
start_timer

**Wait for call 0from above**

rdt_rcv(rcvpkt)
_____
$\Lambda$

rdt_rcv(rcvpkt) &&
( corrupt(rcvpkt) ||
isACK(rcvpkt,1) )
_____
$\Lambda$

**Wait for ACK0**

timeout
_____
udt_send(sndpkt)
start_timer

rdt_rcv(rcvpkt)
&& notcorrupt(rcvpkt)
&& isACK(rcvpkt,1)
_____
stop_timer

rdt_rcv(rcvpkt)
&& notcorrupt(rcvpkt)
&& isACK(rcvpkt,0)
_____
stop_timer

**Wait for ACK1**

timeout
_____
udt_send(sndpkt)
start_timer

rdt_rcv(rcvpkt) &&
( corrupt(rcvpkt) ||
isACK(rcvpkt,0) )
_____
$\Lambda$

**Wait for call 1 from above**

rdt_rcv(rcvpkt)
_____
$\Lambda$

rdt_send(data)
_____
sndpkt = make_pkt(1, data, checksum)
udt_send(sndpkt)
start_timer

# rdt3.0 in action

**sender**   **receiver**

send pkt0 → pkt0 → rcv pkt0
send ACK0
ACK0 ←

rcv ACK0
send pkt1 → pkt1 → rcv pkt1
send ACK1
ACK1 ←

rcvACK1
send pkt0 → pkt0 → rcv pkt0
send ACK0
ACK0 ←

(a) operation with no loss

**sender**   **receiver**

send pkt0 → pkt0 → rcv pkt0
send ACK0
ACK0 ←

rcv ACK0
send pkt1 → pkt1 → X(loss)

timeout
resend pkt1 → pkt1 → rcv pkt1
send ACK1
ACK1 ←

rcvACK1
send pkt0 → pkt0 → rcv pkt0
send ACK0
ACK0 ←

(b) lost packet

# rdt3.0 in action



(c) lost ACK

(d) premature timeout

# 3.4.2 Pipelined Reliable Data Transfer

Performance of rdt3.0

- rdt3.0 works, but performance stinks
- example: 1 Gbps link, 15 ms e-e prop. delay, 1KB packet:

$$T_{transmit} = \frac{L \text{ (packet length in bits)}}{R \text{ (transmission rate, bps)}} = \frac{8kb/pkt}{10^{**}9 \text{ b/sec}} = 8 \text{ microsec}$$

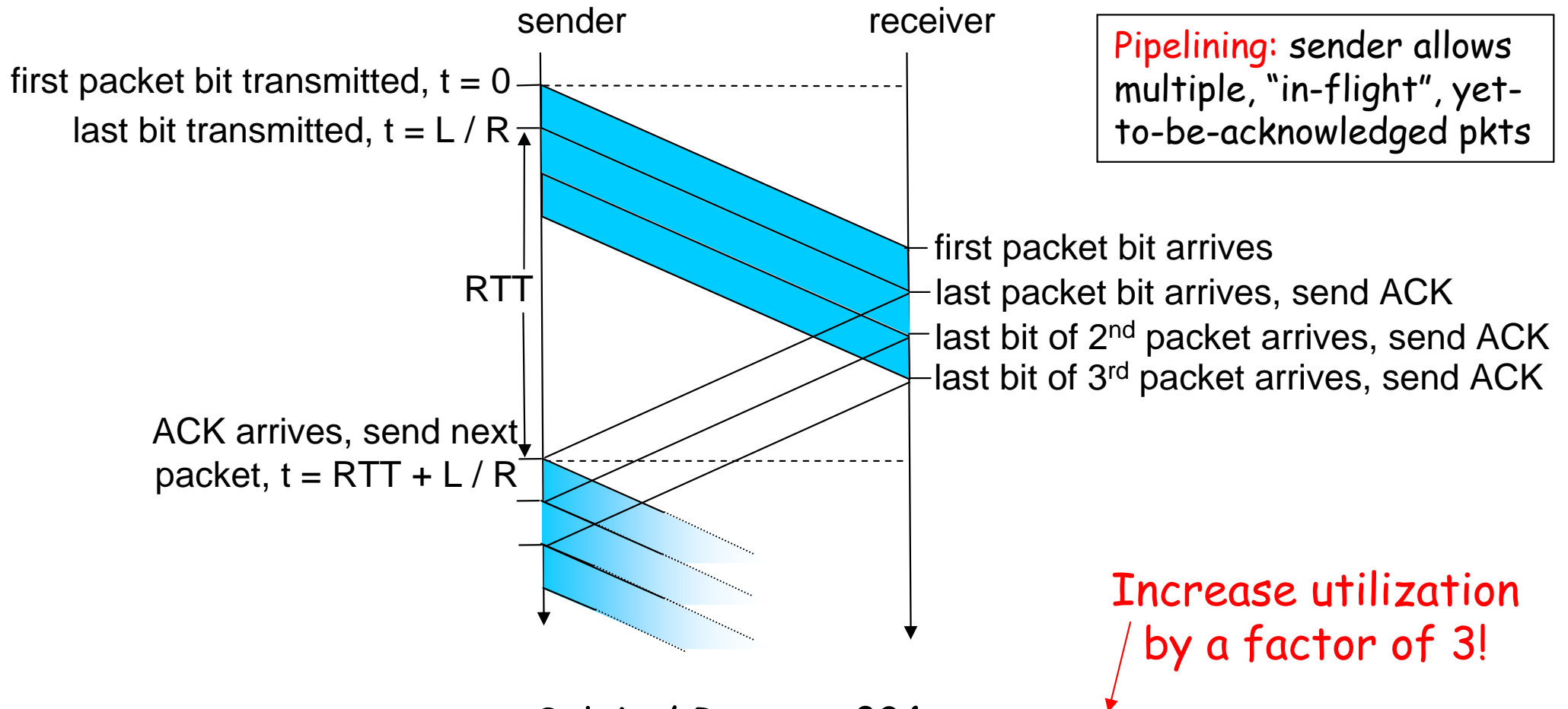$$U_{sender} = \frac{L / R}{RTT + L / R} = \frac{.008}{30.008} = 0.00027$$

- $U_{sender}$: utilization – fraction of time sender busy sending
- 1KB pkt every 30 msec -> 33kB/sec thruput over 1 Gbps link
- network protocol limits use of physical resources!

# rdt3.0: stop-and-wait operation



$$U_{sender} = \frac{L/R}{RTT + L/R} = \frac{.008}{30.008} = 0.00027$$

# Pipelined protocols: increased utilization



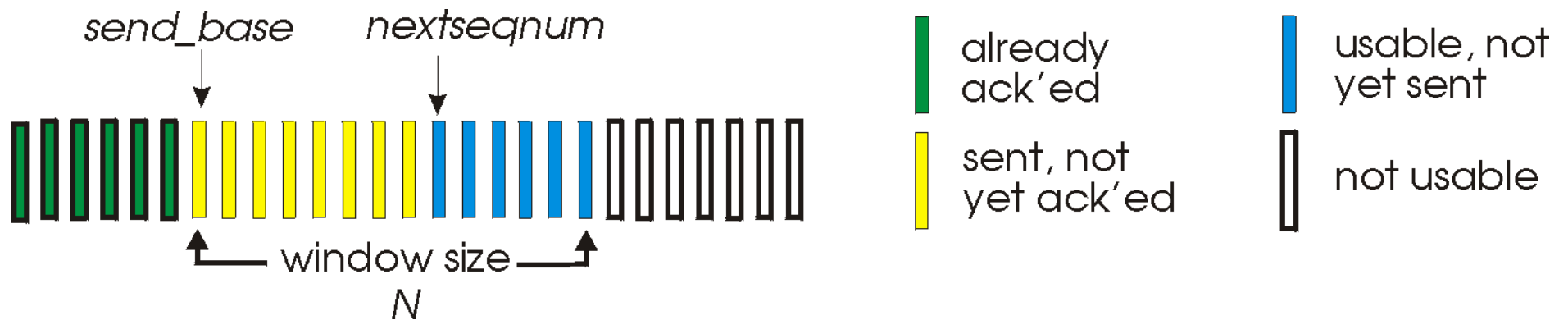sender                                   receiver

first packet bit transmitted, t = 0
last bit transmitted, t = L / R

Pipelining: sender allows multiple, "in-flight", yet-to-be-acknowledged pkts

RTT

first packet bit arrives
last packet bit arrives, send ACK
last bit of 2$^{nd}$ packet arrives, send ACK
last bit of 3$^{rd}$ packet arrives, send ACK

ACK arrives, send next packet, t = RTT + L / R

Increase utilization by a factor of 3!

$$U_{sender} = \frac{3 * L / R}{RTT + L / R} = \frac{.024}{30.008} = 0.0008$$

□ Two generic forms of pipelined protocols: *go-Back-N, selective repeat*

# 3.4.3 Go-Back-N

## Sender:

□ k-bit seq # in pkt header
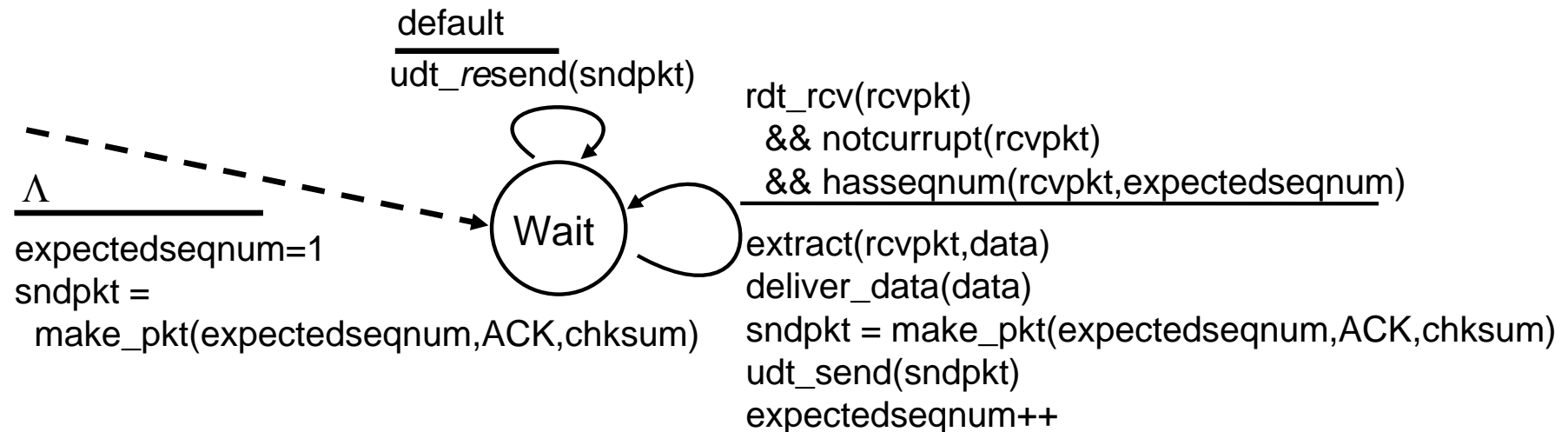
□ "window" of up to N, consecutive unack'ed pkts allowed



□ ACK(n): ACKs all pkts up to, including seq # n - "cumulative ACK"
  ○ may deceive duplicate ACKs (see receiver)
□ timer for each in-flight pkt
□ *timeout(n):* retransmit pkt n and all higher seq # pkts in window
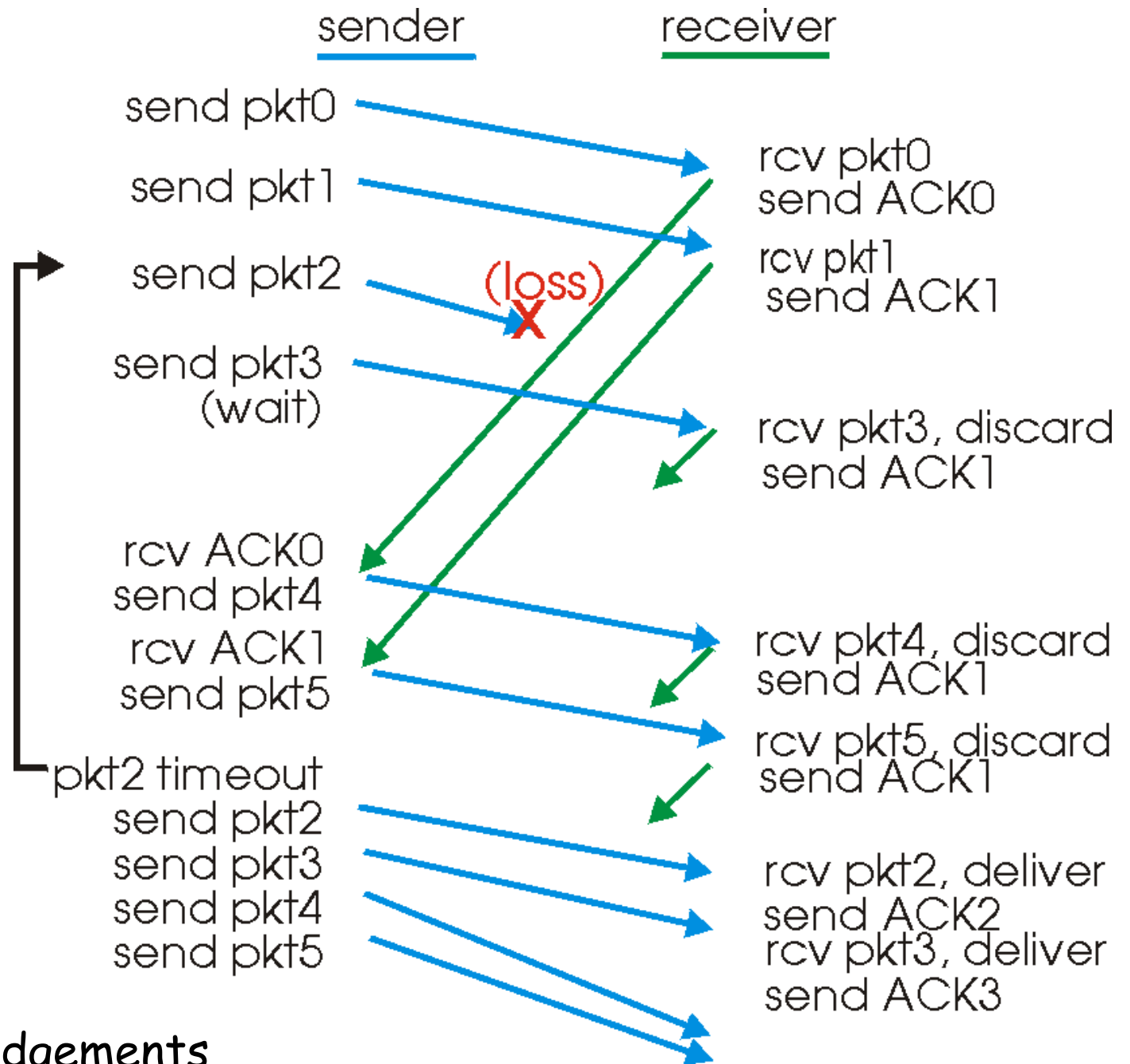
# GBN: sender extended FSM

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)
```

$\Lambda$
_____
base=1
nextseqnum=1

Wait

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

rdt_rcv(rcvpkt)
  && corrupt(rcvpkt)
_____

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

# GBN: receiver extended FSM

default
_____
udt_*re*send(sndpkt)

Λ
_____
expectedseqnum=1
sndpkt =
  make_pkt(expectedseqnum,ACK,chksum)

**Wait**

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

ACK-only: always send ACK for correctly-received pkt
  with highest *in-order* seq #

- may generate duplicate ACKs
- need only remember `expectedseqnum`

□ out-of-order pkt:

- discard (don't buffer) -> no receiver buffering!
- Re-ACK pkt with highest in-order seq #

# GBN in action

sender                          receiver

send pkt0 → rcv pkt0
                                send ACK0
send pkt1 → rcv pkt1
                                send ACK1
send pkt2        (loss) ✗
send pkt3
(wait) → rcv pkt3, discard
                                send ACK1
rcv ACK0
send pkt4 → rcv pkt4, discard
rcv ACK1                        send ACK1
send pkt5 → rcv pkt5, discard
                                send ACK1
pkt2 timeout
send pkt2 → rcv pkt2, deliver
send pkt3 → send ACK2
send pkt4    rcv pkt3, deliver
send pkt5    send ACK3

- sequence numbers
- cumulative acknowledgements
- checksums
- timeout/retransmit operation

# 3.4.4 Selective Repeat

- Go-Back-N: a single packet error can cause GBN to retransmit a large number of packets
- receiver *individually* acknowledges all correctly received pkts
  - buffers pkts, as needed, for eventual in-order delivery to upper layer
- sender only resends pkts for which ACK not received
  - sender timer for each unACKed pkt
- sender window
  - N consecutive seq #'s
  - again limits seq #s of sent, unACKed pkts

# Selective repeat: sender, receiver windows

send_base    nextseqnum

| | already ack'ed | | usable, not yet sent |
| sent, not yet ack'ed | | not usable |

window size — N

(a) sender view of sequence numbers

| | out of order (buffered) but already ack'ed | | acceptable (within window) |
| Expected, not yet received | | not usable |

window size — N

rcv_base

(b) receiver view of sequence numbers

# Selective repeat

## sender

**data from above :**

- if next available seq # in window, send pkt

**timeout(n):**

- resend pkt n, restart timer

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

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

## receiver

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

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

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

- ACK(n)

**otherwise:**

- ignore

# Selective repeat in action

pkt0 sent
`0 1 2 3` 4 5 6 7 8 9

pkt1 sent
`0 1 2 3` 4 5 6 7 8 9

pkt2 sent
`0 1 2 3` 4 5 6 7 8 9

X
(loss)

pkt3 sent, window full
`0 1 2 3` 4 5 6 7 8 9

ACK0 rcvd, pkt4 sent
0 `1 2 3 4` 5 6 7 8 9

ACK1 rcvd, pkt5 sent
0 1 `2 3 4 5` 6 7 8 9

pkt2 TIMEOUT, pkt2 resent
0 1 `2 3 4 5` 6 7 8 9

ACK3 rcvd, nothing sent
0 1 `2 3 4 5` 6 7 8 9

pkt0 rcvd, delivered, ACK0 sent
0 `1 2 3 4` 5 6 7 8 9

pkt1 rcvd, delivered, ACK1 sent
0 1 `2 3 4 5` 6 7 8 9

pkt3 rcvd, buffered, ACK3 sent
0 1 `2 3 4 5` 6 7 8 9

pkt4 rcvd, buffered, ACK4 sent
0 1 `2 3 4 5` 6 7 8 9

pkt5 rcvd, buffered, ACK5 sent
0 1 `2 3 4 5` 6 7 8 9

pkt2 rcvd, pkt2,pkt3,pkt4,pkt5
delivered, ACK2 sent
0 1 2 3 4 5 `6 7 8 9`

# Selective repeat: dilemma

Example:

□ seq #'s: 0, 1, 2, 3

□ window size=3

□ receiver sees no difference in two scenarios!

□ incorrectly passes duplicate data as new in (a)

Q: what relationship between seq # size and window size?

# Chapter 3 outline

# 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
- *send & receive buffers*

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

socket door

application writes data

socket door

application reads data

TCP send buffer

TCP receive buffer

segment

# TCP segment structure

32 bits

URG: urgent data
(generally not used)

ACK: ACK #
valid

PSH: push data now
(generally not used)

RST, SYN, FIN:
connection estab
(setup, teardown
commands)

Internet
checksum
(as in UDP)

| source port # | dest port # |
|---|---|
| sequence number | |
| acknowledgement number | |

| head len | not used | U A P R S F | Receive window |
|---|---|---|---|
| checksum | | | Urg data pnter |

Options (variable length)

application
data
(variable length)

counting
by bytes
of data
(not segments!)

# bytes
rcvr willing
to accept

# TCP seq. #'s and ACKs

Seq. #'s:
- byte stream "number" of first byte in segment's data

ACKs:
- 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

Piggyback: ACK of one direction is carried in a segment carrying data of reverse-direction

Host A      Host B

User types 'C'

Seq=42, ACK=79, data = 'C'

host ACKs receipt of 'C', echoes back 'C'

Seq=79, ACK=43, data = 'C'

host ACKs receipt of echoed 'C'

Seq=43, ACK=80

time

simple telnet scenario

# TCP Round Trip Time and Timeout

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

Q: 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**

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

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

# Example RTT estimation:

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

# TCP Round Trip Time and Timeout

## Setting the timeout

- **EstimtedRTT** plus "safety margin"
  - large variation in **EstimatedRTT** -> larger safety margin
- first estimate of how much SampleRTT deviates from EstimatedRTT:

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

(typically, β = 0.25)
```

Then set timeout interval:

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

# TCP reliable data transfer

- TCP creates rdt service on top of IP's unreliable service
- Pipelined segments
- Cumulative acks
- TCP uses *single* retransmission timer (vs. individual timer for each segment)

- Retransmissions are triggered by:
  - timeout events
  - duplicate acks
- Initially consider simplified TCP sender:
  - ignore duplicate acks
  - ignore flow control, congestion control

# 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 acknowledges previously unacked segments
  - update what is known to be acked
  - start timer if there are outstanding segments

```
NextSeqNum = InitialSeqNum
SendBase = InitialSeqNum

loop (forever) {
  switch(event)

  event: data received from application above
      create TCP segment with sequence number NextSeqNum
      if (timer currently not running)
          start timer
      pass segment to IP
      NextSeqNum = NextSeqNum + length(data)

  event: timer timeout
      retransmit not-yet-acknowledged segment with
          smallest sequence number
      start timer

  event: ACK received, with ACK field value of y
      if (y > SendBase) {
          SendBase = y
          if (there are currently not-yet-acknowledged segments)
              start timer
      }

  }  /* end of loop forever */
```

# TCP sender (simplified)

Comment:
• SendBase-1: last cumulatively ack'ed byte
Example:
• SendBase-1 = 71; y= 73, so the rcvr wants 73+ ; y > SendBase, so that new data is acked

# TCP: retransmission scenarios



Seq=92, 8 bytes data

ACK=100

X loss

timeout

Seq=92, 8 bytes data

ACK=100

SendBase = 100

time

lost ACK scenario

Seq=92, 8 bytes data

Seq=100, 20 bytes data

ACK=100

ACK=120

Seq=92, 8 bytes data

ACK=120

Seq=92 timeout

Seq=92 timeout

Sendbase = 100
SendBase = 120

SendBase = 120

time

premature timeout

# TCP retransmission scenarios (more)



Host A

Host B

Seq=92, 8 bytes data

ACK=100

Seq=100, 20 bytes data

X loss

timeout

SendBase = 120

ACK=120

time

Cumulative ACK scenario

# 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 duplicate ACK, 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 |

# Fast Retransmit

□ Time-out period often relatively long:
  ○ long delay before resending lost packet (Consider LAN cases)

□ Detect lost segments via duplicate ACKs.

  ○ Sender often sends many segments back-to-back

  ○ If segment is lost, there will likely be many duplicate ACKs.

□ If sender receives 3 ACKs for the same data, it supposes that segment after ACKed data was lost:

  ○ fast retransmit: resend segment before timer expires

# Fast retransmit algorithm:

event: ACK received, with ACK field value of y
        if (y > SendBase) {
            SendBase = y
            if (there are currently not-yet-acknowledged segments)
                start timer
        }
        else {
            increment count of dup ACKs received for y
            if (count of dup ACKs received for y = 3) {
                resend segment with sequence number y
            }
        }

a duplicate ACK for
already ACKed segment

fast retransmit

# TCP Flow Control

□ **receive side of TCP connection has a receive buffer:**

  ○ app process may be slow at reading from the buffer

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

□ **Receive window:** give the sender an idea how much free buffer space available at the receiver

# TCP Flow control: how it works



(Suppose TCP receiver discards out-of-order segments)

- spare room in buffer

```
= RcvWindow
= RcvBuffer-[LastByteRcvd -
  LastByteRead]
```

- Rcvr advertises spare room by including value of `RcvWindow` in segments
  - Initial RcvWindow=RcvBuffer
- Sender limits unACKed data to `RcvWindow`
  - guarantees receive buffer doesn't overflow
- 1-byte zero window
  - Can send 1 byte data when RcvWindow is 0

# TCP Connection Management

<span style="color:red">Recall:</span> TCP sender, receiver establish "connection" before exchanging data segments

☐ initialize TCP variables: seq. #s, buffers, flow control info (e.g. `RcvWindow`)

## Three way handshake:

<span style="color:red">Step 1:</span> client host sends TCP SYN segment to server
- specifies initial seq # (isn)
- no data

<span style="color:red">Step 2:</span> server host receives SYN, replies with SYNACK segment
- server allocates buffers
- specifies server initial seq. #

<span style="color:red">Step 3:</span> client receives SYNACK, replies with ACK segment, which may contain data

**Client host**                    **Server host**

Connection request ——  SYN=1, seq=client_isn

—— Connection granted

SYN=1, seq=server_isn, ack=client_isn+1

ACK ——

SYN=0, seq=client_isn+1, ack=server_isn+1

Time                                Time

# TCP Connection Management (cont.)

## Closing a connection:

client closes socket:
```
clientSocket.close();
```

**Step 1:** client end system sends TCP FIN control segment to server

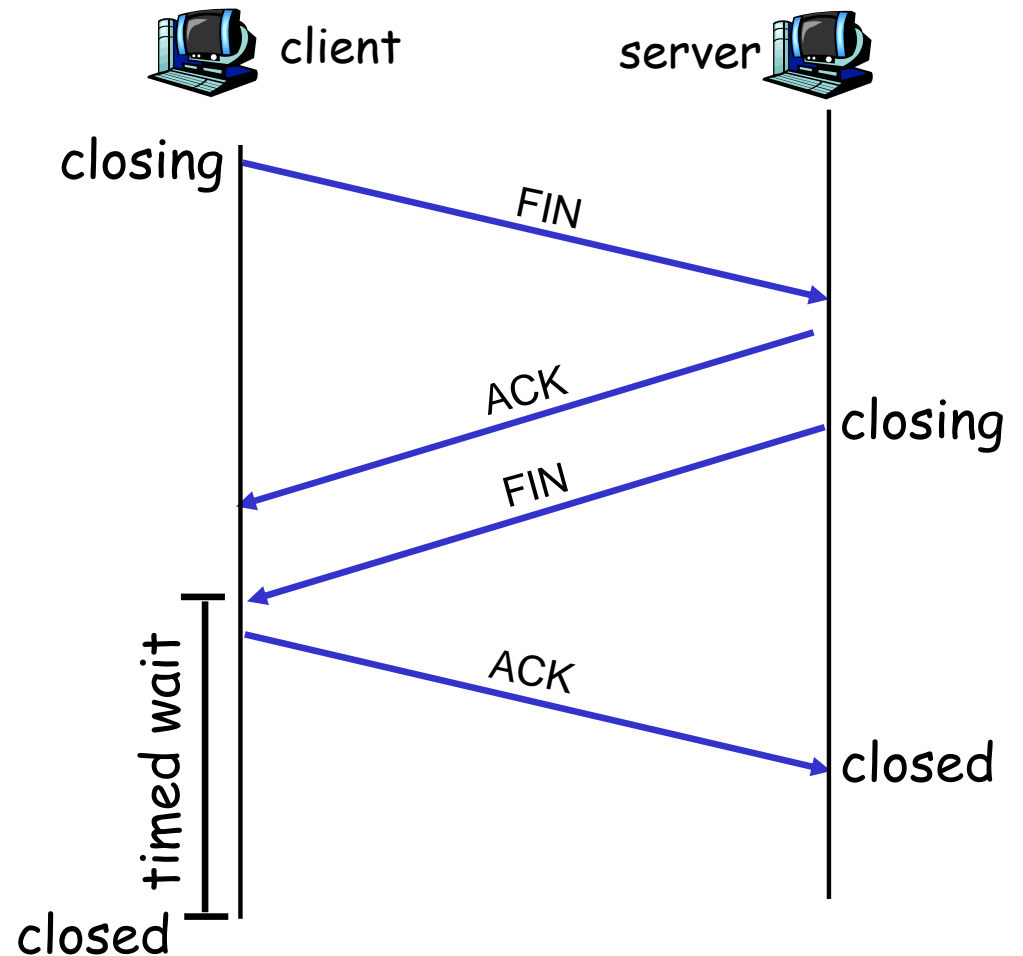**Step 2:** server receives FIN, replies with ACK. Closes connection, sends FIN.

client     server

close    FIN

ACK    close

FIN

timed wait    ACK

closed

# TCP Connection Management (cont.)

**Step 3:** client receives FIN, replies with ACK.

   ○ Enters "timed wait" - will respond with ACK to received FINs

**Step 4:** server, receives ACK. Connection closed.

**Note:** with small modification, can handle simultaneous FINs.

# TCP Connection Management (cont)



TCP client lifecycle

TCP server lifecycle

CLOSED
client application initiates a TCP connection
send SYN
SYN_SENT
receive SYN & ACK send ACK
ESTABLISHED
client application initiates close connection
send FIN
FIN_WAIT_1
receive ACK send nothing
FIN_WAIT_2
receive FIN send ACK
TIME_WAIT
wait 30 seconds

CLOSED
server application creates a listen socket
LISTEN
receive SYN send SYN & ACK
SYN_RCVD
receive ACK send nothing
ESTABLISHED
receive FIN send ACK
CLOSE_WAIT
send FIN
LAST_ACK
receive ACK send nothing

# Chapter 3 outline

# Principles of Congestion Control

**Congestion:**

☐ informally: "too many sources sending too much data too fast for *network* to handle"

☐ different from flow control!

☐ manifestations:

○ lost packets (buffer overflow at routers)

○ long delays (queueing in router buffers)

☐ a top-10 problem!

# Approaches towards congestion control

Two broad approaches towards congestion control:

## End-end congestion control:

- no explicit feedback from network
- congestion inferred from end-system observed loss, delay
- approach taken by TCP

## Network-assisted congestion control:

- routers provide feedback to end systems
  - single bit indicating congestion (SNA, DECbit, TCP/IP ECN, ATM)
  - explicit rate sender should send at

# Chapter 3 outline

# TCP Congestion Control

end-end control (no network assistance)

How does sender limit transmission rate?

□ transmission limit:

**LastByteSent-LastByteAcked**

$$\leq \text{CongWin}$$

□ Roughly,

$$\text{rate} = \frac{\text{CongWin}}{\text{RTT}} \text{ Bytes/sec}$$

□ **CongWin** is dynamic, function of perceived network congestion

How does sender perceive congestion?

□ loss event = timeout *or* 3 duplicate acks

□ TCP sender reduces rate (**CongWin**) after loss event

How to change? three mechanisms:

  ○ AIMD
  ○ slow start
  ○ conservative after timeout events

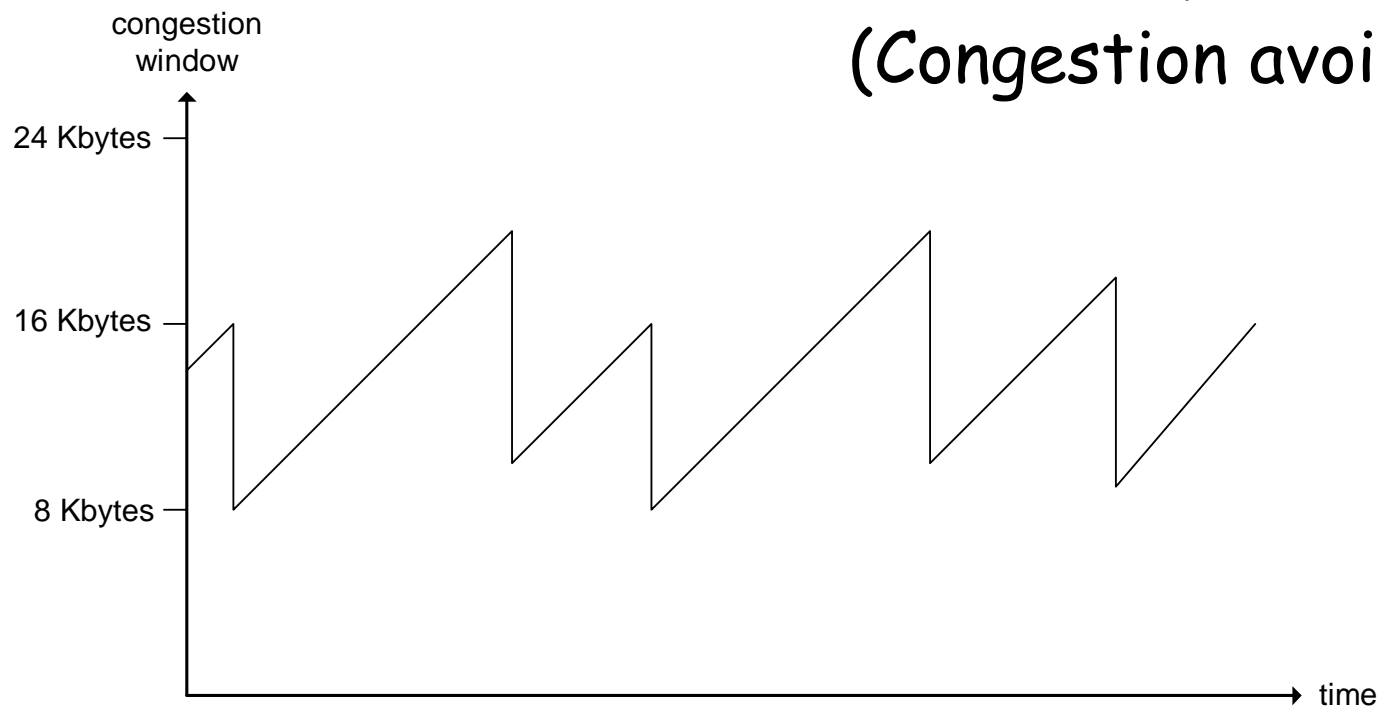# TCP AIMD

<span style="color:red">multiplicative decrease:</span>
cut `CongWin` in half after loss event

<span style="color:red">additive increase:</span>
increase `CongWin` by 1 MSS every RTT in the absence of loss events: *probing*
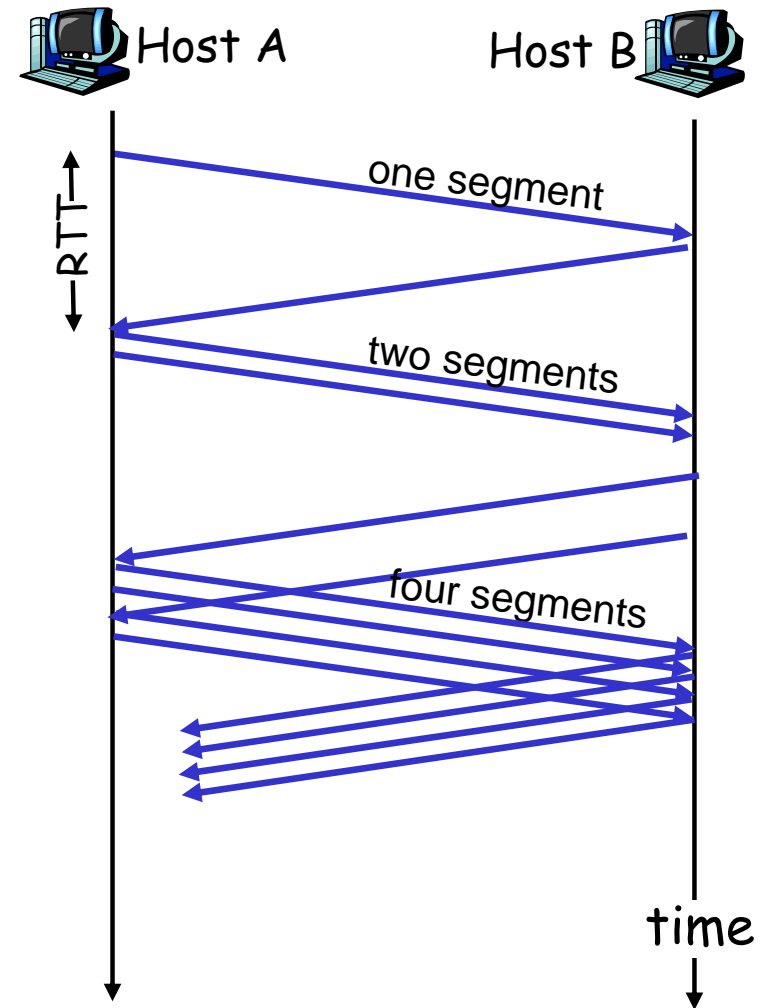
(Congestion avoidance)

congestion window

24 Kbytes —

16 Kbytes —

8 Kbytes —

time

<span style="color:red">Long-lived TCP connection</span>

# TCP Slow Start

□ **When connection begins,**
  **`CongWin` = 1 MSS**

  ○ Example: MSS = 500
    bytes & RTT = 200 msec

  ○ initial rate = 20 kbps

□ **available bandwidth may**
  **be >> MSS/RTT**

  ○ desirable to quickly ramp
    up to respectable rate

□ When connection begins,
  increase rate
  exponentially fast until
  first loss event

# TCP Slow Start (more)

□ When connection begins, increase rate exponentially until first loss event:

○ double `CongWin` every RTT

○ done by incrementing `CongWin` for every ACK received

□ <span style="color:red">Summary:</span> initial rate is slow but ramps up exponentially fast

Host A          Host B

RTT

one segment

two segments

four segments

time

# Refinement

□ After 3 dup ACKs:
   - ○ `CongWin` is cut in half
   - ○ window then grows linearly

□ But after timeout event:
   - ○ `CongWin` instead set to 1 MSS;
   - ○ window then grows exponentially
   - ○ to a threshold, then grows linearly

**Philosophy:**

• 3 dup ACKs indicates network capable of delivering some segments
• timeout before 3 dup ACKs is "more alarming"
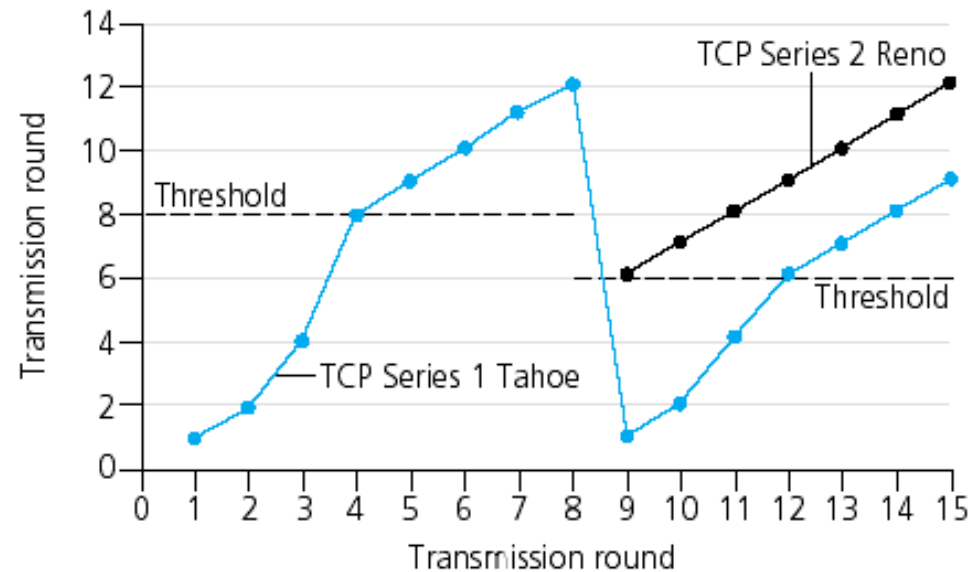
# Refinement (more)

Q: When should the exponential increase switch to linear?

A: When `CongWin` gets to 1/2 of its value before timeout.

## Implementation:

□ Variable Threshold

□ At loss event, Threshold is set to 1/2 of CongWin just before loss event

# Summary: TCP Congestion Control

□ When `CongWin` is below `Threshold`, sender in slow-start phase, window grows exponentially.

□ When `CongWin` is above `Threshold`, sender is in congestion-avoidance phase, window grows linearly.

□ When a triple duplicate ACK occurs, `Threshold` set to `CongWin/2` and `CongWin` set to `Threshold`.

□ When timeout occurs, `Threshold` set to `CongWin/2` and `CongWin` is set to 1 MSS.

# TCP sender congestion control

| Event | State | TCP Sender Action | Commentary |
|---|---|---|---|
| ACK receipt for previously unacked data | Slow Start (SS) | CongWin = CongWin + MSS, If (CongWin > Threshold) set state to "Congestion Avoidance" | Resulting in a doubling of CongWin every RTT |
| ACK receipt for previously unacked data | Congestion Avoidance (CA) | CongWin = CongWin+MSS * (MSS/CongWin) | Additive increase, resulting in increase of CongWin by 1 MSS every RTT |
| Loss event detected by triple duplicate ACK | SS or CA | Threshold = CongWin/2, CongWin = Threshold, Set state to "Congestion Avoidance" | Fast recovery, implementing multiplicative decrease. CongWin will not drop below 1 MSS. |
| Timeout | SS or CA | Threshold = CongWin/2, CongWin = 1 MSS, Set state to "Slow Start" | Enter slow start |
| Duplicate ACK | SS or CA | Increment duplicate ACK count for segment being acked | CongWin and Threshold not changed |