**TRƯỜNG ĐẠI HỌC CÔNG NGHIỆP TP. HỒ CHÍ MINH**
**KHOA CÔNG NGHỆ THÔNG TIN**

INDUSTRIAL
UNIVERSITY
OF HOCHIMINH CITY
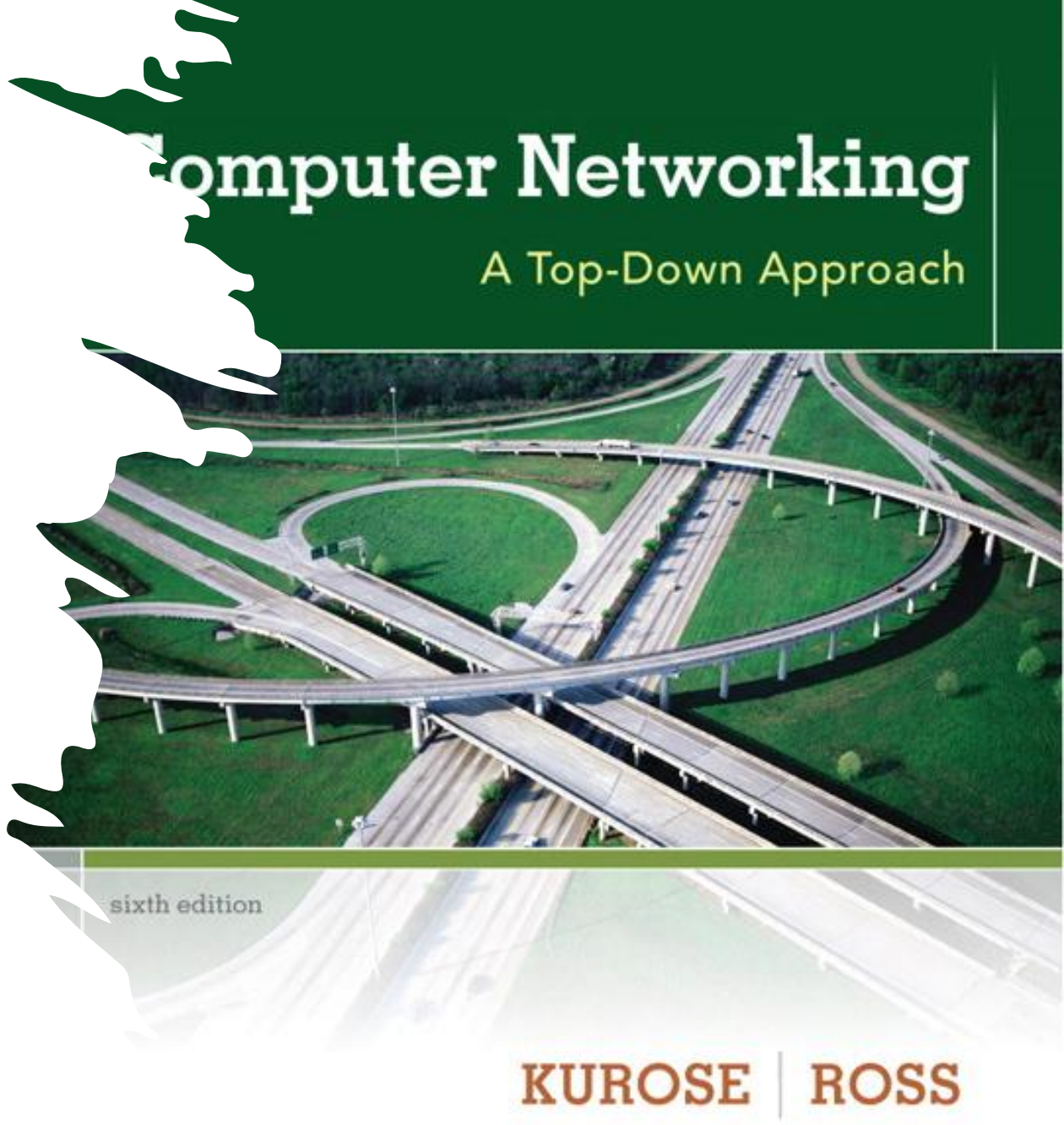
# Chapter 3: Transport Layer

*Computer Networking: A Top-Down Approach*
6th edition
Jim Kurose, Keith Ross
Addison-Wesley

Computer Networking
A Top-Down Approach

sixth edition

KUROSE | ROSS
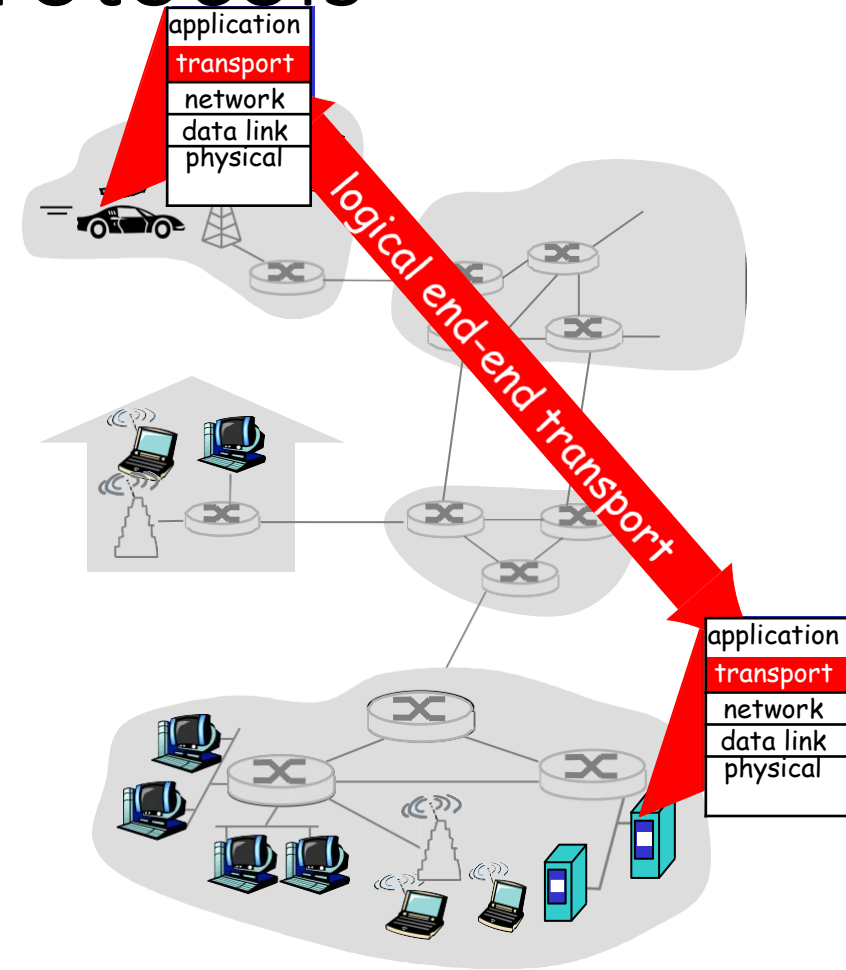
# Chapter 3: Transport Layer

## Our goals:

❖ understand principles behind transport layer services:

  ▪ multiplexing/demultiplexing
  ▪ reliable data transfer
  ▪ flow control
  ▪ congestion control

❖ learn about transport layer protocols in the Internet:

  ▪ UDP: connectionless transport
  ▪ TCP: connection-oriented transport
  ▪ TCP congestion control

# Chapter 3 outline

1. Transport-layer services
2. Multiplexing and demultiplexing
3. Connectionless transport: UDP
4. Principles of reliable data transfer

5. Connection-oriented transport: TCP
   - segment structure
   - reliable data transfer
   - flow control
   - connection management

# Transport services and protocols

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

application
transport
network
data link
physical

logical end-end transport
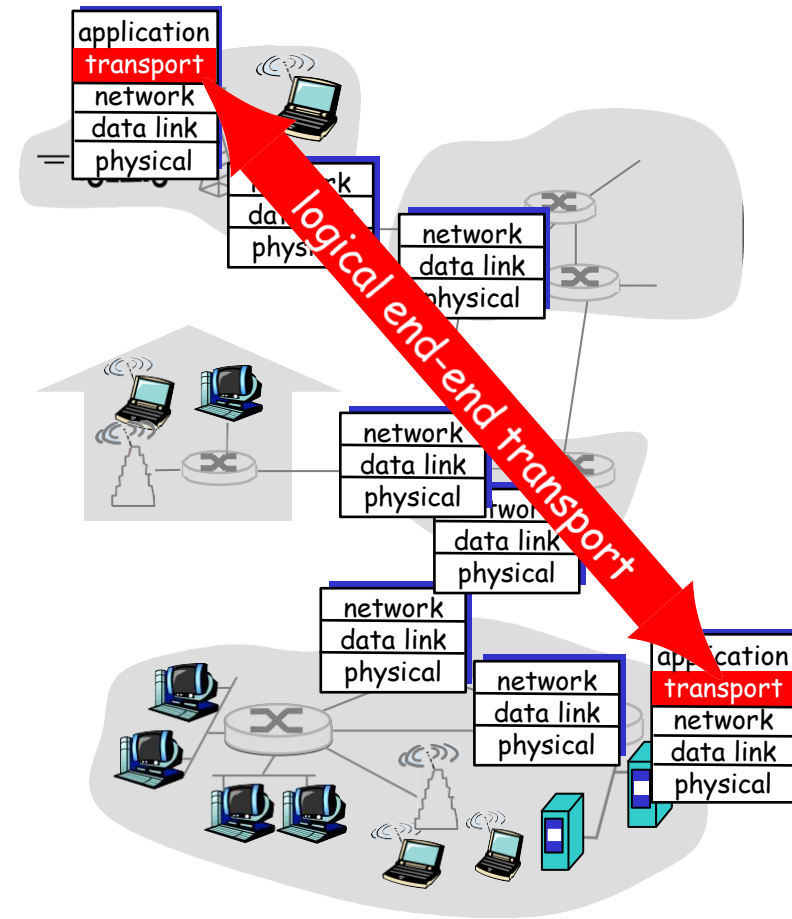
application
transport
network
data link
physical

# Transport vs. network layer

❖ **network layer:** logical communication between hosts

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

# Internet transport-layer protocols
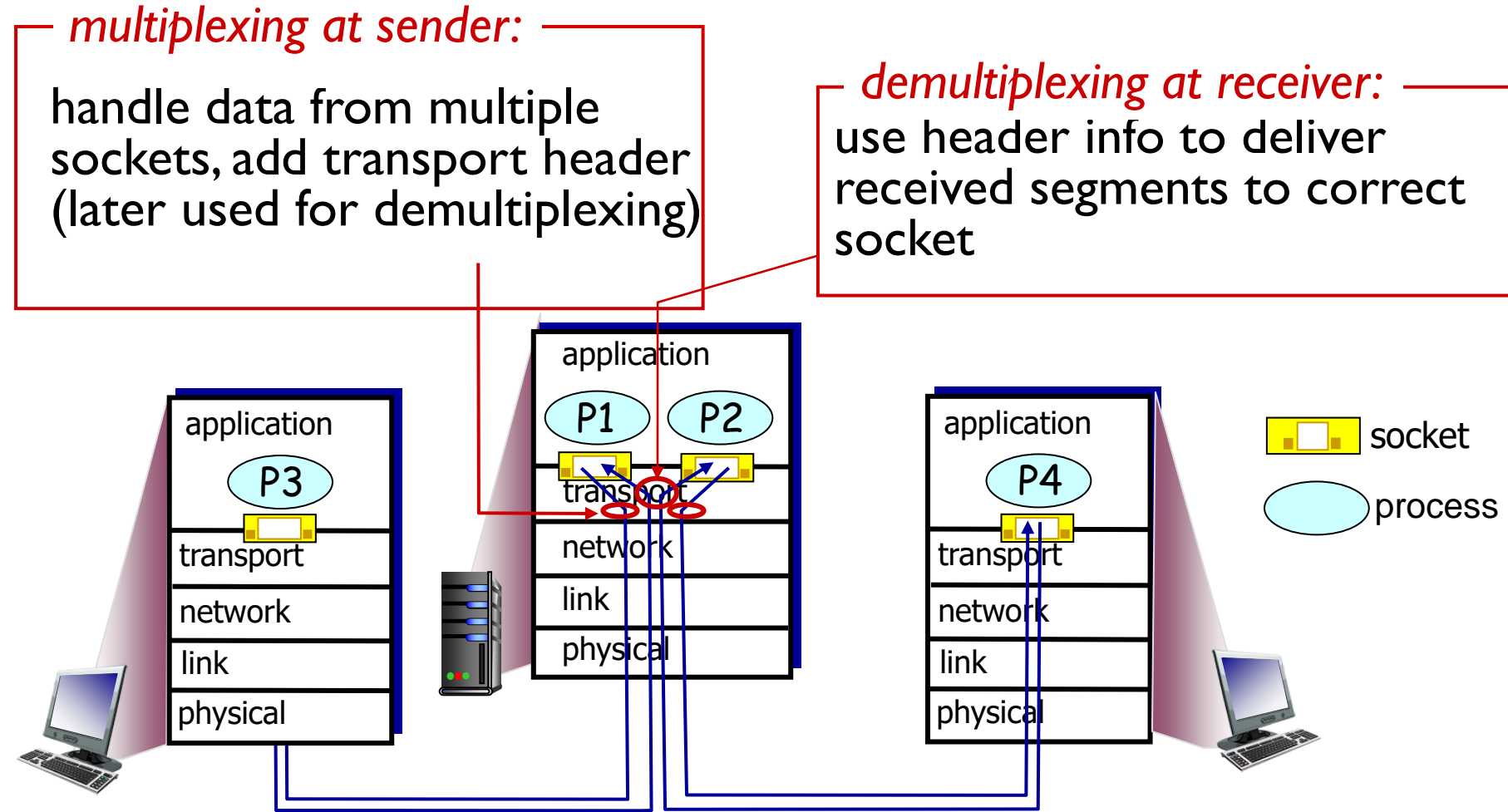
❖ reliable, in-order delivery (TCP)
- congestion control
- flow control
- connection setup

❖ unreliable, unordered delivery: UDP
- no-frills extension of "best-effort" IP

❖ services not available:
- delay guarantees
- bandwidth guarantees

# Chapter 3 outline

1. Transport-layer services
2. <span style="color:red">Multiplexing and demultiplexing</span>
3. Connectionless transport: UDP
4. Principles of reliable data transfer

5. Connection-oriented transport: TCP
   - segment structure
   - reliable data transfer
   - flow control
   - connection management
6. Principles of congestion control
7. TCP congestion control

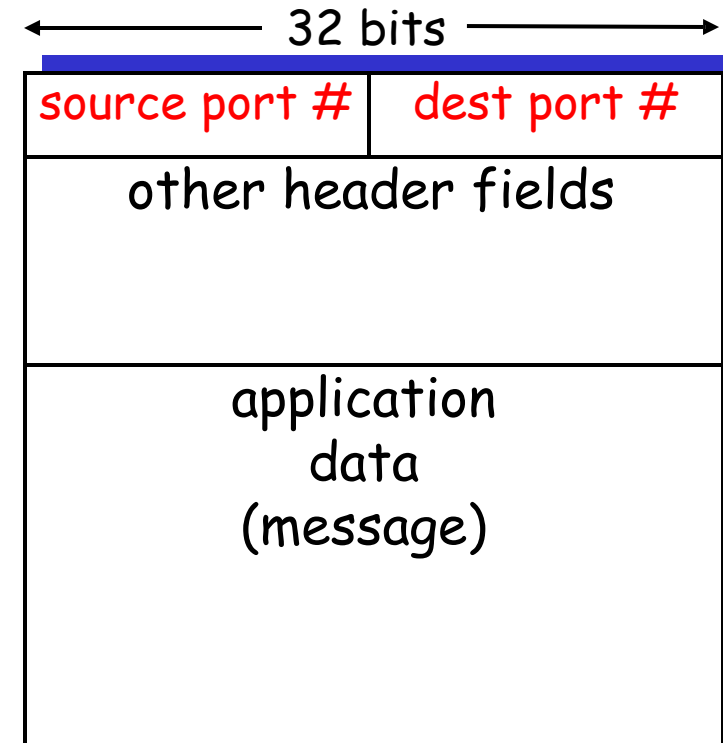# Multiplexing/demultiplexing



*multiplexing at sender:*

handle data from multiple sockets, add transport header (later used for demultiplexing)

*demultiplexing at receiver:*

use header info to deliver received segments to correct socket

socket

process

# 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

❖ **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

❖ *recall:* created socket has host-local port #:

```
DatagramSocket mySocket1
= new DatagramSocket(12534);
```

❖ *recall:* when creating datagram to send into UDP socket, must specify
- destination IP address
- destination port #

---

❖ when host receives UDP segment:
- checks destination port # in segment
- directs UDP segment to socket with that port #

➡ IP datagrams with *same dest. port #,* but different source IP addresses and/or source port numbers will be directed to *same socket* at dest
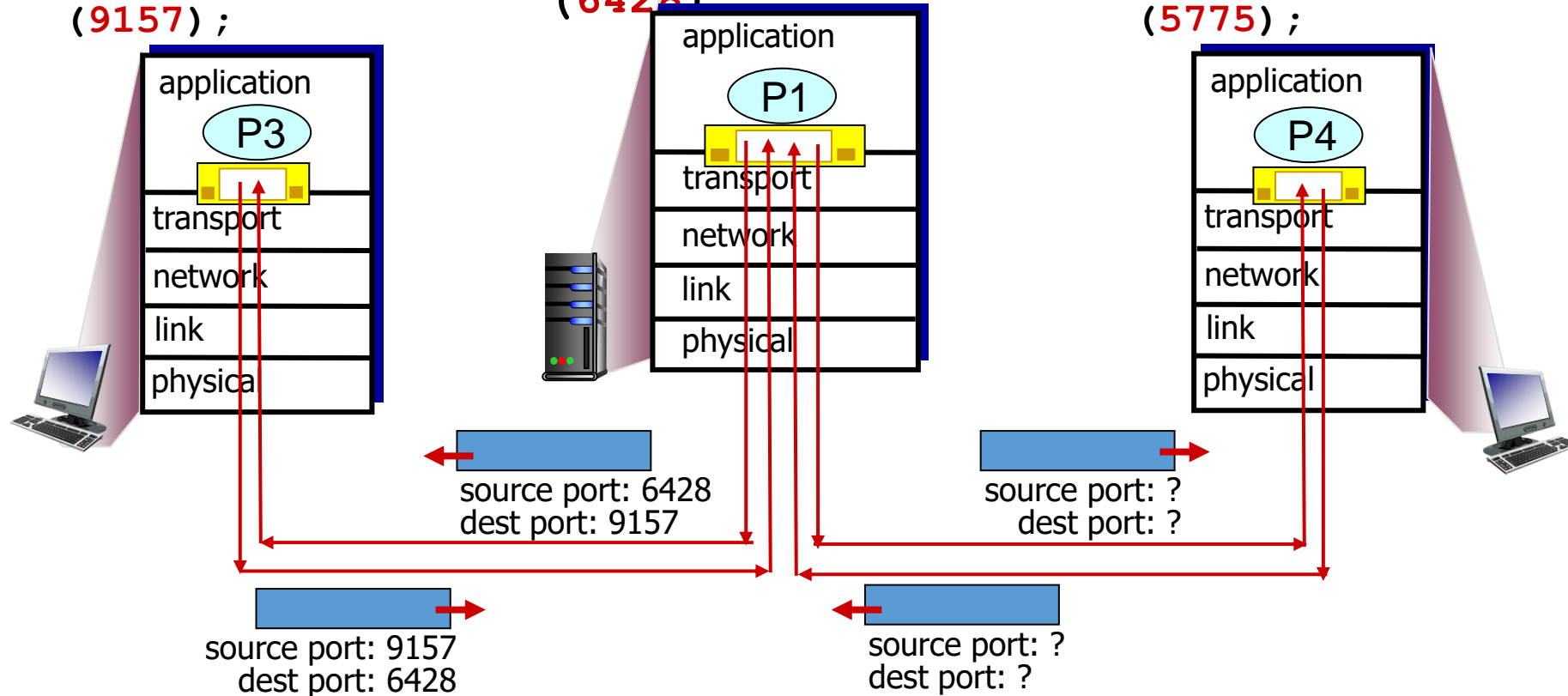
# Connectionless demux (cont)

# 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.)



source IP,port: B,80
dest IP,port: A,9157

source IP,port: A,9157
dest IP, port: B,80

source IP,port: C,5775
dest IP,port: B,80

source IP,port: C,9157
dest IP,port: B,80

host: IP address A

server: IP address B

host: IP address C

three segments, all destined to IP address: B,
dest port: 80 are demultiplexed to *different* sockets

# Connection-oriented demux (cont.)

# Chapter 3 outline

1. Transport-layer services
2. Multiplexing and demultiplexing
3. <span style="color:red">Connectionless transport: UDP</span>
4. Principles of reliable data transfer

5. Connection-oriented transport: TCP
   - segment structure
   - reliable data transfer
   - flow control
   - connection management

# 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

<div style="border: 2px solid red; padding: 10px;">

**Why is there a UDP?**

- ❖ no connection establishment (which can add delay)
- ❖ simple: no connection state at sender, receiver
- ❖ small segment header
- ❖ no congestion control: UDP can blast away as fast as desired

</div>

# 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

**<span style="color:red">Goal:</span>** detect "errors" (e.g., flipped bits) in transmitted segment

**<span style="color:red">Sender:</span>**

- 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

**<span style="color:red">Receiver:</span>**

- compute checksum of received segment
- check if computed checksum equals checksum field value:
  - NO - error detected
  - YES - no error detected. <span style="color:blue">But maybe errors nonetheless? More later ....</span>

# Internet Checksum Example

example: add three words 16-bit

```
0110011001100110
0101010101010101
0000111100001111
```

---

sum          1100101011001010

checksum     0011010100110101

*Note:* when adding numbers, a carryout from the most significant bit needs to be added to the result

# Chapter 3 outline

1. Transport-layer services

2. Multiplexing and demultiplexing

3. Connectionless transport: UDP

4. **Principles of reliable data transfer**

5. Connection-oriented transport: TCP
   - segment structure
   - reliable data transfer
   - flow control
   - connection management

# Principles of Reliable data transfer

❖ important in app., transport, link layers

❖ top-10 list of important networking topics!



(a) provided service

❖ characteristics of unreliable channel will determine complexity of reliable data transfer protocol (rdt)

# Principles of Reliable data transfer

❖ important in app., transport, link layers

❖ top-10 list of important networking topics!

application layer

transport layer

sending process

receiver process

data

data

reliable channel

unreliable channel

(a) provided service

(b) service implementation

❖ characteristics of unreliable channel will determine complexity of reliable data transfer protocol (rdt)

# Principles of Reliable data transfer

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



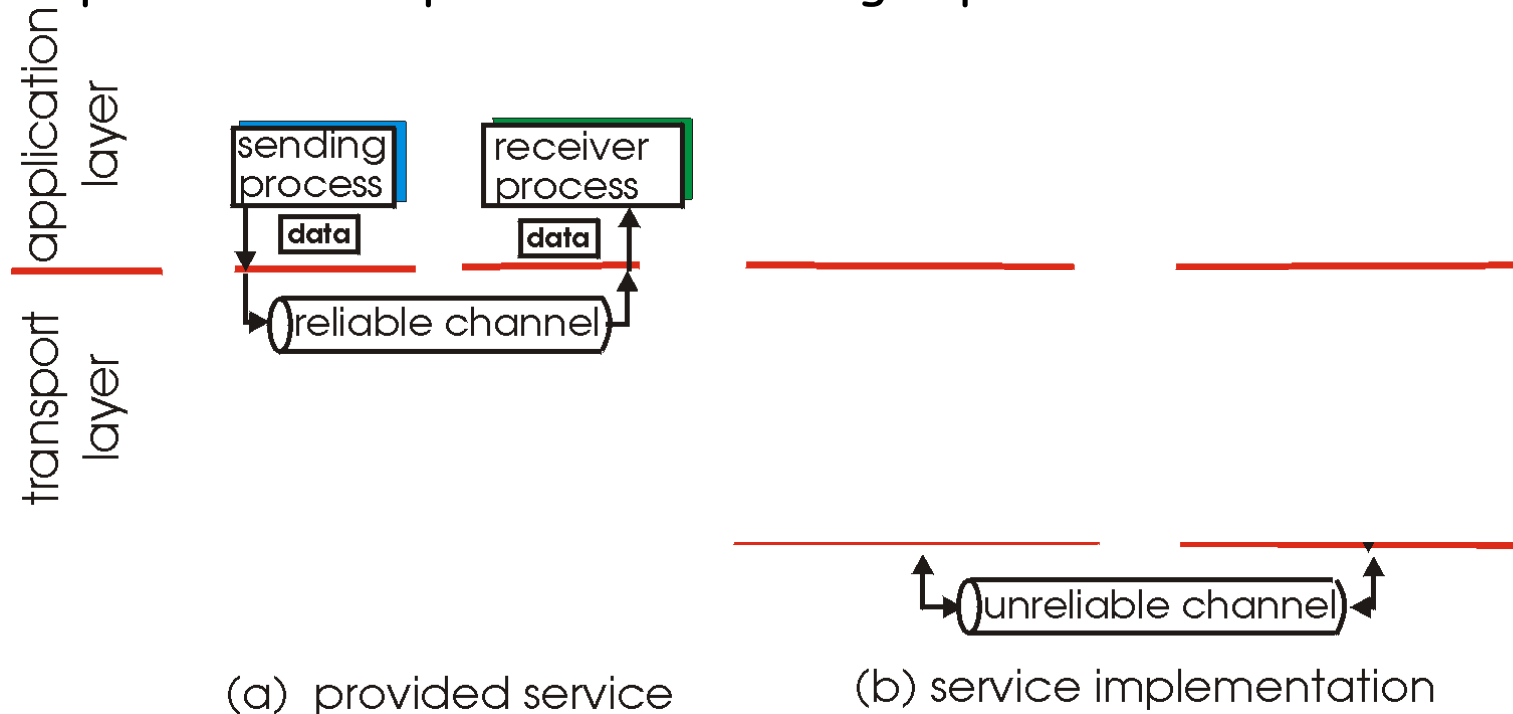(a) provided service          (b) service implementation

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

# 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

**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

state 2

event
actions

# 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:

*How do humans recover from "errors" during conversation?*

# 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:
  ▪ <span style="color:red">acknowledgements (ACKs):</span> receiver explicitly tells sender that pkt received OK
  ▪ <span style="color:red">negative acknowledgements (NAKs):</span> 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 (checksum)
  ▪ receiver feedback: control msgs (ACK,NAK) rcvr->sender

# rdt2.0: FSM specification

rdt_send(data)

sndpkt = make_pkt(data, checksum)

udt_send(sndpkt)

receiver

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

```
Wait  for        Wait for
call from        ACK  or
above            NAK
```

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

**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)

Wait for call from above

Wait for ACK or NAK

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

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)

Wait for call from below

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

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 retransmits current pkt if ACK/NAK corrupted
❖ sender adds *sequence number* to each pkt
❖ 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)
_____

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

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

Wait for
call 0 from
above

Wait for
ACK  or
NAK 0

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

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

Wait for
ACK  or
NAK 1

Wait for
call 1 from
above

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

 rdt_send(data)
_____

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

# 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

## Sender:

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

## Receiver:

- ❖ 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

# Recap

| Rdt Version | Scenario | Features Added |
| --- | --- | --- |
| 1.0 | No error | Nothing |
| 2.0 | Data Bit Error | Checksum, ACK/NAK |
| 2.1 | Data Bit Error<br>ACK/NAK Bit Error | Checksum, ACK/NAK,<br>Sequence Number |
| 3.0 | What else? | |

# rdt3.0: channels with errors and loss

**New assumption:** underlying channel can also lose packets (data or ACKs)
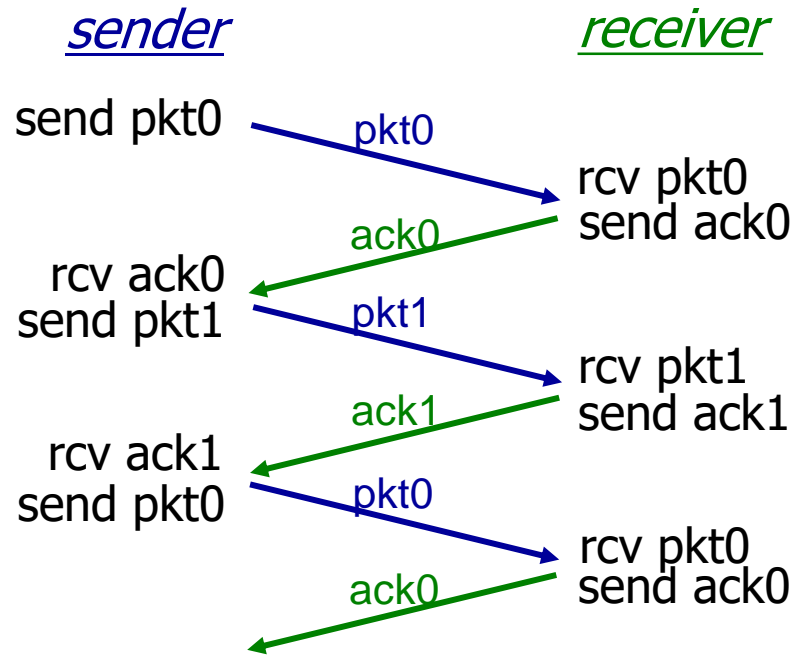
- checksum, seq. #, ACKs, retransmissions will be of help, but not enough
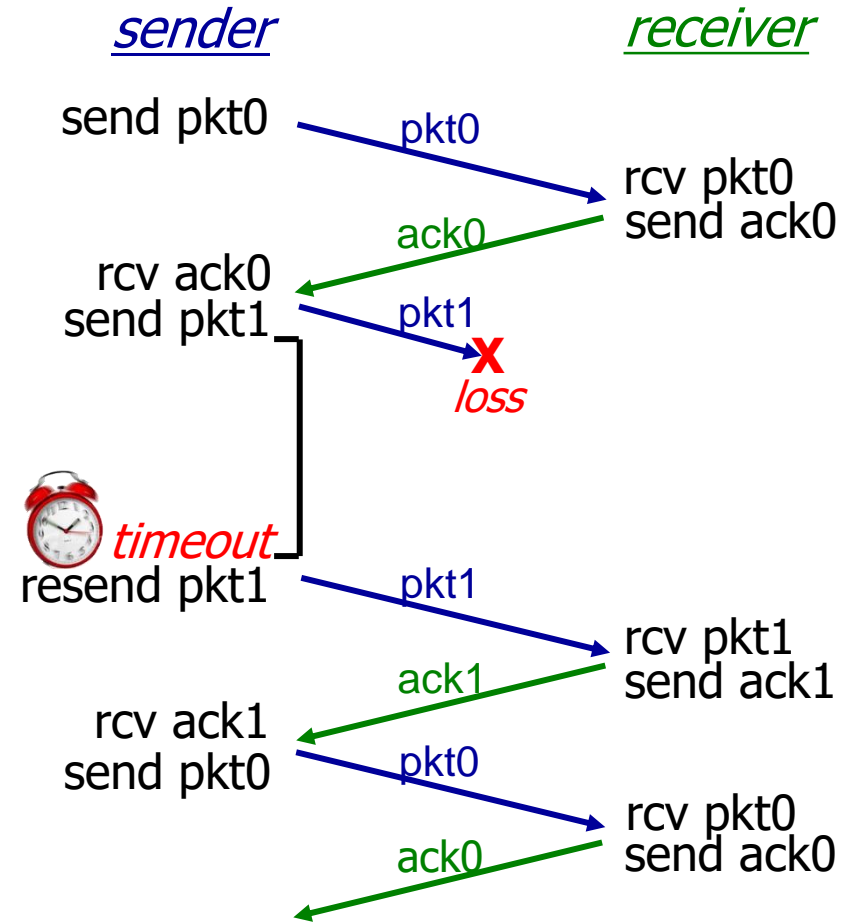
**Approach:** sender waits "reasonable" amount of time for ACK

- ❖ 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
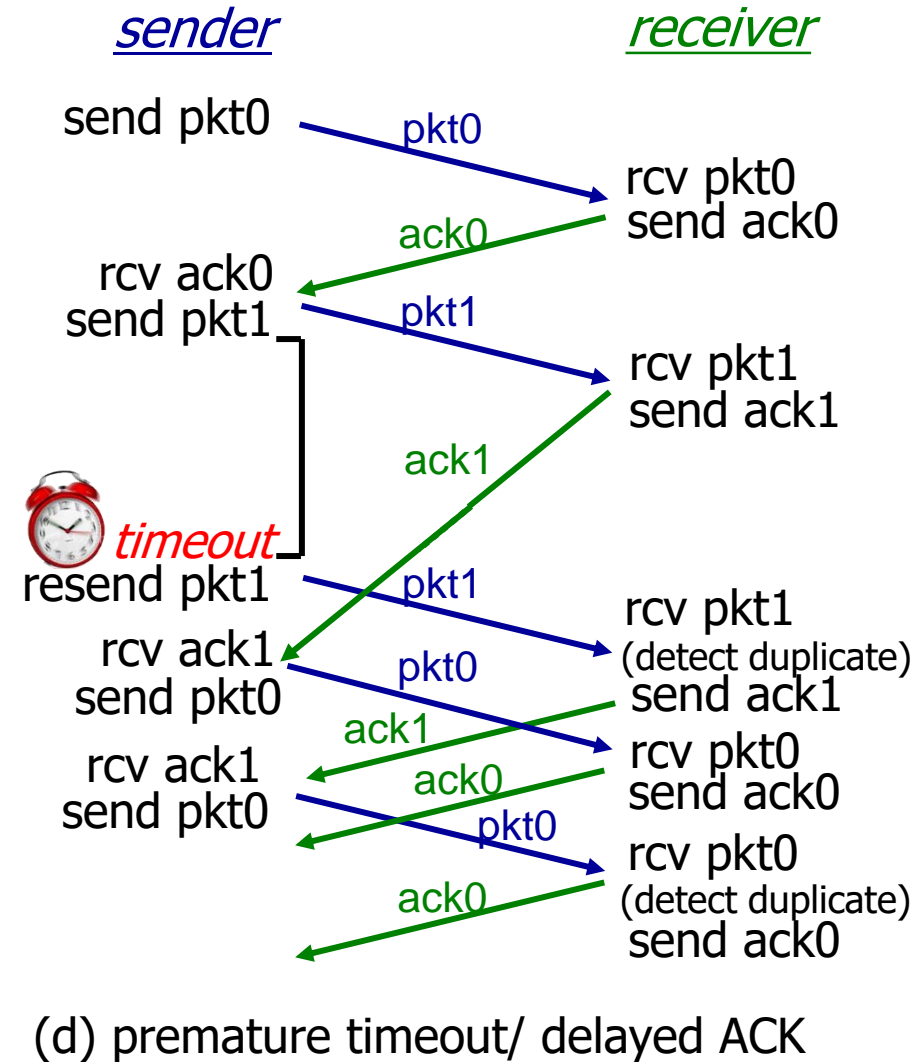- ❖ requires countdown timer

# rdt3.0 in action



(a) no loss

(b) packet loss

# rdt3.0 in action



(c) ACK loss

(d) premature timeout/ delayed ACK

# Recap

| Rdt Version | Scenario | Features Added |
|---|---|---|
| 1.0 | No error | Nothing |
| 2.0 | Data Bit Error | Checksum, ACK/NAK |
| 2.1 | Data Bit Error ACK/NAK Bit Error | Checksum, ACK/NAK, Sequence Number |
| 3.0 | Data Bit Error ACK/NAK Bit Error Packet Loss | Checksum, ACK/NAK, Sequence Number, Timeout |

# Performance of rdt3.0

- ❖ rdt3.0 works, but performance stinks
- ❖ ex: 1 Gbps link, 15 ms prop. delay, 8000 bit packet:
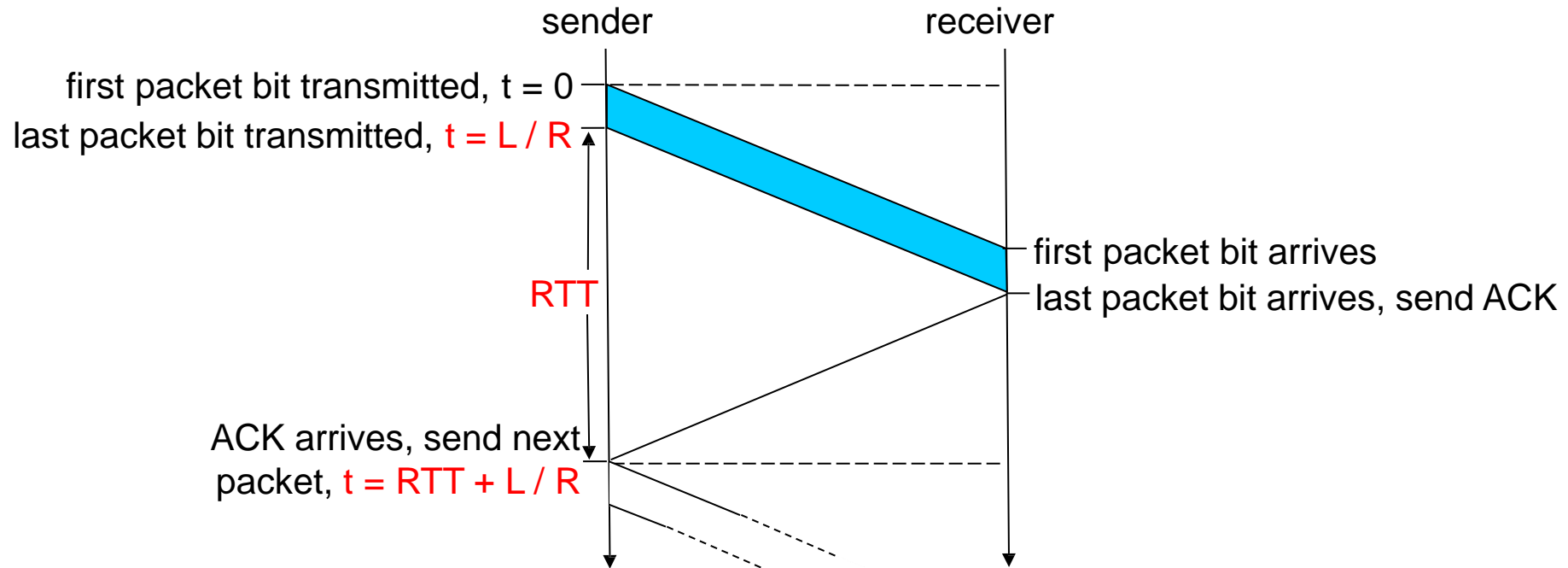
$$d_{trans} = \frac{L}{R} = \frac{8000\,\text{bits}}{10^9\,\text{bps}} = 8\,\text{microseconds}$$

- ▪ $U_{sender}$: utilization – fraction of time sender busy sending

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

- ▪ if RTT=30 msec, 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

sender                                    receiver

first packet bit transmitted, t = 0

last packet bit transmitted, t = L / R

RTT

first packet bit arrives

last packet bit arrives, send ACK

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

**Why send only one packet?
Can we send more?**

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

# Pipelined protocols

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

- range of sequence numbers must be increased
- buffering at sender and/or receiver



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

❖ **two generic forms of pipelined protocols:** go-Back-N, selective repeat

# Pipelining: increased utilization

sender                                    receiver

first packet bit transmitted, t = 0

last bit transmitted, t = L / R

RTT

first packet bit arrives

last packet bit arrives, send ACK

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

last bit of 2^nd packet arrives, send ACK

last bit of 3^rd packet arrives, send ACK

· Increase utilization by a factor of 3!

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

How many packets should be in the pipeline?

# Bandwidth-Delay Product (BDP)

❖ A logical link has a bandwidth of B bits/s and the round trip time is D sec.

❖ Bandwidth-Delay Product or BDP is computed as B*D

❖ Example:

1. Two computers are connected over a WAN and can transfer data at 2Mbps with 100ms round trip time
   - BDP = 2 * 0.1 = 0.2 Mb

2. A satellite link has a bandwidth of 10Mbps and RTT of 1s
   - BDP = 10 * 1 = 2 Mb

3. Two computers from two universities in different continents are connected by a 500Mps link and has RTT of 200ms
   - BDP = 500 * 0.2 = 100 Mb

# BDP Computation (cont'd)

How do you interpret BDP?

Why is BDP interesting?

❖ BDP estimates on the amount of data that the communicate pipe can "store"

- If the amount of unack data (window size) allowed is less than the BDP, sender will be idle at times, leading to lower efficiency

# Pipelined Protocols

- Selective Repeat: big pic
- sender can have up to N unack'ed packets in pipeline
- rcvr sends individual ack for each packet
- sender maintains timer for each unacked packet
  - when timer expires, retransmit only unack'ed packet

- Go-back-N: big picture:
- sender can have up to N unacked packets in pipeline
- rcvr only sends
- cumulative acks
  - doesn't ack packet if
- there's a gap
- sender has timer for oldest unacked packet
  - if timer expires, retransmit all unack'ed packets

# Go-Back-N

**Sender:**
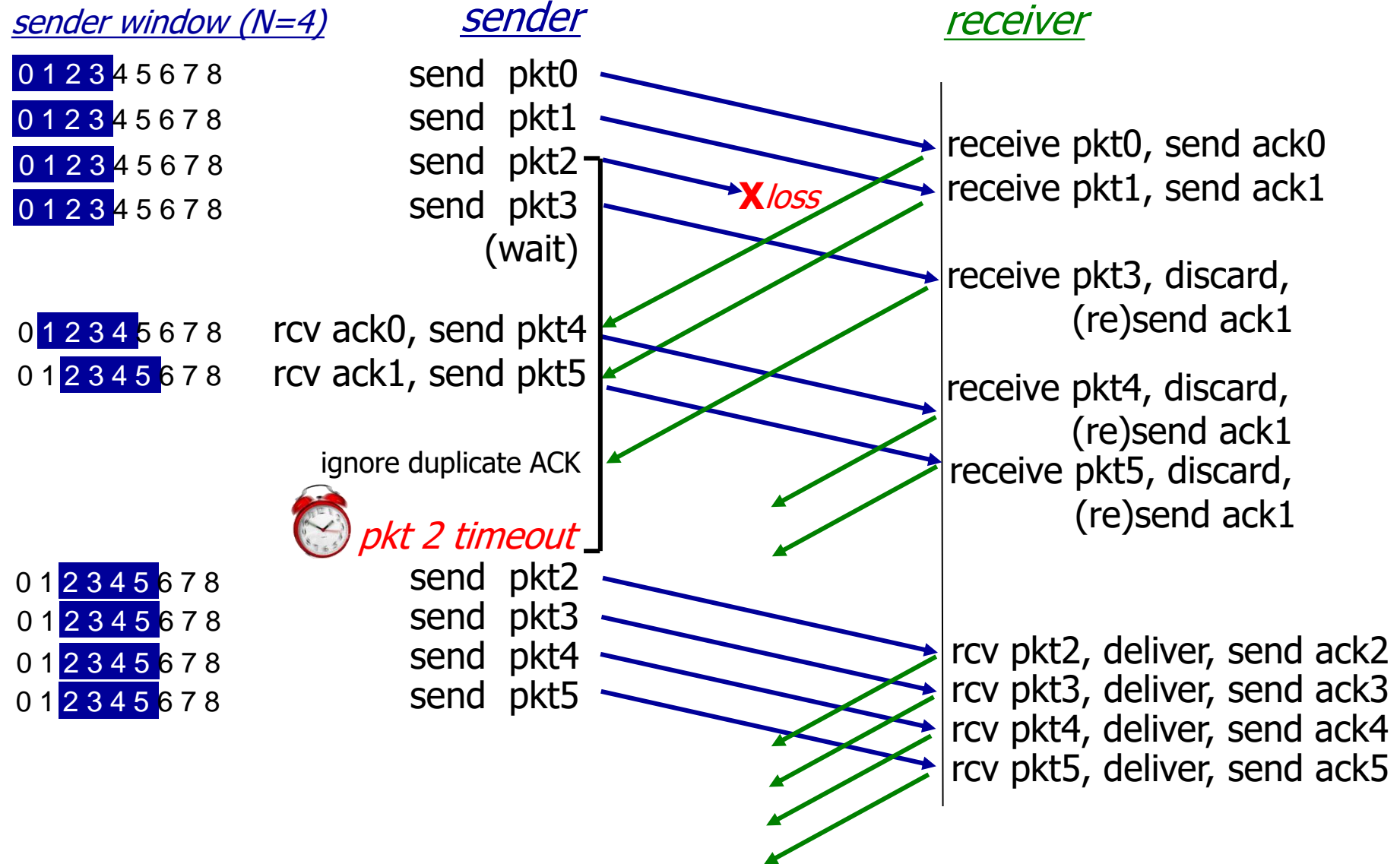
❖ k-bit seq # in pkt header

❖ "window" of up to N, consecutive unack'ed pkts allowed



❖ ACK(n): ACKs all pkts up to, including seq # n - "cumulative ACK"
  ▪ may receive duplicate ACKs (see receiver)

❖ timer for each in-flight pkt

❖ timeout(n): retransmit pkt n and all higher seq # pkts in window

# Go-Back-N



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

# Selective Repeat

- ❖ 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, unACK'ed pkts