

Chapter 3

Transport Layer

Reti di Elaboratori

Corso di Laurea in Informatica

Università degli Studi di Roma "La Sapienza"

Canale A-L

Prof.ssa Chiara Petrioli

Parte di queste slide sono state prese dal materiale associato al libro
Computer Networking: A Top Down Approach, 5th edition.

All material copyright 1996-2009

J.F Kurose and K.W. Ross, All Rights Reserved

Thanks also to Antonio Capone, Politecnico di Milano, Giuseppe Bianchi and
Francesco LoPresti, Un. di Roma Tor Vergata

Chapter 3: Transport Layer

Our goals:

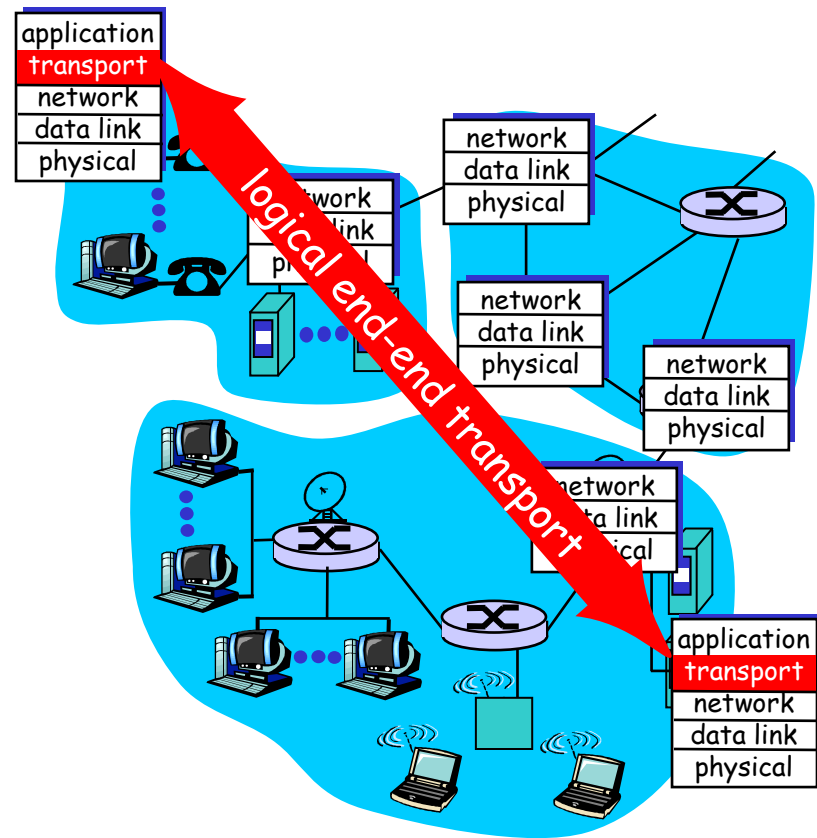
- r understand principles behind transport layer services:
 - m multiplexing/demultiplexing
 - m reliable data transfer
 - m flow control
 - m congestion control
- r learn about transport layer protocols in the Internet:
 - m UDP: connectionless transport
 - m TCP: connection-oriented transport
 - m TCP congestion control

Chapter 3 outline

- r 3.1 Transport-layer services
- r 3.2 Multiplexing and demultiplexing
- r 3.3 Connectionless transport: UDP
- r 3.4 Principles of reliable data transfer
- r 3.5 Connection-oriented transport: TCP
 - m segment structure
 - m reliable data transfer
 - m flow control
 - m connection management
- r 3.6 Principles of congestion control
- r 3.7 TCP congestion control

Transport services and protocols

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



Transport vs. network layer

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

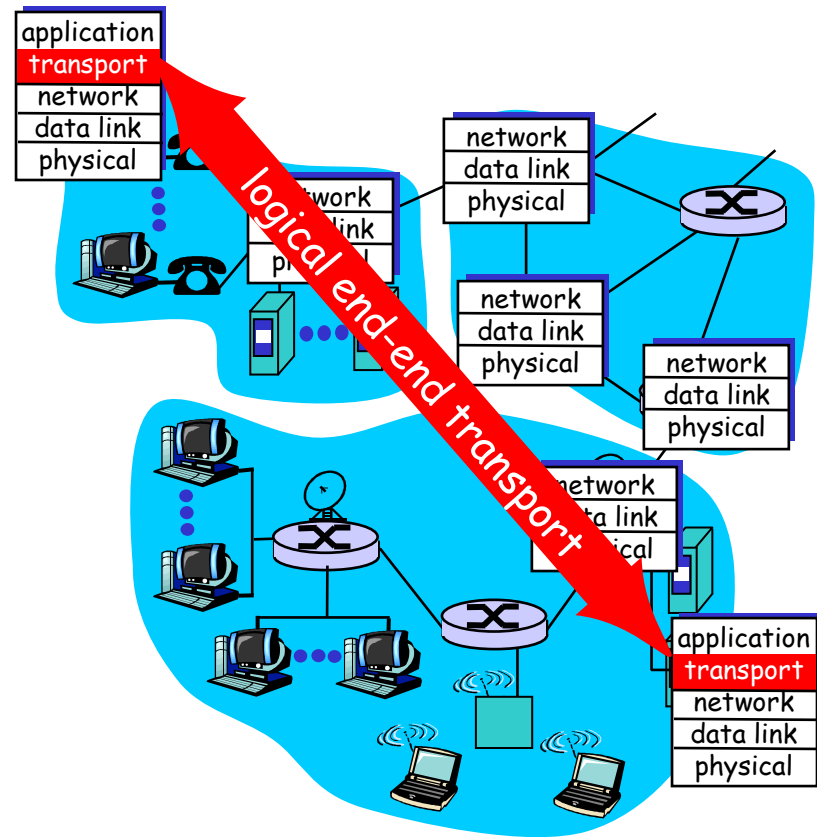
Household analogy:

12 kids sending letters to 12 kids

- r processes = kids
- r app messages = letters in envelopes
- r hosts = houses
- r transport protocol = Ann and Bill
- r network-layer protocol = postal service

Internet transport-layer protocols

- r reliable, in-order delivery (TCP)
 - m congestion control
 - m flow control
 - m connection setup
- r unreliable, unordered delivery: UDP
 - m no-frills extension of "best-effort" IP
- r services not available:
 - m delay guarantees
 - m bandwidth guarantees

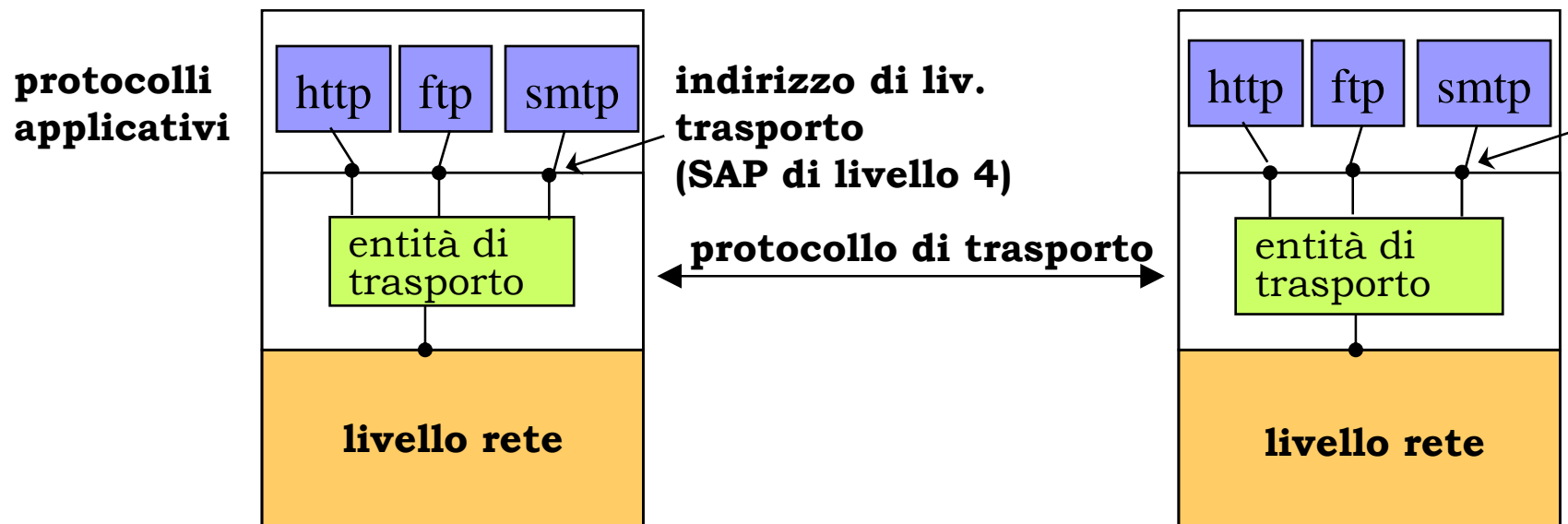


Chapter 3 outline

- r 3.1 Transport-layer services
- r 3.2 Multiplexing and demultiplexing
- r 3.3 Connectionless transport: UDP
- r 3.4 Principles of reliable data transfer
- r 3.5 Connection-oriented transport: TCP
 - m segment structure
 - m reliable data transfer
 - m flow control
 - m connection management
- r 3.6 Principles of congestion control
- r 3.7 TCP congestion control

Servizio di trasporto

- r Più applicazioni possono essere attive su un end system
 - m il livello di trasporto svolge funzioni di multiplexing/demultiplexing
 - m ciascun collegamento logico tra applicazioni è indirizzato dal livello di trasporto



Multiplexing/demultiplexing

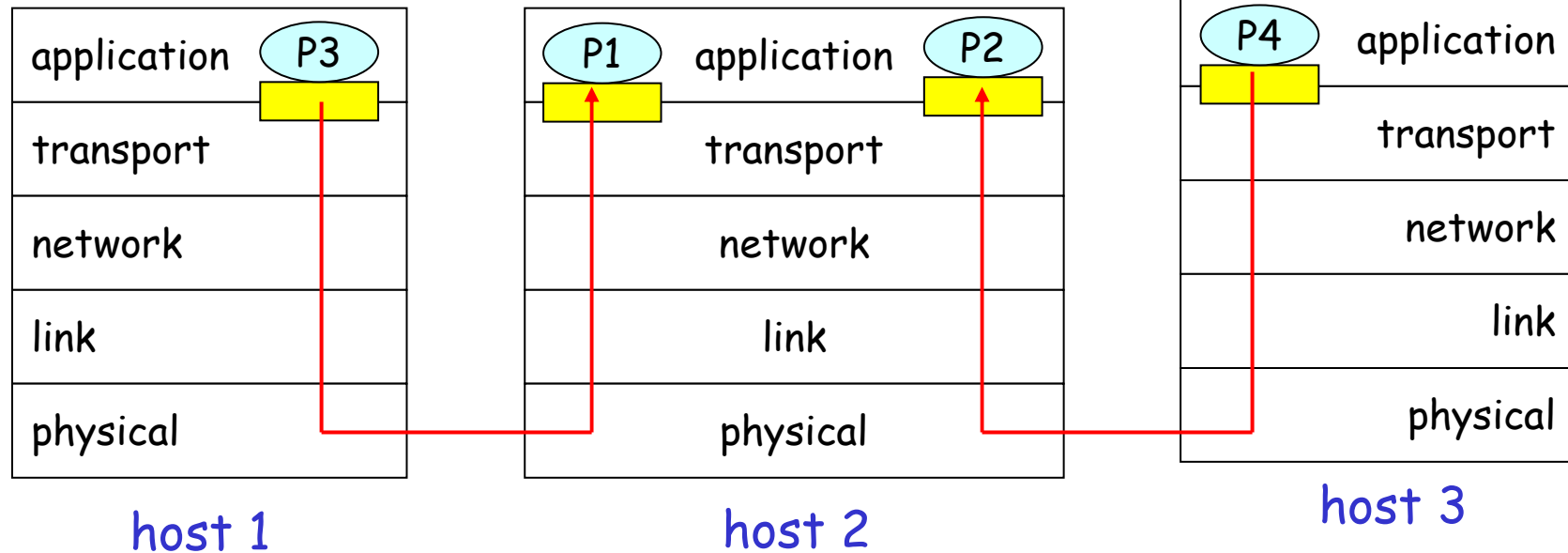
Demultiplexing at rcv host:

delivering received segments to correct socket

Multiplexing at send host:

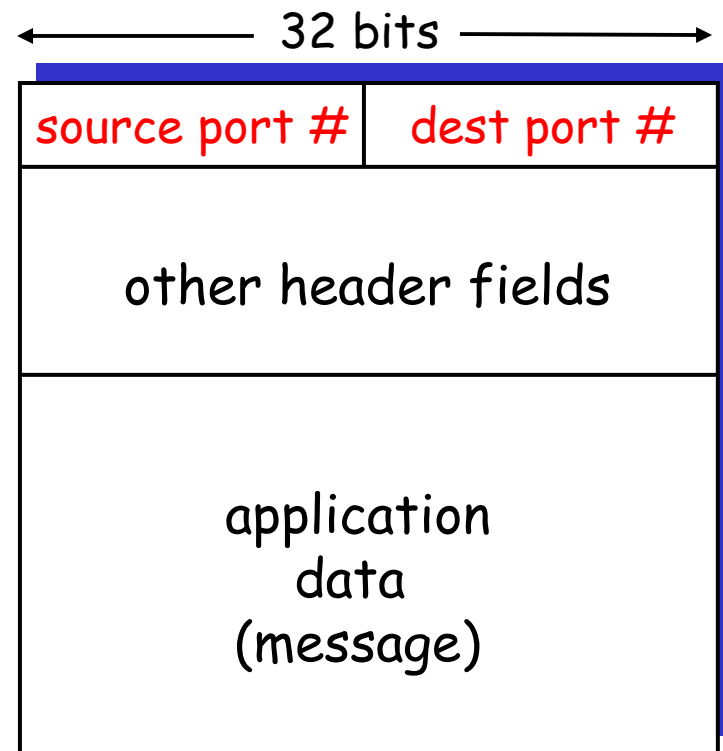
gathering data from multiple sockets, enveloping data with header (later used for demultiplexing)

■ = socket ○ = process



How demultiplexing works

- r **host receives IP datagrams**
 - m each datagram has source IP address, destination IP address
 - m each datagram carries 1 transport-layer segment
 - m each segment has source, destination port number (recall: well-known port numbers for specific applications)
- r **host uses IP addresses & port numbers to direct segment to appropriate socket**



TCP/UDP segment format

Connectionless demultiplexing

- r Create sockets with port numbers:

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

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

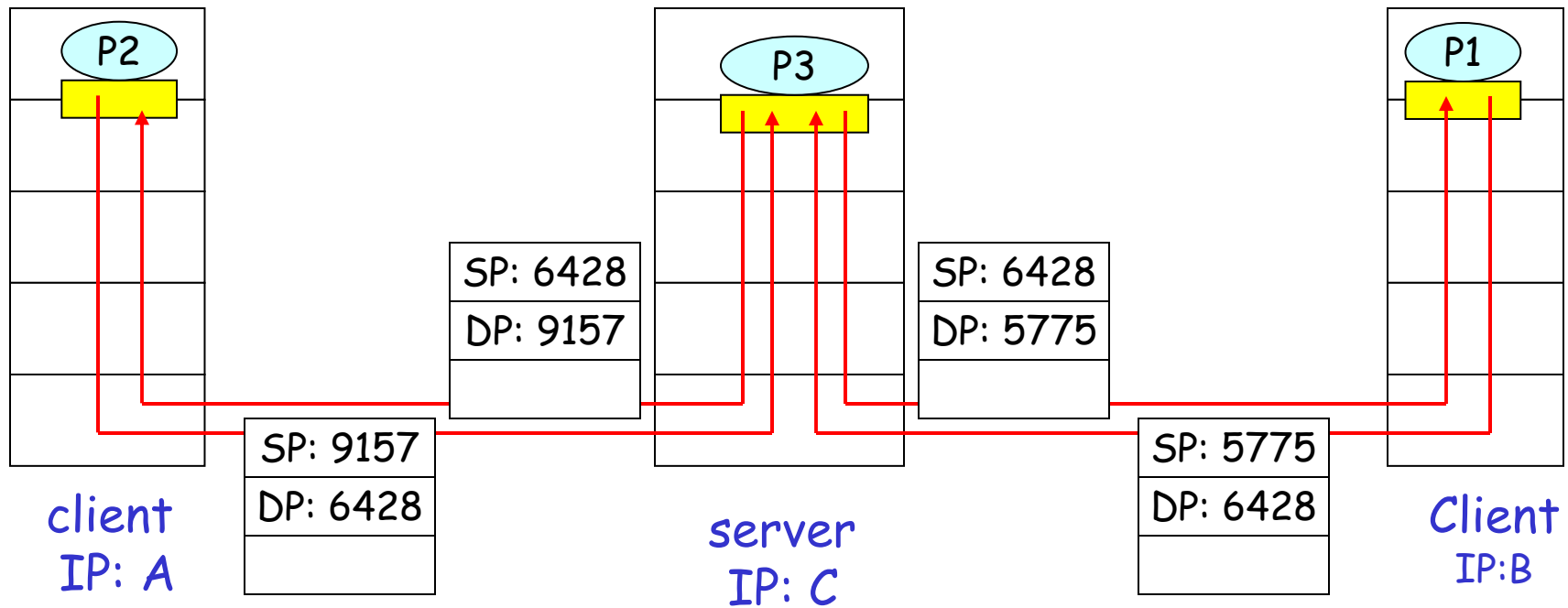
- r UDP socket identified by two-tuple:

(dest IP address, dest port number)

- r When host receives UDP segment:
 - m checks destination port number in segment
 - m directs UDP segment to socket with that port number
- r IP datagrams with different source IP addresses and/or source port numbers directed to same socket

Connectionless demux (cont)

```
DatagramSocket serverSocket = new DatagramSocket(6428);
```

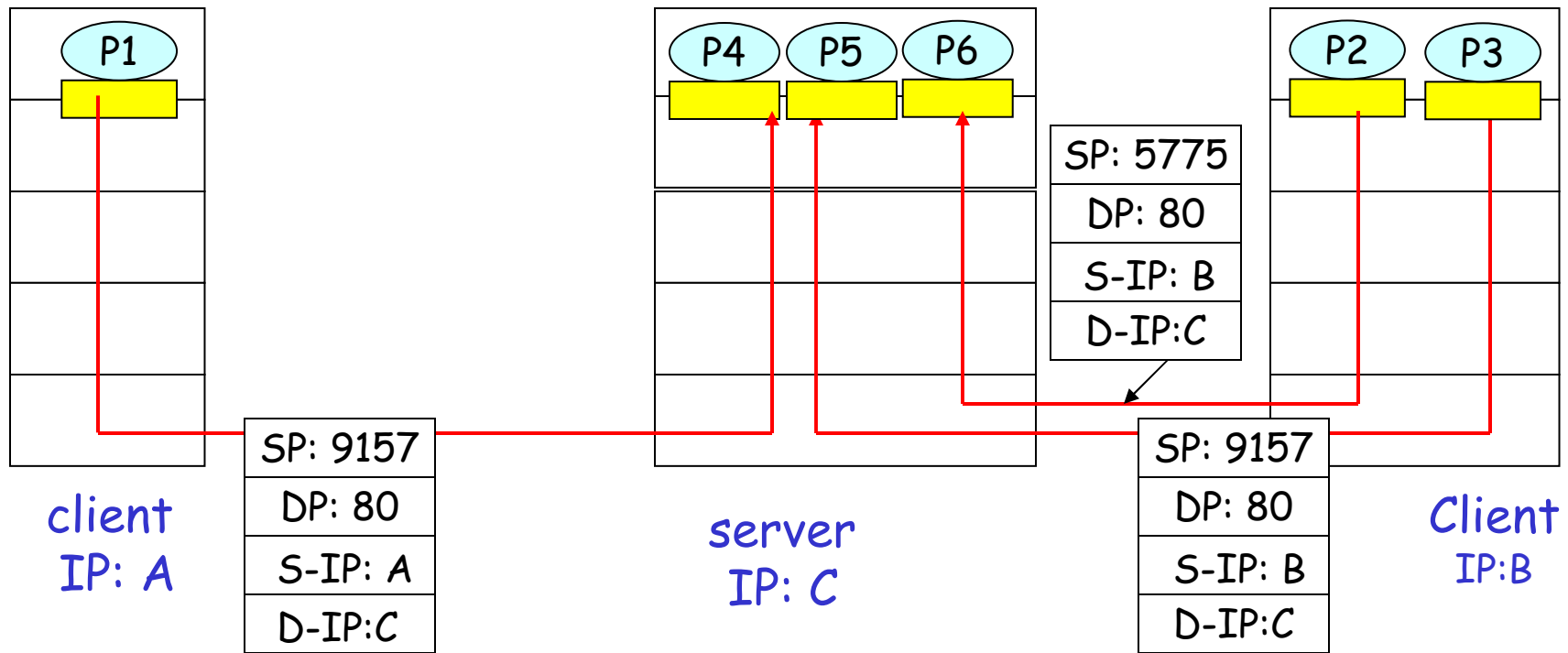


SP provides "return address"

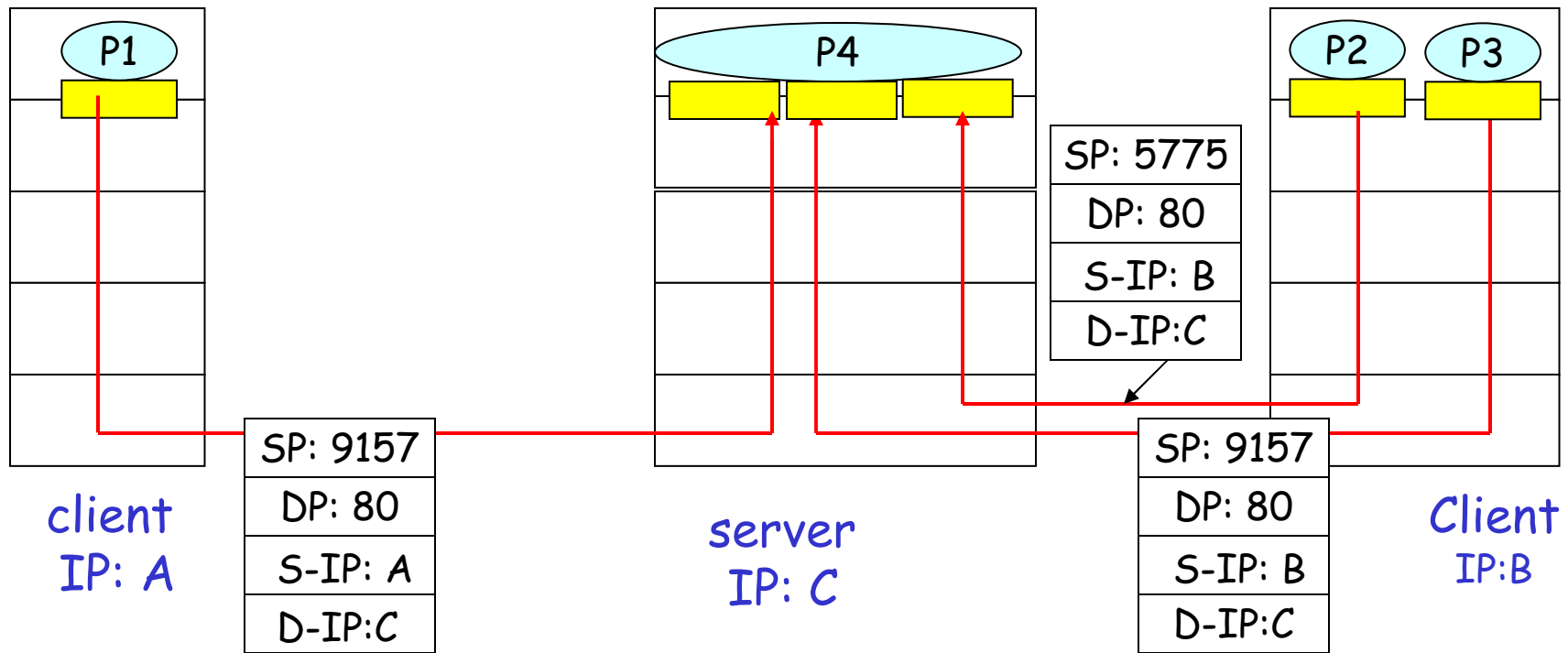
Connection-oriented demux

- r TCP socket identified by 4-tuple:
 - m source IP address
 - m source port number
 - m dest IP address
 - m dest port number
- r recv host uses all four values to direct segment to appropriate socket
- r Server host may support many simultaneous TCP sockets:
 - m each socket identified by its own 4-tuple
- r Web servers have different sockets for each connecting client
 - m non-persistent HTTP will have different socket for each request

Connection-oriented demux (cont)



Connection-oriented demux: Threaded Web Server



Chapter 3 outline

- r 3.1 Transport-layer services
- r 3.2 Multiplexing and demultiplexing
- r 3.3 Connectionless transport: UDP
- r 3.4 Principles of reliable data transfer
- r 3.5 Connection-oriented transport: TCP
 - m segment structure
 - m reliable data transfer
 - m flow control
 - m connection management
- r 3.6 Principles of congestion control
- r 3.7 TCP congestion control

UDP: User Datagram Protocol [RFC 768]

- r "no frills," "bare bones" Internet transport protocol
 - r "best effort" service, UDP segments may be:
 - m lost
 - m delivered out of order to app
- ↓
- reliable transfer over UDP: add reliability at application layer
- m application-specific error recovery!
- r **connectionless:**
 - m no handshaking between UDP sender, receiver
 - m each UDP segment handled independently of others

Why is there a UDP?

- r no connection establishment (which can add delay)
- r **simple: no connection state at sender, receiver**
- r **small segment header (8 byte)**
- r no congestion control: UDP can blast away as fast as desired

often used for streaming multimedia apps

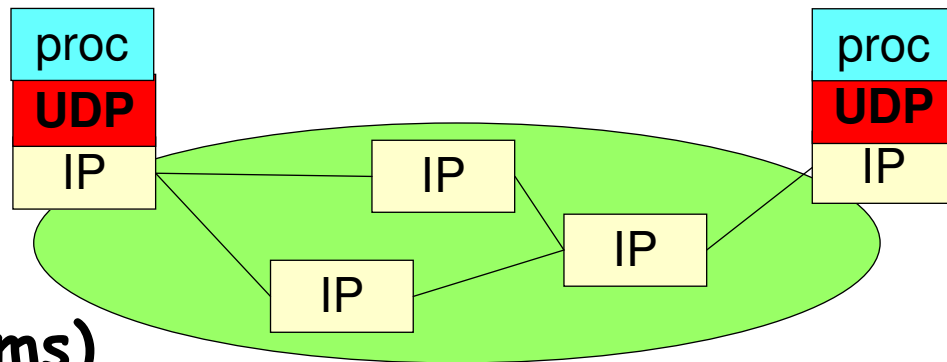
- m loss tolerant
- m rate sensitive

other UDP uses: DNS, SNMP..

UDP Packets

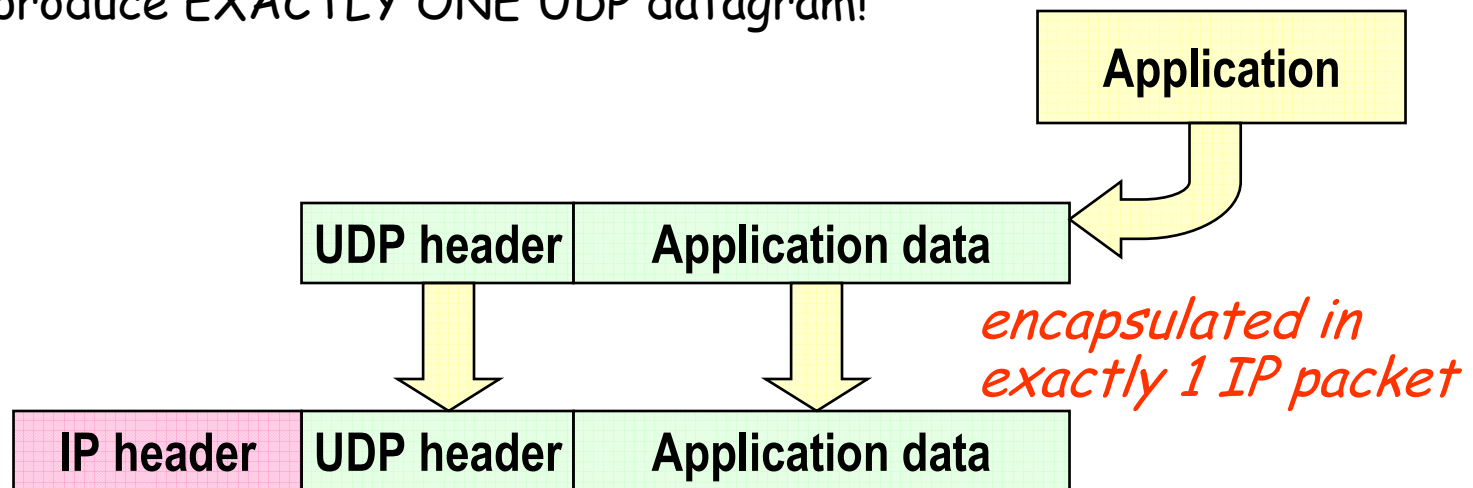
r **Connection-Less**

m (no handshaking)



r **UDP packets (Datagrams)**

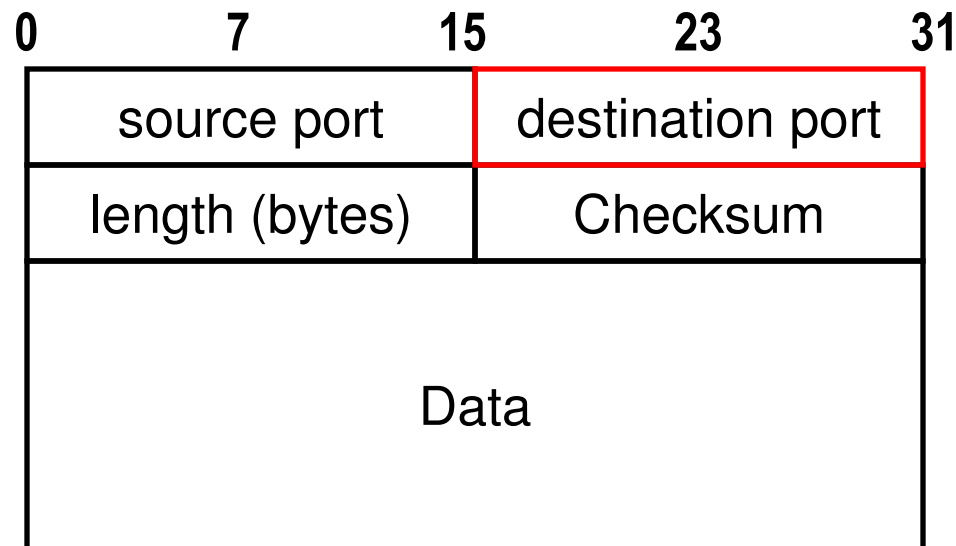
m Each application interacts with UDP transport sw to produce EXACTLY ONE UDP datagram!



This is why, improperly, we use the term UDP packets

UDP datagram format

8 bytes header + variable payload



r **UDP length field**

- m all UDP datagram
- m (header + payload)

r **payload sizes allowed:**

- m Empty
- m even size (bytes)

→ **UDP functions limited to:**

⇒ **addressing**

→ which is the only strictly necessary role of a transport protocol

⇒ **Error checking**

→ which may even be **disabled** for performance

Maximum UDP datagram size

- r 16 bit UDP length field:
 - m Maximum up to $2^{16}-1 = 65535$ bytes
 - m Includes 8 bytes UDP header (max data = 65527)

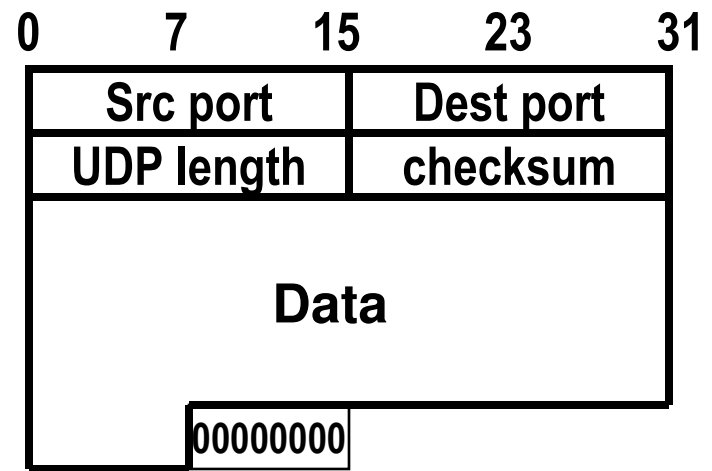
- r But max IP packet size is also 65535
 - m Minus 20 bytes IP header, minus 8 bytes UDP header
 - m Max UDP_data = 65507 bytes!

- r Moreover, most OS impose further limitations!
 - m most systems provide 8192 bytes maximum (max size in NFS)
 - m some OS had (still have?) internal implementation features (bugs?) that limit IP packet size
 - SunOS 4.1.3 had 32767 for max tolerable IP packet transmittable (but 32786 in reception...) - bug fixed only in Solaris 2.2

- r Finally, subnet Maximum Transfer Unit (MTU) limits may fragment datagram - annoying for reliability!
 - m E.g. ethernet = 1500 bytes; PPP on your modem = 576

Error checksum

- r 16 bit checksum field, obtained by:
 - m summing up all 16 bit words in header data and **pseudoheader**, in 1's complement (checksum fields filled with 0s initially)
 - m take 1's complement of result
 - m if result is 0, set it to 11111...11 (65535==0 in 1's complement)
 - m Sender puts checksum value into UDP checksum field

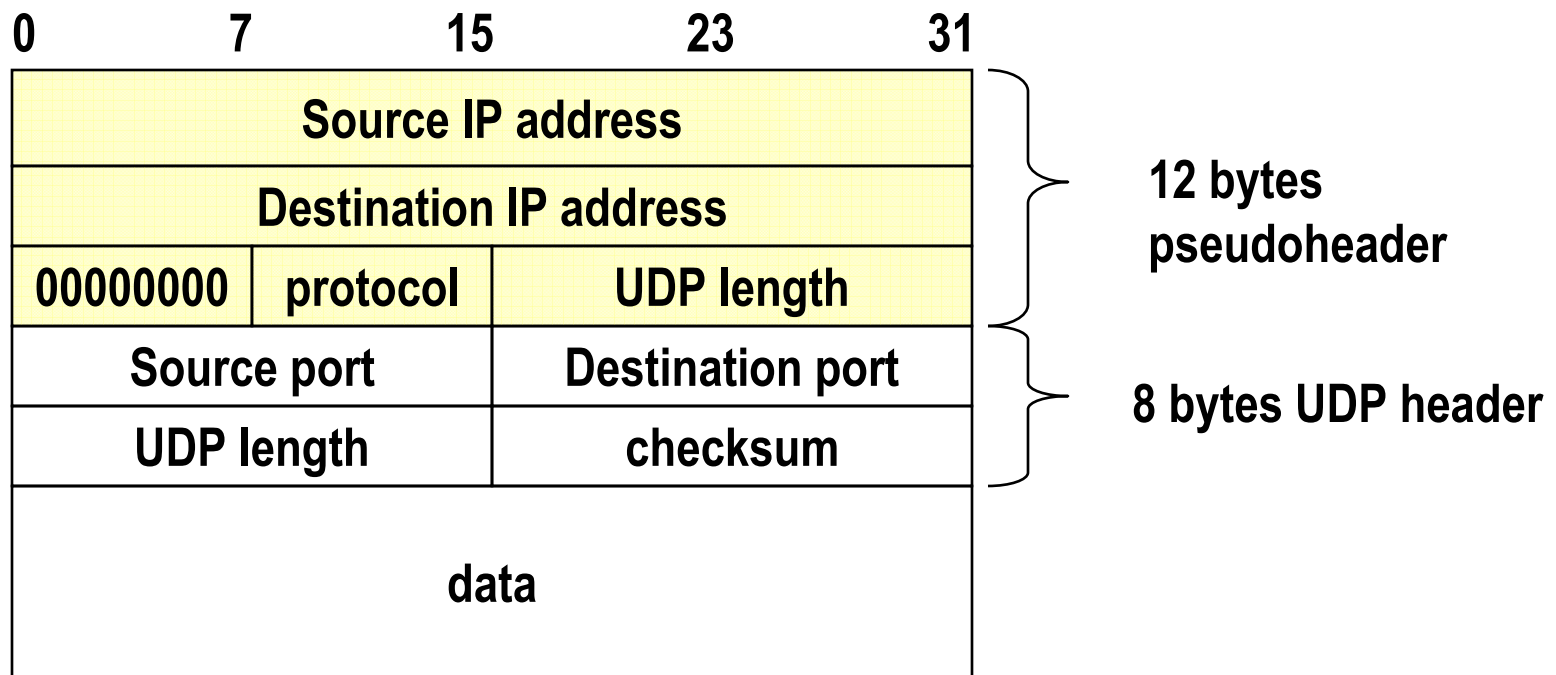


- r at destination:
 - m 1's complement sum should return 0, otherwise error detected
 - m upon error, no action (just packet discard)
- r efficient implementation RFC 1071

- r Zero padding
 - m To multiple of 16 bits
- r checksum disabled
 - m by source, by setting 0 in the checksum field

Pseudo header

- r Is not transmitted!
 - m But it is information available at transmitter and at receiver
 - m intention: double check that packet has arrived at correct destination



Protocol field (TCP=6,UDP=17) necessary, as same checksum calculation used in TCP. UDP length duplicated.

disabling checksum

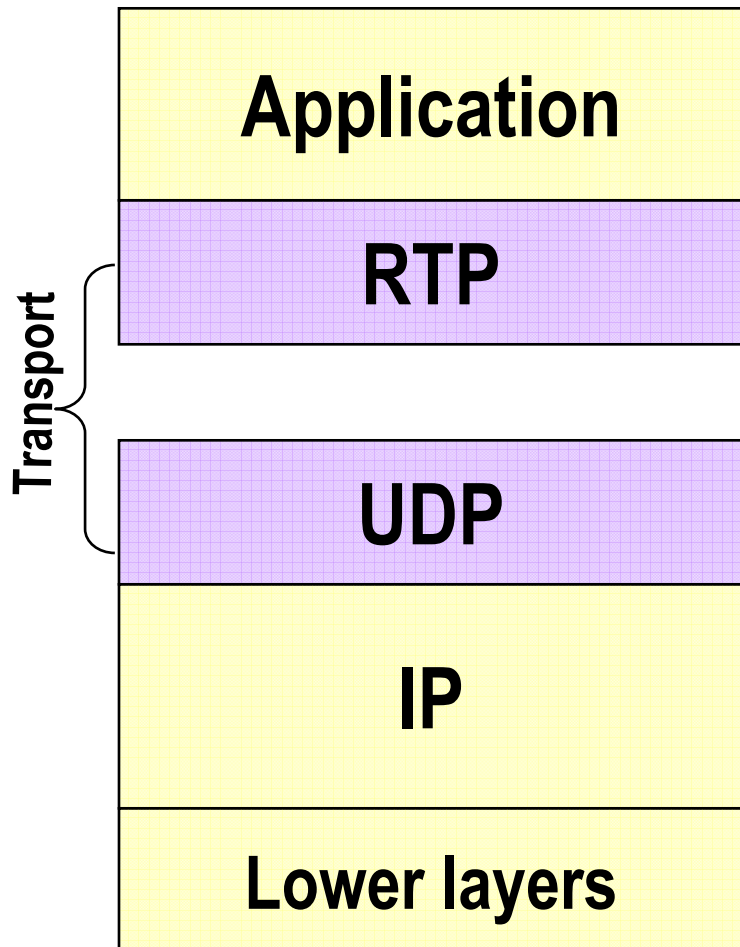
- r In principle never!
 - m Remember that IP packet checksum DOES NOT include packet payload.
- r In practice, often done in NFS
 - m sun was the first, to speed up implementation
- r may be tolerable in LANs under one's control.
- r Definitely dangerous in the wide internet
 - m Exist layer 2 protocols without error checking

UDP: a lightweight protocol

- r No connection establishment
 - m no initial overhead due to handshaking
- r No connection state
 - m greater number of supported connections by a server!
- r Small packet header overhead
 - m 8 bytes only vs 20 in TCP
- r originally intended for simple applications, oriented to short information exchange
 - m DNS
 - m management (e.g. SNMP)
 - m etc
- r No rate limitations
 - m No throttling due to congestion & flow control mechanisms
 - m No retransmission (for certain application loss tolerable)
- r extremely important features for today multimedia applications!
Especially for real time applications which can tolerate some packet loss but require a minimum send rate.

RTP as seen from Application

Be careful: UDP ok for multimedia because it does not provide anything at all (no features = no limits!). Application developers have to provide supplementary transport capabilities at the application layer!



*Solution for audio/video:
Real Time Protocol
(RTP, RFC 1889)*



**SOCKET
INTERFACE**

Application developer integrates RTP into the application by:

- writing code which creates the RTP encapsulating packets;
- sends the RTP packets into a UDP socket interface.

Details of RTP in subsequent courses – unless we are ahead of schedule.

Chapter 3 outline

- r 3.1 Transport-layer services
- r 3.2 Multiplexing and demultiplexing
- r 3.3 Connectionless transport: UDP
- r 3.4 Principles of reliable data transfer
- r 3.5 Connection-oriented transport: TCP
 - m segment structure
 - m reliable data transfer
 - m flow control
 - m connection management
- r 3.6 Principles of congestion control
- r 3.7 TCP congestion control

A MUCH more complex transport for three main reasons

r Connection oriented

m implements mechanisms to setup and tear down a full duplex connection between end points

r Reliable

m implements mechanisms to guarantee error free and ordered delivery of information

r Flow & Congestion controlled

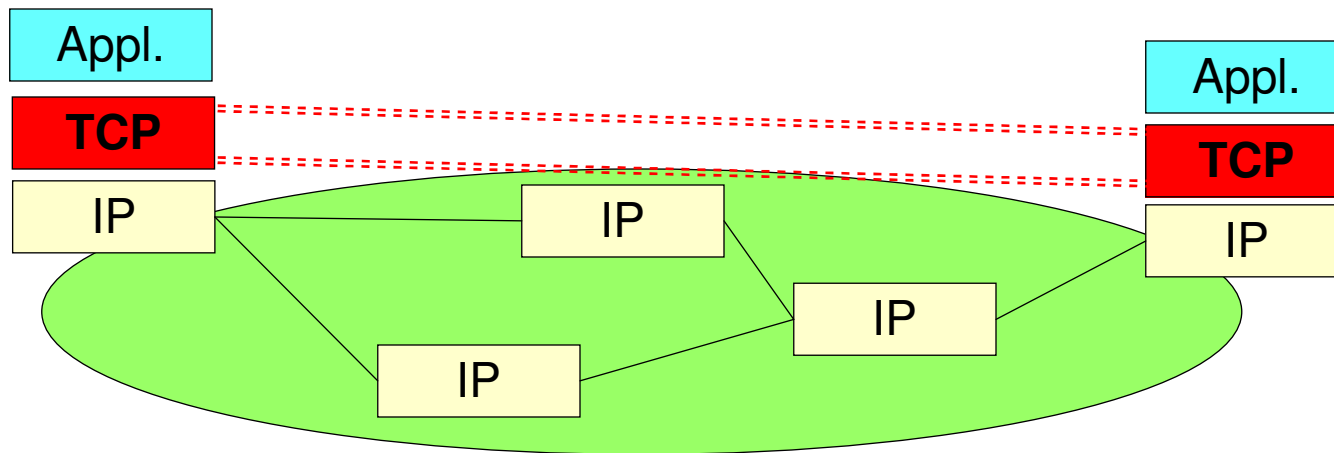
m implements mechanisms to control traffic

TCP services

- r connection oriented
 - m TCP connections
- r *reliable transfer service*
 - m all bytes sent are received

→ TCP functions

- application addressing (ports)
- error recovery (acks and retransmission)
- reordering (sequence numbers)
- flow control
- congestion control



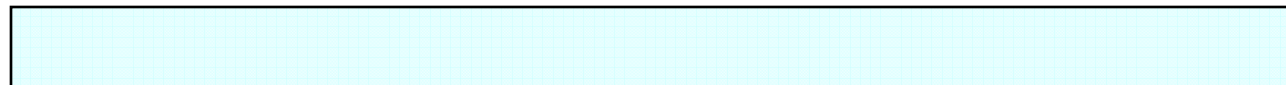
Byte stream service

- r TCP exchange data between applications as a stream of bytes.
- r It does not introduce any data delimiter (an application duty)
 - m source application may enter 10 bytes followed by 1 and 40 (grouped with some semantics)
 - m data is buffered at source, and transmitted
 - m at receiver, may be read in the sequence 25 bytes, 22 bytes and 4 bytes...

Application view

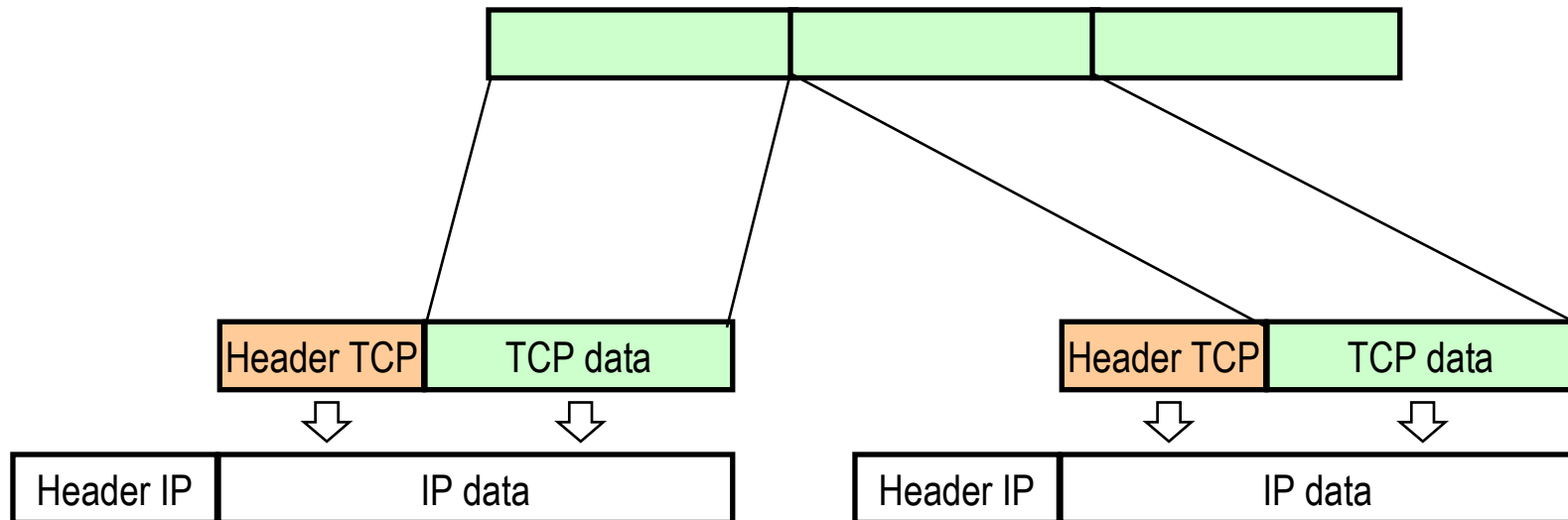


TCP view



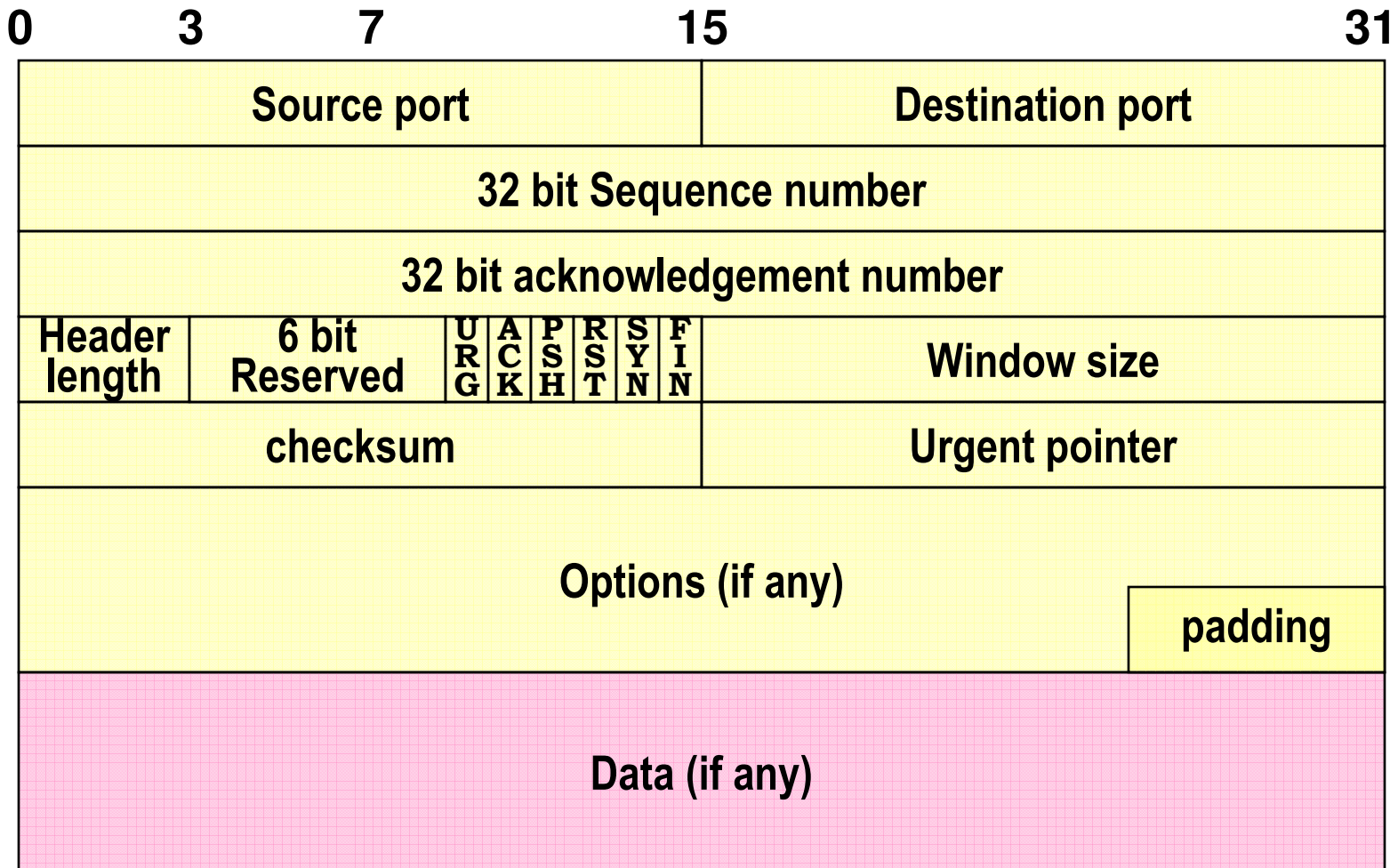
TCP segments

- r Application data broken into segments for transmission
- r segmentation totally up to TCP, according to what TCP considers being the best strategy
- r each segment placed into an IP packet
- r very different from UDP!!



TCP segment format

20 bytes header (minimum)



Source port			Destination port					
32 bit Sequence number								
32 bit acknowledgement number								
Header length	6 bit Reserved	U R G	A C K	P S H	R S T	S Y N	F I N	Window size
checksum				Urgent pointer				

- r Source & destination port + source and destination IP addresses
 - m univocally determine TCP connection

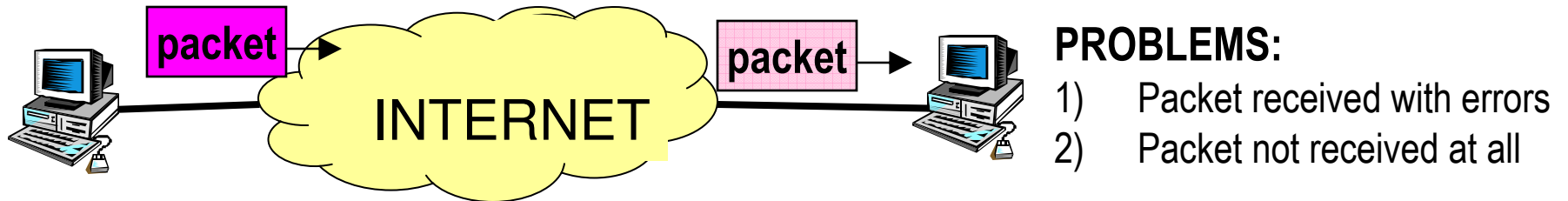
- r checksum as in UDP
 - m same calculation including same pseudoheader

- r no explicit segment length specification

Source port				Destination port				
32 bit Sequence number								
32 bit acknowledgement number								
Header length	6 bit Reserved	URG	ACK	PSH	RST	SYN	FIN	Window size
checksum				Urgent pointer				
Options (if any)							00000000	

- r Header length: 4 bits
 - m specifies the header size ($n \times 4$ byte words) for options
 - m maximum header size: 60 (15×4)
 - m option field size must be multiple of 32bits: zero padding when not.
- r Reserved: 000000 (still today!)

Reliable data transfer: issues



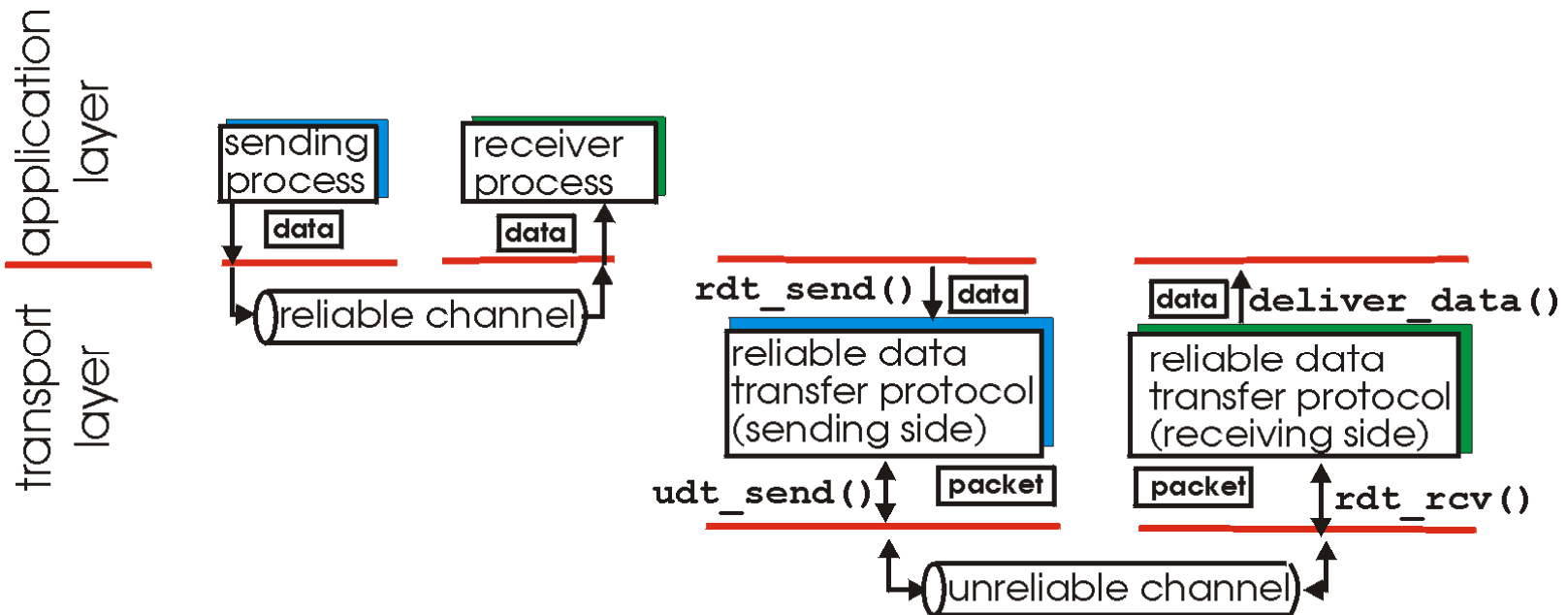
Same problem considered at DATA LINK LAYER

(although it is less likely that a whole packet is lost at data link)

- r mechanisms to guarantee correct reception:
 - m Forward Error Correction (FEC) coding schemes
 - Powerful to correct bits affected by error, not effective in case of packet loss
 - Mostly used at link layer
 - m Error detection (e.g. checksum used in UDP)
 - m Retransmission - issues:
 - ACK
 - NACK
 - TIMEOUT

Principles of Reliable data transfer

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



(a) provided service

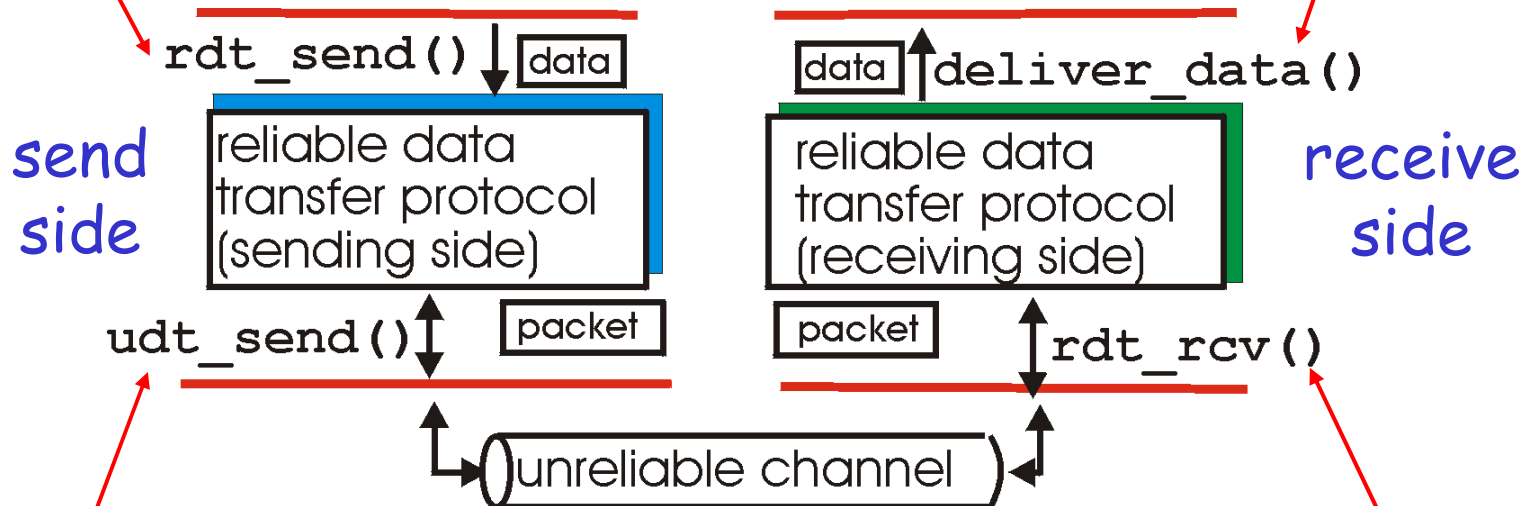
(b) service implementation

- r 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



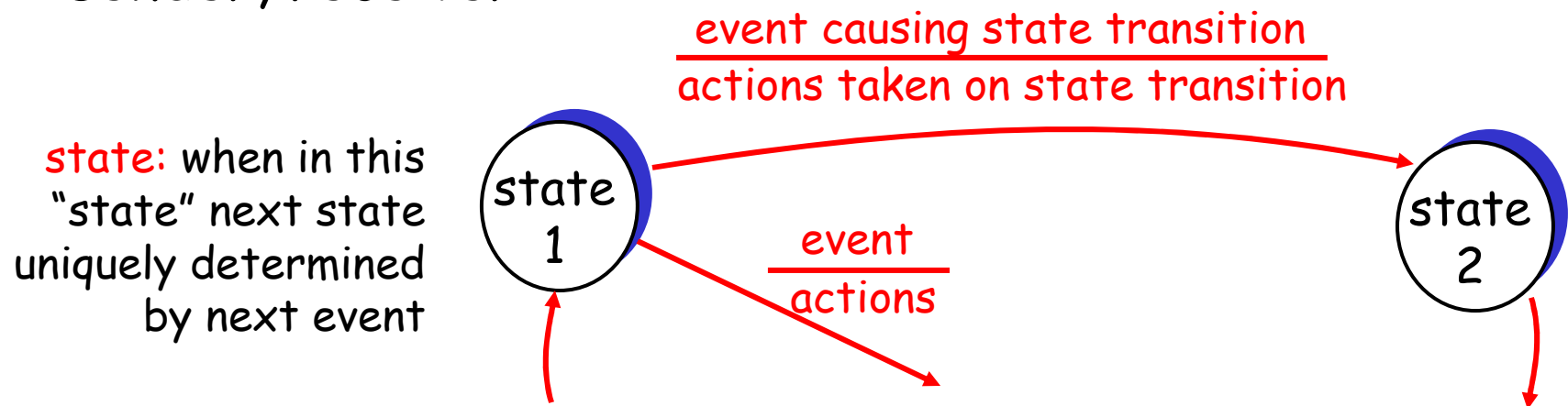
udt_send() : called by rdt, to transfer packet over unreliable channel to receiver

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

Reliable data transfer: getting started

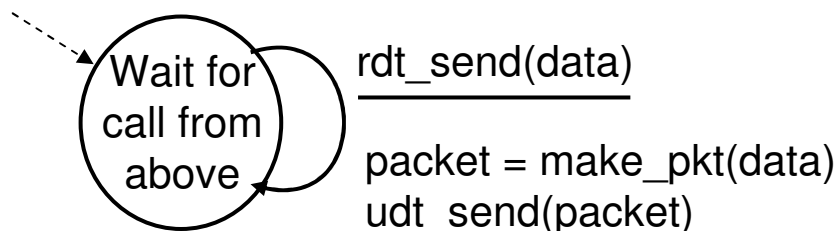
We'll:

- r incrementally develop sender, receiver sides of reliable data transfer protocol (rdt)
- r consider only unidirectional data transfer
 - m but control info will flow on both directions!
- r use finite state machines (FSM) to specify sender, receiver

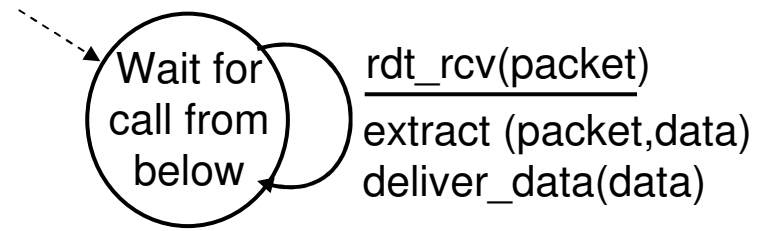


Rdt1.0: reliable transfer over a reliable channel

- r underlying channel perfectly reliable
 - m no bit errors
 - m no loss of packets (→no congestion, no buffer overflows)
- r separate FSMs for sender, receiver:
 - m sender sends data into underlying channel
 - m receiver read data from underlying channel



sender

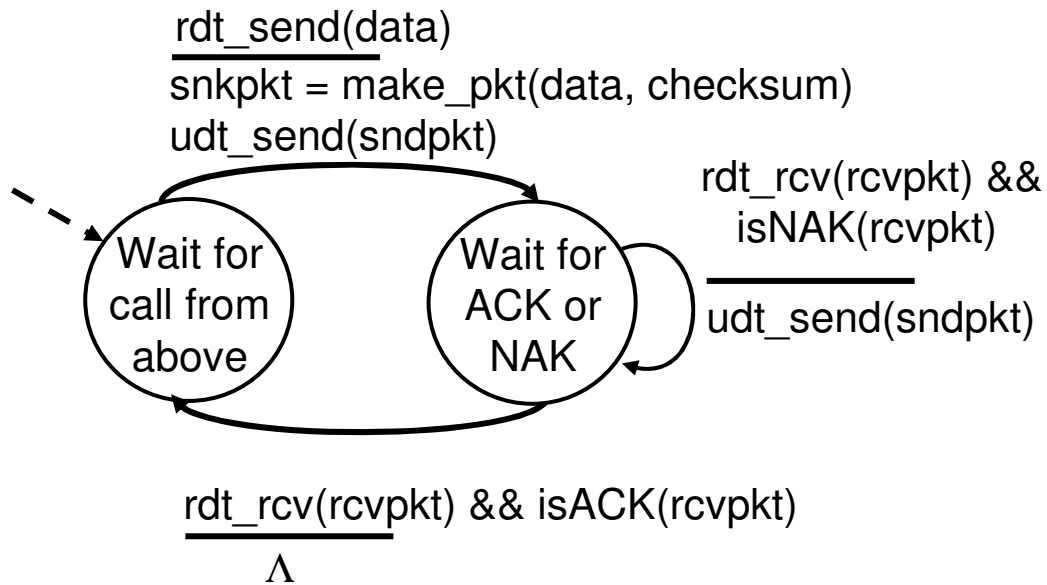


receiver

Rdt2.0: channel with bit errors

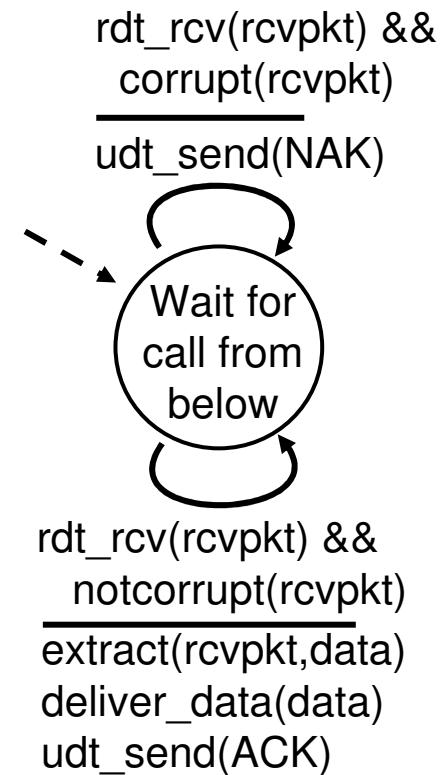
- r underlying channel may flip bits in packet
 - m recall: UDP checksum to detect bit errors
- r **Still no loss!!**
- r *the question: how to recover from errors:*
 - m *acknowledgements (ACKs):* receiver explicitly tells sender that pkt received OK
 - m *negative acknowledgements (NAKs):* receiver explicitly tells sender that pkt had errors
 - m sender retransmits pkt on receipt of NAK
 - m human scenarios using ACKs, NAKs?
- r new mechanisms in rdt2.0 (beyond rdt1.0):
 - m error detection
 - m receiver feedback: control msgs (ACK,NAK) rcvr->sender

rdt2.0: FSM specification

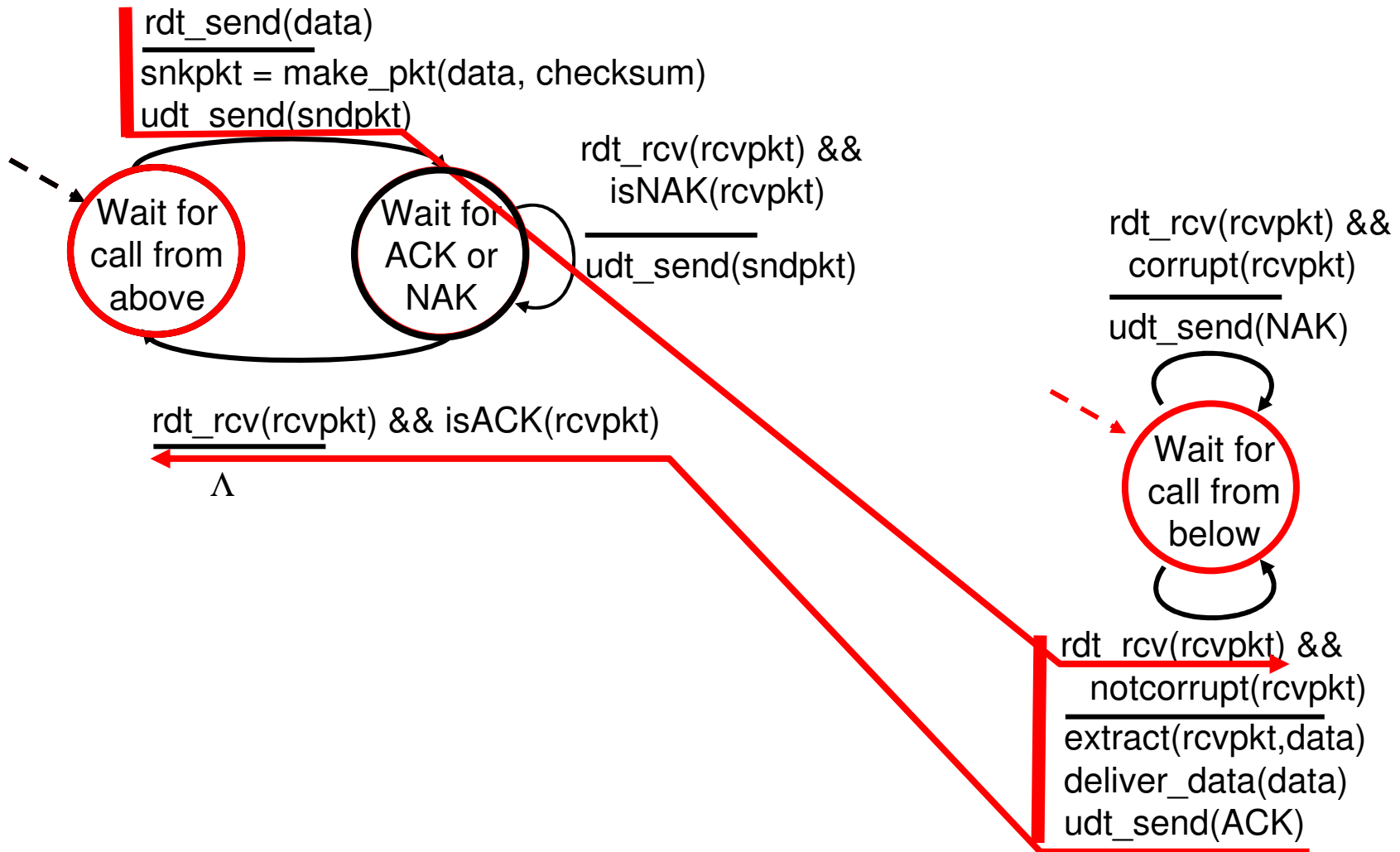


sender

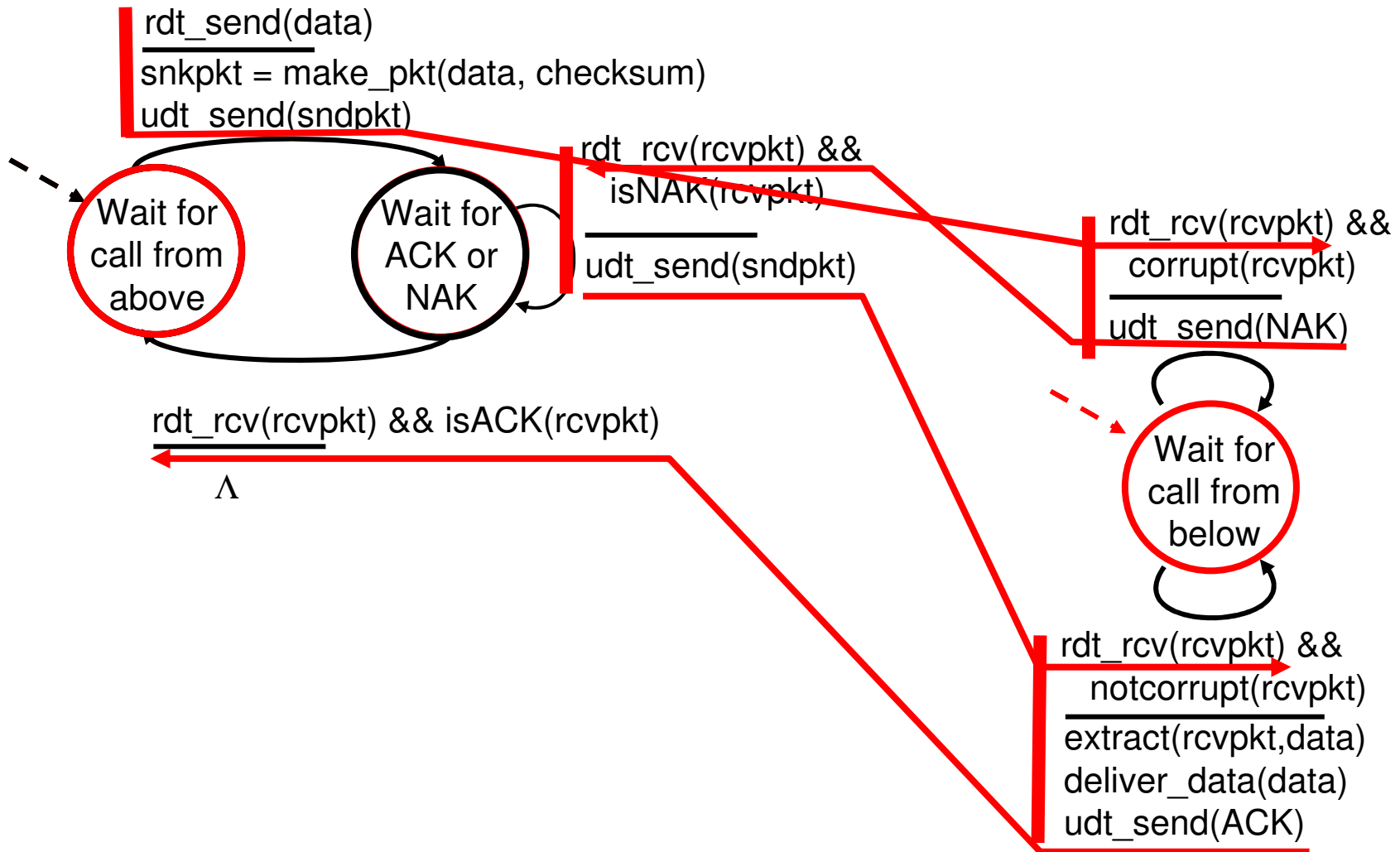
receiver



rdt2.0: operation with no errors



rdt2.0: error scenario



rdt2.0 has a fatal flaw!

What happens if ACK/NAK corrupted?

- r sender doesn't know what happened at receiver!
- r can't just retransmit: possible duplicate

What to do?

- r sender ACKs/NAKs receiver's ACK/NAK? What if sender ACK/NAK lost?
- r retransmit, but this might cause retransmission of correctly received pkt!

Handling duplicates:

- r sender adds *sequence number* to each pkt
- r sender retransmits current pkt if ACK/NAK garbled
- r receiver discards (doesn't deliver up) duplicate pkt

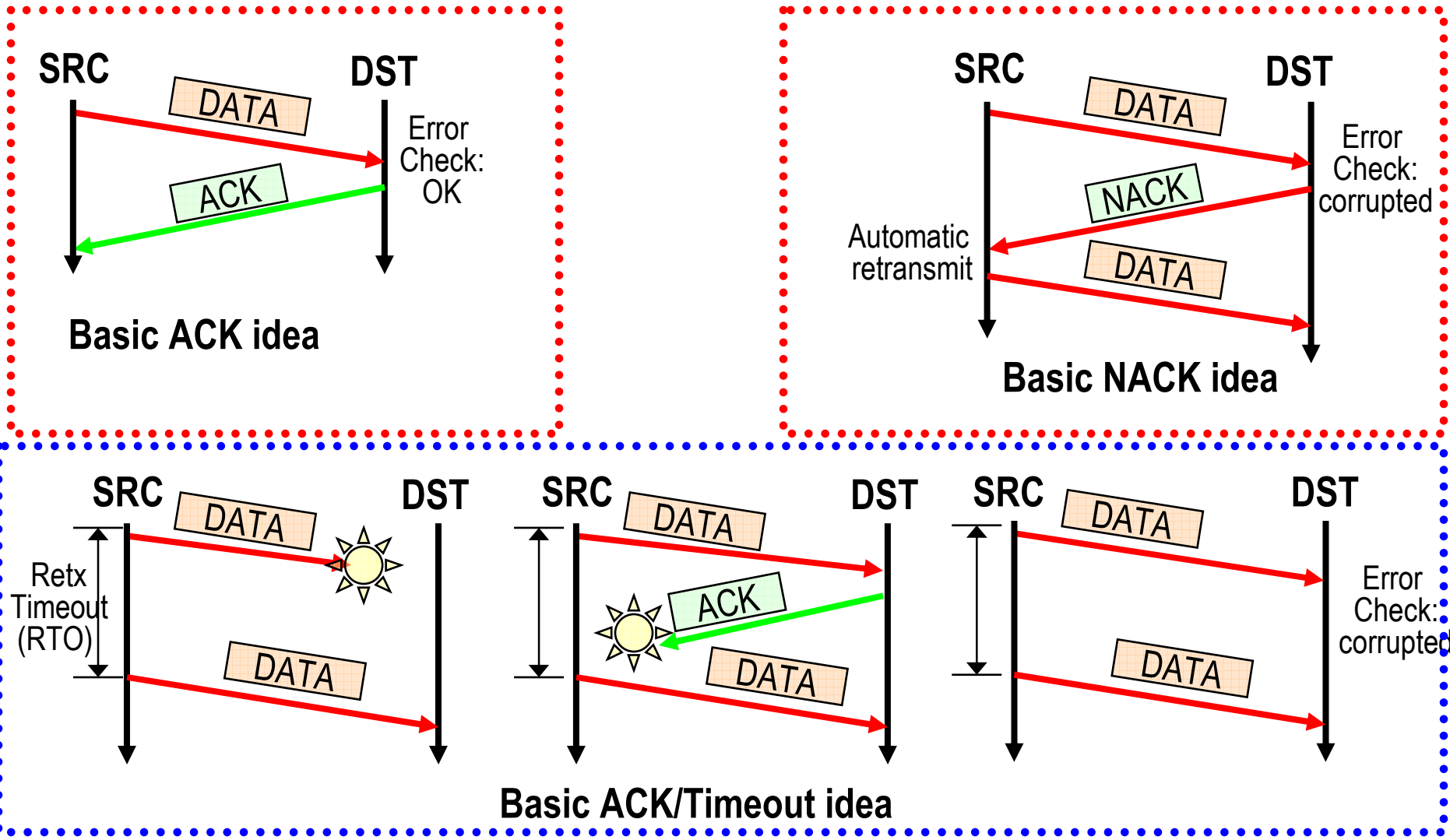
stop and wait

Sender sends one packet, then waits for receiver response

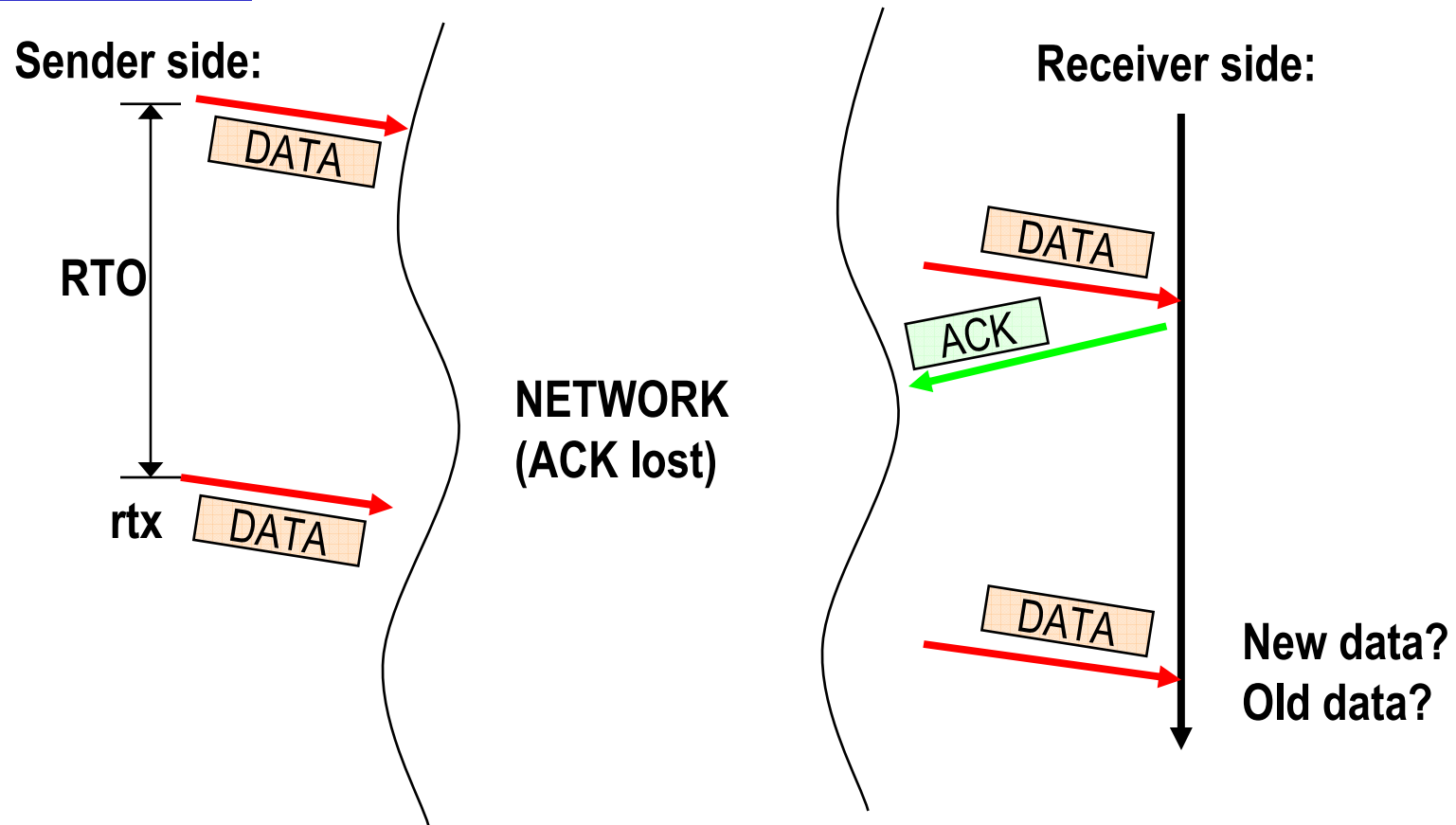
Retransmission scenarios

referred to as ARQ schemes (Automatic Retransmission reQuest)

COMPONENTS: a) error checking at receiver; b) feedback to sender; c) retx

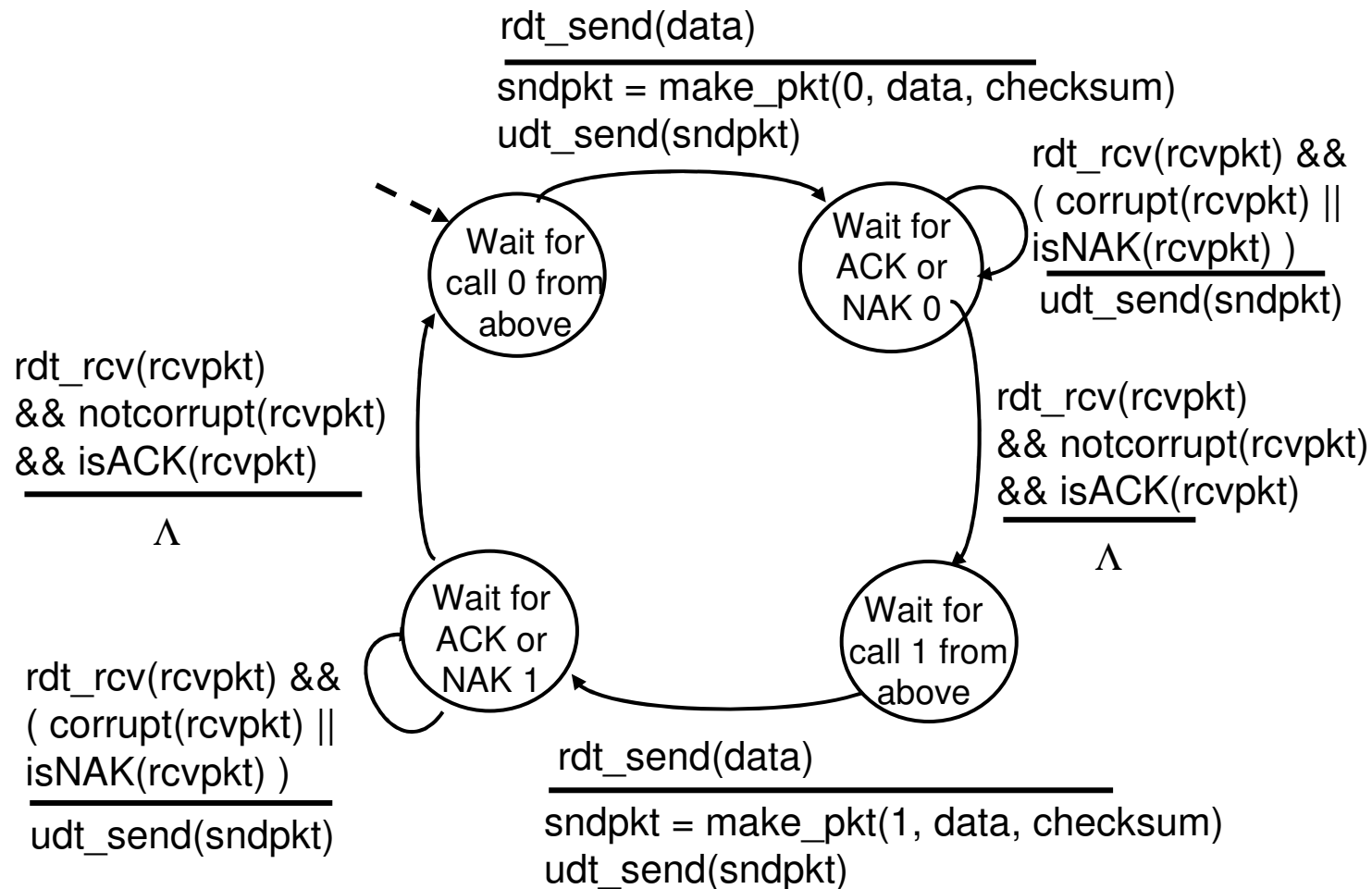


Why sequence numbers? (on data)

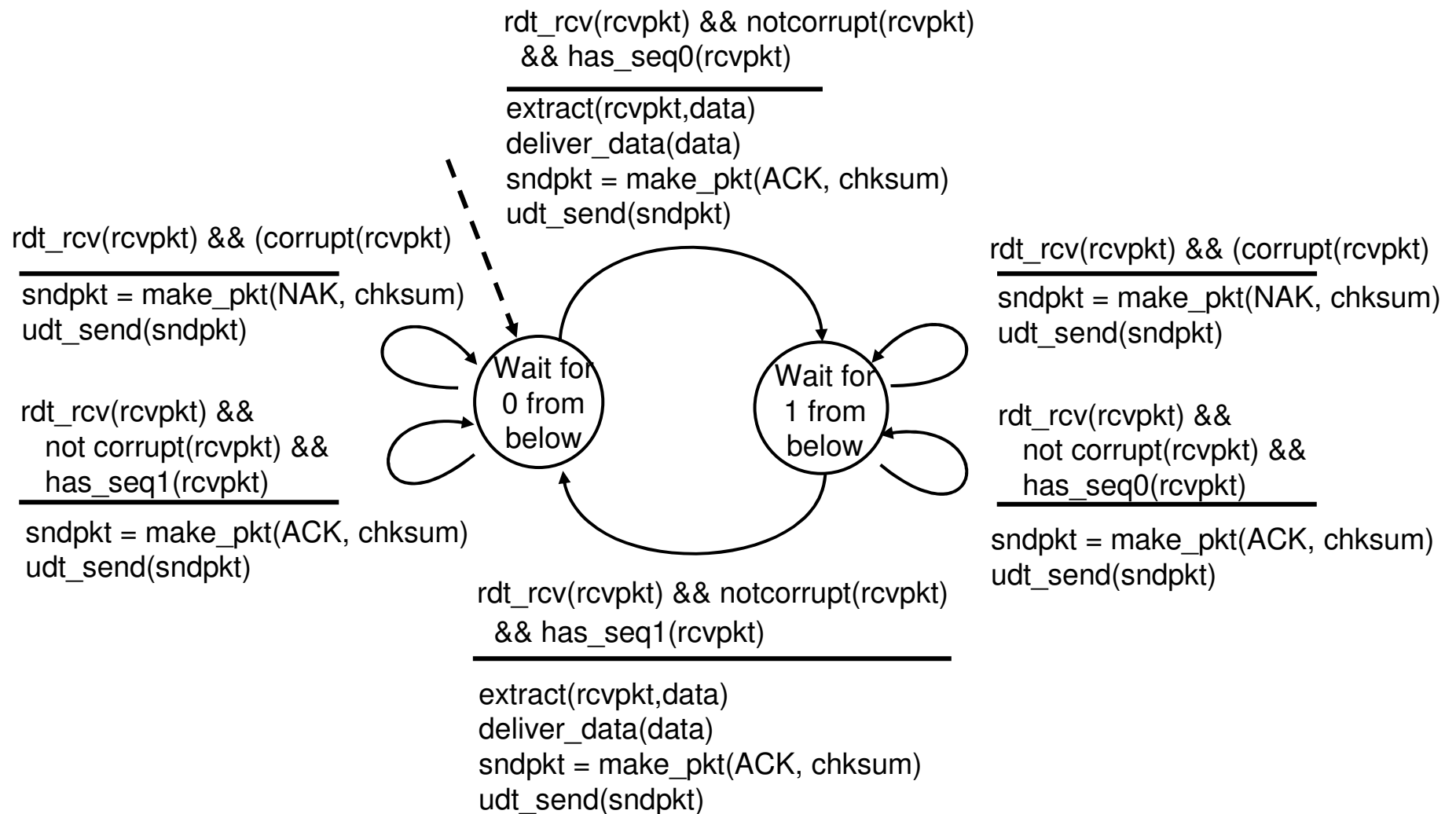


***Need to univocally "label" all packets circulating
in the network between two end points.
1 bit (0-1) enough for Stop-and-wait***

rdt2.1: sender, handles garbled ACK/NAKs



rdt2.1: receiver, handles garbled ACK/NAKs



rdt2.1: discussion

Sender:

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

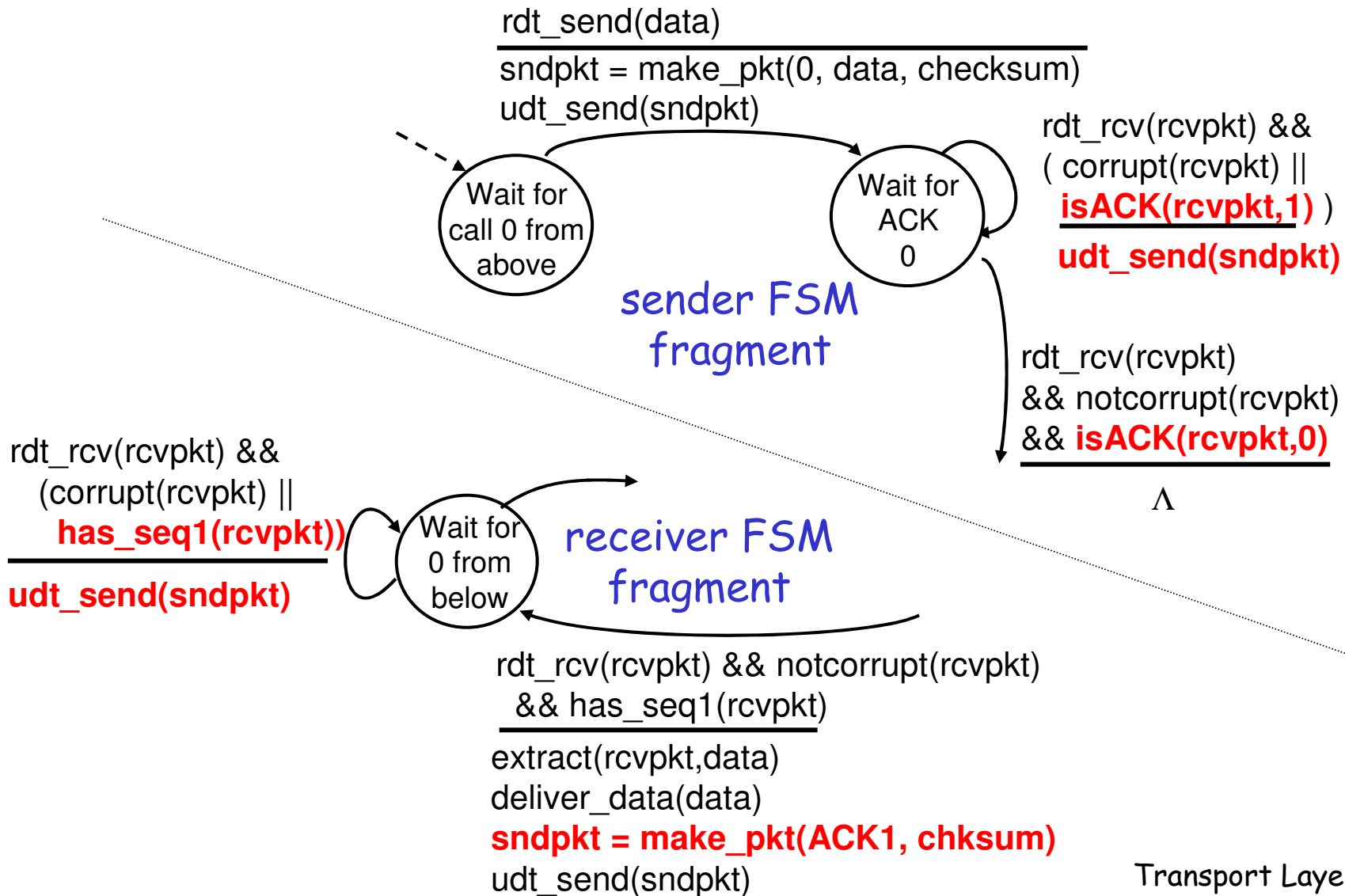
Receiver:

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

rdt2.2: a NAK-free protocol

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

rdt2.2: sender, receiver fragments



rdt3.0: channels with errors *and* loss

New assumption:

underlying channel can also lose packets (data or ACKs)

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

Q: how to deal with loss?

- m sender waits until certain data or ACK lost, then retransmits
- m yuck: drawbacks?

Approach: sender waits

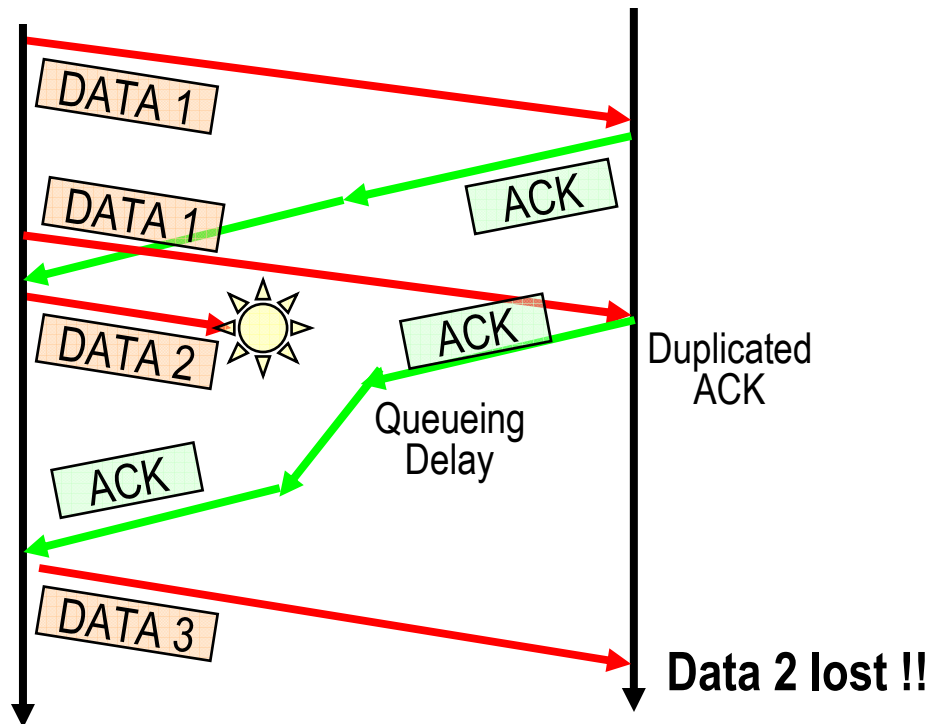
“reasonable” amount of time for ACK

- r retransmits if no ACK received in this time
- r if pkt (or ACK) just delayed (not lost):
 - m retransmission will be duplicate, but use of seq. #'s already handles this
 - m receiver must specify seq # of pkt being ACKed
- r requires countdown timer

Why sequence numbers? (on ack)

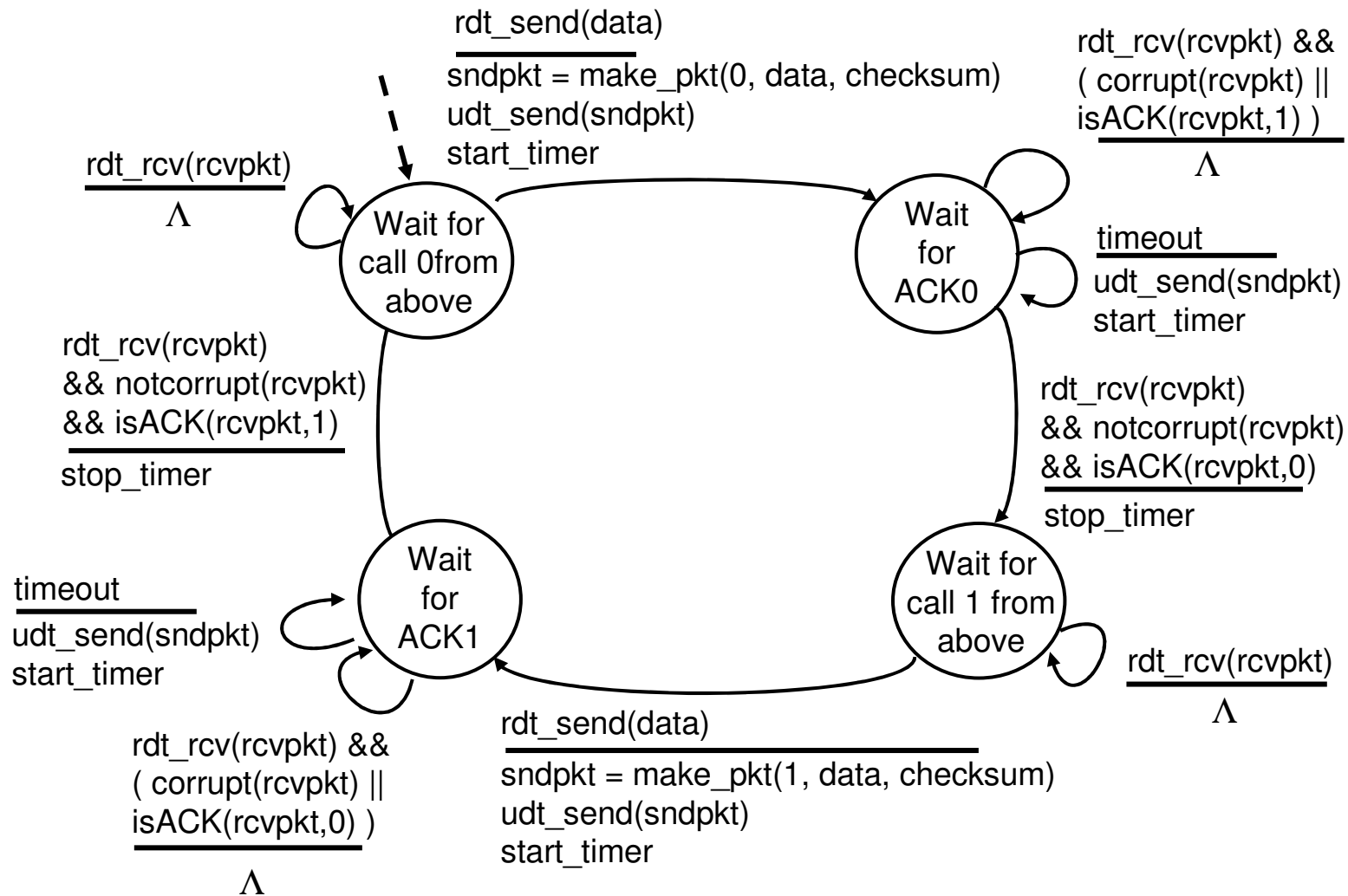
Sender side:

Receiver side:

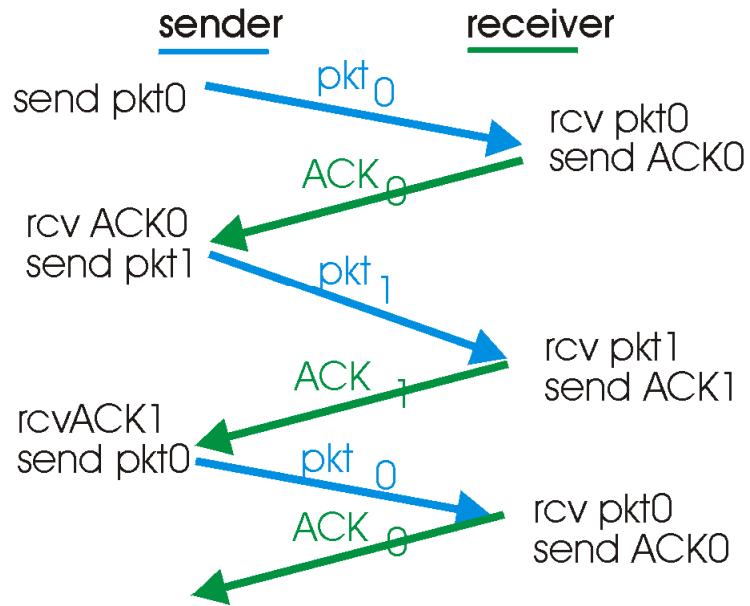


*With pathologically critical network (as the Internet!)
also need to univocally "label" all acks circulating
in the network between two end points.
1 bit (0-1) enough for Stop-and-wait ?*

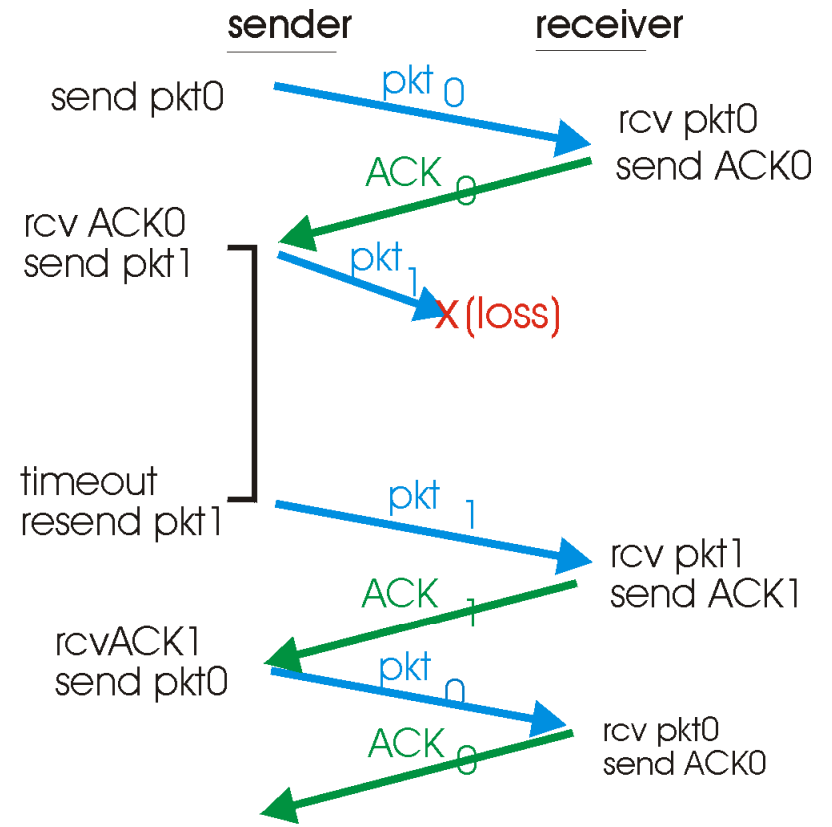
rdt3.0 sender



rdt3.0 in action

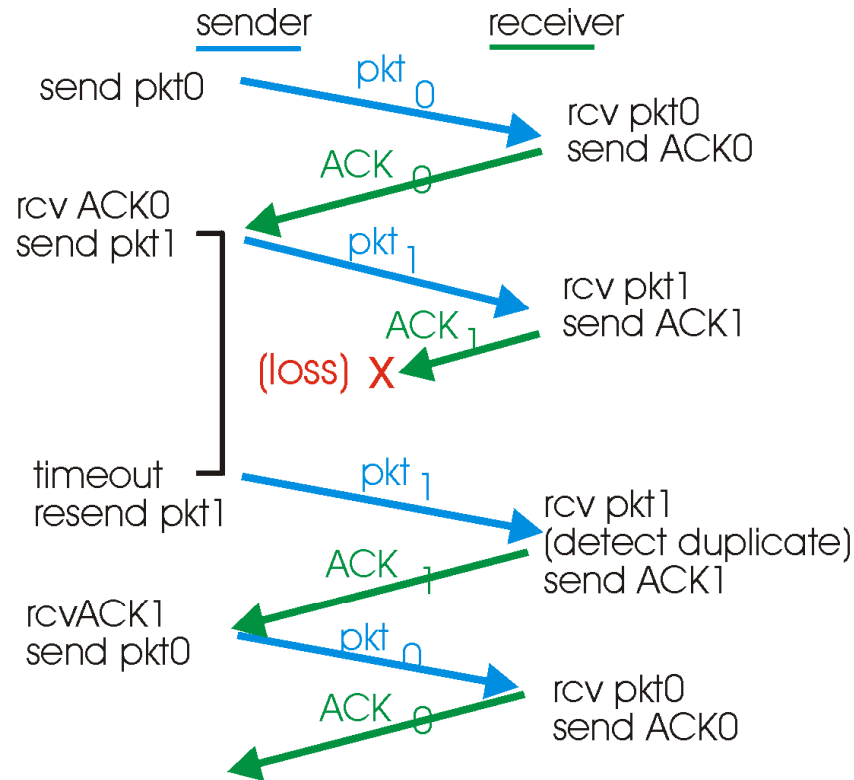


(a) operation with no loss

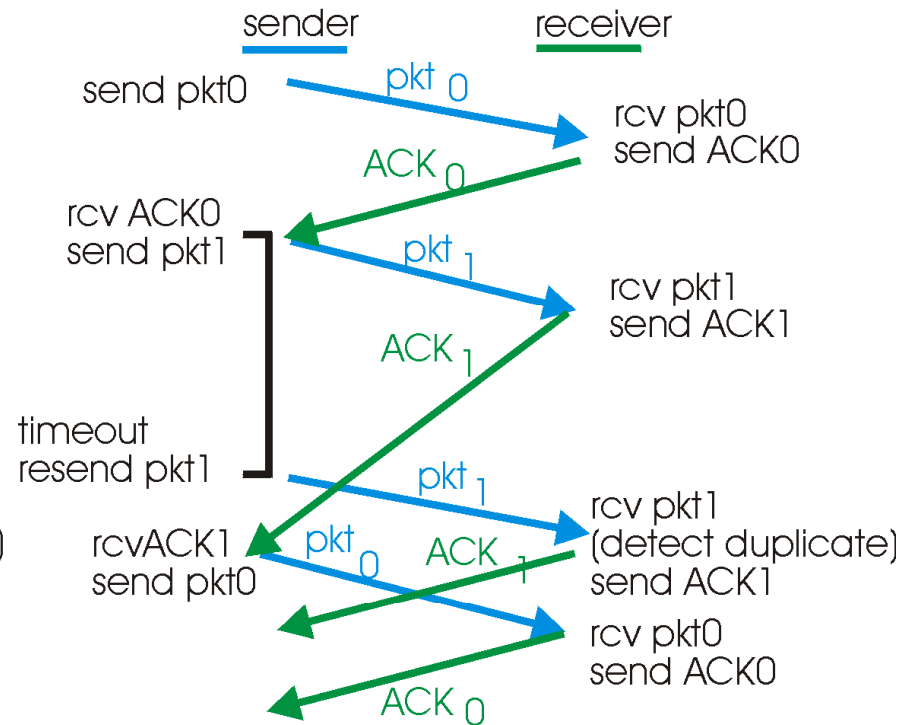


(b) lost packet

rdt3.0 in action



(c) lost ACK



(d) premature timeout

Performance of rdt3.0

r rdt3.0 works, but performance stinks

r example: 1 Gbps link, 15 ms e-e prop. delay, 1KB packet:

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

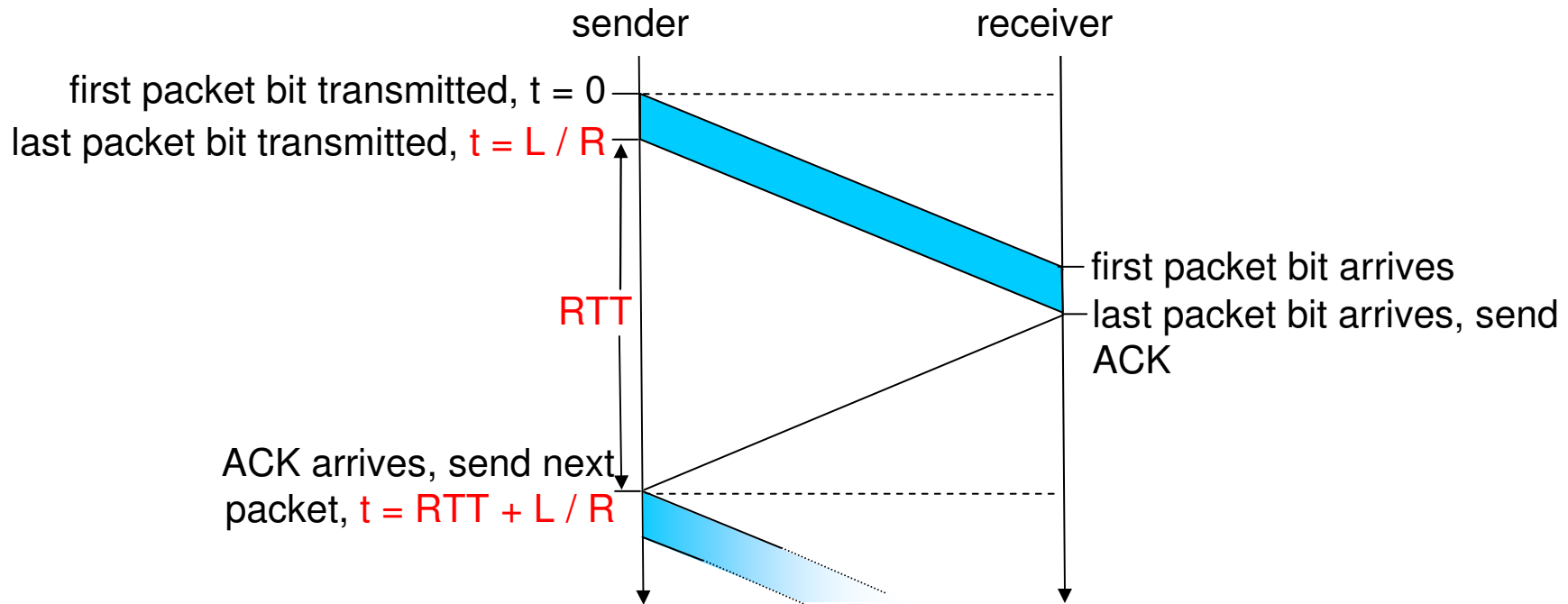
$$U_{\text{sender}} = \frac{L / R}{RTT + L / R} = \frac{.008}{30.008} = 0.00027$$

m U_{sender} : **utilization** - fraction of time sender busy sending

m 1KB pkt every 30 msec -> 33kB/sec throuput over 1 Gbps link

m network protocol limits use of physical resources!

rdt3.0: stop-and-wait operation

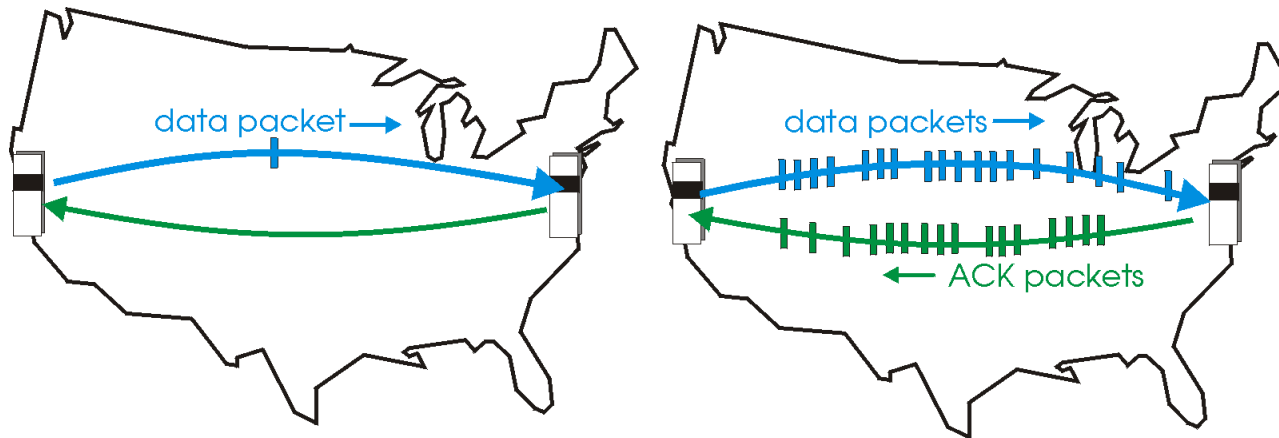


$$U_{\text{sender}} = \frac{L / R}{RTT + L / R} = \frac{.008}{30.008} = 0.00027$$

Pipelined protocols

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

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



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

(b) a pipelined protocol in operation

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