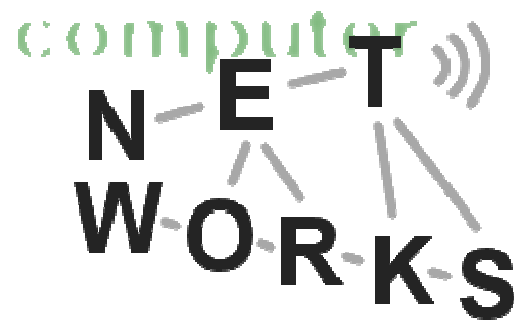# Transport Layer – Part I

## Telematics, Winter 2009/2010

# Chapter 4: The Transport Layer

5: Application Layer

4: Transport Layer

3: Network Layer

2: Link Layer

1: Physical Layer
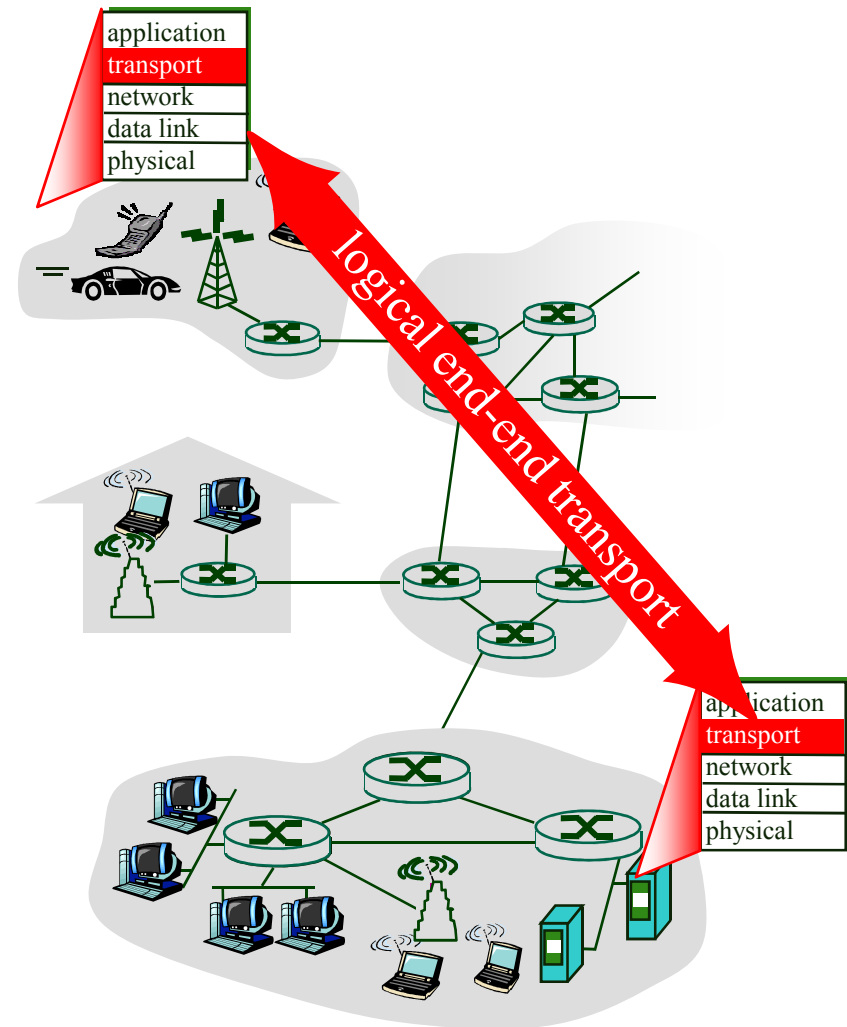
# Chapter 4: The Transport Layer

## Our goals:

o understand principles behind transport layer services:

- o multiplexing/demultiplexing
- o reliable data transfer
- o flow control
- o congestion control

o learn about transport layer protocols in the Internet:

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

# Transport Layer

- 3.1 Transport-layer services

- 3.2 Multiplexing and demultiplexing

- 3.3 Connectionless transport: UDP

- 3.4 Principles of reliable data transfer

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

- 3.6 Principles of congestion control

- 3.7 TCP congestion control

# Transport services and protocols

o   provide *logical communication* between app processes running on different hosts

o   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

o   more than one transport protocol available to apps
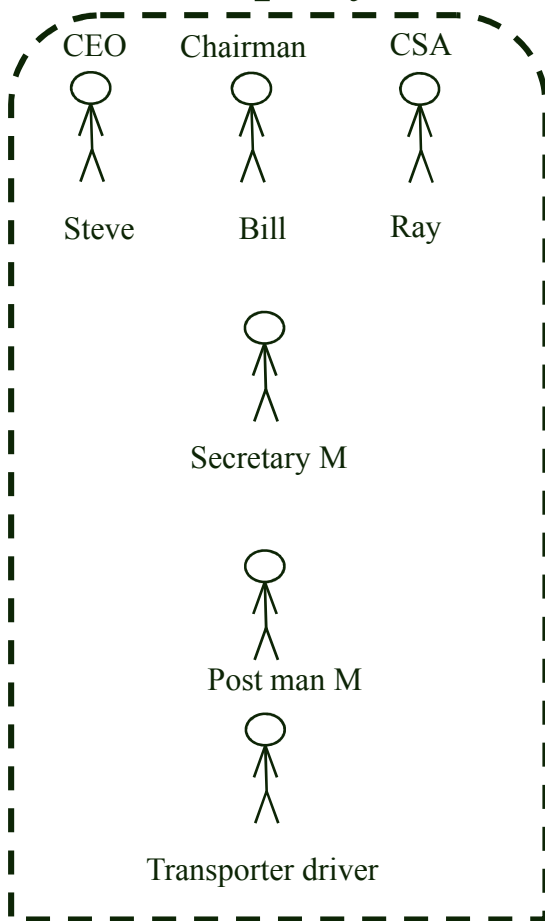-   Internet: TCP and UDP

# Transport vs. network layer

o *network layer:* logical communication between hosts

o *transport layer:* logical communication between processes

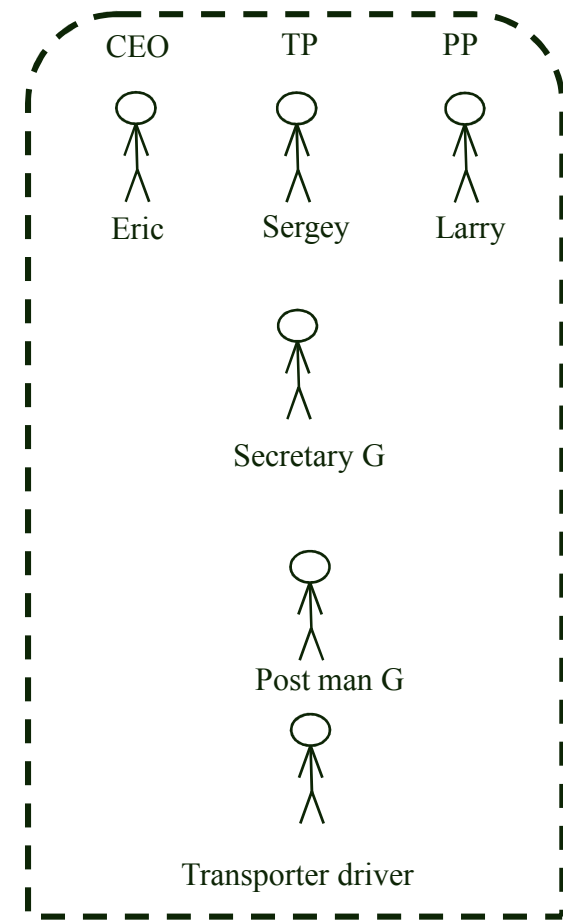    o relies on & enhances, network layer services

# Transport Protocol: Analogy

## Company M

CEO    Chairman    CSA

Steve     Bill     Ray

Secretary M

Post man M

Transporter driver

CEO = Chief Executive Officer
CSA = Chief Software Architect
TP = Technology President
PP = Products President

## Company G

CEO    TP    PP

Eric    Sergey    Larry

Secretary G

Post man G

Transporter driver

Road

# Transport Protocol: Analogy

## Company M

CEO     Chairman     CSA

Steve     Bill     Ray

Secretary M

Post man M

Transporter driver

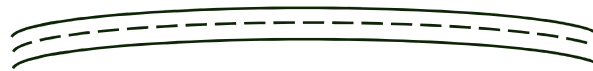CEO = Chief Executive Officer
CSA = Chief Software Architect
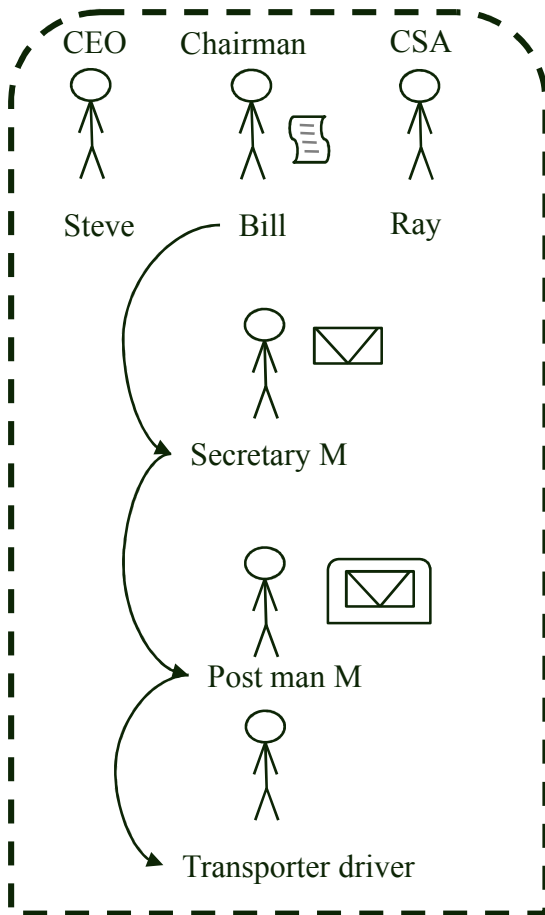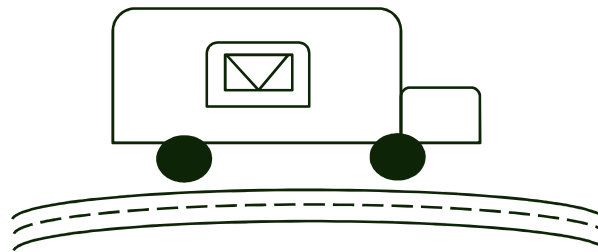TP = Technology President
PP = Products President

## Company G

TP     CEO     PP

Sergey     Eric     Larry

Secretary G

Post man G

Transporter driver

Road

# Transport Protocol: Analogy

o *Postal service (Network Layer):* logical communication between company buildings.

o *Secretary service (Transport Layer):* logical communication between employees of G und M.

   o relies on & enhances, postal services

# Transport Protocol: Analogy

Company

Host

Emp.1    Emp.2    Emp.3

Secretary

Post man

Lorry driver

Road

5: App. Layer (processes)

4: Transport Layer (protocols)

3: Network Layer (protocol)

2: Link Layer (protocols)

1: Physical Layer (medium)

# Internet transport-layer protocols



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

# Excursus: Sockets

## Socket API

o  introduced in BSD4.1 UNIX, 1981

o  explicitly created, used, released by apps

o  client/server paradigm

o  two types of transport service via socket API:

  o  unreliable datagram

  o  reliable, byte stream-oriented
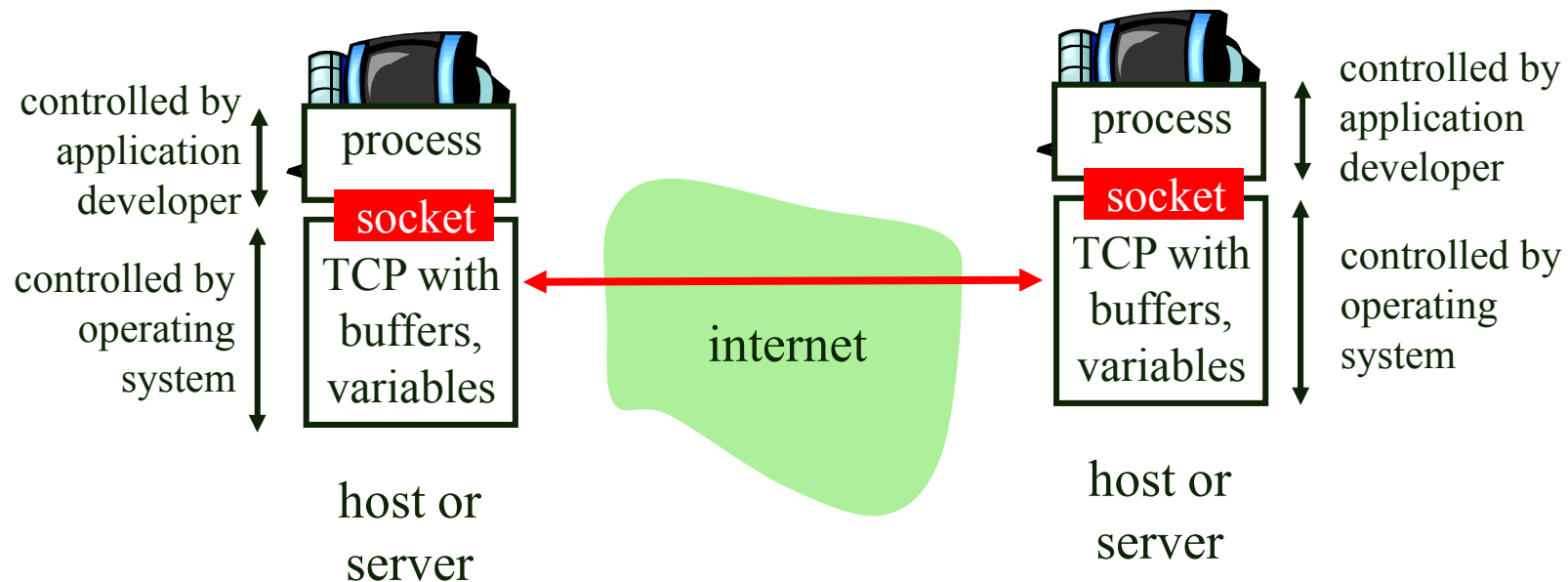
socket

a *host-local*, *application-created*, *OS-controlled* interface (a "door") into which application process can both send and receive messages to/from another application process

# Excursus: Socket programming *with TCP*

Socket: a door between application process and end-end-transport protocol (UDP or TCP)

TCP service: reliable transfer of **bytes** from one process to another

# Excursus: Socket programming *with TCP*

**Client must contact server**

o server process must first be running

o server must have created socket (door) that welcomes client's contact

**Client contacts server by:**

o creating client-local TCP socket

o specifying IP address, port number of server process

o When client creates socket: client TCP establishes connection to server TCP

o When contacted by client, server TCP creates new socket for server process to communicate with client

    o allows server to talk with multiple clients

    o source port numbers used to **distinguish** clients

application viewpoint

*TCP provides reliable, in-order transfer of bytes ("pipe") between client and server*

# Transport Layer

- 3.1 Transport-layer services

- 3.2 Multiplexing and demultiplexing

- 3.3 Connectionless transport: UDP

- 3.4 Principles of reliable data transfer

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

- 3.6 Principles of congestion control

- 3.7 TCP congestion control

# Multiplexing/demultiplexing

Multiplexing at send host:

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

Demultiplexing at rcv host:

delivering received segments
to correct socket

■ ■ = socket      ⬭ = process

| application ICQ 1025 | ICQ 1025 application Skype 1026 | Skype 1026 application |
| transport | transport | transport |
| network | network | network |
| link | link | link |
| physical | physical | physical |

host 1            host 2            host 3

# How demultiplexing works

o host receives IP datagrams

- o each datagram has source IP address, destination IP address
- o each datagram carries 1 transport-layer segment
- o each segment has source, destination port number

o 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

o Create sockets with port numbers:

```
DatagramSocket clientSocket =
    new DatagramSocket();
```

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

o UDP socket identified by  two-tuple:

(dest IP address, dest port number)

o When host receives UDP segment:

- o checks destination port number in segment
- o directs UDP segment to socket with that port number

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

# Connectionless demux (cont)

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

P2

P1

SP provides "return address"

SP: 6428
DP: 5775

server
IP: C

Client
IP:B

SP: 5775
DP: 6428

# Connection-oriented demux

o TCP socket identified by 4-tuple:

- o source IP address
- o source port number
- o dest IP address
- o dest port number

o recv host uses all four values to direct segment to appropriate socket

o Server host may support many simultaneous TCP sockets:

- o each socket identified by its own 4-tuple

o Web servers have different sockets for each connecting client

# Connection-oriented demux (cont)

client
IP: A

server
IP: C

Client
IP:B

P1

P4  P5  P6

P2  P3

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

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

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

# Connection-oriented demux (cont)

client
IP: A

server
IP: C

Client
IP:B

P1

P4 (Apache)

P2    P3

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

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

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

# Transport Layer

- 3.1 Transport-layer services

- 3.2 Multiplexing and demultiplexing

- 3.3 Connectionless transport: UDP

- 3.4 Principles of reliable data transfer

- 3.5 Connection-oriented transport: TCP

  - segment structure

  - reliable data transfer

  - flow control

  - connection management

- 3.6 Principles of congestion control

- 3.7 TCP congestion control

# The Problem with TCP

o TCP offers a **reliable** and **easy** to use transport protocol to programmers.

  o Congestion control
  o Retransmissions etc.

o However congestion control imposes transmission-rate **constraints**.

o If a traffic jam is detected on a path, sender **decreases** sending rate "dramatically".

o **Problem**: One cannot "switch" off functions of TCP ex. Congestion control.

# UDP: User Datagram Protocol [RFC 768]

o "no frills," "bare bones" Internet transport protocol

o "best effort" service, UDP segments may be:
  - o lost
  - o delivered out of order to app

o *connectionless:*
  - o no handshaking between UDP sender, receiver
  - o each UDP segment handled **independently** of others

Why is there a UDP?

o no connection establishment (which can add delay)

o simple: no connection state (buffers & parameters) at sender, receiver

o small segment header (8 bytes v.s. 20 bytes)

o no congestion control & retransmission: UDP can blast away as fast as desired (e.g. used by VOIP)

# UDP: more

o often used for streaming multimedia apps
  - o loss tolerant
  - o rate sensitive
o other UDP uses
  - o DNS
  - o SNMP
o reliable transfer over UDP: add reliability at application layer
  - o application-specific error recovery!
  - o ex.  ACK/NAK, retransmissions (non-trivial).

32 bits

Length, in bytes of UDP segment, including header

| source port # | dest port # |
|---|---|
| length | checksum |

Application data (message)

ex.
- DNS query
- audio sample

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 ….

# UPD checksum example

o Lets take the word "hi" (8bit ASCII)

o Convert it to binary

    o h = 01101000

    o i = 01101001

o Add both words

```
  01101000  (h)
+ 01101001 (i)
  ─────────
  11010001 (h+i)
```

o UDP checksum works with 16 Bit words, but we use 8 Bits for simplicity

o The 1s complement is obtained by inverting ones to zeros and vice versa.

o 11010001 -> 00101110 (checksum)

# UPD checksum example

- Check (unaltered bits):

  01101000  (h)

  + 01101001 (i)

  11010001 (h+i)

  + 00101110 (checksum)

  **11111111 (OK)**

- Check (altered bits):

  01101000  (h)

  + 01101001 (i)

  110100**11** (h+i)

  + 00101110 (checksum)

  **100000001 (NOK!)**

UDP segment

32 bits

| source port # | dest port # |
|---|---|
| length | 00101110 |

01101000 01101001
(h)          (i)

# UDP checksum
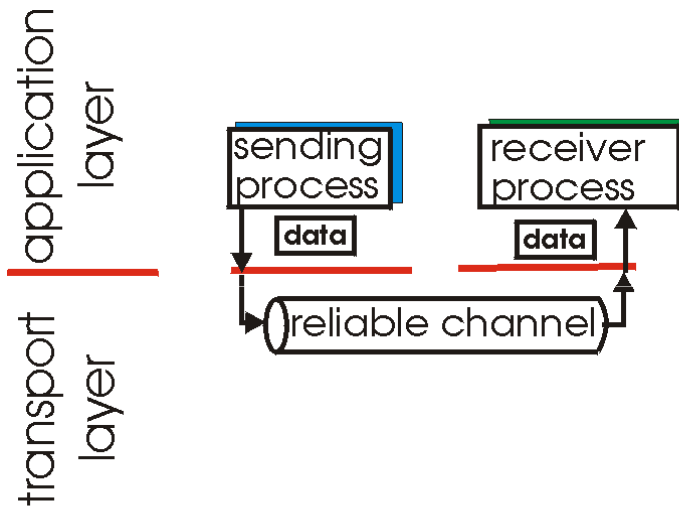
o Why error detection in the first place?

o Link Layer provides CRC! (Ethernet)

o No guarantee for:
  o link-to-link reliability (e.g. non ethernet)
  o memory error detection on routers

o IP is designed to run on any layer 2 protocol (ethernet, PPP, 802.11, 802.16).

o End-to-end error detection is safety measure

o UPD does not recover from errors (discard/warning)

# Transport Layer

# Principles of Reliable data transfer

o   important in app., transport, link layers

o   top-10 list of important networking topics!

application layer

sending process

receiver process

data

data

reliable channel

transport layer

(a)  provided service

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

# Principles of Reliable data transfer

o important in app., transport, link layers

o top-10 list of important networking topics!



(a) provided service          (b) service implementation

o 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
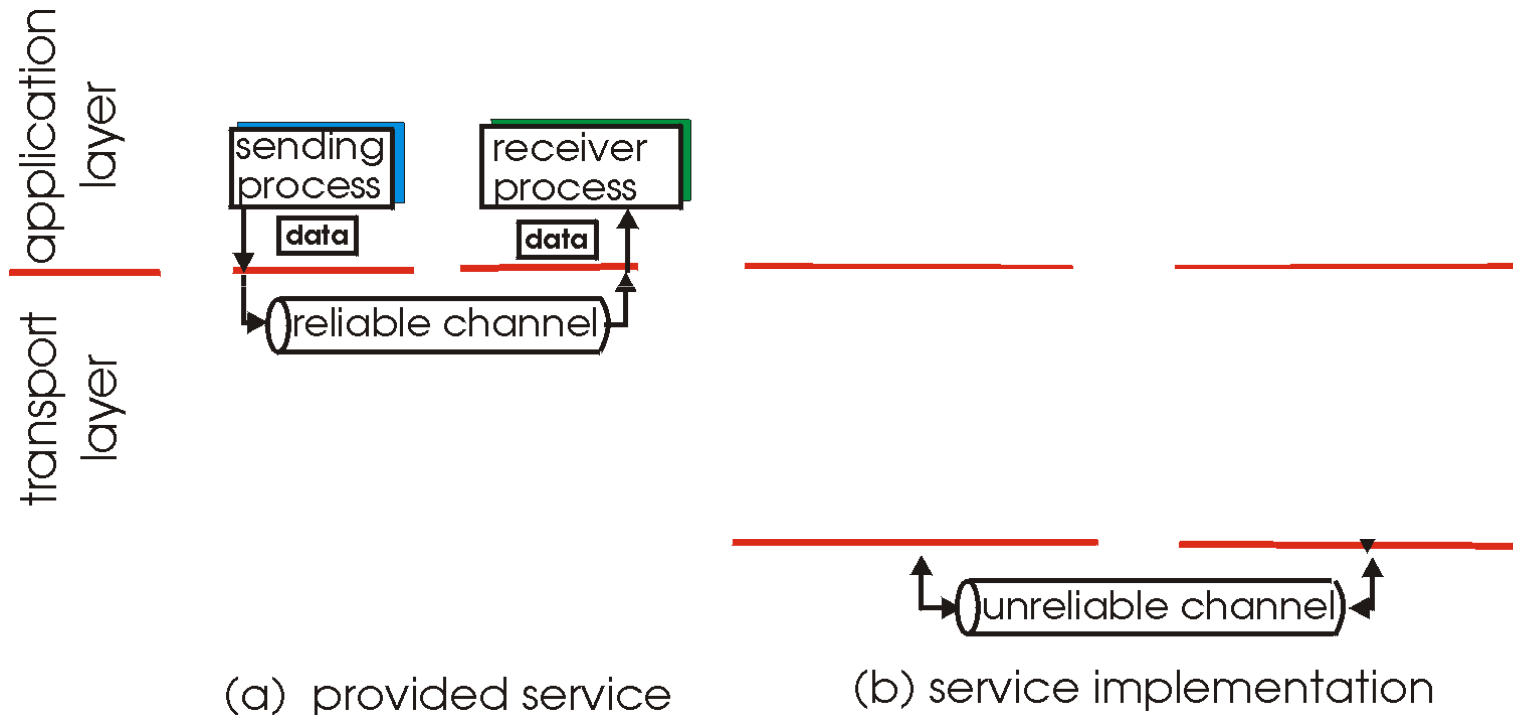
transport layer

sending process

receiver process

data

data

reliable channel

(a) provided service

rdt_send() data

data deliver_data()

reliable data transfer protocol (sending side)

reliable data transfer protocol (receiving side)

udt_send() packet

packet rdt_rcv()

unreliable channel

(b) service implementation

- characteristics of unreliable channel will determine complexity of reliable data transfer protocol (rdt)
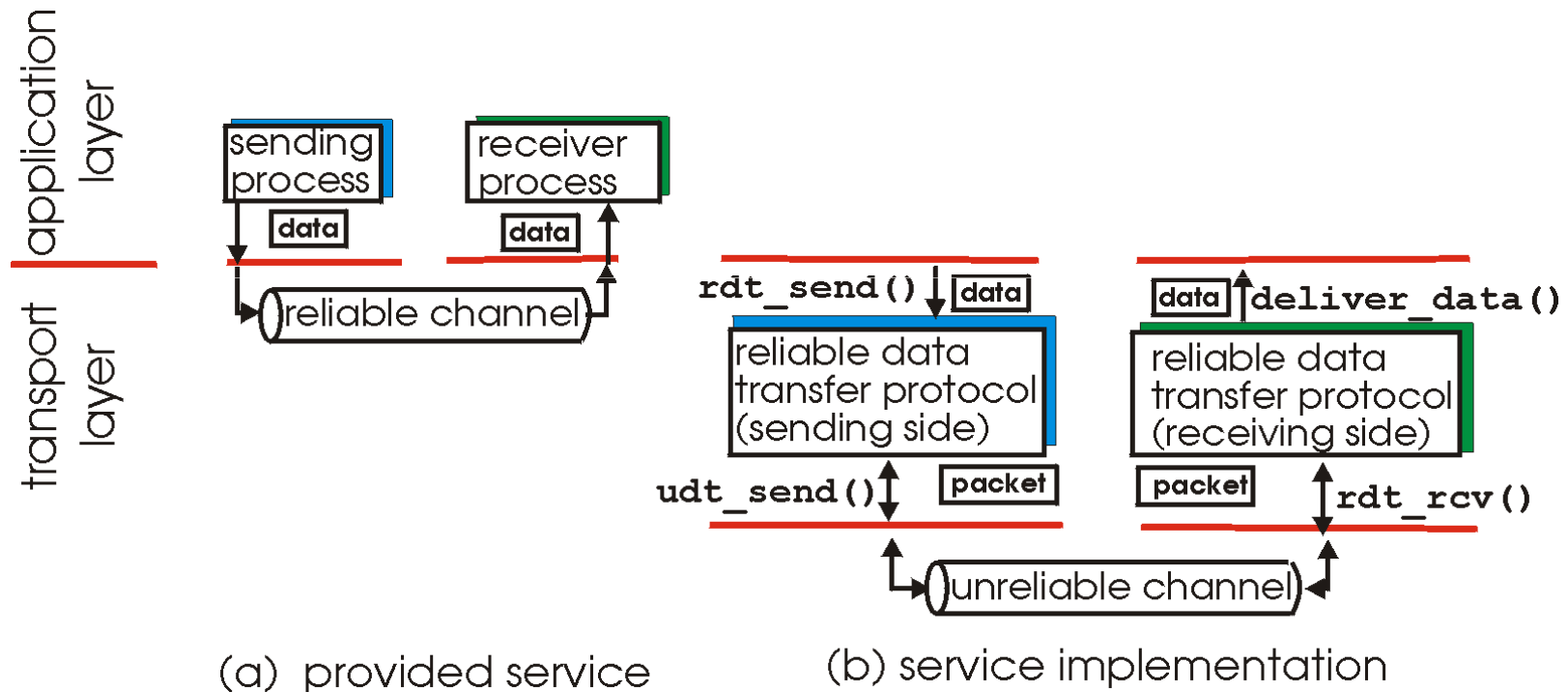
# Reliable data transfer: getting started

`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

`rdt_send()` ↓ data

data ↑ `deliver_data()`

send side

reliable data transfer protocol (sending side)

reliable data transfer protocol (receiving side)

receive side

`udt_send()` ↕ packet

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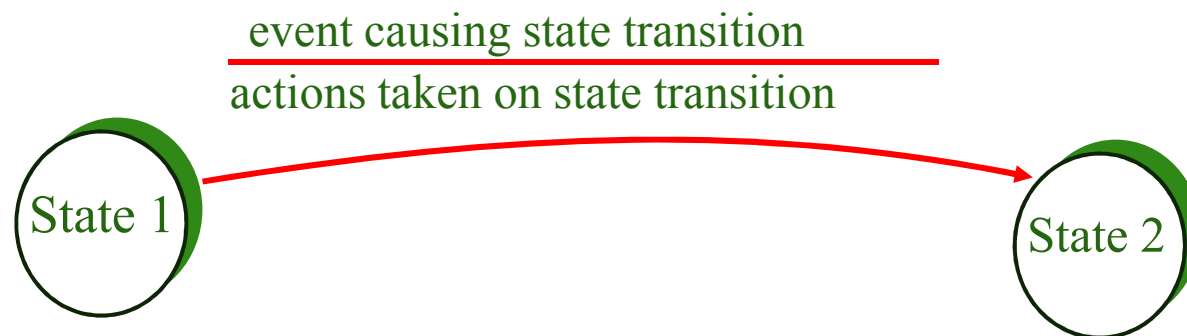
o incrementally develop sender, receiver sides of reliable data transfer protocol (rdt)

o consider only unidirectional data transfer

- o but control info will flow on both directions!

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

o Use generic term "packet" rather than "segment"

# Finite State Machine

o FSM is a model of behavior composed of a finite number of

   o states

   o transitions between states on events

   o actions taken upon events

o Necessary to define the behavior of our protocol, prior to implementation

event causing state transition

actions taken on state transition

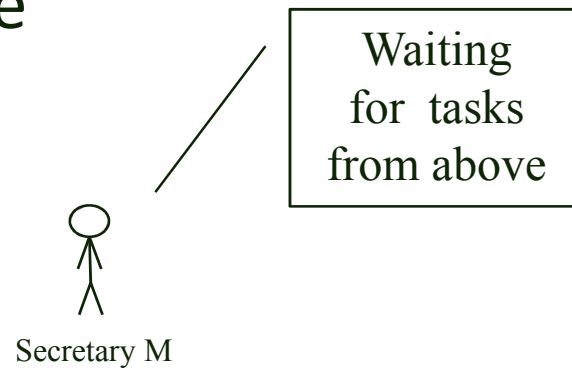State 1                     State 2

# Rdt1.0: reliable transfer over a reliable channel

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

- o We will first look at an analogy with the secretary then the state machines.

# Rdt1.0: reliable transfer over a reliable channel (Analogy)

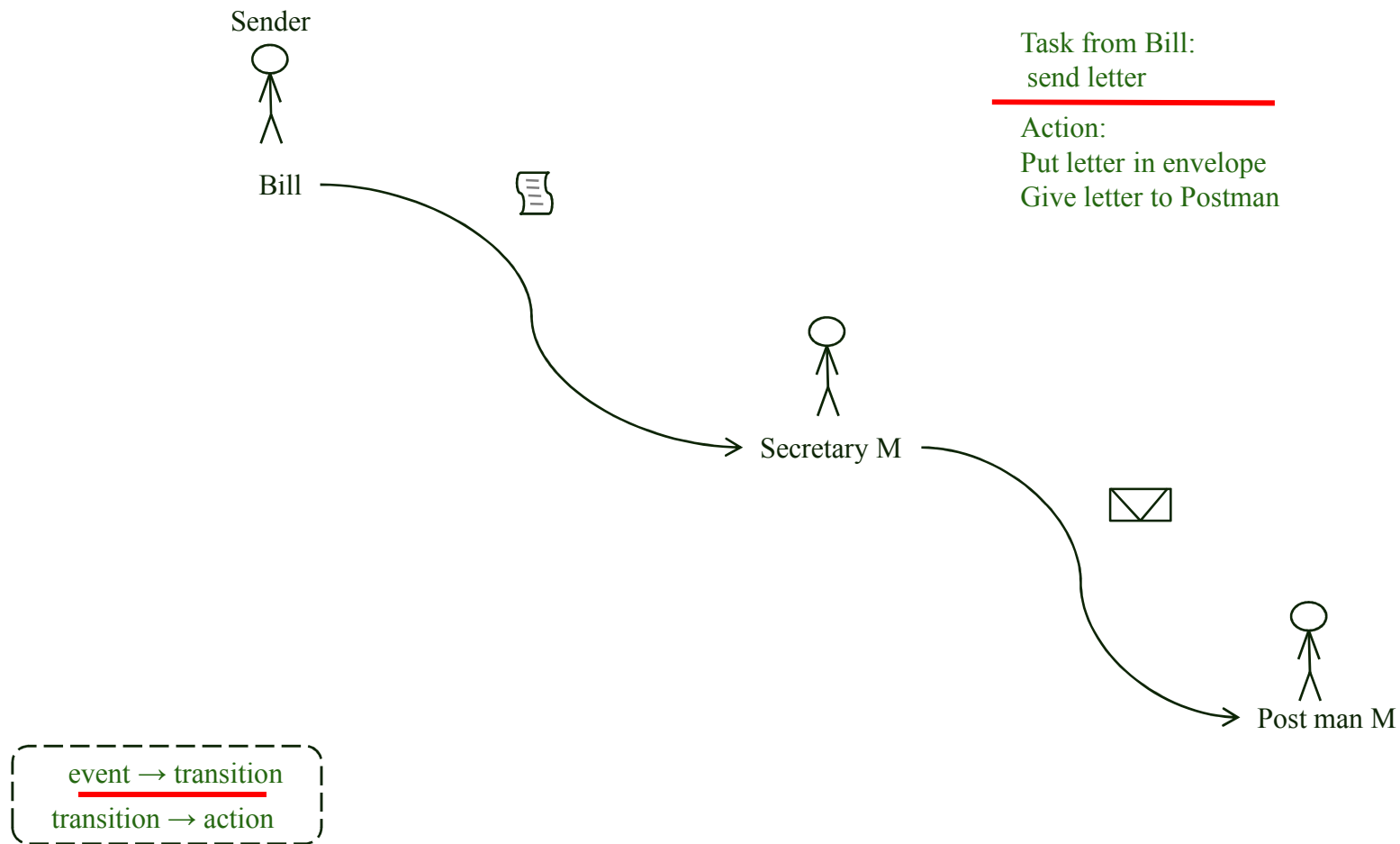o The secretary from
   our previous example
   has one state

o He waits for tasks
   from his boss

o Task is sending letters

Waiting
for tasks
from above

Secretary M

# Rdt1.0: reliable transfer over a reliable channel (Analogy)

Sender

Bill

Task from Bill:
 send letter

Action:
Put letter in envelope
Give letter to Postman

Secretary M

Post man M

event → transition

transition → action

# Rdt1.0: reliable transfer over a reliable channel (Analogy)

o The secretary goes back to his state, waiting for more tasks.

Waiting for tasks from above

Secretary M

# Rdt1.0: reliable transfer over a reliable 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

event → transition
-------------------
transition → action

Λ = no event/action

N E T W O R K S

# Rdt2.0: channel with bit errors

○ underlying channel may flip bits in packet

    ○ checksum to detect bit errors     00101110

○ *the* **question**: how to recover from errors?

○ **Analogy**:

    ○ Imagine you dictate phone number over cell phone to friend.

    ○ Bad reception may scramble your voice.



You          Friend

0176

OK

1234

Not OK, repeat please

1234

OK

# Rdt2.0: channel with bit 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

- **A**utomatic **R**epeat re**Q**uest type of protocol (**ARQ**)

# rdt2.0: FSM specification

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

receiver

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

Wait for
call from
above

Wait for
ACK or
NAK

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

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

sender

Wait for
call from
below

event → transition
transition → action

Λ = no event/action

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?

o sender doesn't know what happened at receiver!

o can't just retransmit: possible duplicate

## Handling duplicates:

o sender retransmits current pkt if ACK/NAK garbled

o sender adds *sequence number* to each pkt

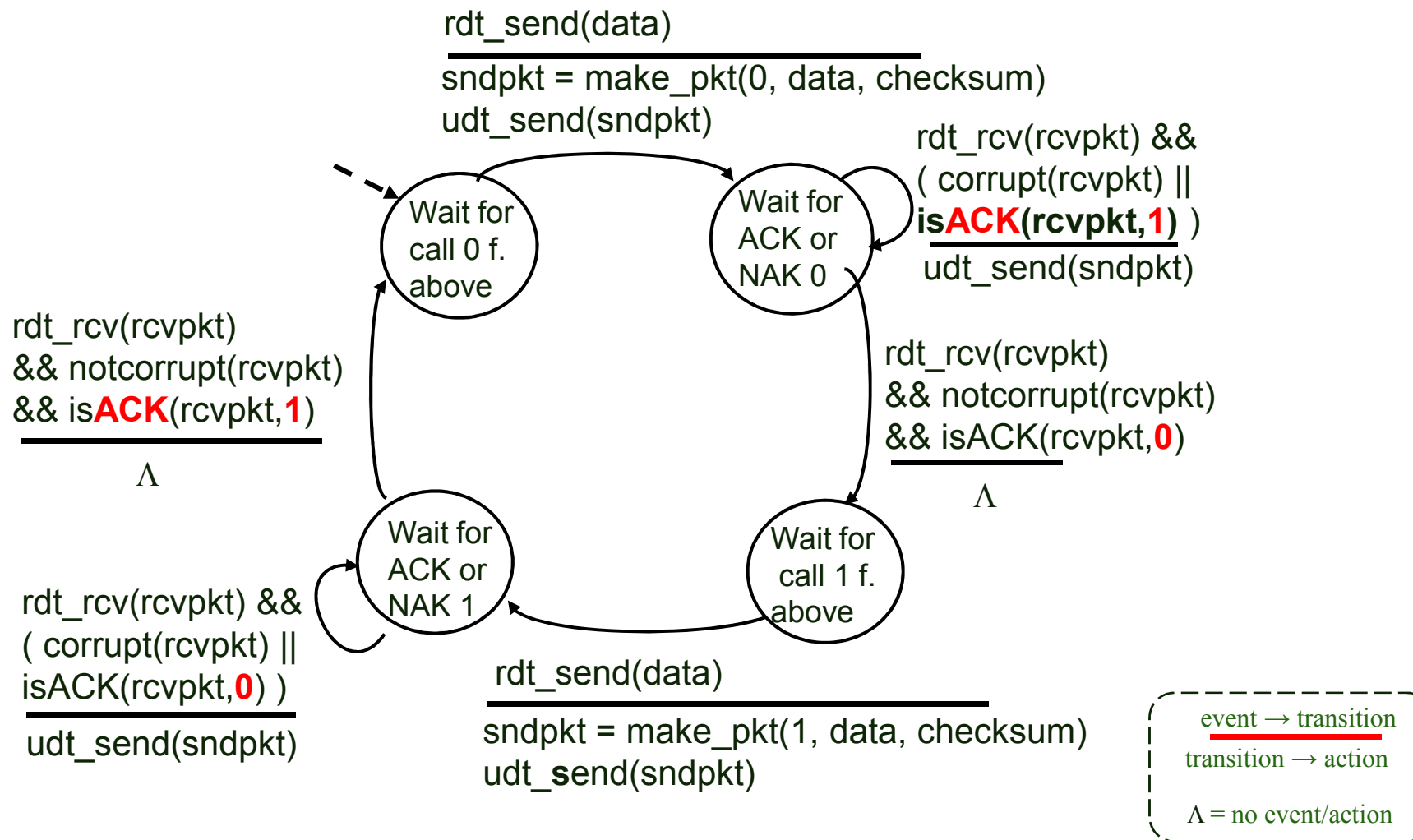o receiver discards (doesn't deliver up) duplicate pkt

## Using only ACK + Sequence:

o We can **discard** NAK packets, by using only ACK + Seq.#

o duplicate ACK at sender results in same action as NAK: *retransmit current pkt*

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

# rdt2.2: sender, handles garbled ACKs

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 f.
above

Wait for
ACK or
NAK 0

rdt_rcv(rcvpkt)
&& notcorrupt(rcvpkt)
&& is**ACK**(rcvpkt,**1**)
_____
Λ

rdt_rcv(rcvpkt)
&& notcorrupt(rcvpkt)
&& isACK(rcvpkt,**0**)
_____
Λ

Wait for
ACK or
NAK 1

Wait for
call 1 f.
above

rdt_rcv(rcvpkt) &&
( corrupt(rcvpkt) ||
isACK(rcvpkt,**0**) )
_____
udt_send(sndpkt)

rdt_send(data)
_____

sndpkt = make_pkt(1, data, checksum)
udt_**s**end(sndpkt)

event → transition
_____
transition → action

Λ = no event/action

# rdt2.2: receiver, handles garbled ACKs

rdt_rcv(rcvpkt) && notcorrupt(rcvpkt)
  && has_seq0(rcvpkt)
_____
extract(rcvpkt,data)
deliver_data(data)
sndpkt = make_pkt(ACK,0, chksum)
udt_send(sndpkt)

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

**Wait for 0 from below**

**Wait for 1 from below**

rdt_rcv(rcvpkt) &&
  (corrupt(rcvpkt) ||
  has_seq0(rcvpkt)
_____
udt_send(sndpkt)

rdt_rcv(rcvpkt) && notcorrupt(rcvpkt)

  && has_seq1(rcvpkt)
_____
extract(rcvpkt,data)
deliver_data(data)
sndpkt = make_pkt(ACK, 1,chksum)
udt_send(sndpkt)

event → transition
_____
transition → action

Λ = no event/action

# rdt2.2: discussion

## Sender:

o seq # added to pkt

o two seq. #'s (0,1) will suffice.  Why?

o must check if received ACK corrupted

o twice as many states

  o state must "remember" whether "current" pkt has 0 or 1 seq. #

## Receiver:

o must check if received packet is duplicate

  o state indicates whether 0 or 1 is expected pkt seq #

o note: receiver can *not* know if its last ACK received OK at sender

# rdt: What do we have so far?

- ## rdt 1.0
  - simple transfer over reliable channel (unrealistic)
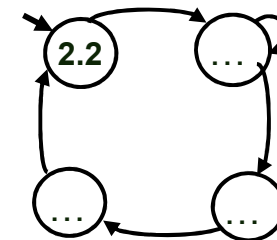
- ## rdt 2.0
  - bit error prone channel (more realistic)
  - checksum (data), ACK/NAK, retransmit
  - **but what if ACK corrupt**?

- ## rdt 2.2
  - checksum (data & **ACK**)
  - retransmit if ACK corrupt
  - **but what if data OK, but ACK corrupt? -> duplicate**
  - introduce sequence numbers (more states)
  - slimed down: discard NAK by introducing seq. in ACK
  - **but what if channel looses packets?**
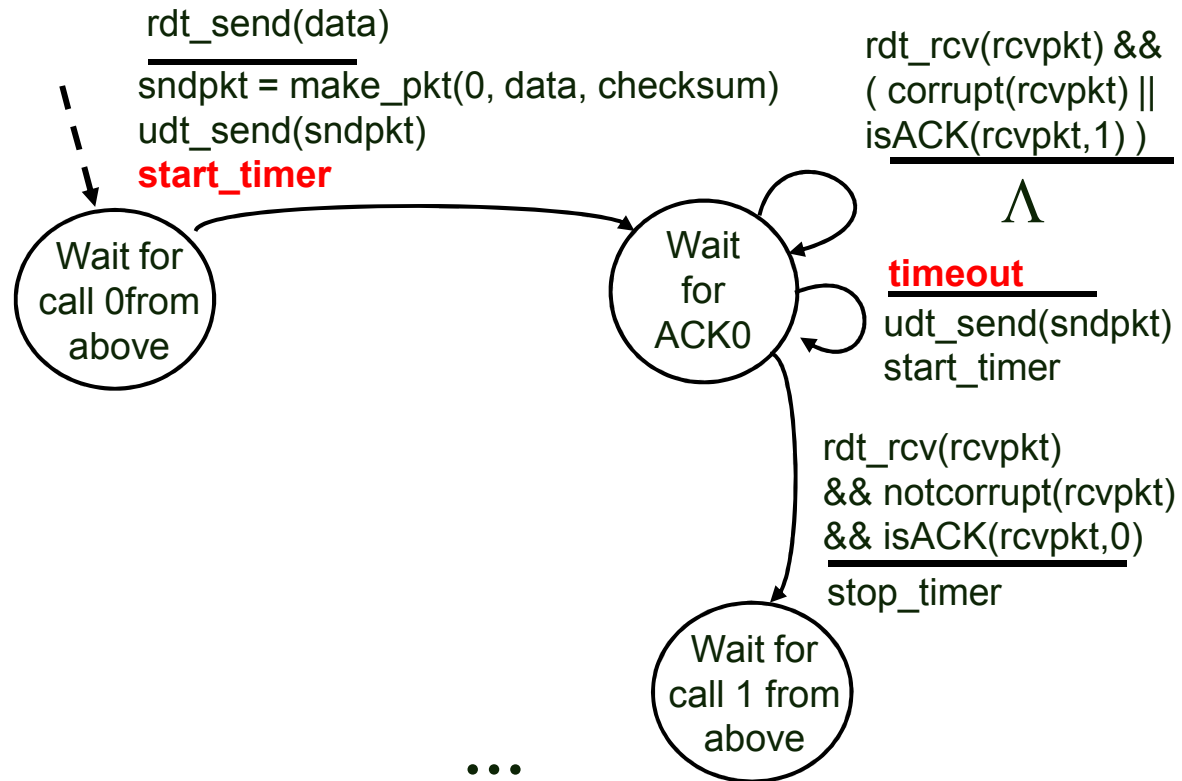
# rdt3.0: channels with errors *and* loss

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

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

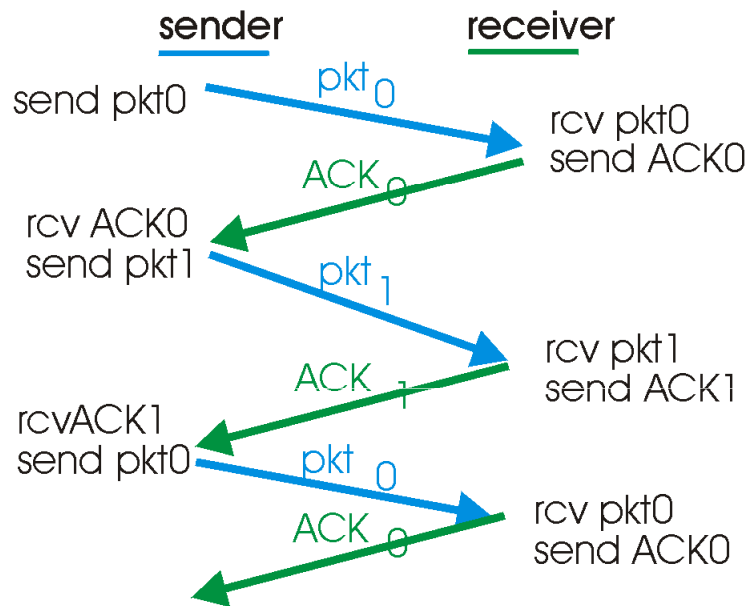Approach: sender waits "reasonable" amount of time for ACK

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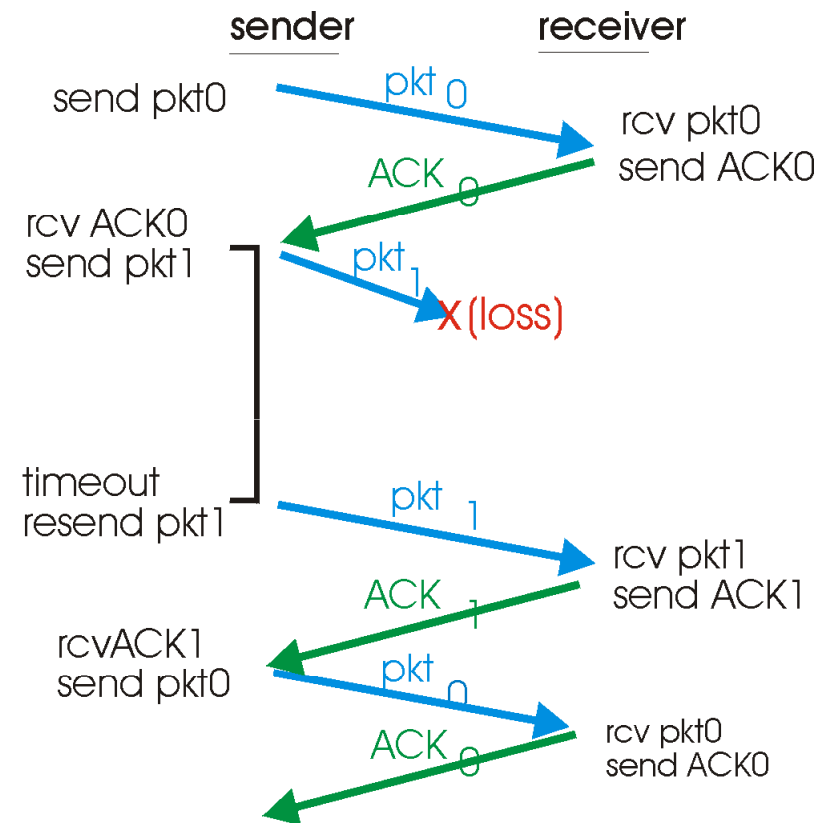- o requires countdown timer

# 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) &&
( corrupt(rcvpkt) ||
isACK(rcvpkt,1) )
---

$\Lambda$

**Wait for ACK0**

**timeout**
---
udt_send(sndpkt)
start_timer

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

**Wait for call 1 from above**

$\cdots$

event $\rightarrow$ transition
---
transition $\rightarrow$ action

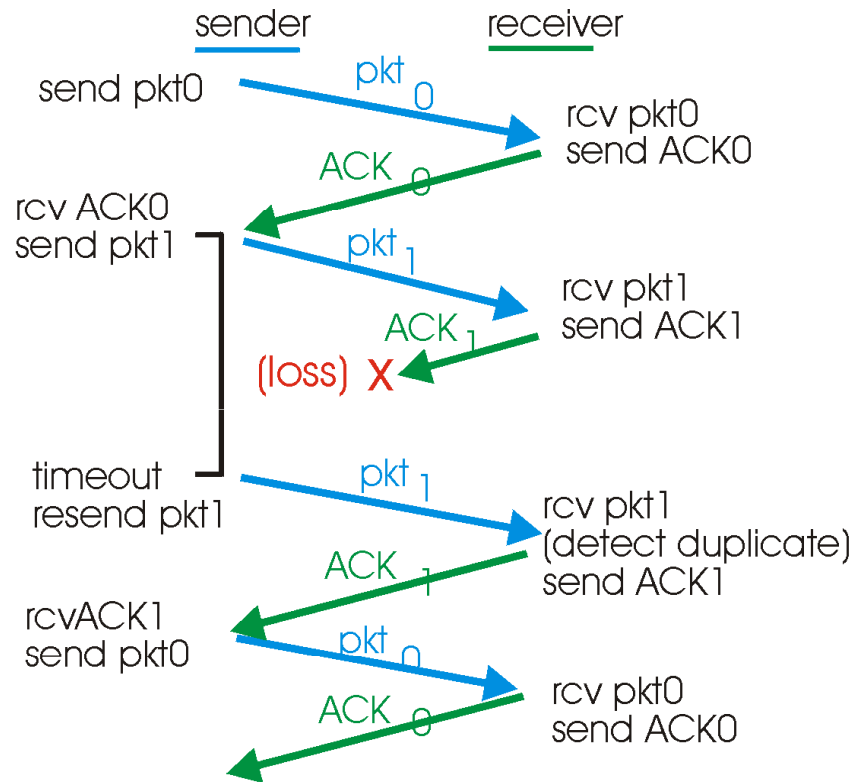$\Lambda$ = no event/action

# rdt3.0 in action



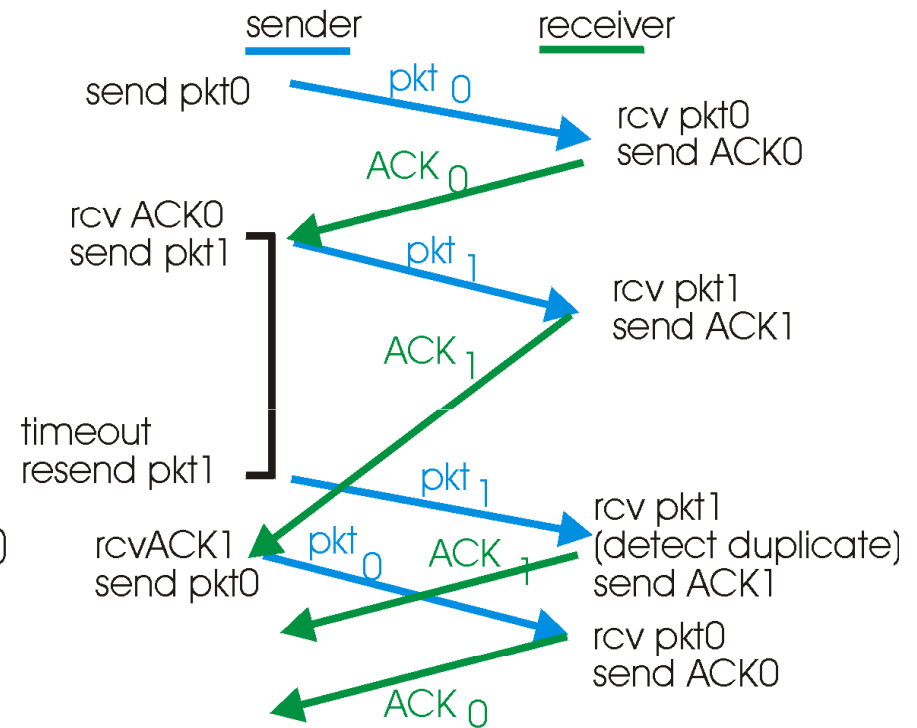(a) operation with no loss



(b) lost packet

# rdt3.0 in action



(c) lost ACK

(d) premature timeout

# Performance of rdt3.0

o rdt3.0 works, but performance stinks

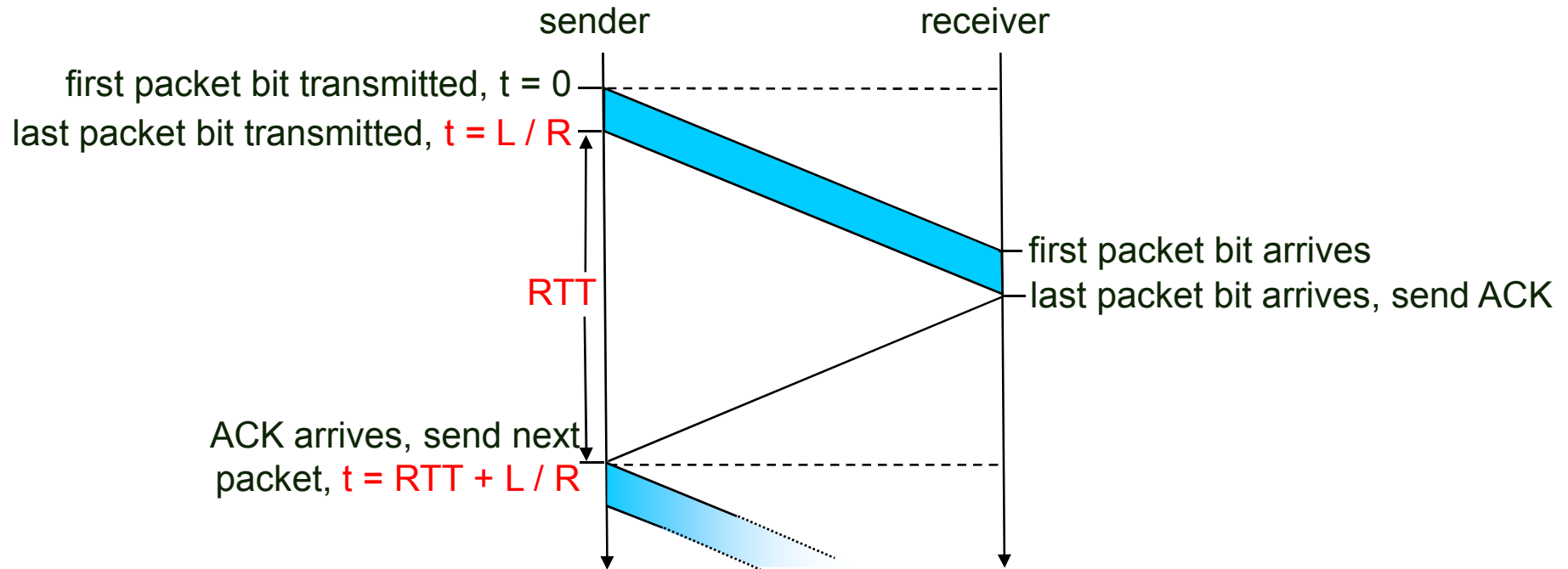o ex: 1 Gbps link, 15 ms prop. delay, 8000 bit packet:

$$d_{trans} = \frac{L}{R} = \frac{8000\,\text{bits}}{10^9\,\text{bps}} = 8\,\text{microseconds}$$

○ U $_{sender}$: utilization – fraction of time sender busy sending

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

○ 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



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

RTT

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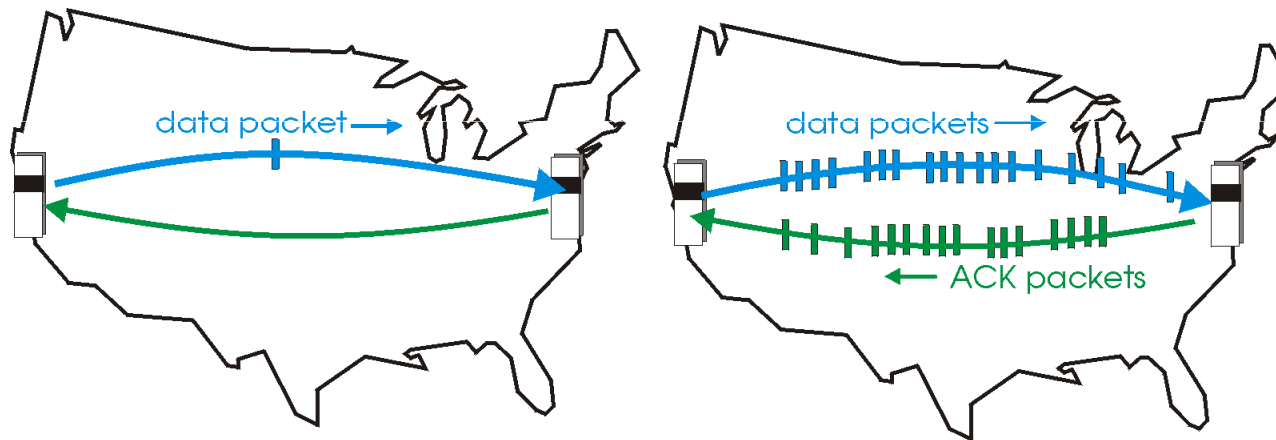
ACK arrives, send next
packet, t = RTT + L / R

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

# 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



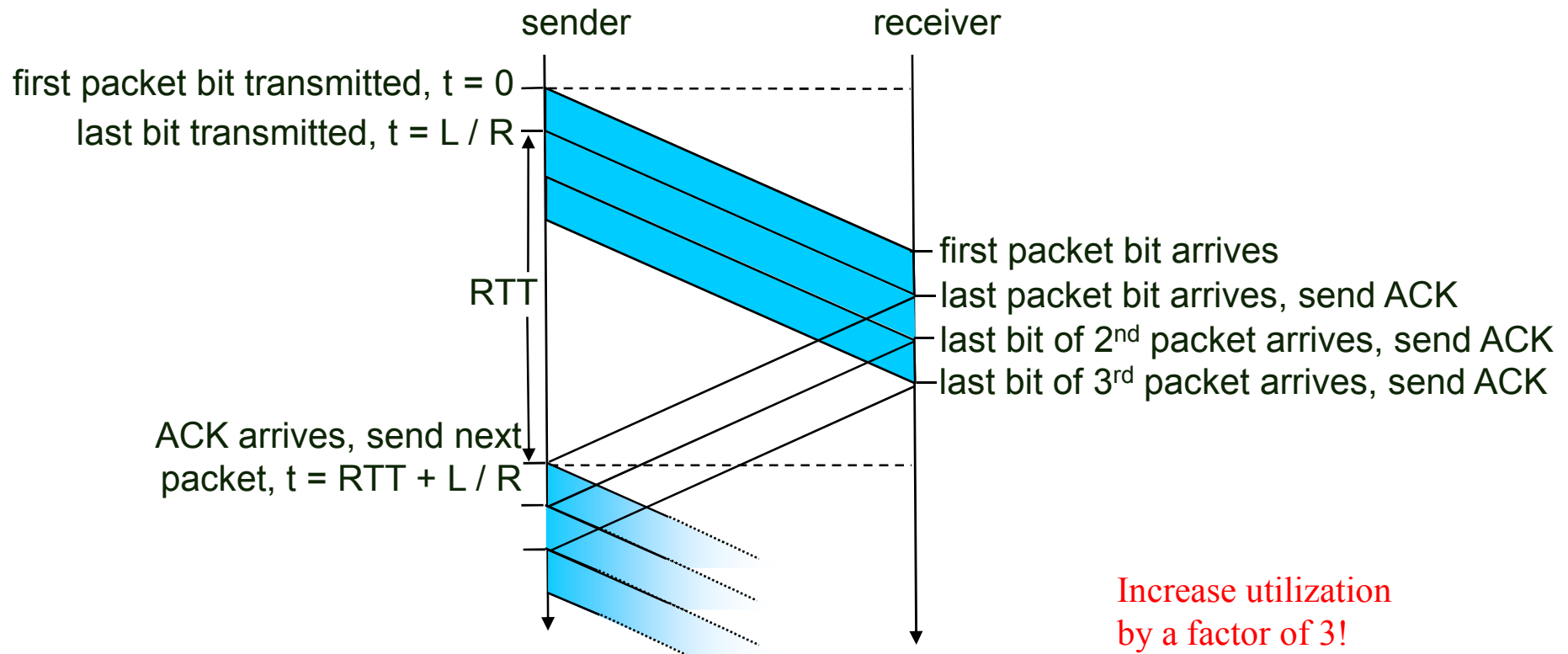(a) a stop-and-wait protocol in operation     (b) a pipelined protocol in operation

- 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

last bit of 2nd packet arrives, send ACK

last bit of 3rd packet arrives, send ACK

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

Increase utilization by a factor of 3!

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

# Pipelining Protocols

## Go-back-N: big picture:

o Sender can have up to N unacked packets in pipeline

o Rcvr only sends cumulative acks
  - o Doesn't ack packet if there's a gap

o Sender has timer for oldest unacked packet
  - o If timer expires, retransmit all unacked packets

## Selective Repeat: big pic

o Sender can have up to N unacked packets in pipeline

o Rcvr acks individual packets

o Sender maintains timer for each unacked packet
  - o When timer expires, retransmit only unack packet

# Go-Back-N (GBN) Demonstration

o Protocol Demo (Link)

o http://media.pearsoncmg.com/aw/aw_kurose_net work_2/applets/go-back-n/go-back-n.html

# Chapter 4: Summary

o principles behind transport layer services:

    o multiplexing, demultiplexing

    o reliable data transfer

Next:

    o flow control

    o congestion control

o instantiation and implementation in the Internet

    o UDP

    o TCP

# Thank you

Any questions?