Mic!

Chapters
3.4 & 3.5

# Principles of Reliable Transport

INTERNATIONAL EDITION

Computer Networking
Top-Down Approach

James F. Kurose • Keith W. Ross

ALWAYS LEARNING          PEARSON

Geoff Coulson
Week 15 Lecture 1

# Preamble: we are being very selective

- we'll focus on *principles of reliability* than on the details of TCP
  - although we will look at TCP as well
- we'll be using slides from Kurose and Ross
  - I hereby acknowledge their copyright!
  - indication as we go of topics we're omitting from the book
  - but it's very well worthwhile reading and understanding the omitted sections!

# Chapter 3 outline (from the book)

3.1 transport-layer services

3.2 multiplexing and demultiplexing

3.3 connectionless transport: UDP

3.4 principles of reliable data transfer (omit: selective resend)

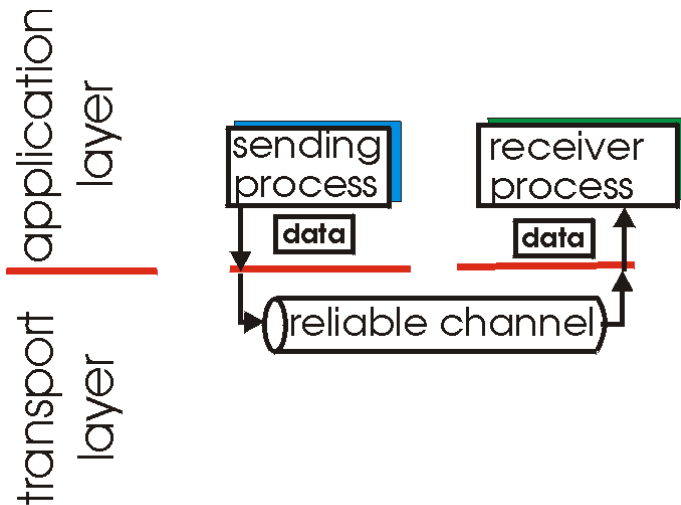3.5 connection-oriented transport: TCP
- segment structure
- reliable data transfer
- flow control
- connection management

3.6 principles of congestion control

3.7 TCP congestion control
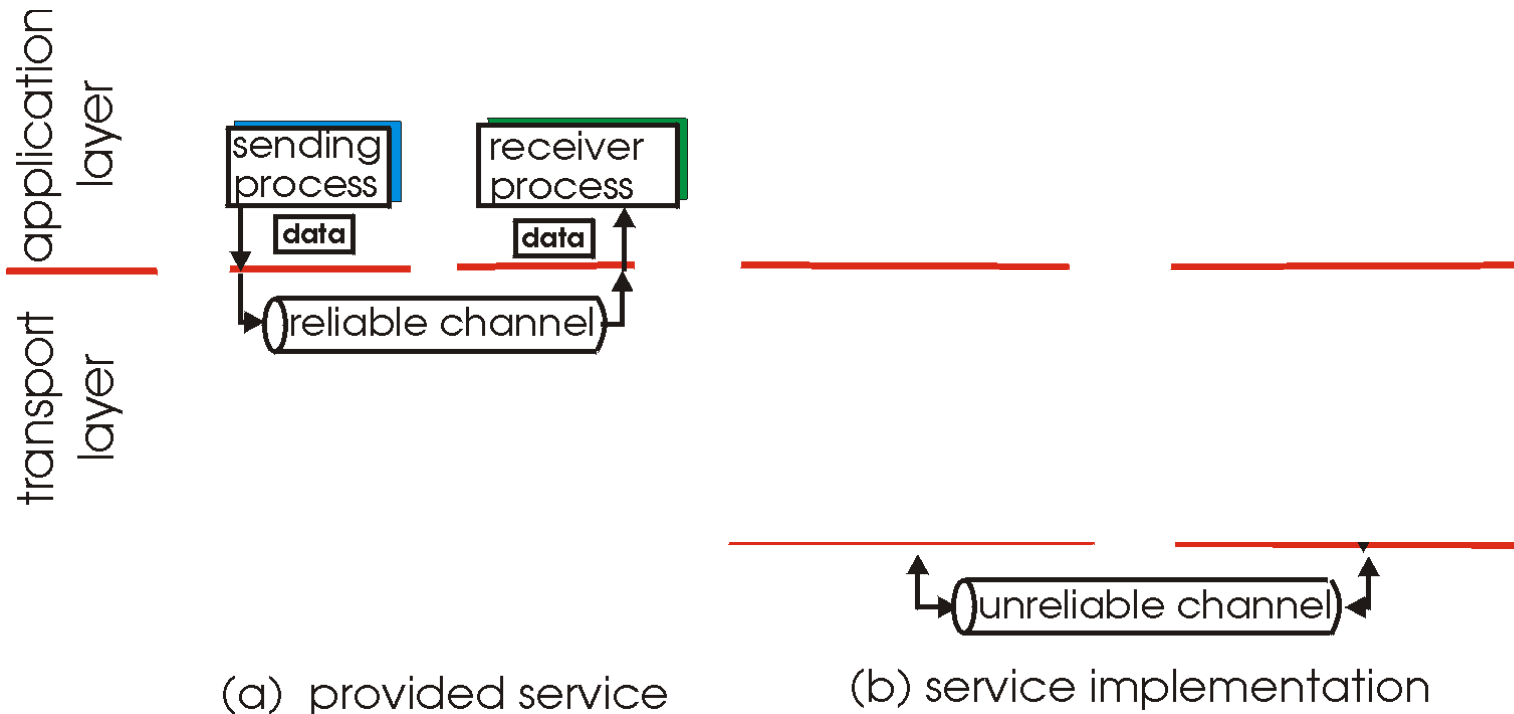
# Principles of reliable data transfer

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



(a) provided service

# Principles of reliable data transfer

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



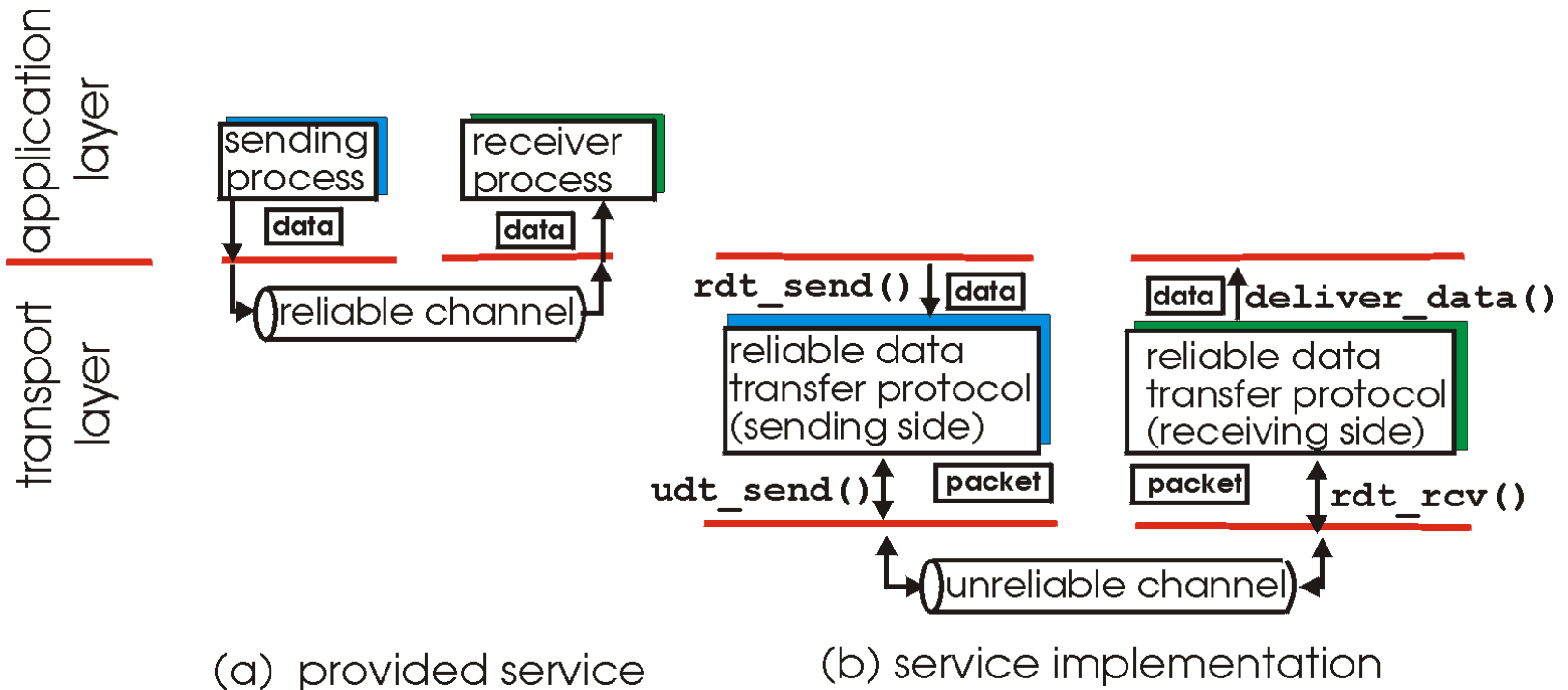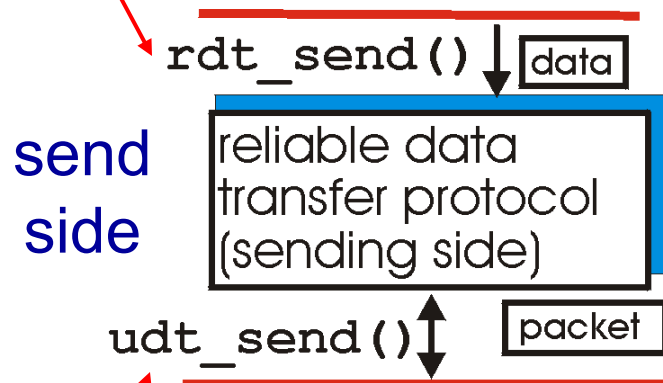(a) provided service

(b) service implementation

# Principles of reliable data transfer

- **important in application, transport, link layers**
  - top-10 list of important networking topics!



(a) provided service          (b) service implementation

# Reliable data transfer: getting started

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

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

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:

- 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 is uniquely determined by next event

state 1

event...
actions...

state 2

# rdt1.0: reliable transfer over a reliable channel

- **underlying channel assued to be perfectly reliable**
  - no bit errors
  - no loss of packets
- **we define separate FSMs for sender, receiver:**
  - sender sends data into underlying channel
  - receiver reads data from underlying channel

Wait for call from above

rdt_send(data)
_____

packet = make_pkt(data)
udt_send(packet)

Wait for call from below

rdt_rcv(packet)
_____
extract (packet,data)
deliver_data(data)

**sender**

**receiver**

# rdt2.0: channel with bit errors

- this time, underlying channel may flip bits in packet
  - employ *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

- this time, underlying channel may flip bits in packet
  - employ *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
  - feedback: control msgs (ACK,NAK) from receiver to sender

# rdt2.0: FSM specification

rdt_send(data)
_____
sndpkt = 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

**sender**

**receiver**

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

**Wait for call from below**

rdt_rcv(rcvpkt) &&
**not**corrupt(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)
_____
$\Lambda$

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

Wait for call from below

rdt_rcv(rcvpkt) && **not**corrupt(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) && **not**corrupt(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!
- sender can't just retransmit: possible **duplicate**

## how can we handle duplicates?

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

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

# rdt2.1: outline

## sender:

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

## receiver:

- twice as many states here, too
- must check if received packet is a duplicate
  - state indicates whether seq # 0 or 1 is expected

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

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

**Wait for call 0 from above**

**Wait for ACK or NAK 0**

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

rdt_rcv(rcvpkt)
&& **not**corrupt(rcvpkt)
&& isACK(rcvpkt)
_____
$\Lambda$

rdt_rcv(rcvpkt)
&& **not**corrupt(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(sndpkt) *

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

* = dup send

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

rdt_rcv(rcvpkt) && **not**corrupt(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) && **not**corrupt(rcvpkt)
&& has_seq1(rcvpkt)
___
extract(rcvpkt,data)
deliver_data(data)
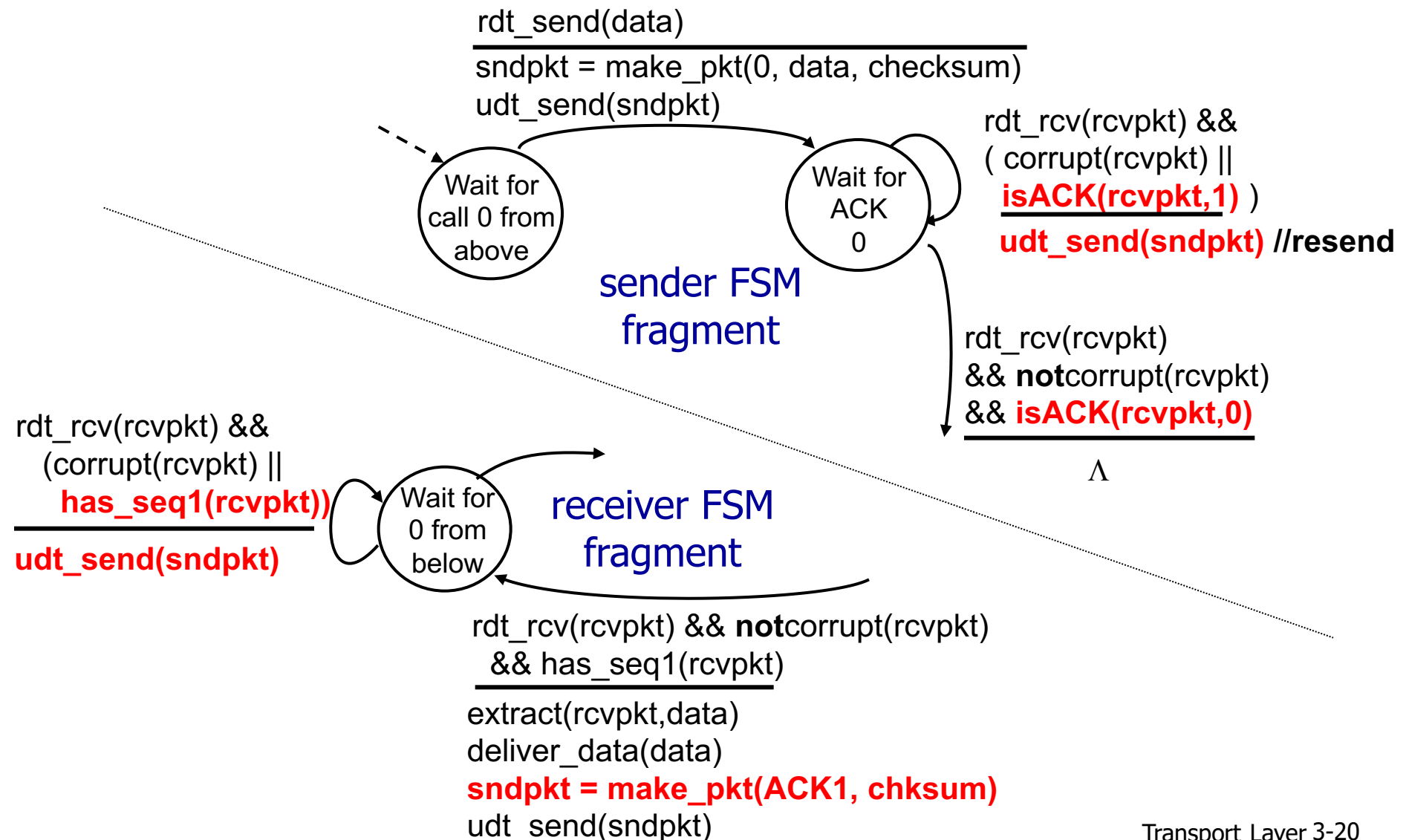sndpkt = make_pkt(ACK, chksum)
udt_send(sndpkt)

* = ACK a dup send

# rdt2.2: a NAK-free protocol

- same functionality as rdt2.1, but uses ACKs only
- instead of a NAK, on receipt of a bad packet, the receiver sends an *ACK for last pkt correctly received*
  - so, receiver ACK must explicitly include seq # of pkt being ACKed
- duplicate ACK received at sender results in same action as NAK: *retransmit current pkt*

Why?
simplifies; eases further extensions; it's what TCP does

# 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) //resend**

( Wait for call 0 from above )

( Wait for ACK 0 )

**sender FSM fragment**

rdt_rcv(rcvpkt)
&& **not**corrupt(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) && **not**corrupt(rcvpkt)
 && has_seq1(rcvpkt)
_____
extract(rcvpkt,data)
deliver_data(data)
**sndpkt = make_pkt(ACK1, chksum)**
udt_send(sndpkt)

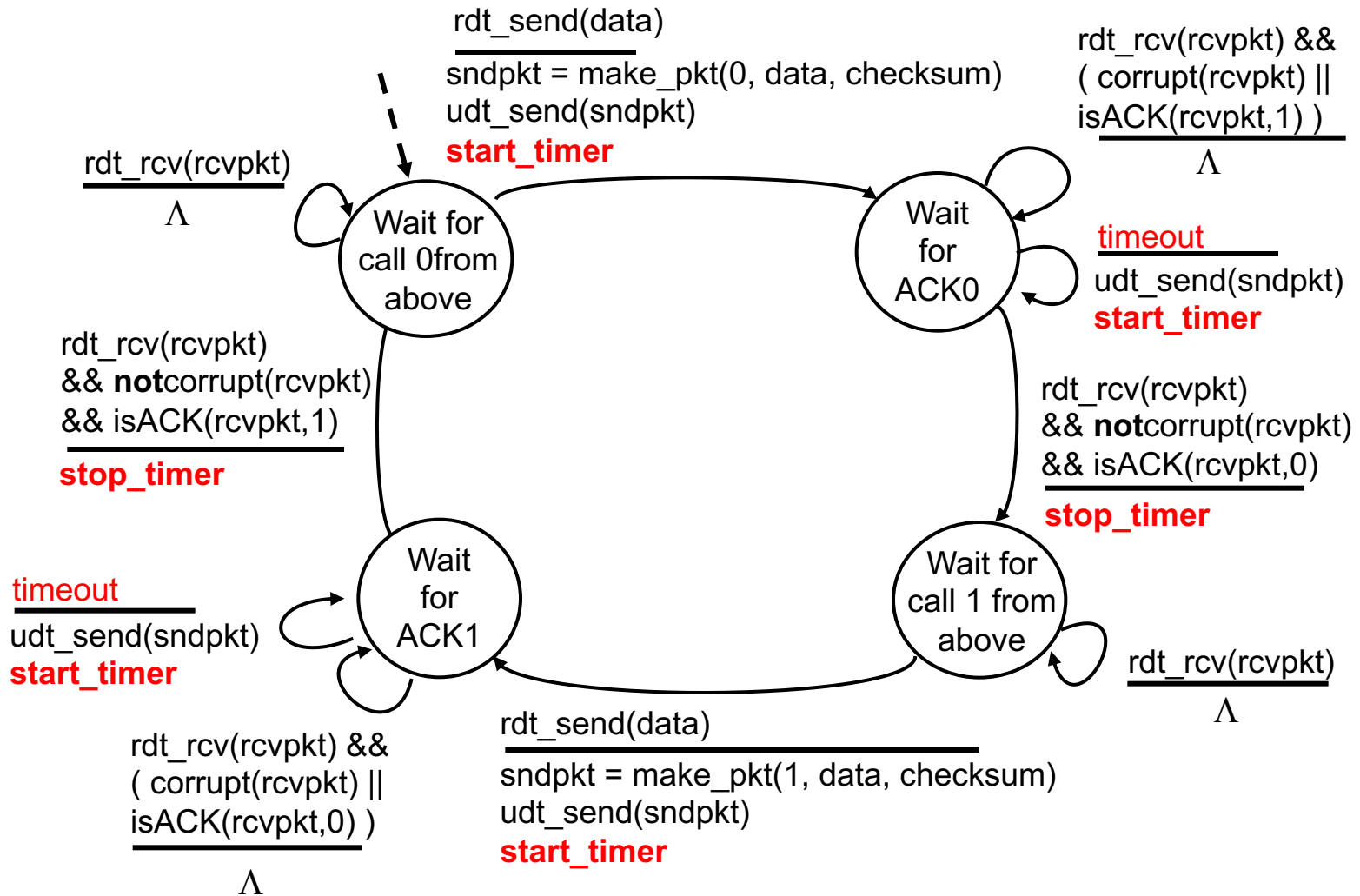# rdt3.0: channels with errors *and* loss

new assumption: as well as corrupting packets, the underlying channel may now also *lose* packets (both data and ACKs)

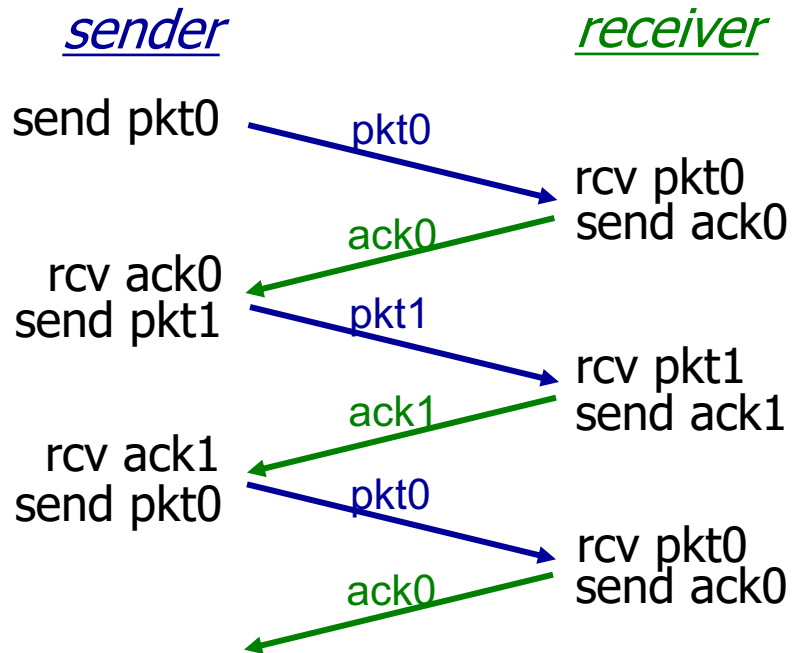- checksum, seq. #, ACKs, retransmissions will be of help ... but not enough

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

- retransmits if no ACK received in this time
- if ACK was just delayed (not lost):
  - retransmission will be received as a duplicate; but seq. #'s already handle this
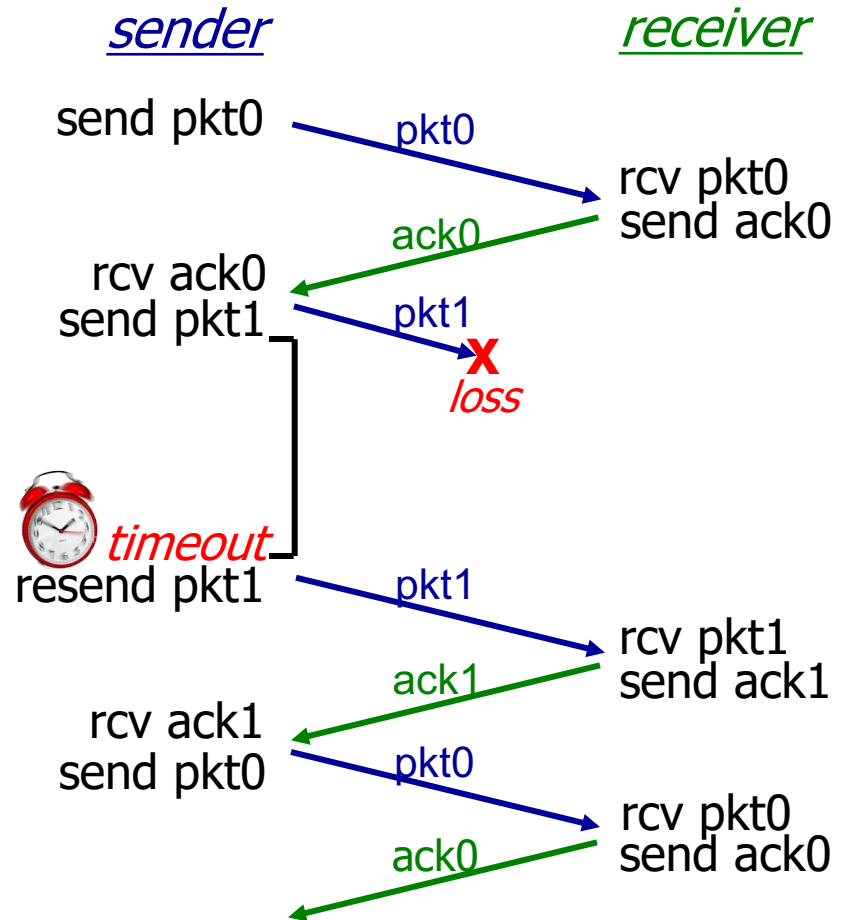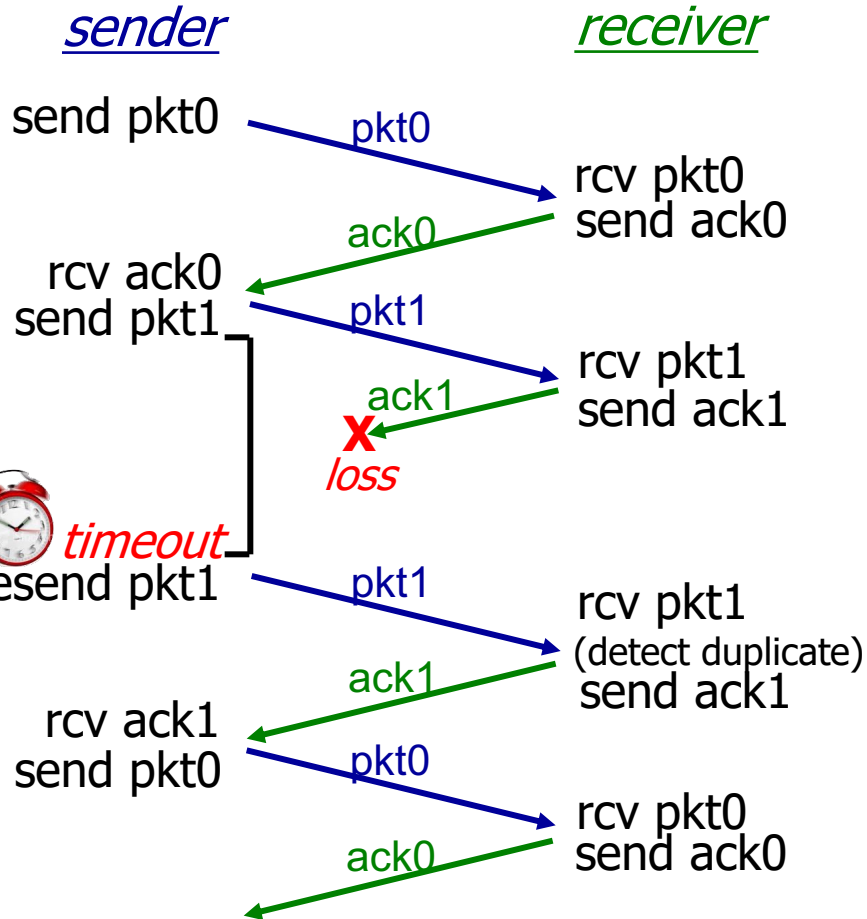- requires *countdown timer*

# rdt3.0 sender



rdt_send(data)

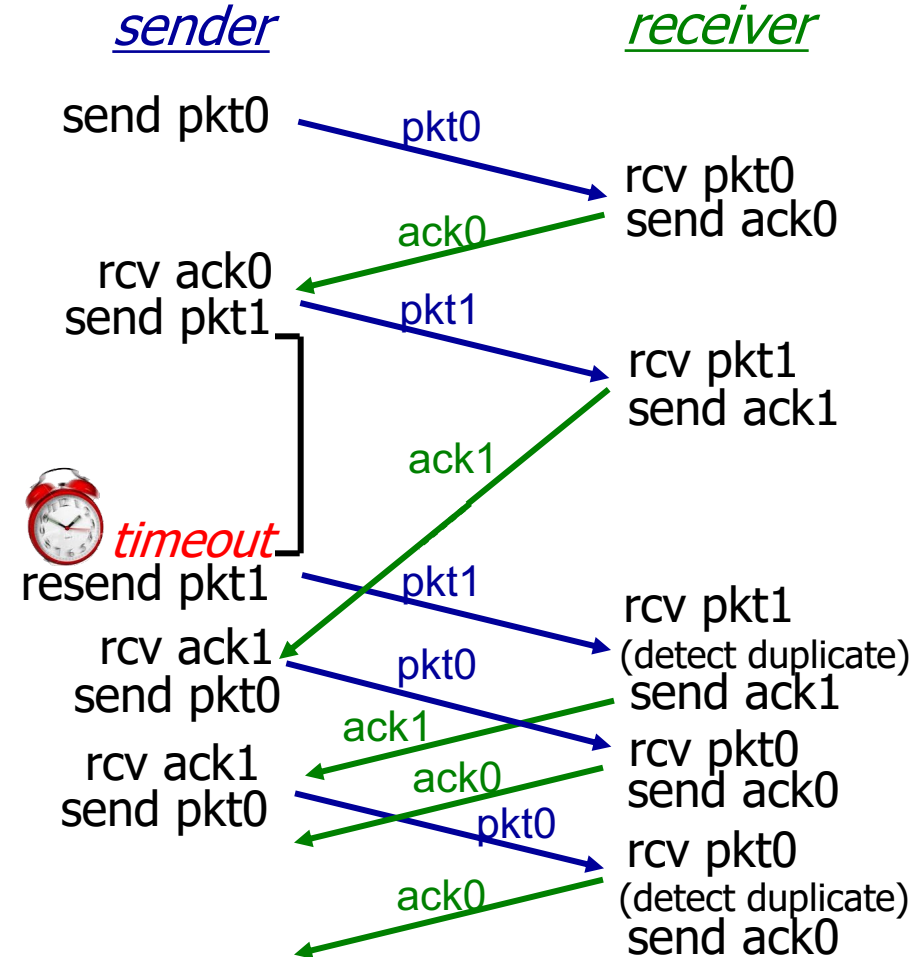sndpkt = make_pkt(0, data, checksum)
udt_send(sndpkt)
**start_timer**

rdt_rcv(rcvpkt) &&
( corrupt(rcvpkt) ||
isACK(rcvpkt,1) )

$\Lambda$

rdt_rcv(rcvpkt)

$\Lambda$

Wait for call 0from above

Wait for ACK0

timeout

udt_send(sndpkt)
**start_timer**

rdt_rcv(rcvpkt)
&& **not**corrupt(rcvpkt)
&& isACK(rcvpkt,1)

**stop_timer**

rdt_rcv(rcvpkt)
&& **not**corrupt(rcvpkt)
&& isACK(rcvpkt,0)

**stop_timer**

timeout

udt_send(sndpkt)
**start_timer**

Wait for ACK1

Wait for call 1 from above

rdt_rcv(rcvpkt)

$\Lambda$

rdt_rcv(rcvpkt) &&
( corrupt(rcvpkt) ||
isACK(rcvpkt,0) )

$\Lambda$

rdt_send(data)

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

# rdt3.0 in action



(a) no packet loss

(b) packet loss

# rdt3.0 in action

**(c) ACK loss**

sender

receiver

send pkt0
→ pkt0 →
rcv pkt0
send ack0

rcv ack0
← ack0 ←
send pkt1
→ pkt1 →
rcv pkt1
send ack1

ack1
X loss

timeout
resend pkt1
→ pkt1 →
rcv pkt1
(detect duplicate)
send ack1

rcv ack1
← ack1 ←
send pkt0
→ pkt0 →
rcv pkt0
send ack0

← ack0 ←

**(d) premature timeout (=delayed ACK)**

sender

receiver

send pkt0
→ pkt0 →
rcv pkt0
send ack0

rcv ack0
← ack0 ←
send pkt1
→ pkt1 →
rcv pkt1
send ack1

timeout
ack1

resend pkt1
→ pkt1 →
rcv pkt1
(detect duplicate)
send ack1

rcv ack1
send pkt0
→ pkt0 →
rcv pkt0
send ack0

rcv ack1
← ack1 ←
send pkt0
← ack0 ←
→ pkt0 →
rcv pkt0
(detect duplicate)
send ack0

← ack0 ←

(please work though this yourself and convince yourself it works as it should!)

# Performance of rdt3.0

- rdt3.0 is **correct**, but performance is terrible
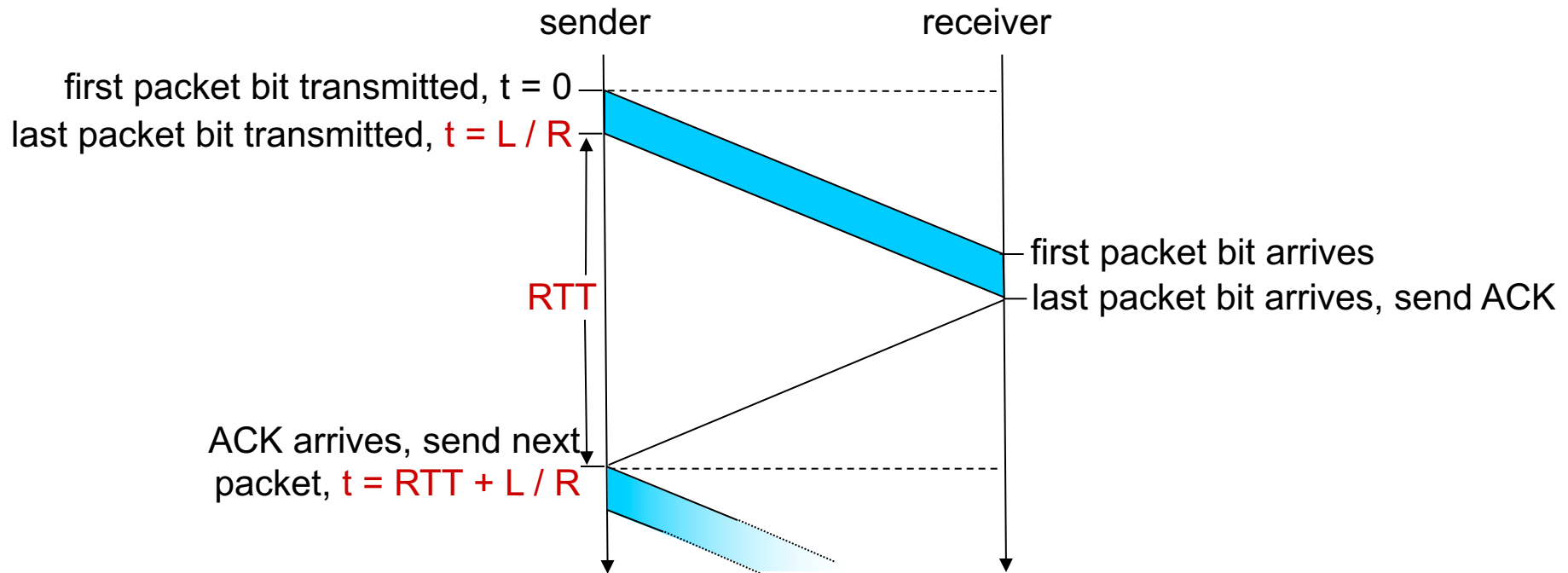- e.g.: 1 Gbps link, 15 ms prop. delay, 1KByte packet:

$$D_{trans} = \frac{L}{R} = \frac{8000 \text{ bits}}{10^9 \text{ bits/sec}} = 8 \text{ microsecs}$$

- $U_{sender}$: *utilization* (i.e., fraction of time sender busy sending)

(n.b. RTT=30ms)

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

- so, 1KByte pkt every 30.008 msec
  - = 33kB/sec (approx) over a 1 Gbps link!
- protocol severely limits potential of network!

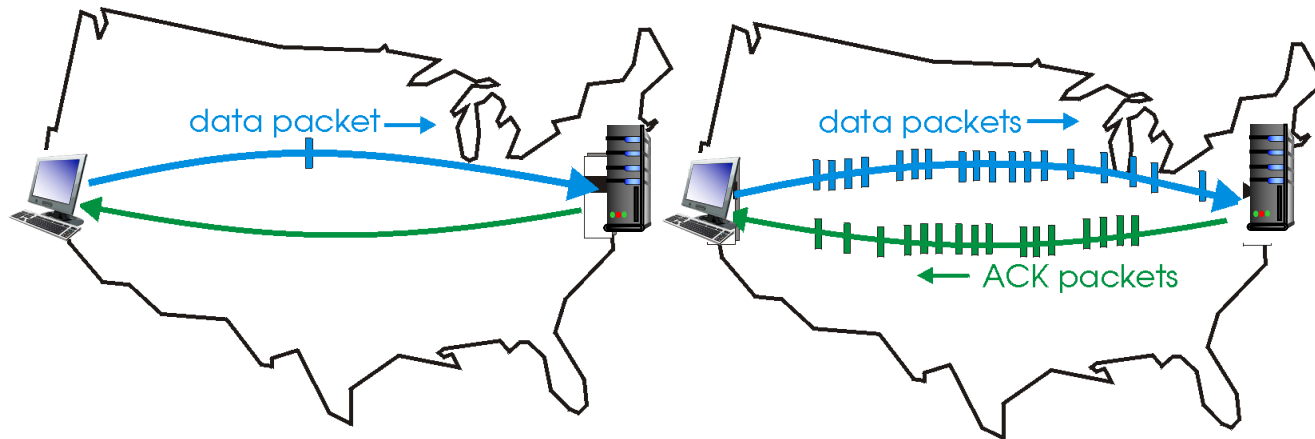# rdt3.0: stop-and-wait performs poorly

sender                    receiver

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

first packet bit arrives
last packet bit arrives, send ACK

RTT

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, yet-to-be-acknowledged, pkts to be "in flight" simultaneously

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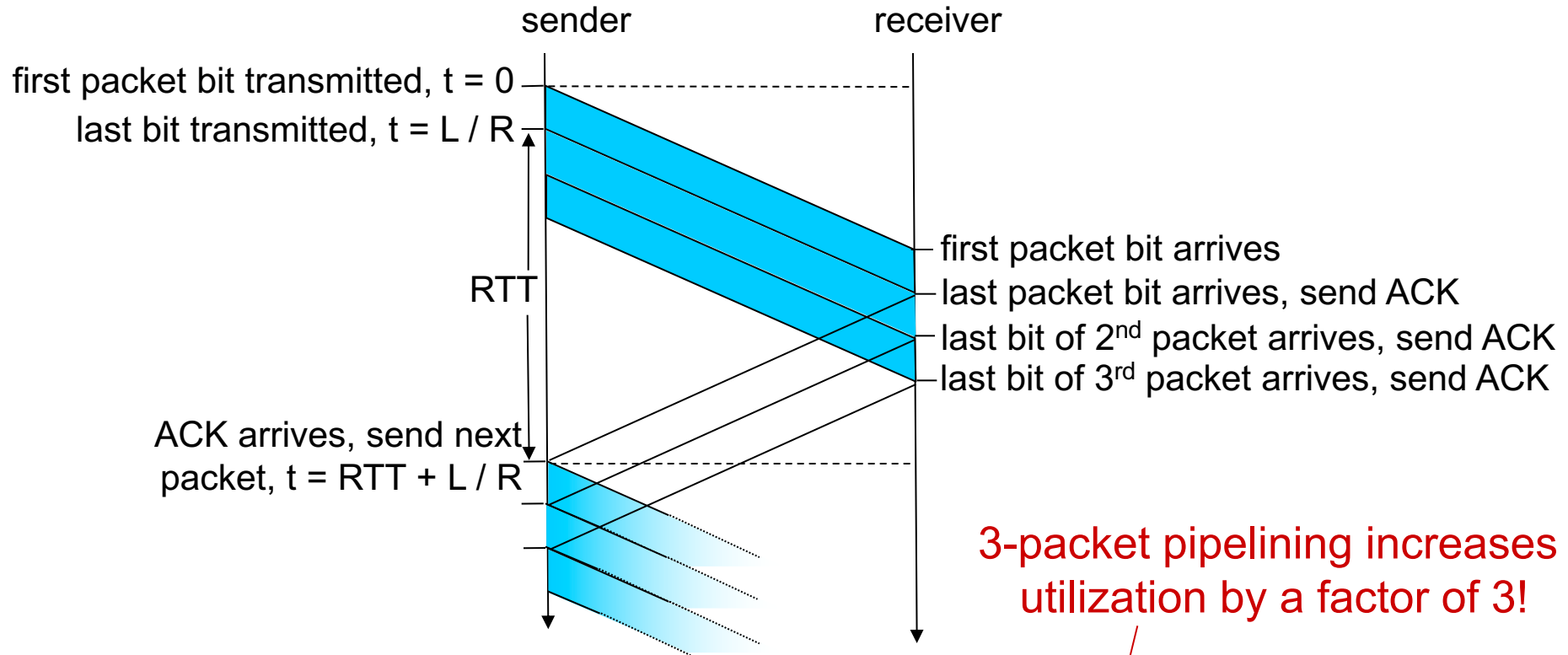- needs buffering at sender and/or receiver



(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 2$^{nd}$ packet arrives, send ACK

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

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

3-packet pipelining increases utilization by a factor of 3!

$$U_{sender} = \frac{3L/R}{RTT + L/R} = \frac{.0024}{30.008} = 0.00081$$

# Pipelined protocols: overview

## Go-back-N:

- sender may push up to N unacked packets into the pipeline ("window")
- receiver only sends *cumulative ack*
  - doesn't ack packet if there's a gap
- sender maintains a timer for oldest unacked packet
  - when timer expires, retransmit *all* unacked packets

## (Selective Repeat:

- 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 that unacked packet)

# Chapter 3 outline

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 (omit: RTT estimation)
- reliable data transfer
- flow control
- connection management

3.6 principles of congestion control

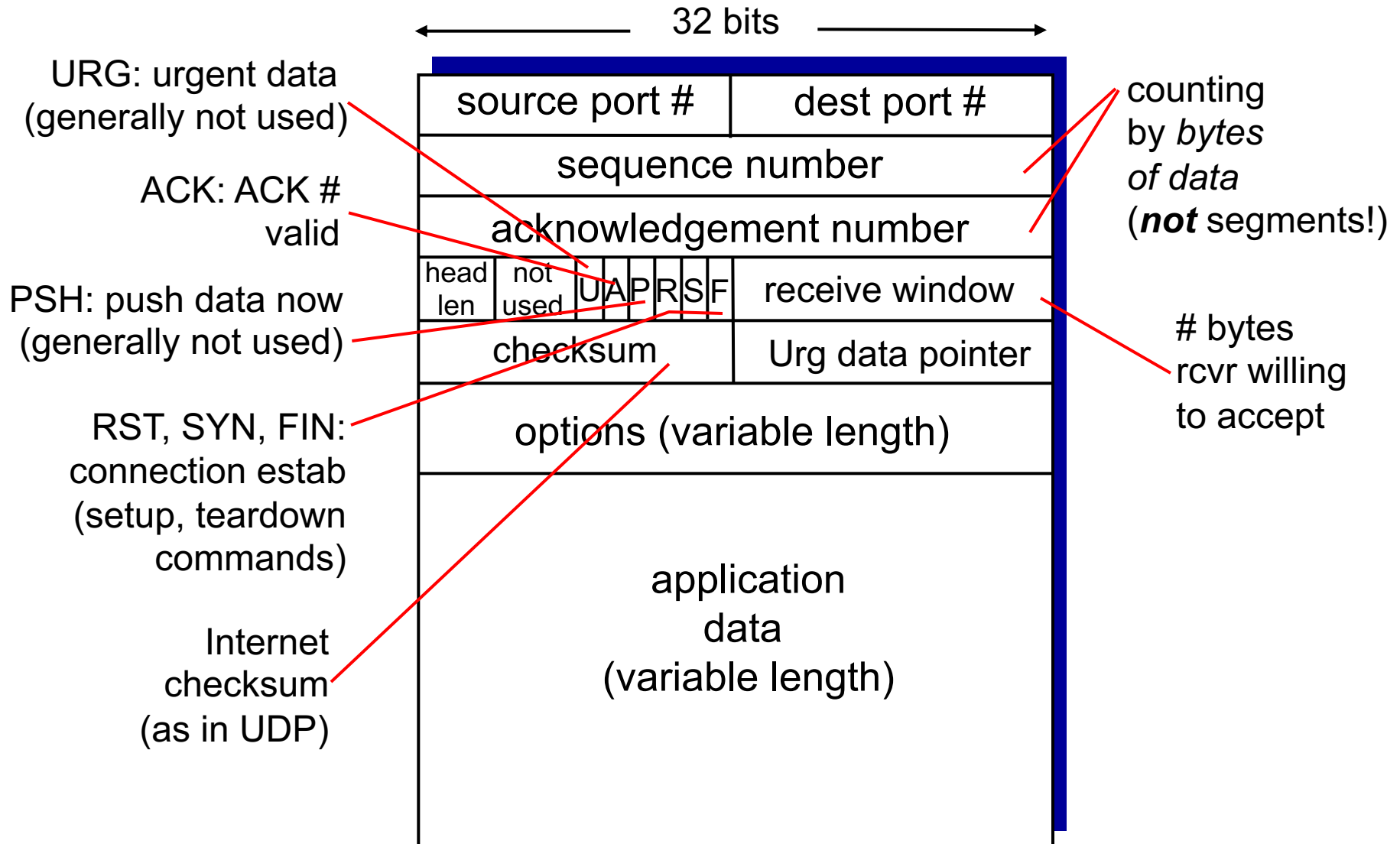3.7 TCP congestion control

# 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

- **full duplex data:**
  - bi-directional data flow in same connection
- **connection-oriented:**
  - handshaking (exchange of control msgs) inits sender, receiver state before data exchange
- **flow controlled:**
  - sender will not overwhelm receiver

# TCP segment structure

32 bits

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

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

options (variable length)

application
data
(variable length)

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)

counting
by *bytes*
of data
(*not* segments!)

# bytes
rcvr willing
to accept

# TCP seq. numbers, ACKs

sequence numbers:
- "index" of the first byte of a segment's data in the ongoing byte stream

acknowledgements:
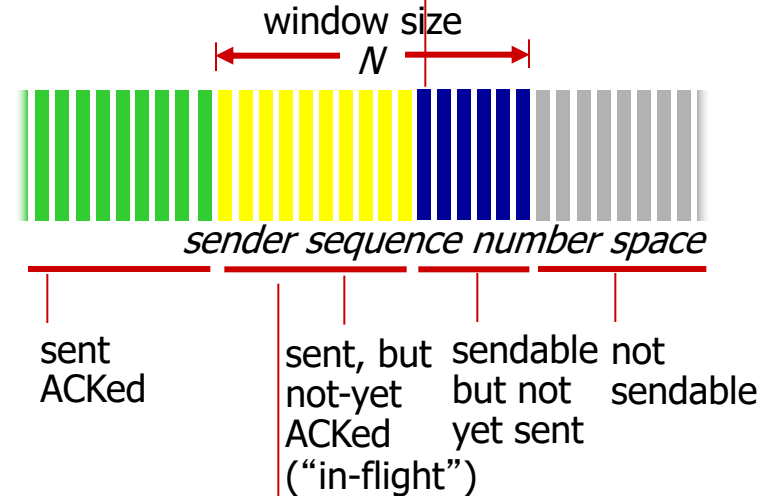- seq # of next byte expected from send side
- cumulative ACK

Q: how does receiver handle out-of-order segments?
- A: TCP spec doesn't say, - up to implementer

outgoing segment from sender

| source port # | dest port # |
|---|---|
| sequence number | |
| acknowledgement number | |
| | rwnd |
| checksum | urg pointer |

window size
$N$

*sender sequence number space*

sent ACKed

sent, but not-yet ACKed ("in-flight")

sendable but not yet sent

not sendable

incoming segment to sender

| source port # | dest port # |
|---|---|
| sequence number | |
| acknowledgement number | |
| A | rwnd |
| checksum | urg pointer |

# TCP round trip time, timeout

Q: how best to set TCP timeout value?

- *too short:* premature timeout, unnecessary retransmissions
- *too long:* slow reaction to segment loss
- certainly set it longer than RTT
  - but RTT varies...

Q: how to estimate RTT?

- **SampleRTT**: measured time from segment transmission until ACK receipt
- **SampleRTT** will vary; ideally, we want a "smoother" estimated RTT average
- so, incorporate several *recent* measurements, not just current **SampleRTT**

# Chapter 3 outline

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
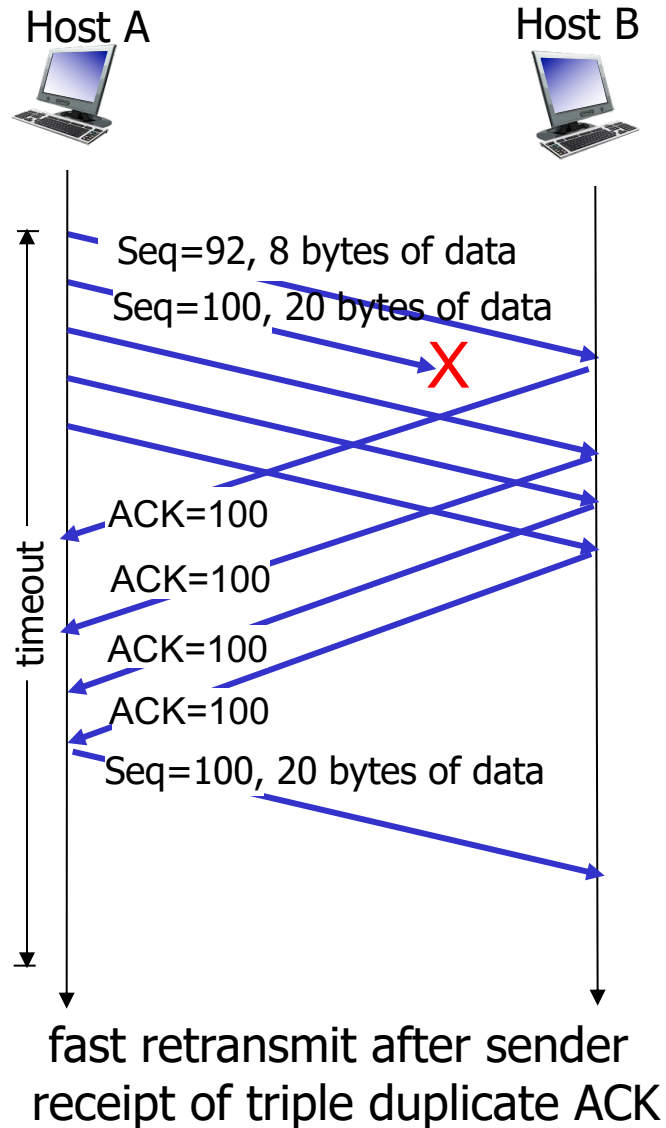
3.7 TCP congestion control

# TCP fast retransmit (a quick example of one TCP optimization)

- time-out period is often relatively "long":
  - so, long delay before resending lost packet
- as we know, we detect corrupt/lost segments via duplicate ACKs
  - sender often sends many segments back-to-back
  - if segments are lost, gaps will likely cause many duplicate ACKs

*TCP fast retransmit*

if sender receives 3 ACKs for same data ("triple duplicate ACKs"),

=> immediately resend unACKed segment with smallest seq #

- likely that unacked segment was lost, so don't wait for timeout

# TCP fast retransmit

Host A                    Host B

Seq=92, 8 bytes of data

Seq=100, 20 bytes of data

X

ACK=100

ACK=100

ACK=100

ACK=100

Seq=100, 20 bytes of data

timeout

fast retransmit after sender
receipt of triple duplicate ACK

# Chapter 3 outline

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 (OMITTED)
- connection management (OMITTED)

3.6 principles of congestion control (OMITTED)

3.7 TCP congestion control (OMITTED)

# Brief word on the omitted topics in section 3.5

- **TCP connection management**
  - establishing the basis for communication – two sides exchange initial state, and each becomes aware that the other is connected
  - employs 3-way (not 2-way) handshake
- **TCP flow control**
  - preventing a sender swamping a receiver
  - use of rwnd field

- **principles of congestion and TCP congestion control**
  - different from flow control – congestion control is about avoiding swamping the network
  - "slow start" mechanism
  - big topic!

# Summary

- we've seen the development of an abstract protocol that ensures 100% reliability in the face of corrupted and/or lost packets
- TCP builds on these ideas and adds many practical features
  - connection management
  - flow and congestion control
  - fast retransmit

- there's a lot to TCP that we have not yet seen!
  - see sections 3.5, 3.6 and 3.7 of the book