

SC250 Computer Networking I

Transport Layer Principles

Prof. Matthias Grossglauser

School of Computer and Communication Sciences
EPFL

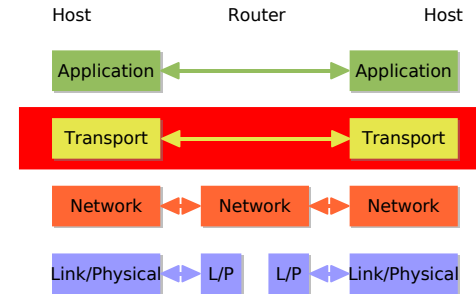
<http://lcawww.epfl.ch>



1

Principles of Reliable Data Transfer

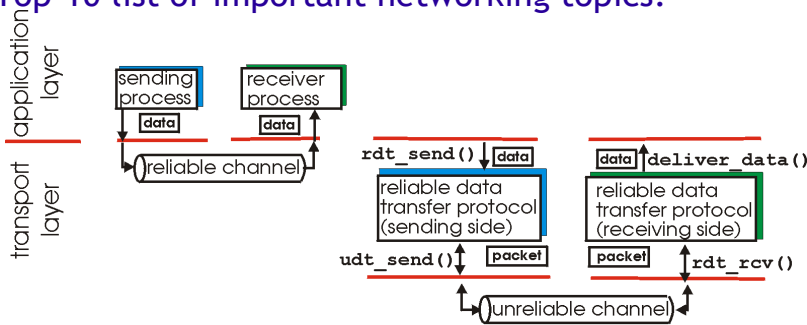
- Reliable transfer
 - Getting data without any alteration (bit errors, loss, duplication,...) from a sender to a receiver
 - Achieving this over an unreliable channel, i.e., a channel that **does** alter information
 - We will need two-way unreliable channel, even to achieve one-way reliable transfer
- Internet:
 - IP: unreliable channel
 - TCP: reliable service, available through socket API



2

Principles of Reliable Data Transfer

- Applies to transport layer (TCP), but also to some link-layer protocols
- Top-10 list of important networking topics!



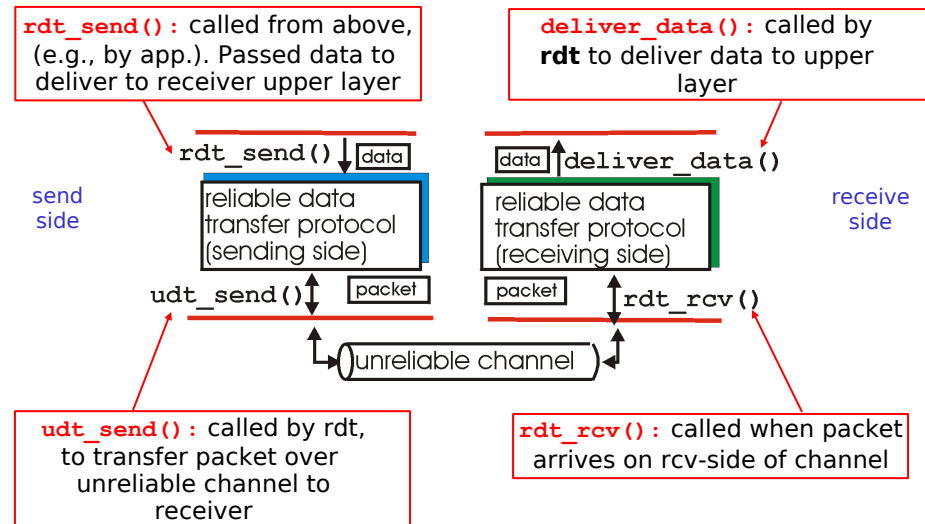
(a) provided service

(b) service implementation

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

3

Reliable Data Transfer: Getting Started



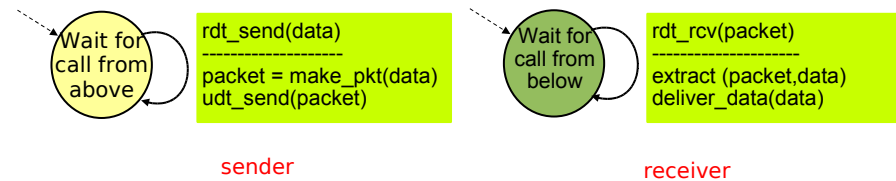
4

Reliable Data Transfer: Getting Started

- 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

RDT1.0: Reliable Transfer over 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



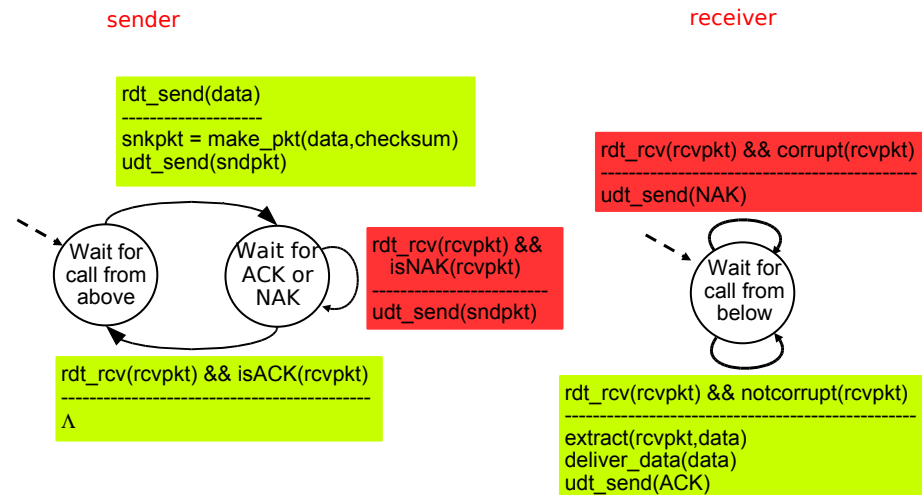
5

6

RDT2.0: Channel with Bit Errors

- Underlying channel may flip bits in packet
 - checksum to detect bit errors: send (m, c) , where $c=h(m)$ is the checksum; receiver checks whether $c=h(m)$
- The question: how to recover from errors:
 - Automatic repeat request (ARQ): control information from receiver back to sender
 - 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, ACK/NAK control messages

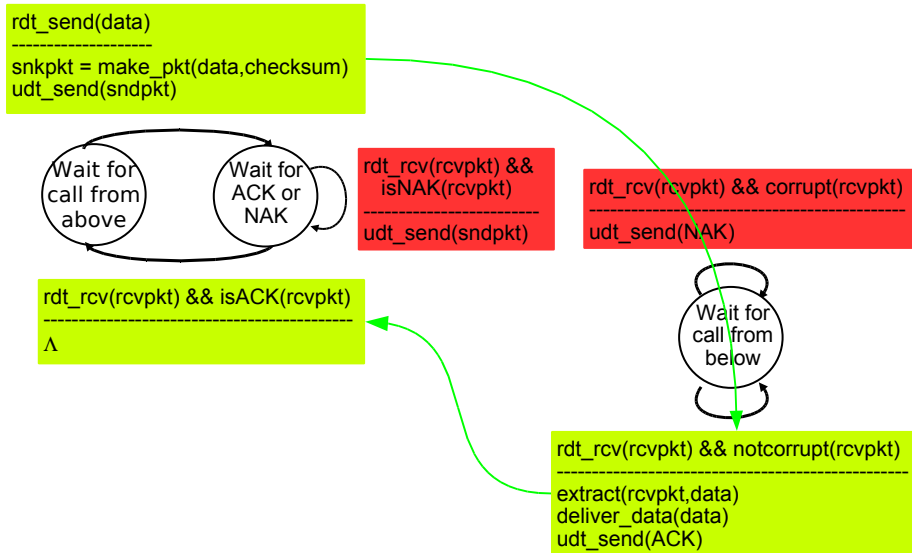
RDT2.0: FSM Specification



7

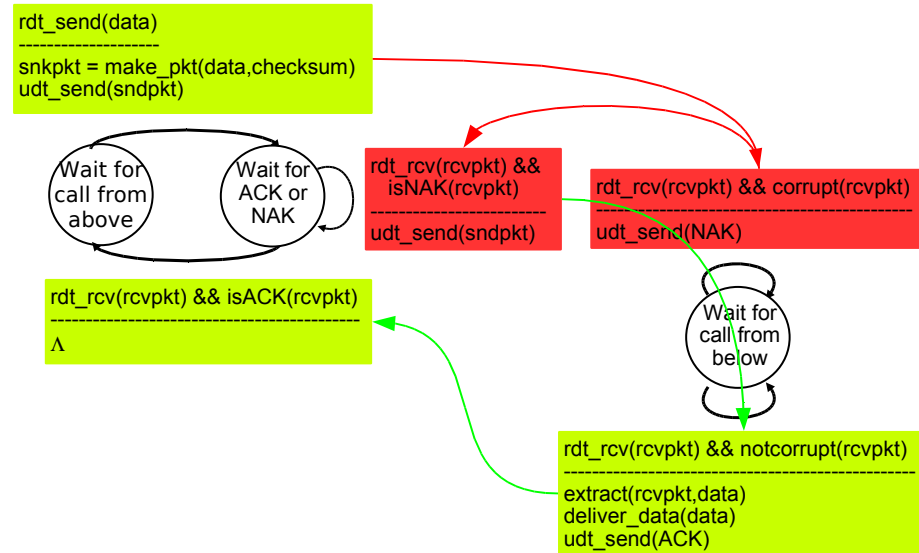
8

RDT2.0: Operation with no Errors



9

RDT2.0: Error Scenario



10

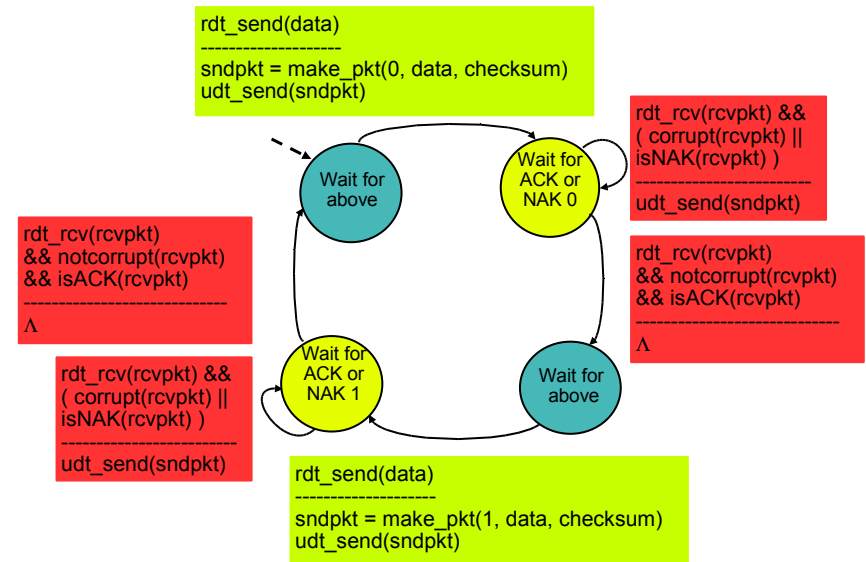
RDT2.0 has a Fatal Flaw!

- What happens if ACK/NAK corrupted?
 - sender doesn't know what happened at receiver!
 - can't just retransmit: possible duplicate
- What to do?
 - sender ACKs/NAKs receiver's ACK/NAK? What if sender ACK/NAK lost?
 - retransmit, but this might cause retransmission of correctly received pkt!
- Handling duplicates:
 - sender adds sequence number to each pkt
 - sender retransmits current pkt if ACK/NAK garbled
 - receiver discards (doesn't deliver up) duplicate pkt

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

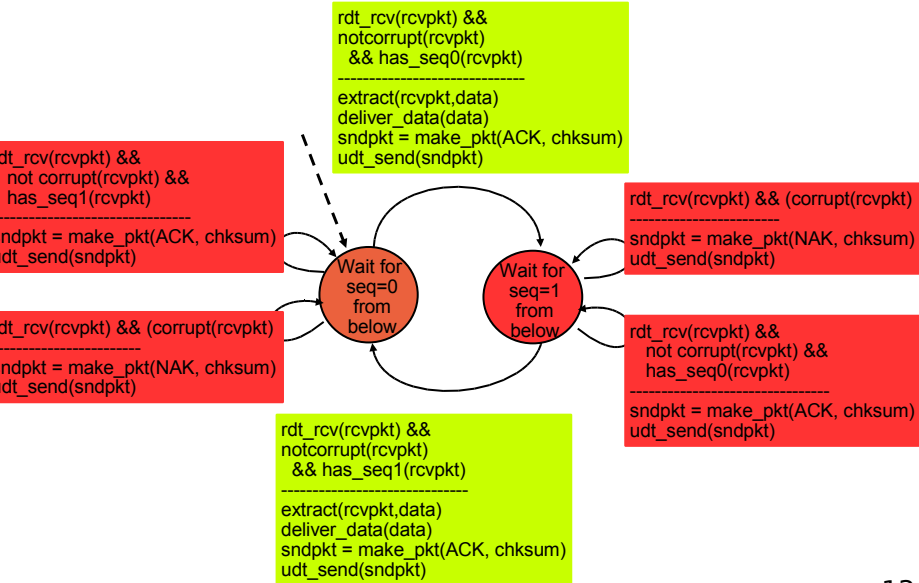
11

RDT2.1 Sender: Handles Garbled ACK/NAKs



12

RDT2.1, Receiver: Handles Garbled ACK/NAKs



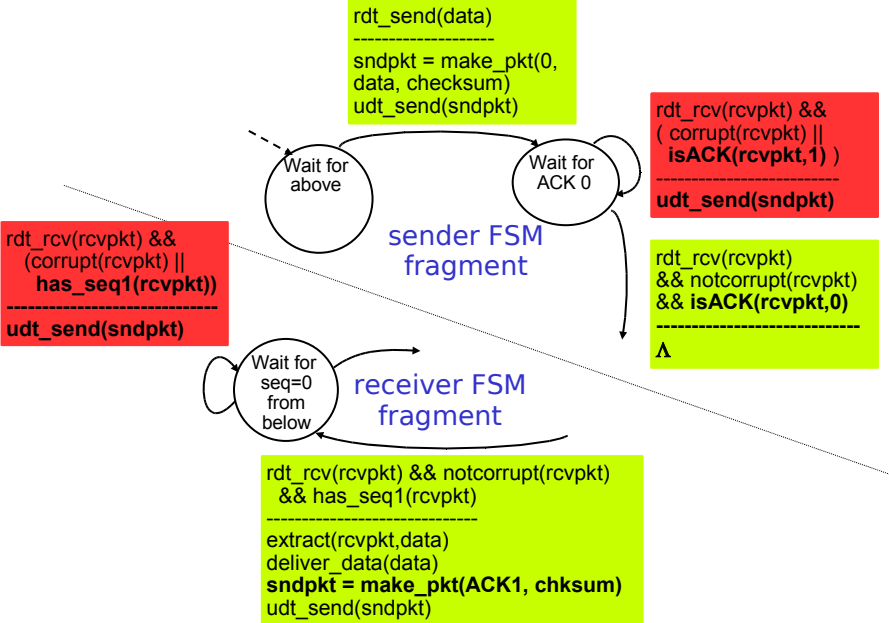
RDT2.1: Discussion

- Sender:
 - Seq # added to pkt
 - Two seq. #'s (0,1) will suffice. Why?
 - Must check if received ACK/NAK corrupted
 - Twice as many states
 - state must “remember” whether “current” pkt has 0 or 1 seq. #
- 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

RDT2.2: A NAK-Free Protocol

- Same functionality as rdt2.1, using ACKs (NAKs) only
- Instead of NAK, receiver sends ACK for last pkt received OK
 - receiver must *explicitly* include seq # of pkt being ACKed
- Duplicate ACK at sender results in same action as NAK: *retransmit current pkt*

RDT2.2 Sender and Receiver Fragments

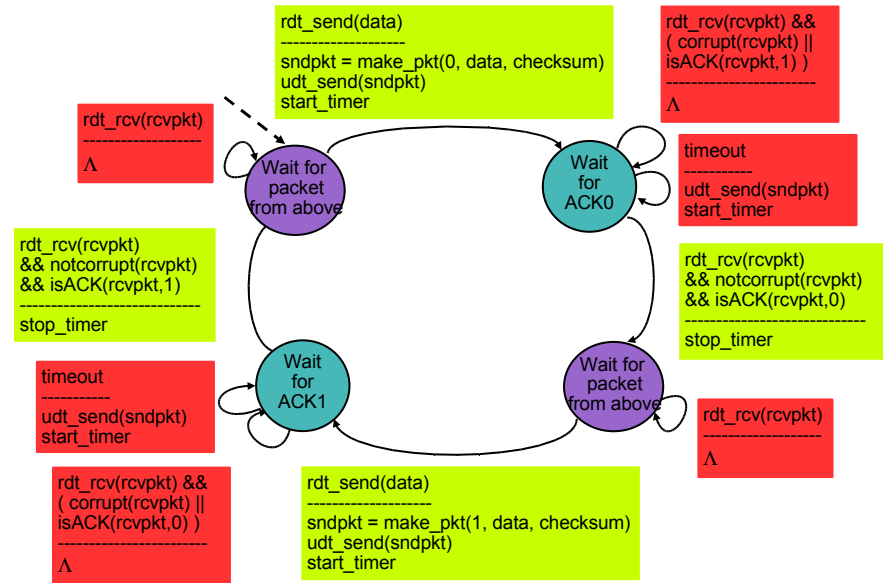


RDT3.0: Channels with Errors and Loss

- New assumption:** underlying channel can also lose packets (data or ACKs)
 - checksum, seq. #, ACKs, retransmissions will be of help, but not enough
- Q: how to deal with loss?**
 - sender waits until certain that data or ACK lost, then retransmits
 - drawbacks?
- 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

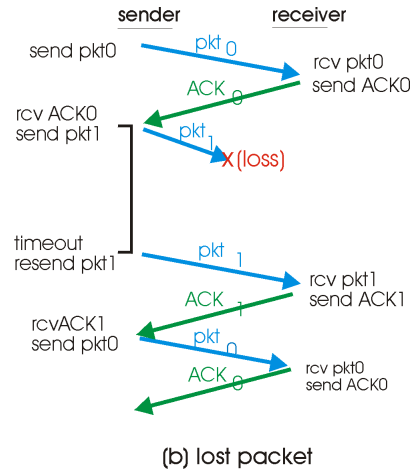
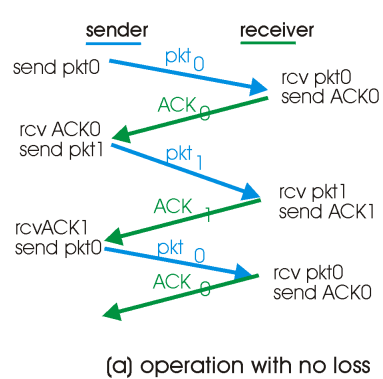
17

RDT3.0 Sender

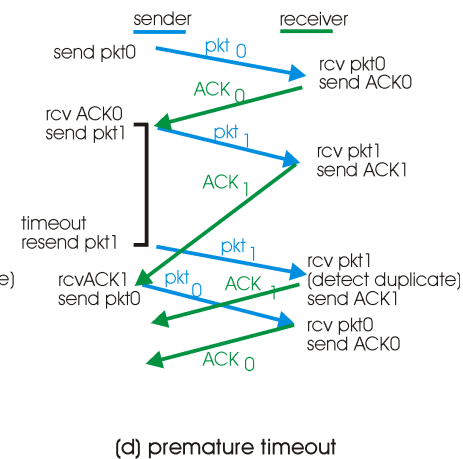
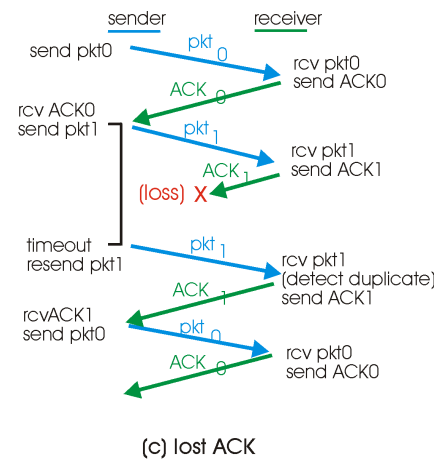


18

RDT3.0 in Action



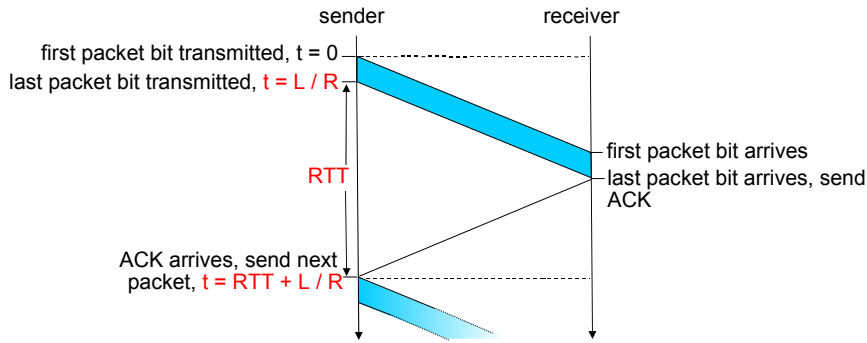
RDT3.0 in Action



19

20

RDT3.0: Stop-and-Wait Operation



$$U_{sender} = \frac{L/R}{L/R + RTT}$$

Performance

- rdt3.0 works, but lousy performance
- example: 1 Gbps link, 15 ms e2e prop. delay, 1KB packet:

$$T_{transmit} = \frac{L \text{ (packet length in bits)}}{R \text{ (transmission rate, bps)}} = \frac{8\text{kb/pkt}}{1\text{e9 b/sec}} = 8 \text{ microsec}$$

$$U_{sender} = \frac{L/R}{RTT + L/R} = \frac{0.008}{30.008} = 0.00027$$

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

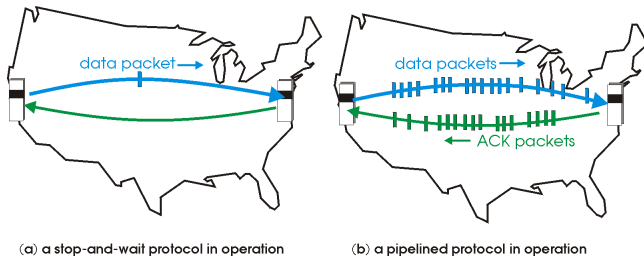
21

22

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



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

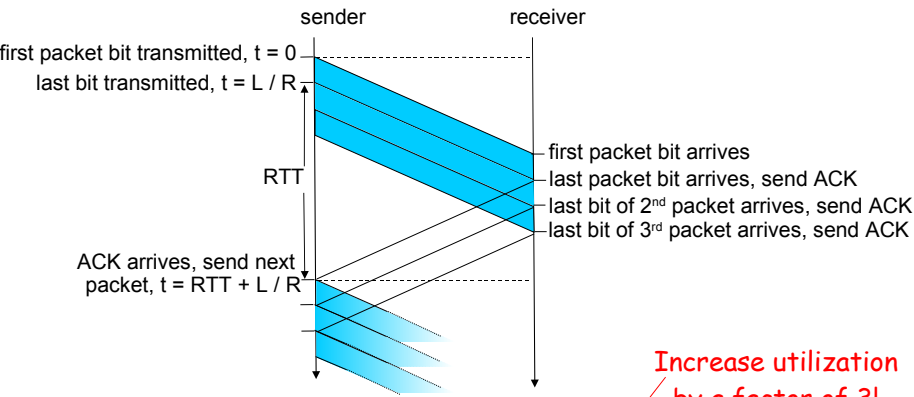
23

Concepts for Pipelined Protocols

- Acknowledgements:**
 - Selective ACK(n), NAK(n):
 - Says “packet n received/not received”
 - Cumulative cACK(n) = {...ACK(n-2),ACK(n-1),ACK(n)}
 - Says “every packet up to and including n received”
- Timer:**
 - Should be larger than RTT, otherwise risk of malfunction
- Sequence number space:**
 - Set of identifiers for packets (stop-and-go: {0,1})
- Sliding window:**
 - Interval of sequence number space, advances upon receipt
 - Budget of in-flight data for pipelined operation
 - Note: sequence # space imposes constraint on window size

24

Pipelining: Increased Utilization

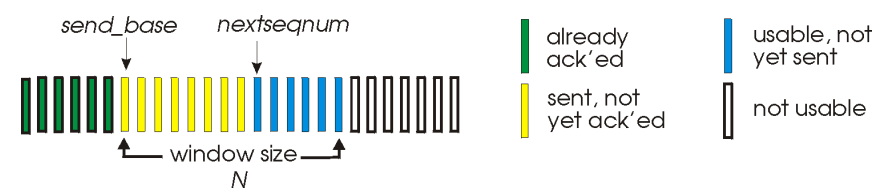


$$U_{sender} = \frac{3L/R}{L/R + RTT}$$

Go-Back-N (GBN)

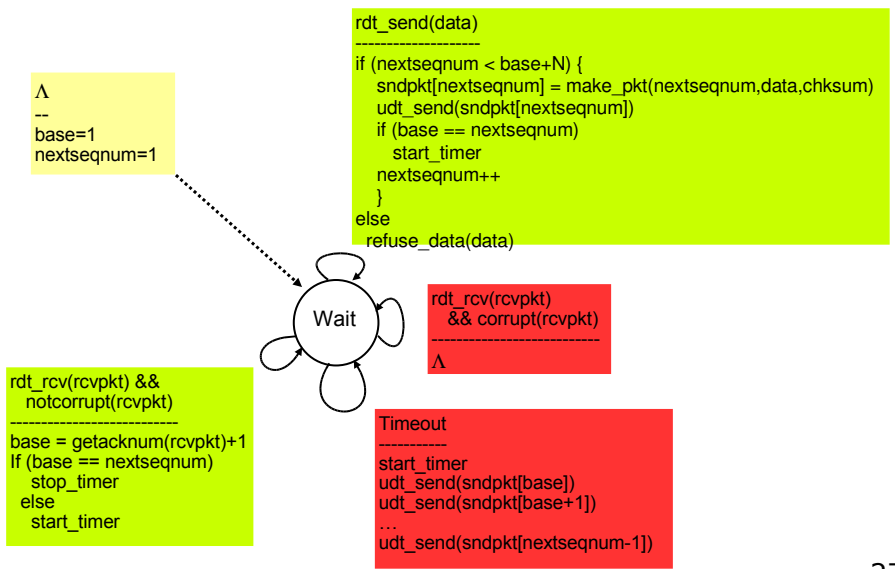
Sender:

- Sequence number (seq #) in packet header
- Window (budget) of up to N consecutive unacknowledged pkts allowed

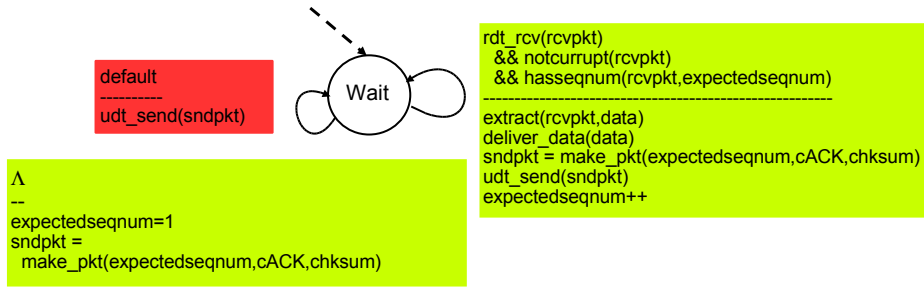


- cACK(n): ACKs all pkts up to, including seq # n - "cumulative ACK"
 - Says: "we are done with everything before (n+1)"
 - may receive duplicate ACKs (see receiver)
- timer for each in-flight pkt (i.e., sent, but not yet received)
- timeout(n): retransmit pkt n and all higher seq # pkts in window

GBN, Sender: Extended FSM

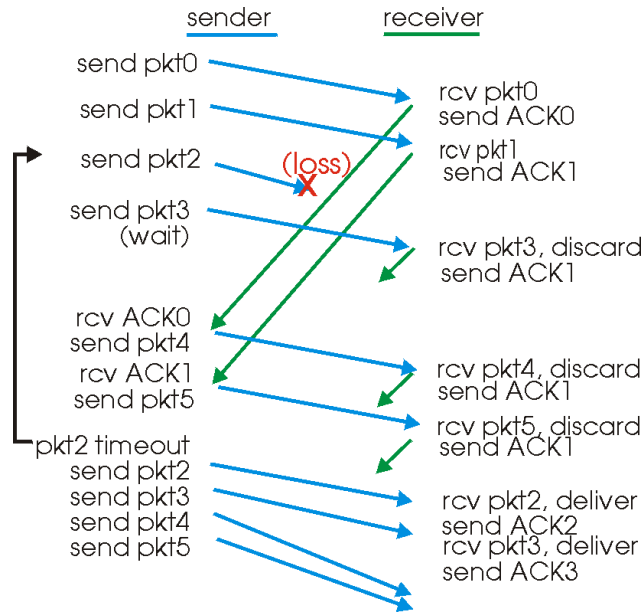


GBN, Receiver: Extended FSM



- ACK-only: always send cACK for correctly-received pkt with highest in-order seq #
 - may generate duplicate cACKs
 - need only remember expectedseqnum
- Out-of-order packets:
 - Discard (don't buffer) -> no receiver buffering!
 - Re-ACK pkt with highest in-order seq #

GBN in Action



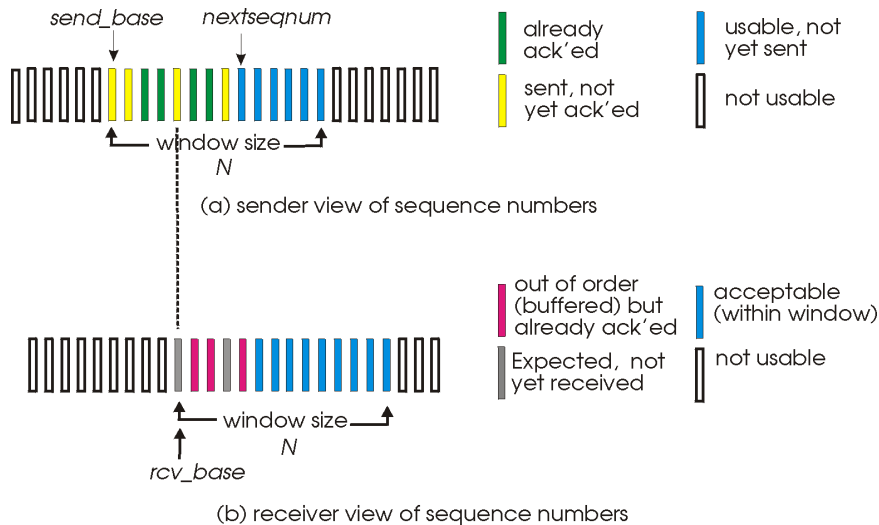
29

Selective Repeat (SR)

- Drawback of GBN:
 - Entire window gets retransmitted due to a single error
 - Especially costly with large window
- 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, unacknowledged pkts

30

Selective Repeat: Sender, Receiver Windows



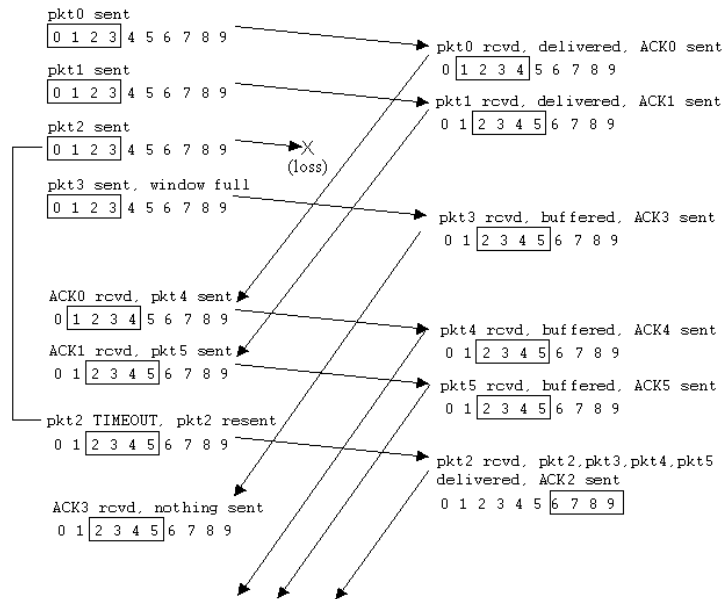
31

Selective Repeat

sender	receiver
<p>data from above :</p> <ul style="list-style-type: none"> if next available seq # in window, send pkt 	<p>pkt n in [rcvbase, rcvbase+N-1]</p> <ul style="list-style-type: none"> send ACK(n) out-of-order: buffer in-order: deliver (also deliver buffered, in-order pkts), advance window to next not-yet-received pkt
<p>timeout(n):</p> <ul style="list-style-type: none"> resend pkt n, restart timer 	<p>pkt n in [rcvbase-N, rcvbase-1]</p> <ul style="list-style-type: none"> ACK(n)
<p>ACK(n) in [sendbase, sendbase+N]:</p> <ul style="list-style-type: none"> mark pkt n as received if n smallest unACKed pkt, advance window base to next unACKed seq # 	<p>otherwise:</p> <ul style="list-style-type: none"> ignore

32

Selective Repeat in Action

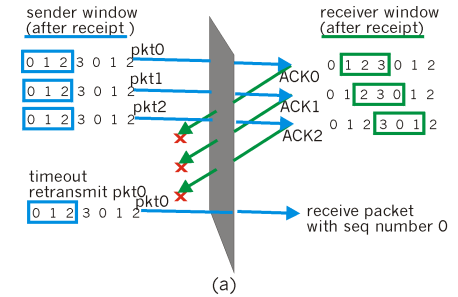


33

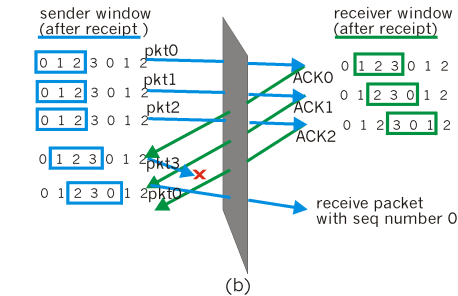
SR: Sizes of Window and Sequence # Space

Example:

- seq #'s: 0, 1, 2, 3
- window size=3
- receiver sees no difference in two scenarios!
- incorrectly passes duplicate data as new in (a)



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



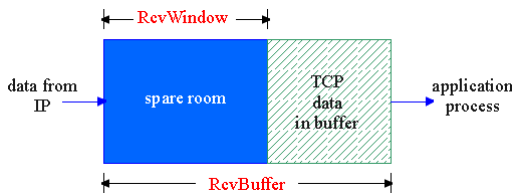
34

Flow Control

- Receive side of TCP connection has a receive buffer:

flow control

sender won't overflow receiver's buffer by transmitting too much, too fast

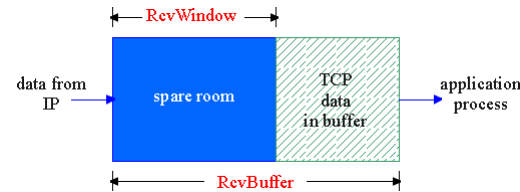


- Speed-matching service: matching the send rate to the receiving app's drain rate

- Application process may be slow at reading from buffer

35

Flow Control: How it Works



- Rcvr advertises spare room by including value of **RcvWindow** in segments
- Sender limits unACKed data to **RcvWindow**
 - guarantees receive buffer doesn't overflow

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

- spare room in buffer = **RcvWindow**
- = **RcvBuffer - [LastByteRcvd - LastByteRead]**

36

Summary: Transport Layer Principles

- Transport layer functionality
 - Types of errors
 - Corrupted packets (bit errors) -> checksums
 - Packet loss (queue overflow) -> sequence numbers
 - Not considered: reordering (both of data and ACK packets)
 - Reliable delivery
 - ARQ (Automatic Repeat Request): ACKs, NAKs, cACKs
 - Pipelining for performance
 - Flow control
 - Matching sender speed to receiver
 - Note: these are all end-to-end issues; another issue is congestion control, i.e., how to share the network
- TCP: main reliable transport protocol in TCP/IP