

Computer Networking

Slide Set 3

Andrew W. Moore

Andrew.Moore@cl.cam.ac.uk

Topic 5 – Transport

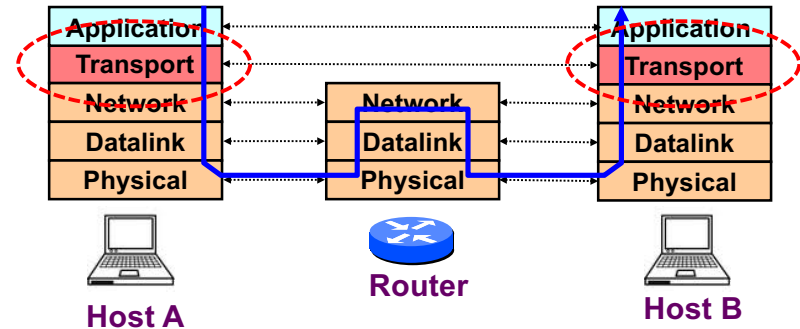
Our goals:

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

2

Transport Layer

- Commonly a layer **at end-hosts**, between the application and network layer



3

Why a transport layer?

- IP packets are addressed to a host but end-to-end communication is between application/processes/tasks at hosts
 - Need a way to decide which packets go to which applications (*more multiplexing*)

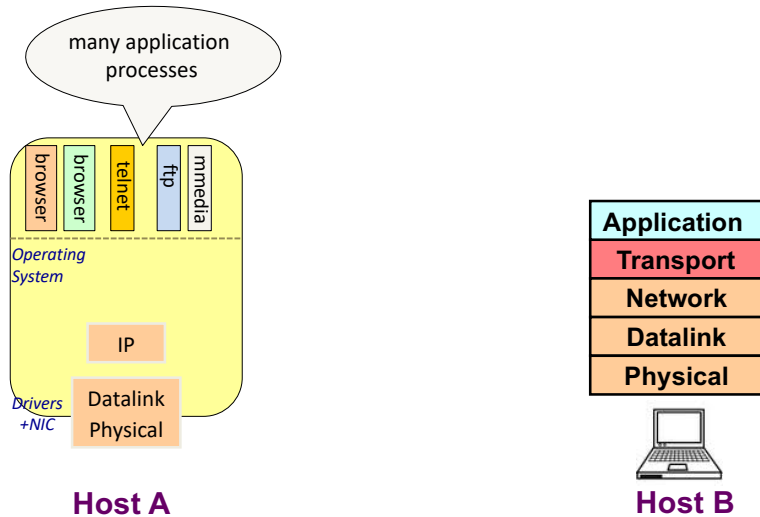
4

Why a transport layer?



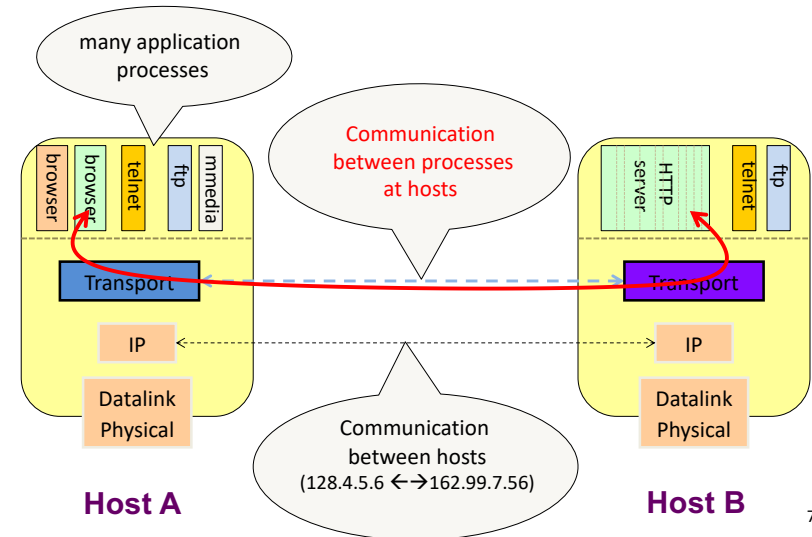
5

Why a transport layer?



6

Why a transport layer?



7

Why a transport layer?

- IP packets are addressed to a host but end-to-end communication is between application processes at hosts
 - Need a way to decide which packets go to which applications (mux/demux)
- IP provides a weak service model (*best-effort*)
 - Packets can be corrupted, delayed, dropped, reordered, duplicated
 - No guidance on how much traffic to send and when
 - Dealing with this is tedious for application developers

8

Role of the Transport Layer

- Communication between application processes
 - Multiplexing between application processes
 - Implemented using *ports*

9

Role of the Transport Layer

- Communication between application processes
- Provide common end-to-end services for app layer [optional]
 - Reliable, in-order data delivery
 - Paced data delivery: flow and congestion-control
 - too fast may overwhelm the network
 - too slow is not efficient

(Just Like Computer Networking Lectures...)

10

Role of the Transport Layer

- Communication between processes
- Provide common end-to-end services for app layer [optional]
- TCP and UDP are the common transport protocols
 - also SCTP, MTCP, SST, RDP, DCCP, ...

11

Role of the Transport Layer

- Communication between processes
- Provide common end-to-end services for app layer [optional]
- TCP and UDP are the common transport protocols
- UDP is a minimalist, no-frills transport protocol
 - only provides mux/demux capabilities

12

Role of the Transport Layer

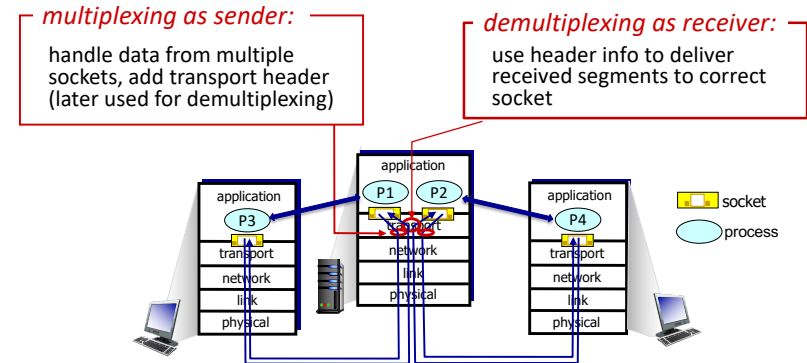
- Communication between processes
- Provide common end-to-end services for app layer [optional]
- TCP and UDP are the common transport protocols
- UDP is a minimalist, no-frills transport protocol
- TCP is the *totus porcus* protocol
 - offers apps a reliable, in-order, byte-stream abstraction
 - with congestion control
 - but **no** performance (delay, bandwidth, ...) guarantees

13

Role of the Transport Layer

- Communication between processes
 - mux/demux from and to application processes
 - implemented using ports

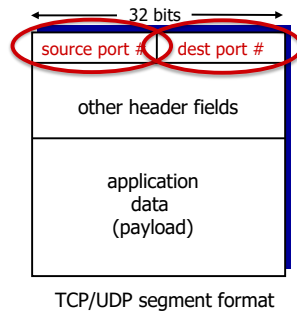
Multiplexing/demultiplexing



14

How demultiplexing Works

- host receives IP datagrams
 - each datagram has source IP address, destination IP address
 - each datagram carries one transport-layer segment
 - each segment has source, destination port number
- host uses *IP addresses & port numbers* to direct segment to appropriate socket



Connectionless demultiplexing

- when creating socket, must specify *host-local* port #:


```
DatagramSocket mySocket1 = new DatagramSocket(12634);
```
- when creating datagram to send into UDP socket, must specify
 - destination IP address
 - destination port #

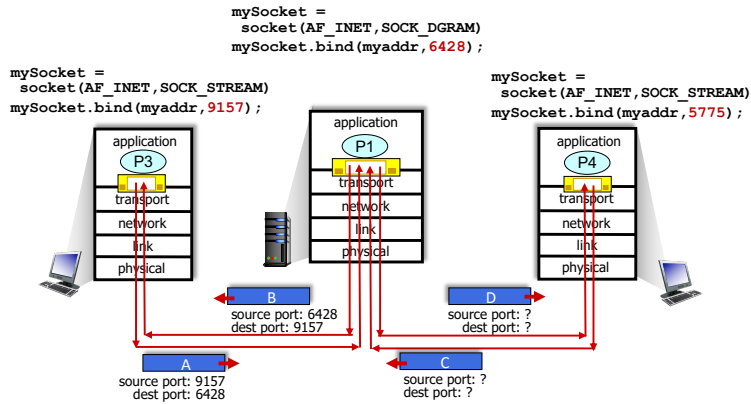
when receiving host receives *UDP* segment:

- checks destination port # in segment
- directs *UDP* segment to socket with that port #



IP/UDP datagrams with *same dest. port #*, but different source IP addresses and/or source port numbers will be directed to *same socket* at receiving host

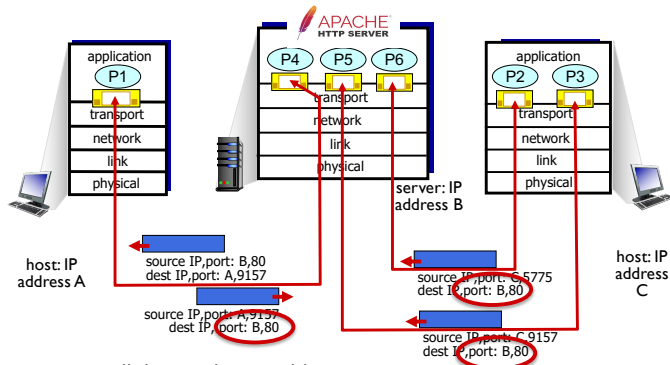
Connectionless demultiplexing: an example



Connection-oriented demultiplexing

- TCP socket identified by 4-tuple:**
 - source IP address
 - source port number
 - dest IP address
 - dest port number
 - server may support many simultaneous TCP sockets:**
 - each socket identified by its own 4-tuple
 - each socket associated with a different connecting client
- slight lie alert.... I should say that a common network tuple has FIVE values
- source IP address
 - source port number
 - dest IP address
 - dest port number AND
 - protocol e.g. TCP (6) or UDP (17)
- demux: receiver uses *all four values (4-tuple)* to direct segment to appropriate socket**

Connection-oriented demultiplexing: example



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

Summary

- Multiplexing, demultiplexing: based on segment, datagram header field values
- UDP:** demultiplexing using destination port number (only)
- TCP:** demultiplexing using 4-tuple: source and destination IP addresses, and port numbers
- Multiplexing/demultiplexing can happen at *any* layer

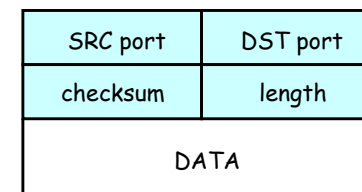
More on Ports

- Separate 16-bit port address space for UDP and TCP
- “Well known” ports (0-1023): everyone agrees which services run on these ports
 - e.g., ssh:22, http:80, https:443
 - helps client know server’s port
- Ephemeral ports (most 1024-65535): dynamically selected: as the source port for a client process

29

UDP: User Datagram Protocol

- Lightweight communication between processes
 - Avoid overhead and delays of ordered, reliable delivery
- UDP described in RFC 768 – (1980!)
 - Destination IP address and port to support demultiplexing
 - Optional error checking on the packet contents
 - (checksum field of 0 means “don’t verify checksum”) **not in IPv6!**
 - ((this idea of optional checksum is removed in IPv6))



30

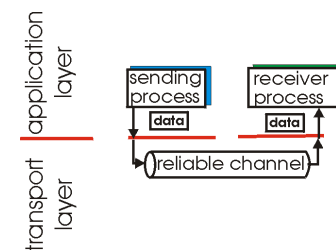
Why a transport layer?

- IP packets are addressed to a host but end-to-end communication is between application processes at hosts
 - Need a way to decide which packets go to which applications (mux/demux)
- IP provides a weak service model (*best-effort*)
 - Packets can be corrupted, delayed, dropped, reordered, duplicated

31

Principles of Reliable data transfer

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



(a) provided service

- In a perfect world, reliable transport is easy

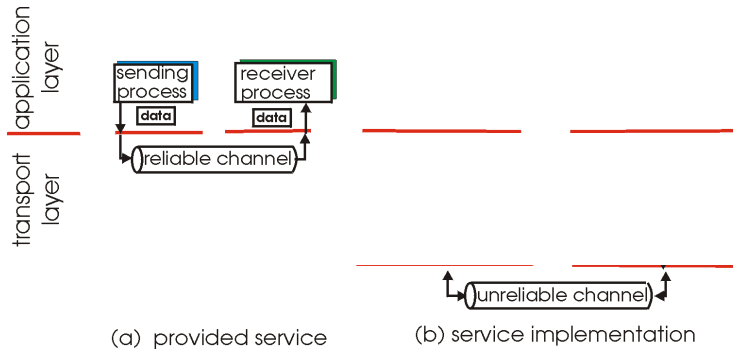
But the Internet default is *best-effort*

- All the bad things best-effort can do
 - a packet is corrupted (bit errors)
 - a packet is lost
 - a packet is delayed (*why?*)
 - packets are reordered (*why?*)
 - a packet is duplicated (*why?*)

32

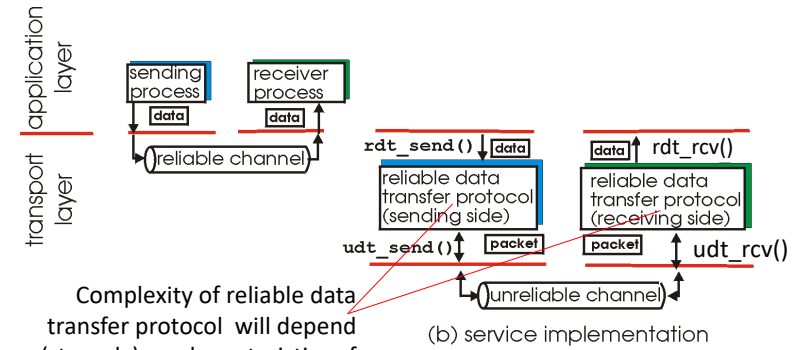
Principles of Reliable data transfer

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



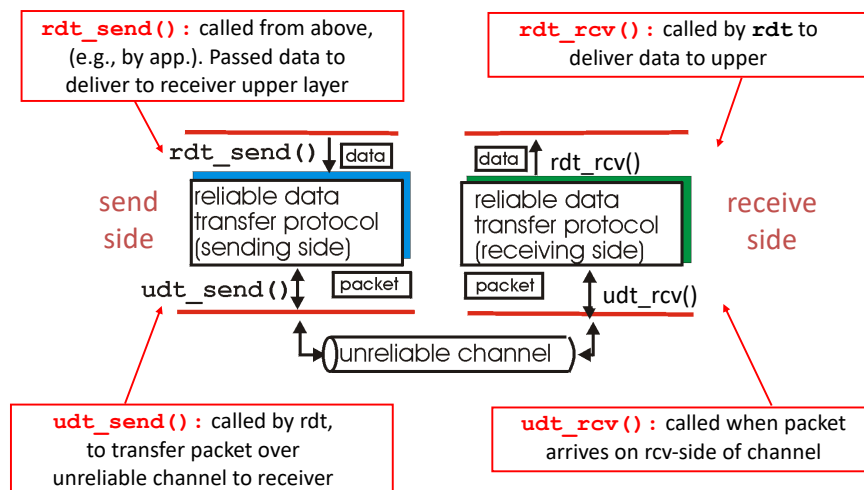
Principles of Reliable data transfer

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



Complexity of reliable data transfer protocol will depend (strongly) on characteristics of unreliable channel (lose, corrupt, reorder data?)

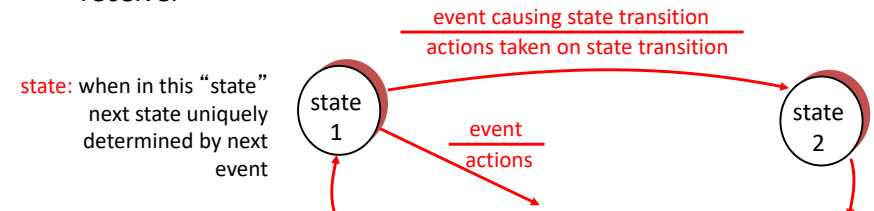
Reliable data transfer: getting started



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



KR state machines – a note.

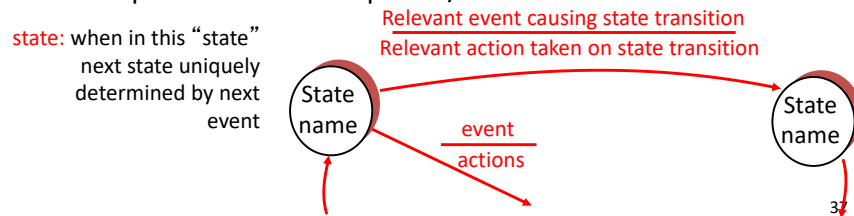
Beware

Kurose and Ross has a confusing/confused attitude to state-machines.

I've attempted to normalise the representation.

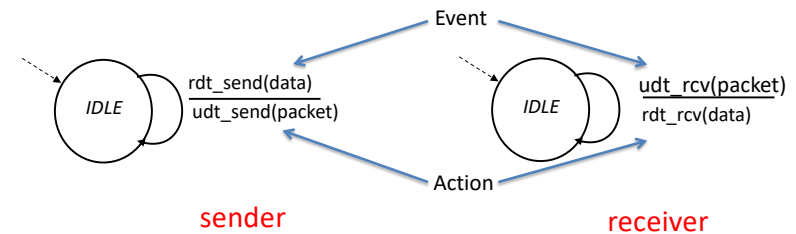
UPSHOT: these slides have differing information to the KR book (from which the RDT example is taken.)

in KR "actions taken" appear wide-ranging, my interpretation is more specific/relevant.



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

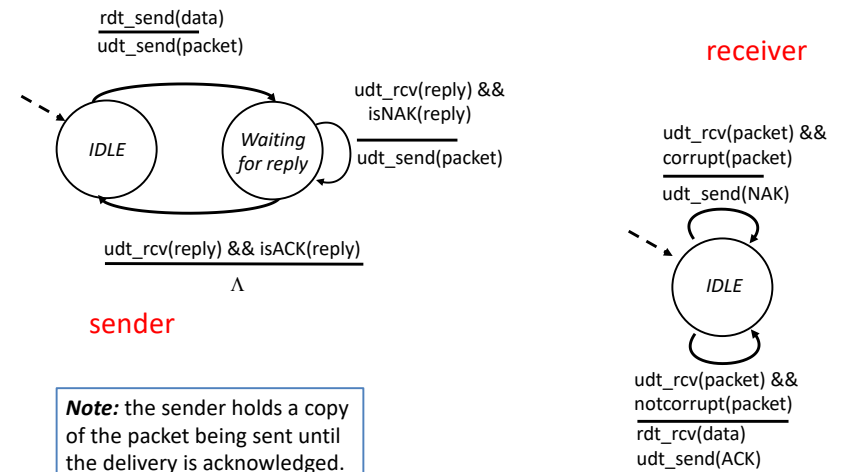


38

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:
 - acknowledgements (ACKs):** receiver explicitly tells sender that packet received is OK
 - negative acknowledgements (NAKs):** receiver explicitly tells sender that packet had errors
 - sender retransmits packet on receipt of NAK
- new mechanisms in **rdt2.0** (beyond **rdt1.0**):
 - error detection
 - receiver feedback: control msgs (ACK,NAK) receiver->sender

rdt2.0: FSM specification

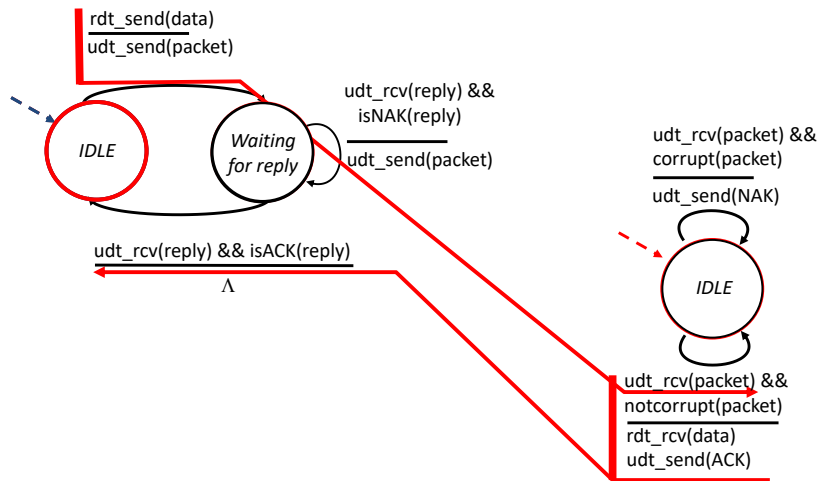


Note: the sender holds a copy of the packet being sent until the delivery is acknowledged.

39

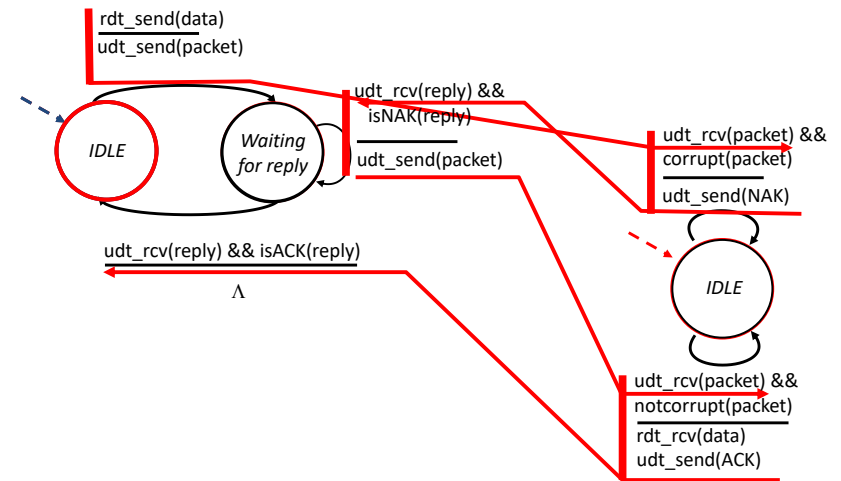
41

rdt2.0: operation with no errors



42

rdt2.0: error scenario



43

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:

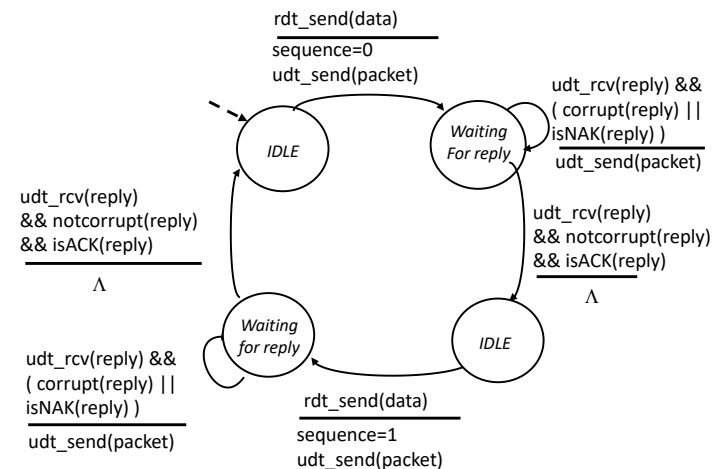
- sender retransmits current packet if ACK/NAK garbled
- sender adds *sequence number* to each packet
- receiver discards (doesn't deliver) duplicate packet

stop and wait

Sender sends one packet, then waits for receiver response

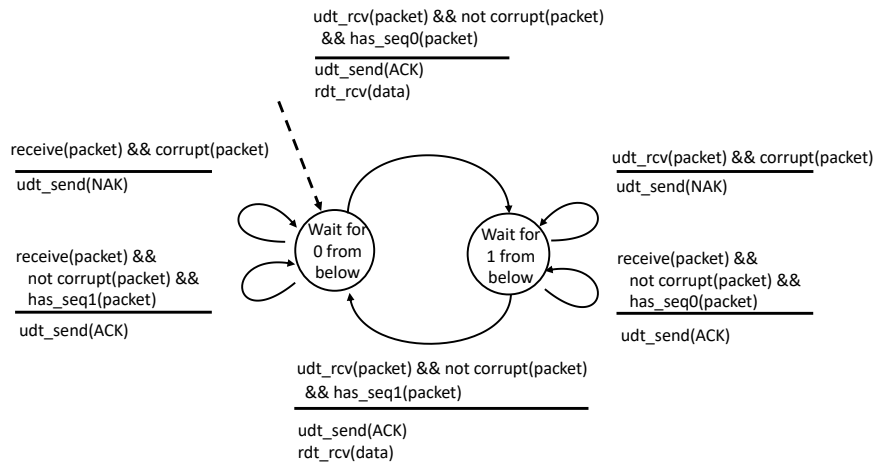
44

rdt2.1: sender, handles garbled ACK/NAKs



46

rdt2.1: receiver, handles garbled ACK/NAKs



47

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 a 0 or 1 sequence number

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

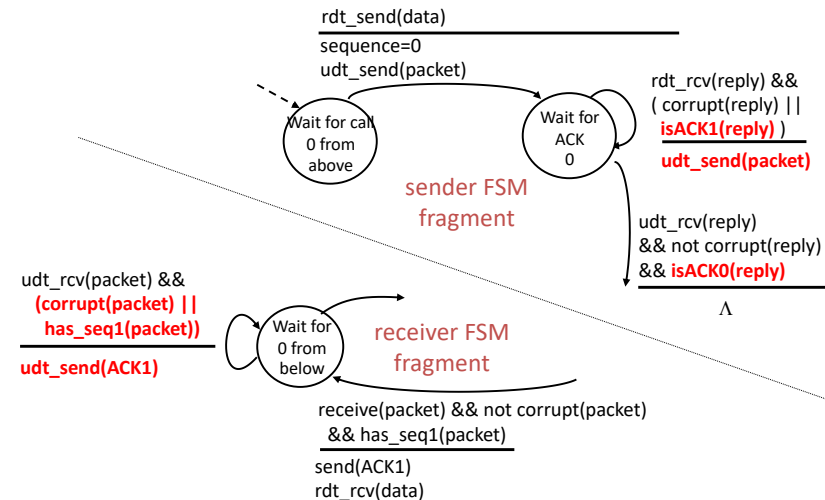
48

rdt2.2: a NAK-free protocol

- same functionality as rdt2.1, using ACKs only
 - receiver must *explicitly* include seq # of pkt being ACKed
- duplicate ACK at sender results in same action as NAK: *retransmit current pkt*

As we will see, TCP uses this approach to be NAK-free

rdt2.2: sender, receiver fragments



49

50

rdt3.0: channels with errors *and* loss

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

- checksum, sequence #s, ACKs, retransmissions will be of help ... but not quite enough

Q: How do *humans* handle lost sender-to-receiver words in conversation?

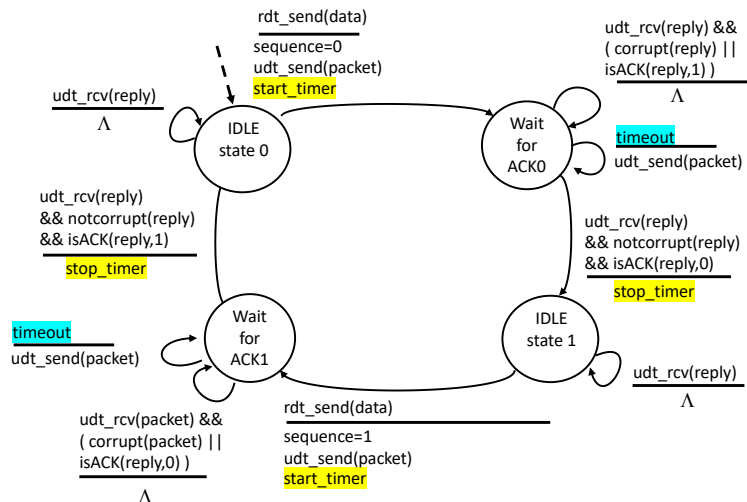
rdt3.0: channels with errors *and* loss

Approach: sender waits “reasonable” amount of time for ACK

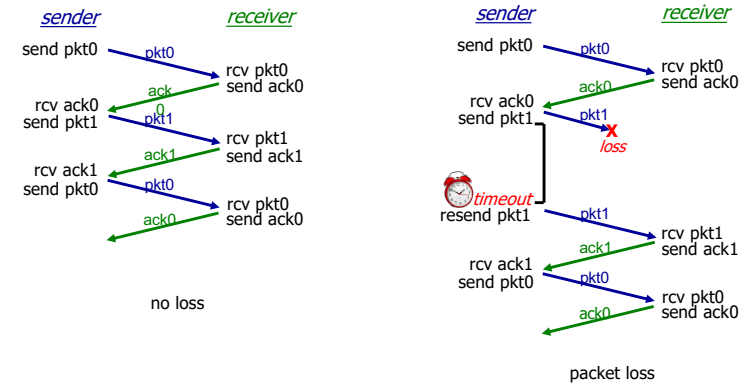
- retransmits if no ACK received in this time
- if pkt (or ACK) just delayed (not lost):
 - retransmission will be duplicate, but seq #s already handles this!
 - receiver must specify seq # of packet being ACKed
- use countdown timer to interrupt after “reasonable” amount of time



rdt3.0 sender

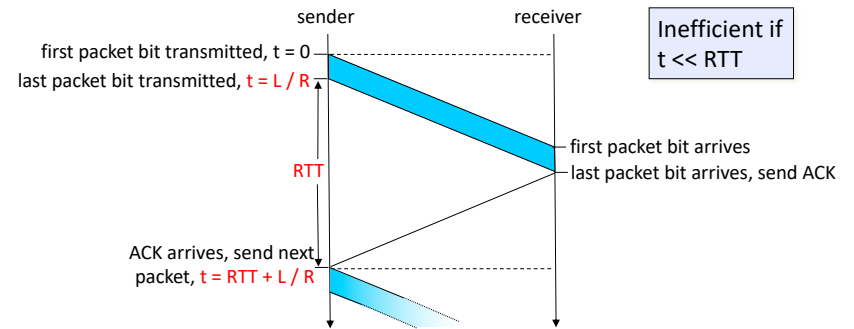
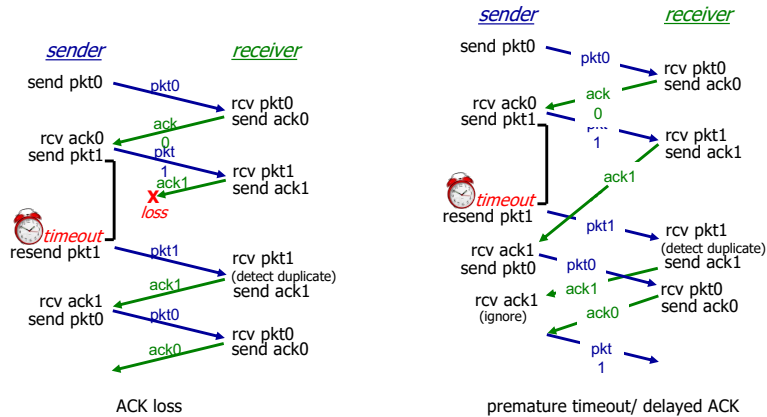


rdt3.0 in action



rdt3.0: stop-and-wait operation

rdt3.0 in action



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

61

Performance of rdt3.0 (stop-and-wait)

- rdt3.0 works, but performance stinks
- ex: 1 Gbps link, 15 ms prop. delay, 8000 bit packet:

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

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

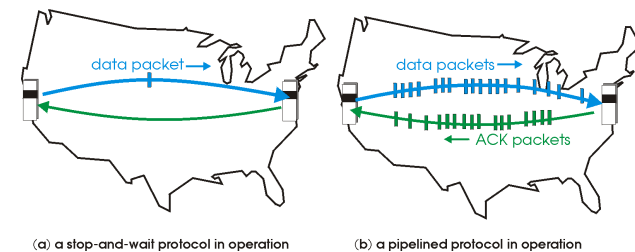
$$U_{\text{sender}} = \frac{L/R}{RTT + L/R} = \frac{.008}{30.008} = 0.00027$$

- 1KB pkt every 30 msec -> 33kB/sec throughput over 1 Gbps link
- The network protocol limits use of physical resources!

Pipelined (Packet-Window) 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

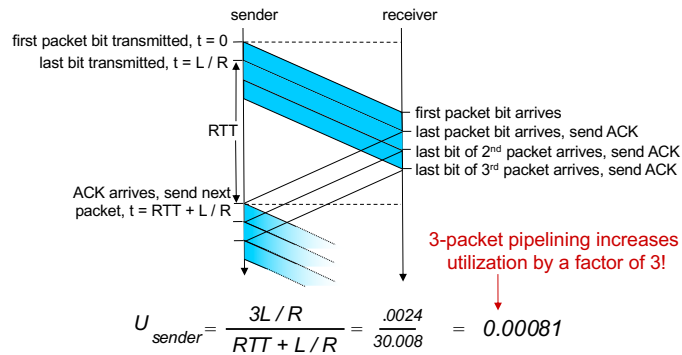


62

63

A Sliding Packet Window

Pipelining: increased utilization



- **window** = set of adjacent sequence numbers
 - The size of the set is the **window size**; assume window size is n
- General idea: send up to n packets at a time
 - Sender can send packets in its window
 - Receiver can accept packets in its window
 - Window of acceptable packets “slides” on successful reception/acknowledgement

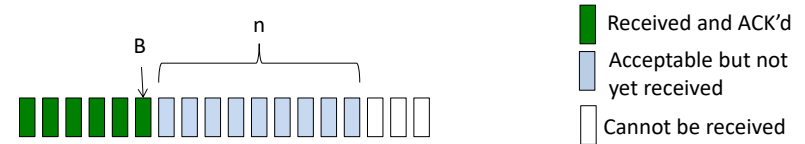
65

Acknowledgements w/ Sliding Window

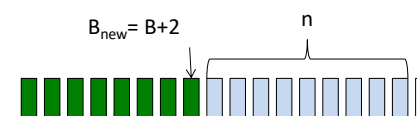
- Two common options
 - cumulative ACKs: ACK carries next in-order sequence number that the receiver expects

Cumulative Acknowledgements (1)

- At receiver



- After receiving $B+1, B+2$



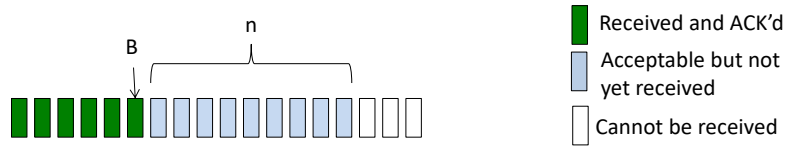
- Receiver sends $\text{ACK}(B_{\text{new}}+1)$

67

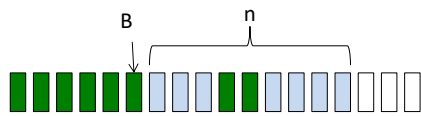
68

Cumulative Acknowledgements (2)

- At receiver



- After receiving B+4, B+5



- Receiver sends ACK(B+???)

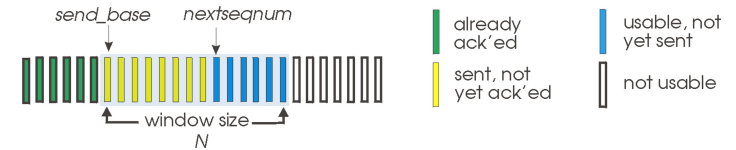
Oh....
how do we
recover?

69

Go-Back-N: sender

- sender: "window" of up to N, consecutive transmitted but unACKed pkts

- k-bit seq # in pkt header



- cumulative ACK:** ACK(n): ACKs all packets up to, including seq # n
 - on receiving ACK(n): move window forward to begin at n+1
- timer for oldest in-flight packet
- timeout(n): retransmit packet n and all higher seq # packets in window

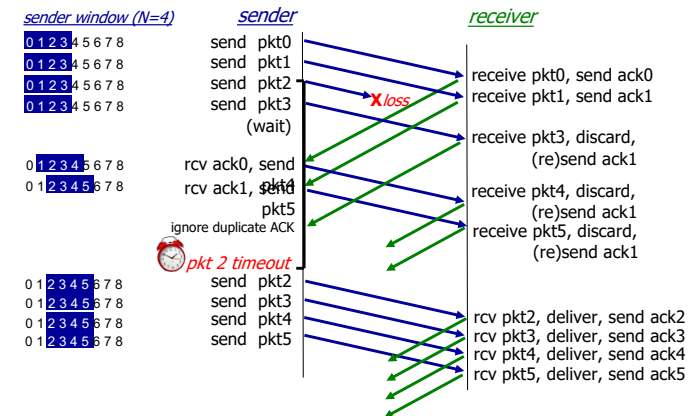
Go-Back-N: receiver

- ACK-only: always send ACK for correctly-received packet so far, with highest *in-order* seq #
 - may generate duplicate ACKs
 - need only remember *rcv_base*
- on receipt of out-of-order packet:
 - can discard (don't buffer) or buffer: an implementation decision
 - re-ACK pkt with highest in-order seq #

Receiver view of sequence number space:



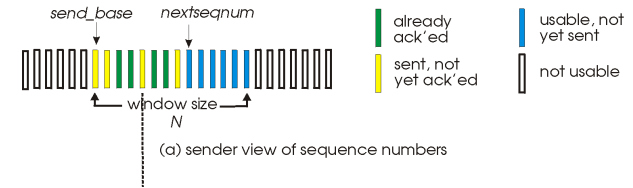
Go-Back-N in action



Selective repeat

- receiver *individually* acknowledges all correctly received packets
 - buffers packets, as needed, for eventual in-order delivery to upper layer
- sender times-out/retransmits individually for unACKed packets
 - sender maintains timer for each unACKed pkt
- sender window
 - N consecutive seq #s
 - limits seq #s of sent, unACKed packets

Selective repeat: sender, receiver windows



Selective repeat: sender and receiver

sender

data from above:

- if next available seq # in window, send packet

timeout(n):

- resend packet n , restart timer

ACK(n) in $[\text{sendbase}, \text{sendbase}+N-1]$:

- mark packet n as received
- if n smallest unACKed packet, advance window base to next unACKed seq #

receiver

packet n in $[\text{rcvbase}, \text{rcvbase}+N-1]$

- send ACK(n)
- out-of-order: buffer
- in-order: deliver (also deliver buffered, in-order packets), advance window to next not-yet-received packet

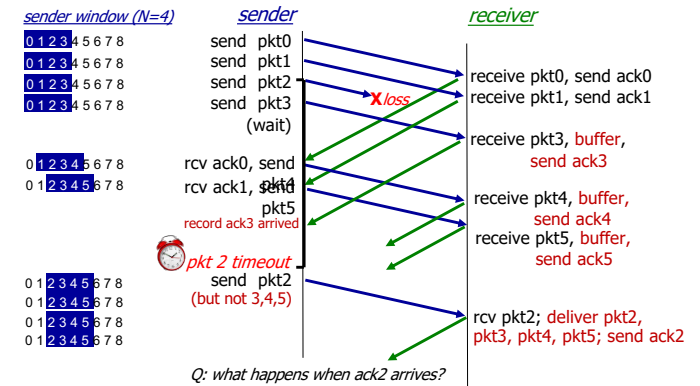
packet n in $[\text{rcvbase}-N, \text{rcvbase}-1]$

- ACK(n)

otherwise:

- ignore

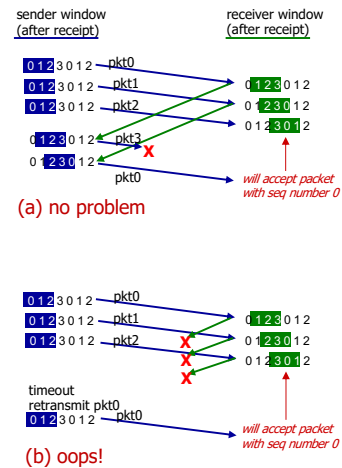
Selective Repeat in action



Selective repeat: a dilemma!

example:

- seq #s: 0, 1, 2, 3 (base 4 counting)
- window size=3



Selective repeat: a dilemma!

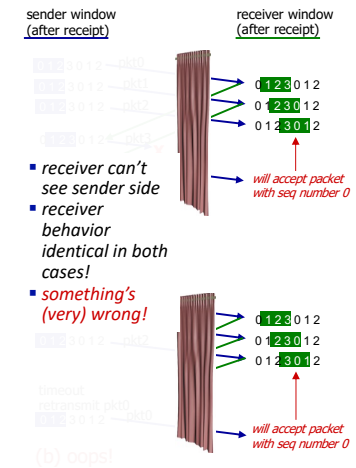
example:

- seq #s: 0, 1, 2, 3 (base 4 counting)
- window size=3

Q: what relationship is needed between sequence # size and window size to avoid problem in scenario (b)?

Solution:

maximum allowable window size =
half the sequence number space.



Observations

- With sliding windows, it is possible to fully utilize a link, provided the window size (n) is large enough. Throughput is $\sim (n/RTT)$
 - Stop & Wait is like $n = 1$.
- Sender has to buffer all unacknowledged packets, because they may require retransmission
- Receiver may be able to accept out-of-order packets, but only up to its buffer limits
- Implementation complexity depends on protocol details (GBN vs. SR)

Recap: components of a solution

- Checksums (for error detection)
- Timers (for loss detection)
- Acknowledgments
 - cumulative
 - selective
- Sequence numbers (duplicates, windows)
- Sliding Windows (for efficiency)
- Reliability protocols use the above to decide when and what to retransmit or acknowledge

What does TCP do?

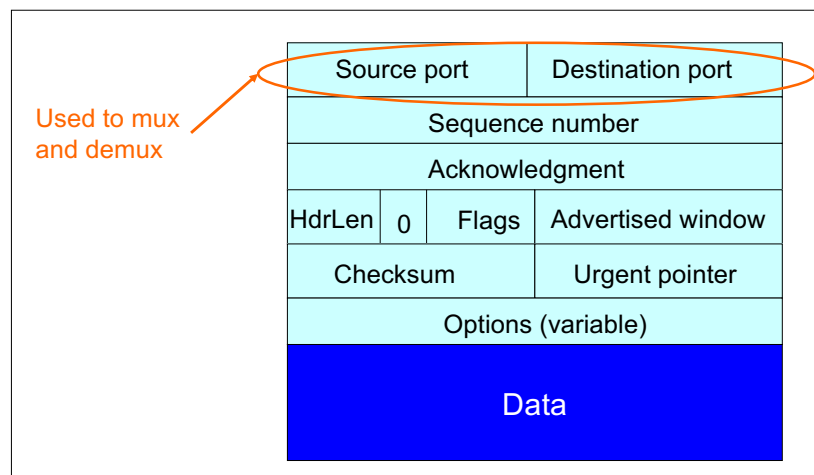
Most of our previous tricks + a few more beside

- Sequence numbers are byte offsets
- Sender and receiver maintain a sliding window
- Receiver sends cumulative acknowledgements (like GBN)
- Sender maintains a single retx. timer
- Receivers do not drop out-of-sequence packets (like SR)
- Introduces **fast retransmit** : optimization that uses duplicate ACKs to trigger early retx
- Introduces timeout estimation algorithms

TCP: overview RFCs: 793,1122, 2018, 5681, 7323

- **point-to-point:**
 - one sender, one receiver
- **reliable, in-order byte stream:**
 - no “message boundaries”
- **full duplex data:**
 - bi-directional data flow in same connection
 - MSS: maximum segment size
- **cumulative ACKs**
- **pipelining:**
 - TCP congestion and flow control set window size
- **connection-oriented:**
 - handshaking (exchange of control messages) initializes sender, receiver state before data exchange
- **flow controlled:**
 - sender will not overwhelm receiver

TCP Header

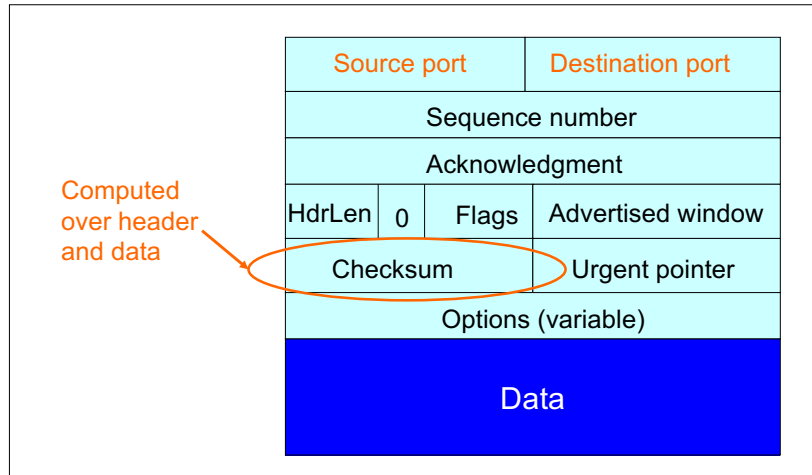


What does TCP do?

Many of our previous ideas, but some key differences

- Checksum

TCP Header



97

What does TCP do?

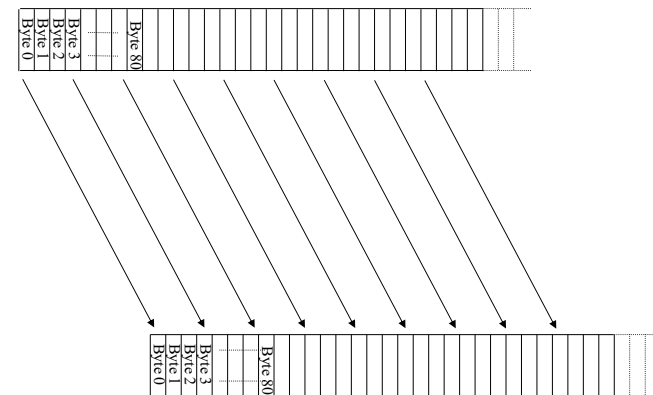
Many of our previous ideas, but some key differences

- Checksum
- Sequence numbers are byte offsets

TCP: Segments and Sequence Numbers

TCP “Stream of Bytes” Service...

Application @ Host A

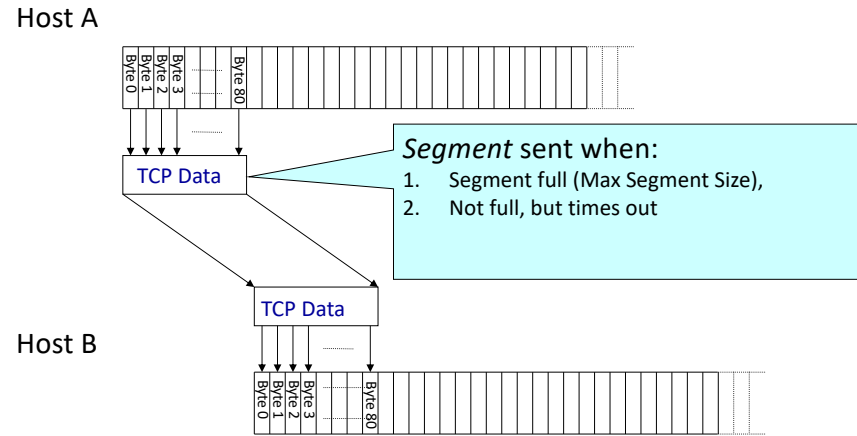


Application @ Host B

99

100

... Provided Using TCP "Segments"



101

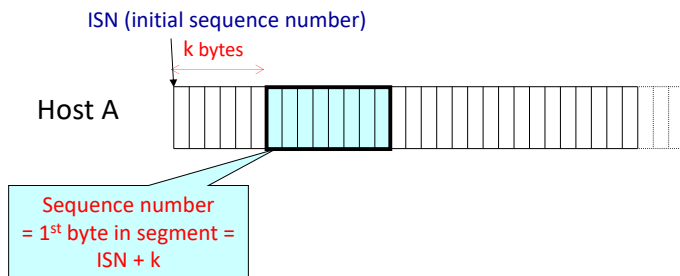
TCP Segment



- IP packet
 - No bigger than Maximum Transmission Unit (MTU)
 - E.g., up to 1500 bytes with Ethernet
- TCP packet
 - IP packet with a TCP header and data inside
 - TCP header ≥ 20 bytes long
- TCP segment
 - No more than Maximum Segment Size (MSS) bytes
 - E.g., up to 1460 consecutive bytes from the stream
 - $MSS = MTU - (IP\ header) - (TCP\ header)$

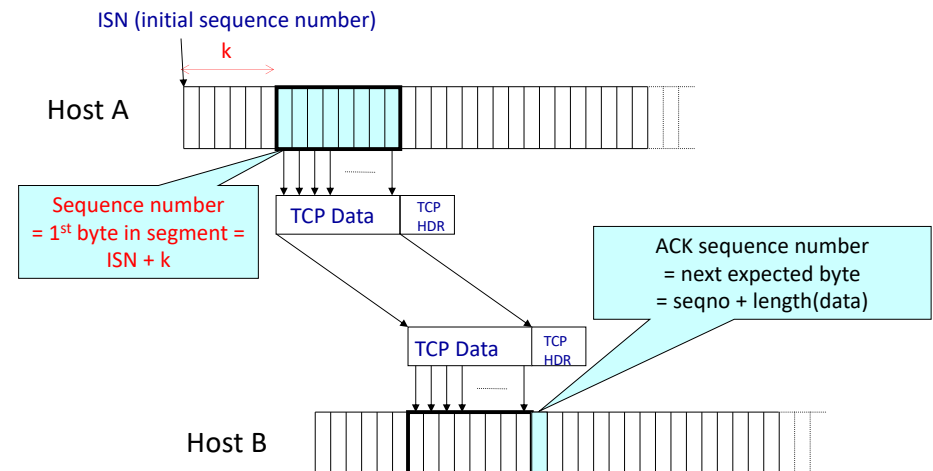
102

Sequence Numbers



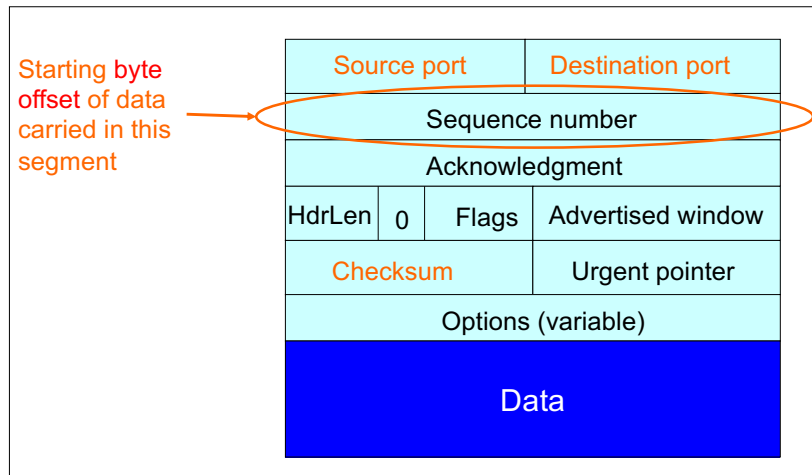
103

Sequence Numbers



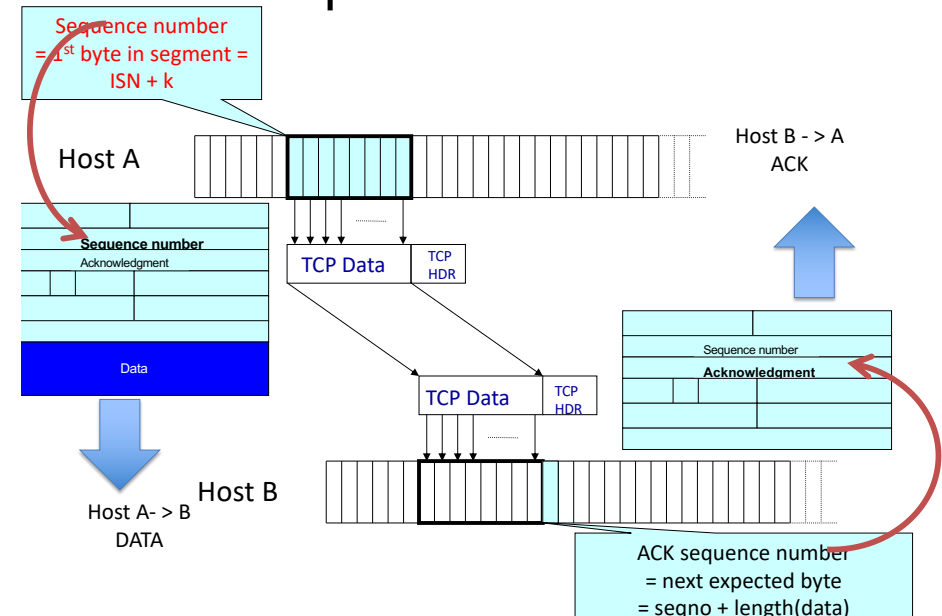
104

TCP Header



105

Sequence Numbers



TCP Sequences and ACKS

TCP is full duplex by default

- two independently flows of sequence numbers

Sequence acknowledgement is given in terms of BYTES (not packets); the window is in terms of bytes.

number of packets = window size (bytes) / Segment Size

Servers and Clients are not Source and Destination

Piggybacking increases efficiency but many flows may only have data moving in one direction

107

What does TCP do?

Most of our previous tricks, but a few differences

- Checksum
- Sequence numbers are byte offsets
- Receiver sends cumulative acknowledgements (like GBN)

ACKing and Sequence Numbers

- Sender sends packet
 - Data starts with sequence number X
 - Packet contains B bytes $[X, X+1, X+2, \dots, X+B-1]$
- Upon receipt of packet, receiver sends an ACK
 - If all data prior to X already received:
 - ACK acknowledges $X+B$ (because that is next expected byte)
 - If highest in-order byte received is Y s.t. $(Y+1) < X$
 - ACK acknowledges $Y+1$
 - Even if this has been ACKed before

109

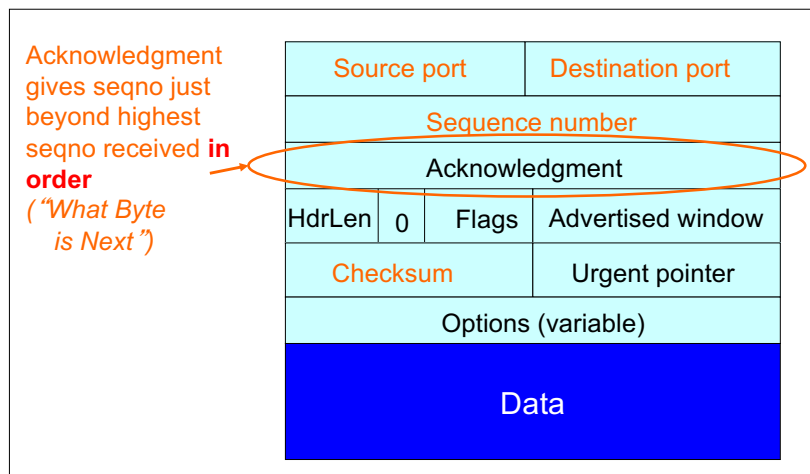
Normal Pattern

- Sender: seqno= X , length= B
- Receiver: ACK= $X+B$
- Sender: seqno= $X+B$, length= B
- Receiver: ACK= $X+2B$
- Sender: seqno= $X+2B$, length= B

- Seqno of next packet is same as last ACK field

110

TCP Header



111

What does TCP do?

Most of our previous tricks, but a few differences

- Checksum
- Sequence numbers are byte offsets
- Receiver sends cumulative acknowledgements (like GBN)
- Receivers **can** buffer out-of-sequence packets (like SR)

112

Loss with cumulative ACKs

- Sender sends packets with 100B and seqnos.:
 - 100, 200, 300, 400, 500, 600, 700, 800, 900, ...
- Assume the fifth packet (seqno 500) is lost, but no others
- Stream of ACKs will be:
 - 200, 300, 400, 500, 500, 500, 500, ...

113

What does TCP do?

Most of our previous tricks, but a few differences

- Checksum
- Sequence numbers are byte offsets
- Receiver sends cumulative acknowledgements (like GBN)
- Receivers may not drop out-of-sequence packets (like SR)
- Introduces **fast retransmit**: optimization that uses duplicate ACKs to trigger early retransmission

114

Loss with cumulative ACKs

- “Duplicate ACKs” are a sign of an isolated loss
 - The lack of ACK progress means 500 hasn’t been delivered
 - Stream of ACKs means some packets are being delivered
- Therefore, could trigger resend upon receiving k duplicate ACKs
 - TCP uses k=3
- But response to loss is trickier....

115

Loss with cumulative ACKs

- Two choices:
 - Send missing packet and increase W by the number of dup ACKs
 - Send missing packet, and wait for ACK to increase W
- Which should TCP do?

116

What does TCP do?

Most of our previous tricks, but a few differences

- Checksum
- Sequence numbers are byte offsets
- Receiver sends cumulative acknowledgements (like GBN)
- Receivers do not drop out-of-sequence packets (like SR)
- Introduces fast retransmit: optimization that uses duplicate ACKs to trigger early retransmission
- Sender maintains a single retransmission timer (like GBN) and retransmits on timeout

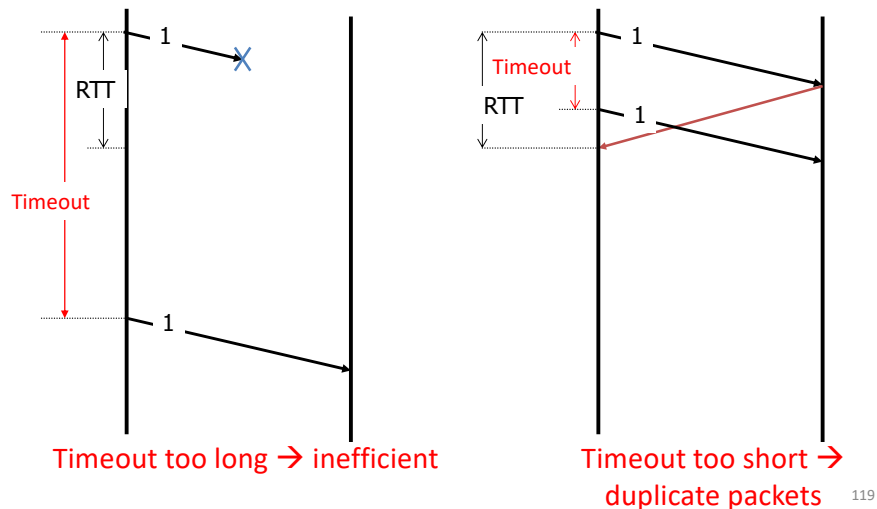
117

Retransmission Timeout

- If the sender hasn't received an ACK by timeout, retransmit the first packet in the window
- How do we pick a timeout value?

118

Timing Illustration



119

Retransmission Timeout

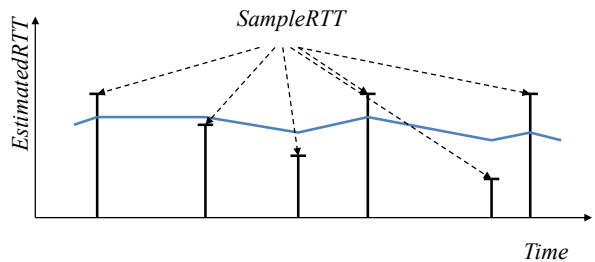
- If haven't received ack by timeout, retransmit the first packet in the window
- How to set timeout?
 - Too long: connection has low throughput
 - Too short: retransmit packet that was just delayed
- Solution: make timeout proportional to RTT
- But how do we measure RTT?

120

RTT Estimation

- Use exponential averaging of RTT samples

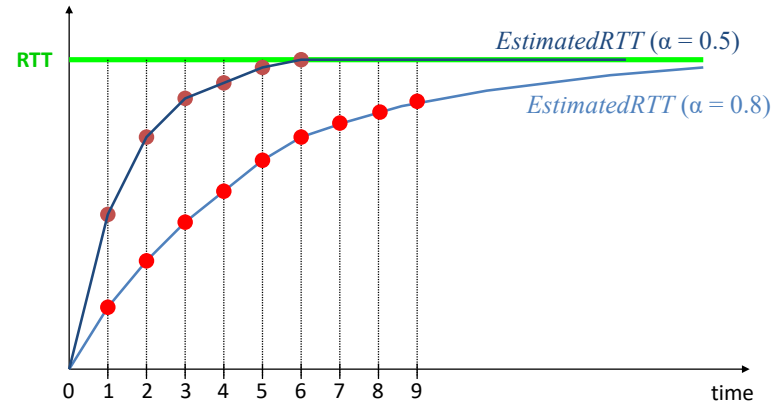
$SampleRTT = AckRcvdTime - SendPacketTime$
 $EstimatedRTT = \alpha \times EstimatedRTT + (1 - \alpha) \times SampleRTT$
 $0 < \alpha \leq 1$



121

Exponential Averaging Example

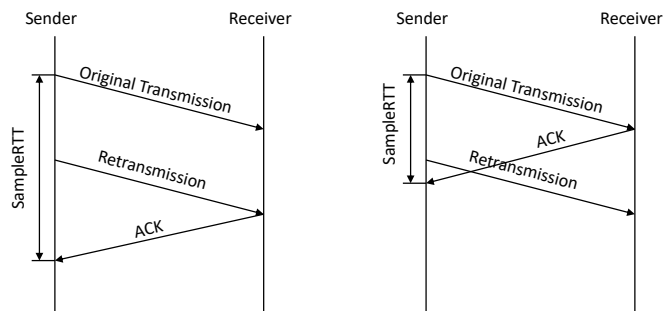
$EstimatedRTT = \alpha * EstimatedRTT + (1 - \alpha) * SampleRTT$
 Assume RTT is constant $\rightarrow SampleRTT = RTT$



122

Problem: Ambiguous Measurements

- How do we differentiate between the real ACK, and ACK of the retransmitted packet?



123

Karn/Partridge Algorithm

Discard junk measures

- Measure *SampleRTT* only for original transmissions
 - Once a segment has been retransmitted, do not use it for any further measurements
- Computes EstimatedRTT using $\alpha = 0.875$
- Timeout value (RTO) = $2 \times EstimatedRTT$
- Employs **exponential backoff**
 - Every time RTO timer expires, set $RTO \leftarrow 2 \cdot RTO$
 - (Up to maximum ≥ 60 sec)
 - Every time new measurement comes in (= successful original transmission), collapse RTO back to $2 \times EstimatedRTT$

124

Jacobson/Karels Algorithm

Add a safety margin

- Problem: need to better capture variability in RTT
 - Directly measure **deviation**
- Deviation = | SampleRTT – EstimatedRTT |
- EstimatedDeviation: exponential average of Deviation
- RTO = EstimatedRTT + 4 x EstimatedDeviation

126

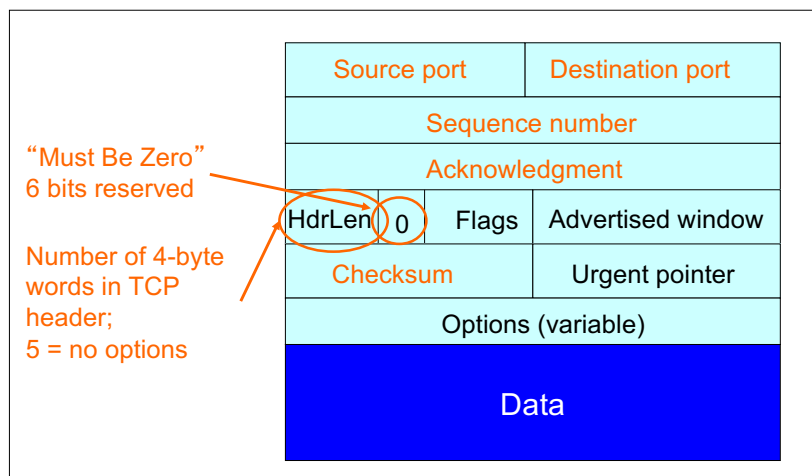
What does TCP do?

Most of our previous ideas, but some key differences

- Checksum
- Sequence numbers are byte offsets
- Receiver sends cumulative acknowledgements (like GBN)
- Receivers do not drop out-of-sequence packets (like SR)
- Introduces fast retransmit: optimization that uses duplicate ACKs to trigger early retransmission
- Sender maintains a single retransmission timer (like GBN) and retransmits on timeout

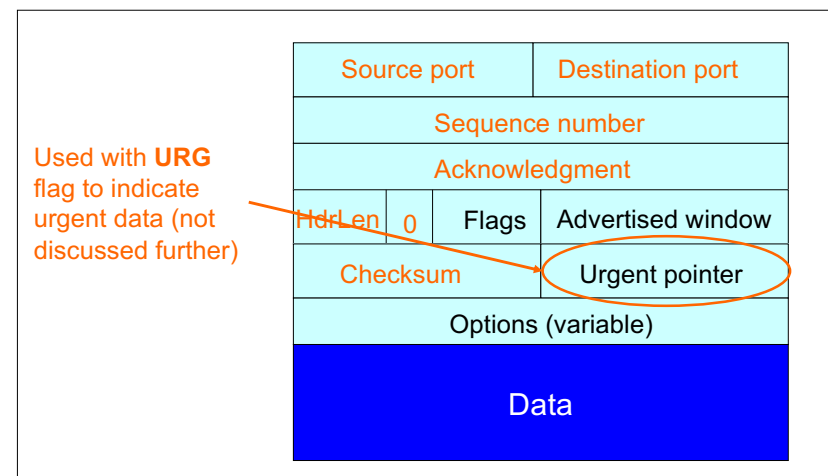
128

TCP Header: What's left?



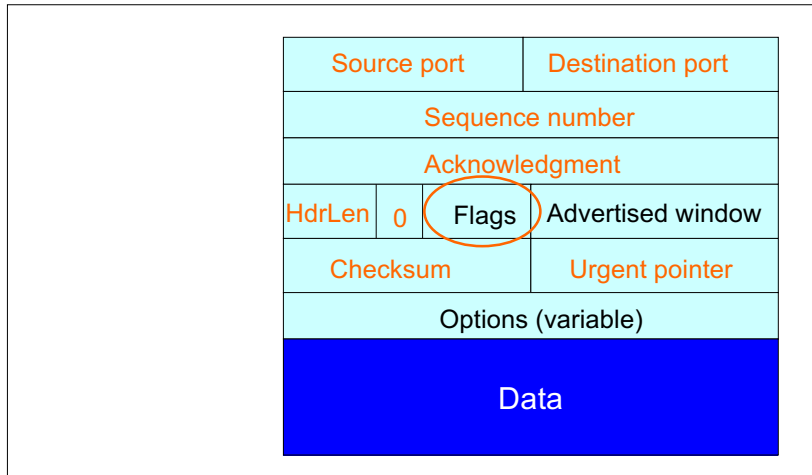
129

TCP Header: What's left?



130

TCP Header: What's left?



131

TCP Connection Establishment and Initial Sequence Numbers

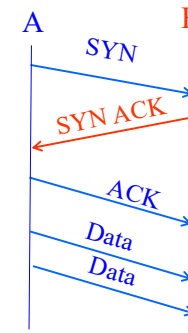
132

Initial Sequence Number (ISN)

- Sequence number for the very first byte
- Why not just use ISN = 0?
- Practical issue
 - IP addresses and port #s uniquely identify a connection
 - Eventually, though, these port #s do get **used again**
 - ... small chance an old packet is **still in flight**
- TCP therefore **requires** changing ISN
- Hosts exchange ISNs when they establish a connection

133

Establishing a TCP Connection

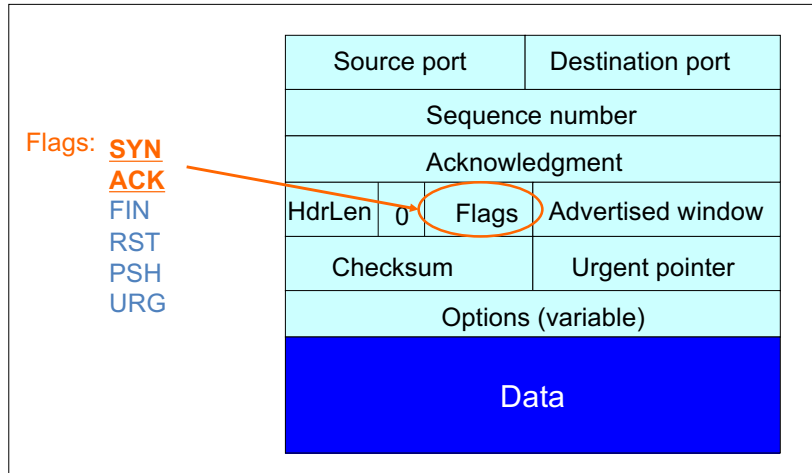


Each host tells its ISN to the other host.

- Three-way handshake to establish connection
 - Host A sends a **SYN** (open; "synchronize sequence numbers") to host B
 - Host B returns a SYN acknowledgment (**SYN ACK**)
 - Host A sends an **ACK** to acknowledge the SYN ACK

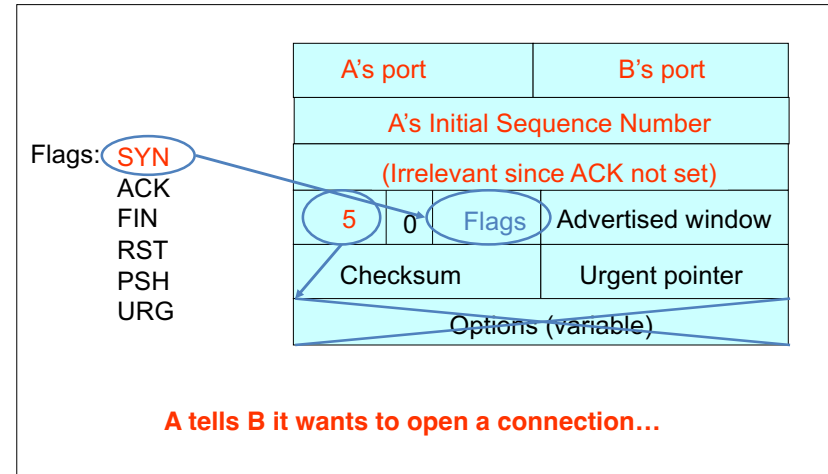
134

TCP Header



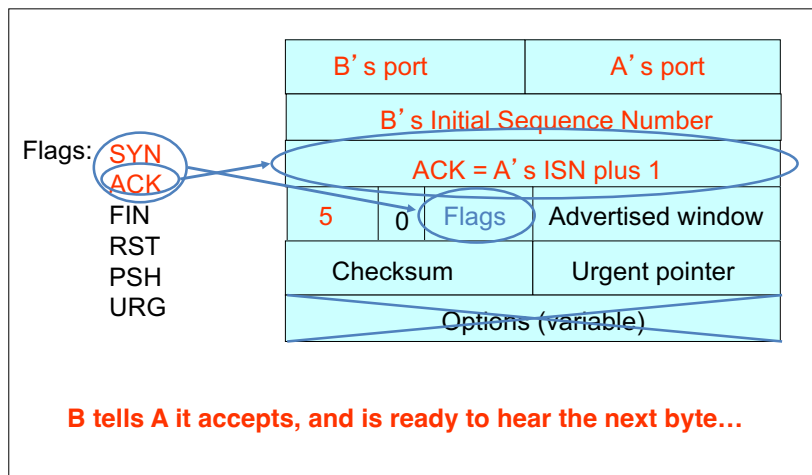
135

Step 1: A's Initial SYN Packet



136

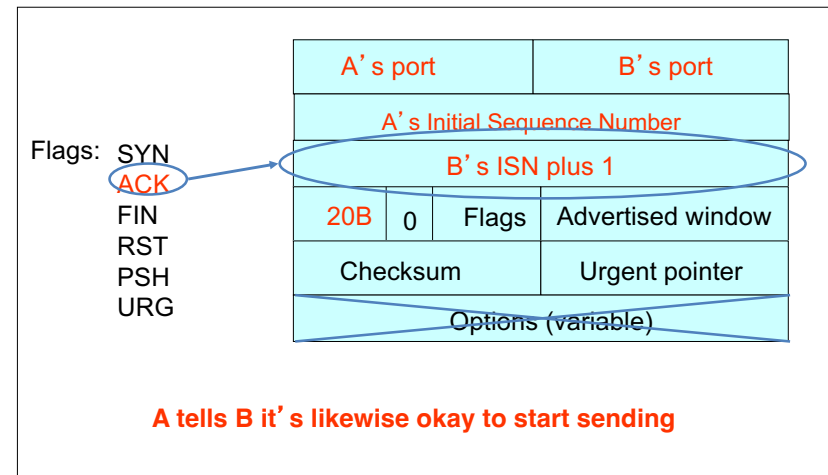
Step 2: B's SYN-ACK Packet



... upon receiving this packet, A can start sending data

137

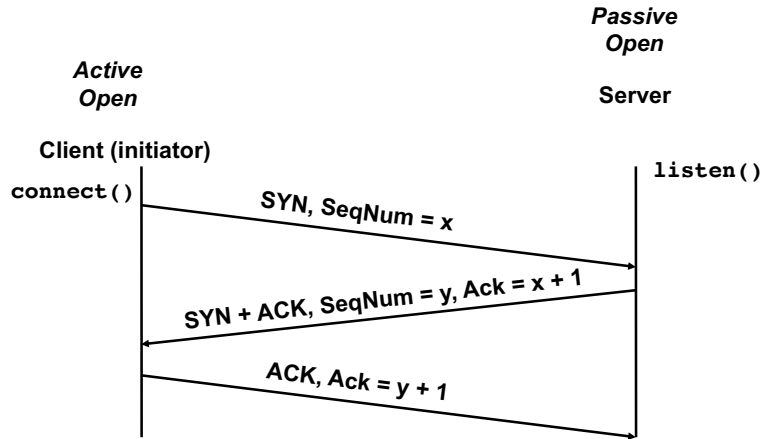
Step 3: A's ACK of the SYN-ACK



... upon receiving this packet, B can start sending data

138

Timing Diagram: 3-Way Handshaking



139

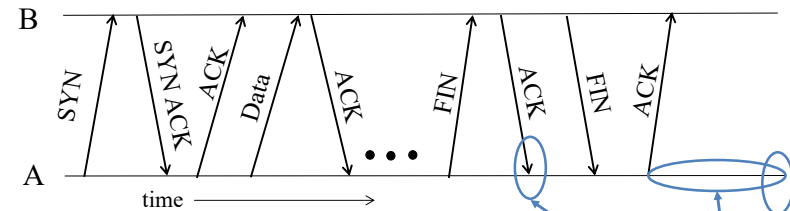
What if the SYN Packet Gets Lost?

- Suppose the SYN packet gets lost
 - Packet is lost inside the network, or:
 - Server **discards** the packet (e.g., it's too busy)
- Eventually, no SYN-ACK arrives
 - Sender sets a **timer** and **waits** for the SYN-ACK
 - ... and retransmits the SYN if needed
- How should the TCP sender set the timer?
 - Sender has **no idea** how far away the receiver is
 - Hard to guess a reasonable length of time to wait
 - **SHOULD** (RFCs 1122 & 2988) use default of **3 seconds**
 - Some implementations instead use 6 seconds

140

Tearing Down the Connection

Normal Termination, One Side At A Time

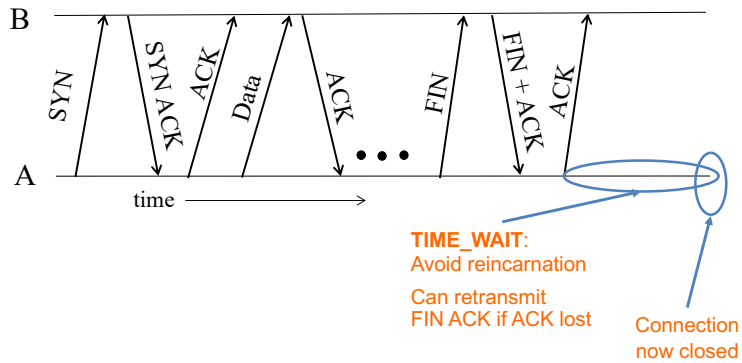


- Finish (**FIN**) to close and receive remaining bytes
 - **FIN** occupies one byte in the sequence space
 - Other host acks the byte to confirm
 - Closes A's side of the connection, but **not B's**
 - Until B likewise sends a **FIN**
 - Which A then acks
- Connection now **half-closed**
- Connection now **closed**
- TIME_WAIT:**
Avoid reincarnation
B will retransmit FIN if ACK is lost

141

142

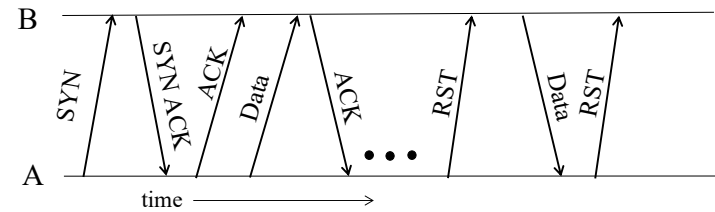
Normal Termination, Both Together



- Same as before, but B sets **FIN** with their ack of A's **FIN**

143

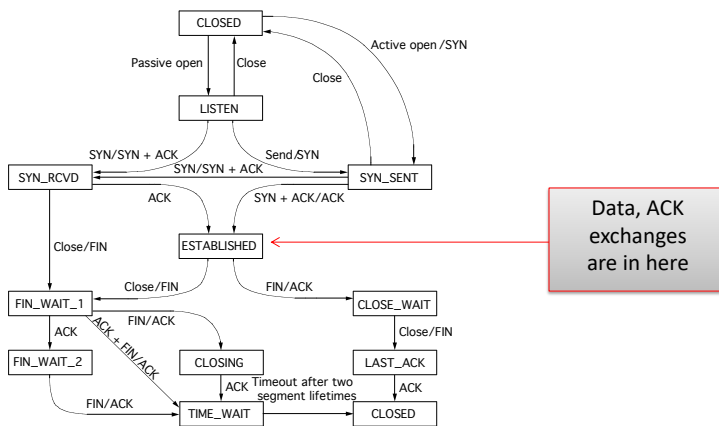
Abrupt Termination



- A sends a RESET (**RST**) to B
 - E.g., because application process on A **crashed**
- **That's it**
 - B does **not** ack the **RST**
 - Thus, **RST** is **not** delivered **reliably**
 - And: any data in flight is **lost**
 - But: if B sends anything more, will elicit **another RST**

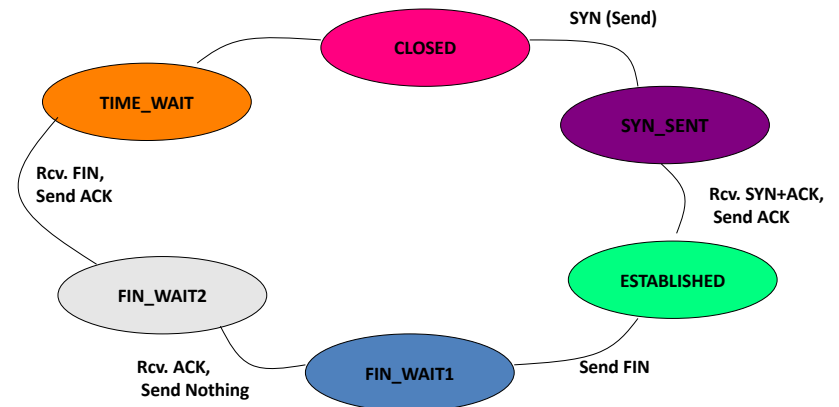
144

TCP State Transitions



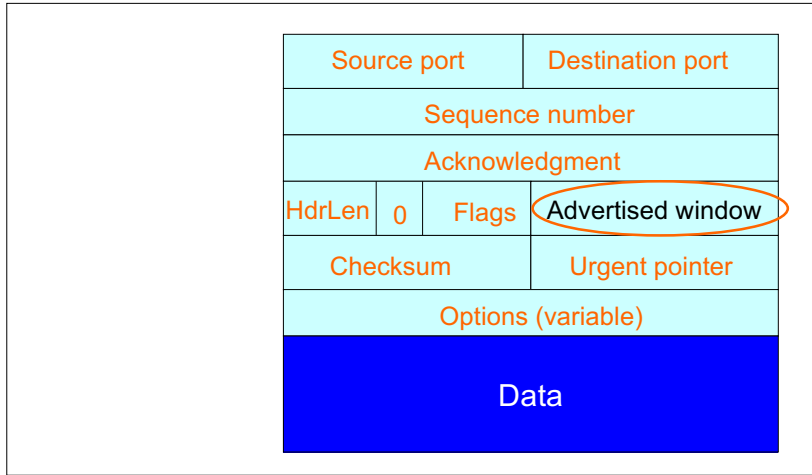
145

An Simpler View of the Client Side



146

TCP Header



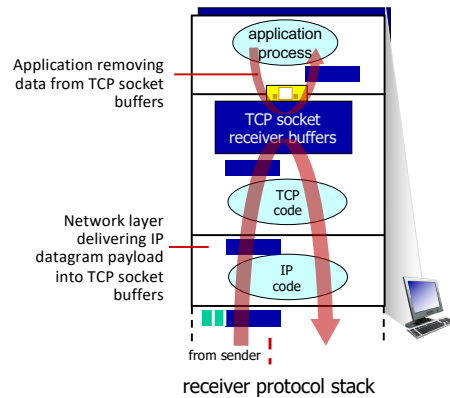
147

- What does TCP do?
 - ARQ windowing, set-up, tear-down
- Flow Control in TCP

148

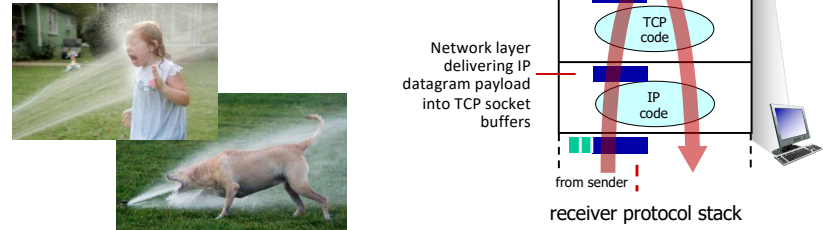
TCP flow control

Q: What happens if network layer delivers data faster than application layer removes data from socket buffers?



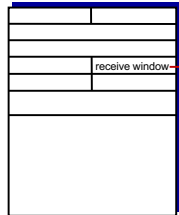
TCP flow control

Q: What happens if network layer delivers data faster than application layer removes data from socket buffers?

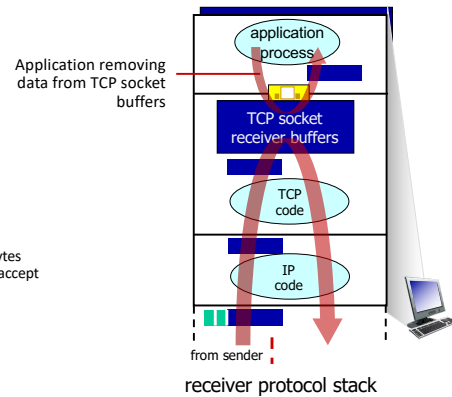


TCP flow control

Q: What happens if network layer delivers data faster than application layer removes data from socket buffers?



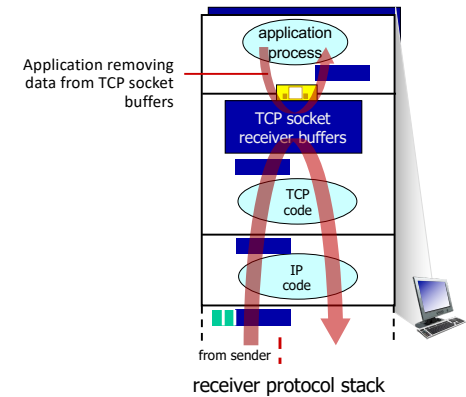
flow control: # bytes receiver willing to accept



TCP flow control

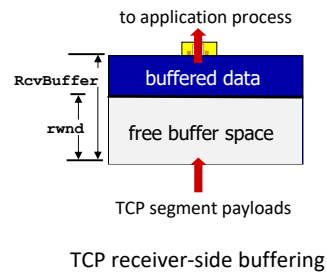
Q: What happens if network layer delivers data faster than application layer removes data from socket buffers?

flow control
receiver controls sender, so sender won't overflow receiver's buffer by transmitting too much, too fast



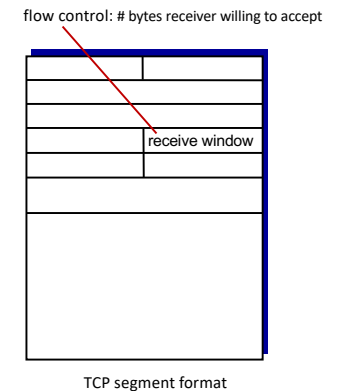
TCP flow control

- TCP receiver “advertises” free buffer space in **rwnd** field in TCP header
 - **RcvBuffer** size set via socket options (typical default is 4096 bytes)
 - many operating systems autoadjust **RcvBuffer**
- sender limits amount of unACKed (“in-flight”) data to received **rwnd**
- guarantees receive buffer will not overflow



TCP flow control

- TCP receiver “advertises” free buffer space in **rwnd** field in TCP header
 - **RcvBuffer** size set via socket options (typical default is 4096 bytes)
 - many operating systems autoadjust **RcvBuffer**
- sender limits amount of unACKed (“in-flight”) data to received **rwnd**
- guarantees receive buffer will not overflow



TCP segment format

- What does TCP do?
 - ARQ windowing, set-up, tear-down
- Flow Control in TCP
- Congestion Control in TCP

We have seen:

- **Flow control:** adjusting the sending rate to keep from overwhelming a slow *receiver*

Now lets attend...

- **Congestion control:** adjusting the sending rate to keep from overloading the *network*

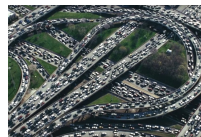
165

166

Principles of congestion control

Congestion:

- informally: “too many sources sending too much data too fast for *network* to handle”
- manifestations:
 - long delays (queueing in router buffers)
 - packet loss (buffer overflow at routers)
- different from flow control!
- a top-10 problem!



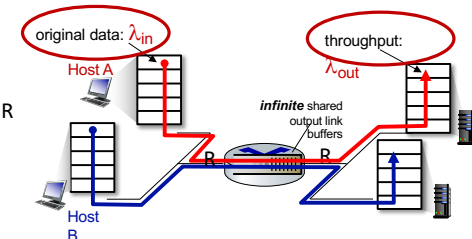
congestion control: too many senders, sending too fast

flow control: one sender too fast for one receiver

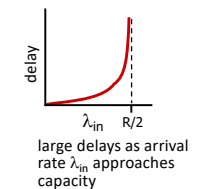
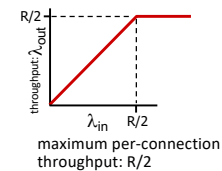
Causes/costs of congestion: scenario 1

Simplest scenario:

- one router, infinite buffers
- input, output link capacity: R
- two flows
- no retransmissions needed

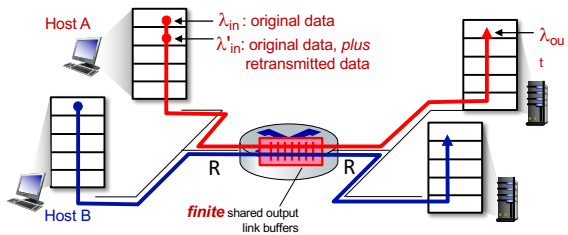


Q: What happens as arrival rate λ_{in} approaches $R/2$?



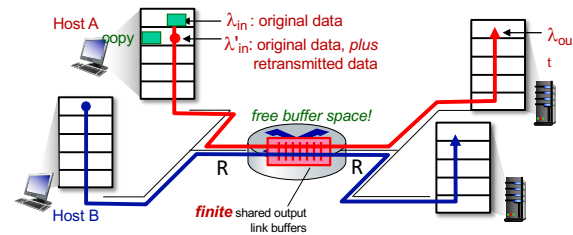
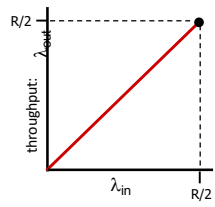
Causes/costs of congestion: scenario 2

- one router, *finite* buffers
- sender retransmits lost, timed-out packet
 - application-layer input = application-layer output: $\lambda_{in} = \lambda_{out}$
 - transport-layer input includes *retransmissions*: $\lambda'_{in} \geq \lambda_{in}$



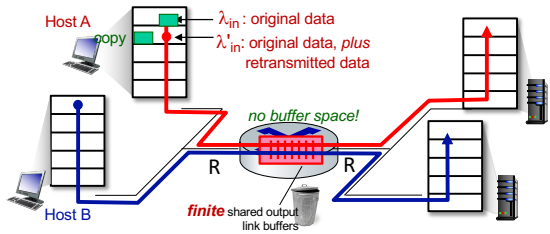
Causes/costs of congestion: scenario 2

- Idealization: *perfect* knowledge
- sender sends only when router buffers available



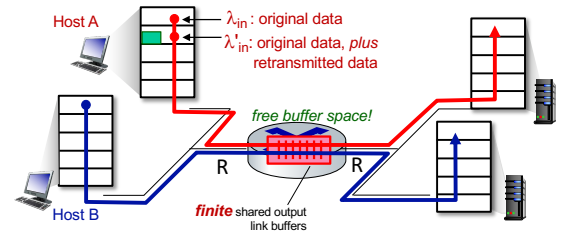
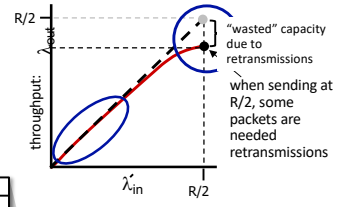
Causes/costs of congestion: scenario 2

- Idealization: *some* perfect knowledge
- packets can be lost (dropped at router) due to full buffers
 - sender knows when packet has been dropped: only resends if packet *known* to be lost



Causes/costs of congestion: scenario 2

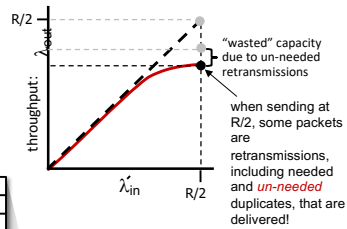
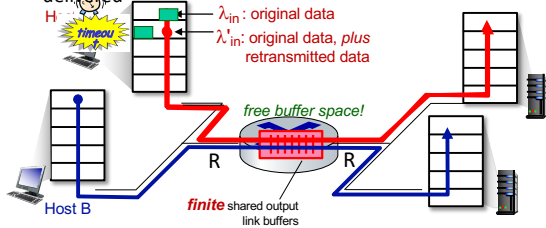
- Idealization: *some* perfect knowledge
- packets can be lost (dropped at router) due to full buffers
 - sender knows when packet has been dropped: only resends if packet *known* to be lost



Causes/costs of congestion: scenario 2

Realistic scenario: *un-needed duplicates*

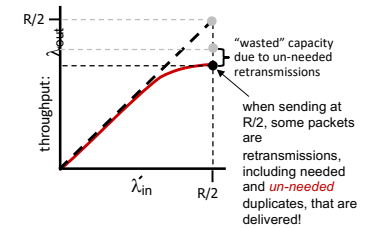
- packets can be lost, dropped at router due to full buffers – requiring retransmissions
- but sender times can time out prematurely, sending *two* copies, *both* of which are delivered



Causes/costs of congestion: scenario 2

Realistic scenario: *un-needed duplicates*

- packets can be lost, dropped at router due to full buffers – requiring retransmissions
- but sender times can time out prematurely, sending *two* copies, *both* of which are delivered



"costs" of congestion:

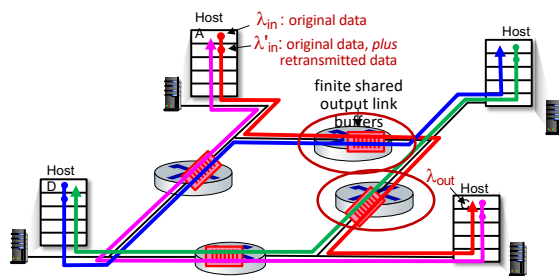
- more work (retransmission) for given receiver throughput
- unneeded retransmissions: link carries multiple copies of a packet
 - decreasing maximum achievable throughput

Causes/costs of congestion: scenario 3

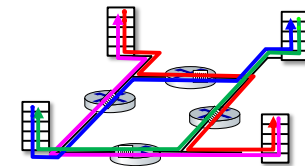
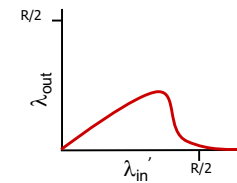
- four senders
- multi-hop paths
- timeout/retransmit

Q: what happens as λ_{in} and λ'_{in} increase ?

A: as red λ'_{in} increases, all arriving blue pkts at upper queue are dropped, blue throughput $\rightarrow 0$



Causes/costs of congestion: scenario 3

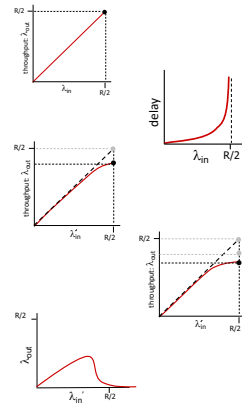


another "cost" of congestion:

- when packet dropped, any upstream transmission capacity and buffering used for that packet was wasted!

Causes/costs of congestion: insights

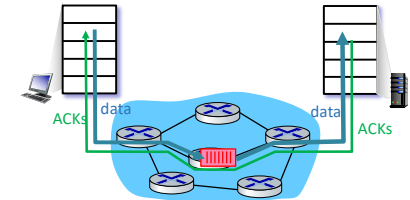
- throughput can never exceed capacity
- delay increases as capacity approached
- loss/retransmission decreases effective throughput
- un-needed duplicates further decreases effective throughput
- upstream transmission capacity / buffering wasted for packets lost downstream



Approaches towards congestion control

End-end congestion control:

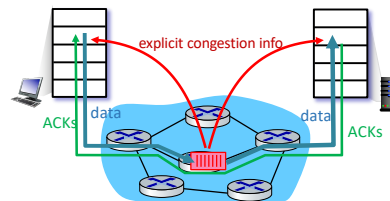
- no explicit feedback from network
- congestion *inferred* from observed loss, delay
- approach taken by TCP



Approaches towards congestion control

Network-assisted congestion control:

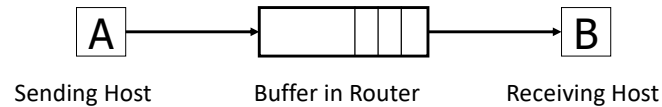
- routers provide *direct* feedback to sending/receiving hosts with flows passing through congested router
- may indicate congestion level or explicitly set sending rate
- TCP ECN, ATM, DECBbit protocols



Three Issues to Consider

- Discovering the available (bottleneck) bandwidth
- Adjusting to variations in bandwidth
- Sharing bandwidth between flows

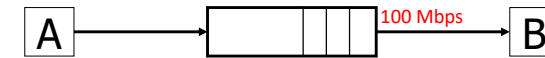
Abstract View



- Ignore internal structure of router and model it as having a single queue for a particular input-output pair

183

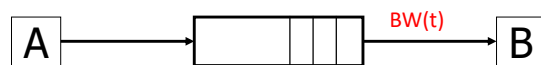
Discovering available bandwidth



- Pick sending rate to match bottleneck bandwidth
 - Without any *a priori* knowledge
 - Could be gigabit link, could be a modem

184

Adjusting to variations in bandwidth



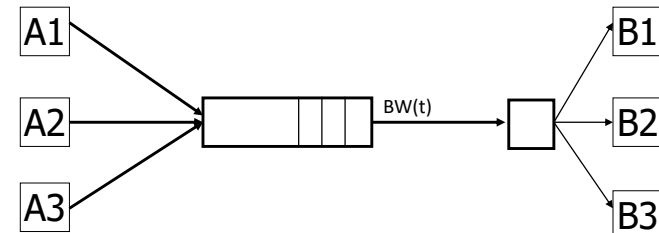
- Adjust rate to match **instantaneous** bandwidth
 - Assuming you have rough idea of bandwidth

185

Multiple flows and sharing bandwidth

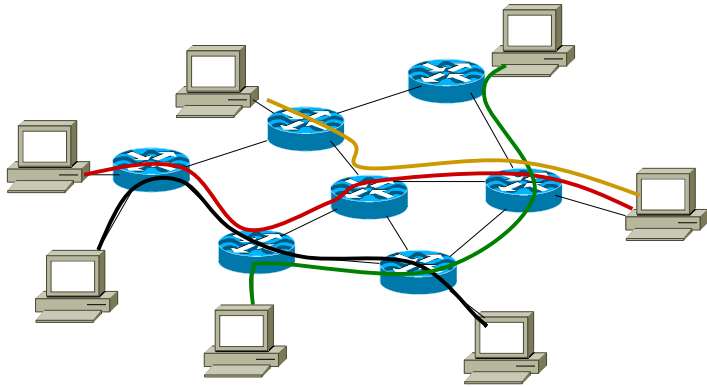
Two Issues:

- Adjust total sending rate to match bandwidth
- Allocation of bandwidth between flows



186

Reality

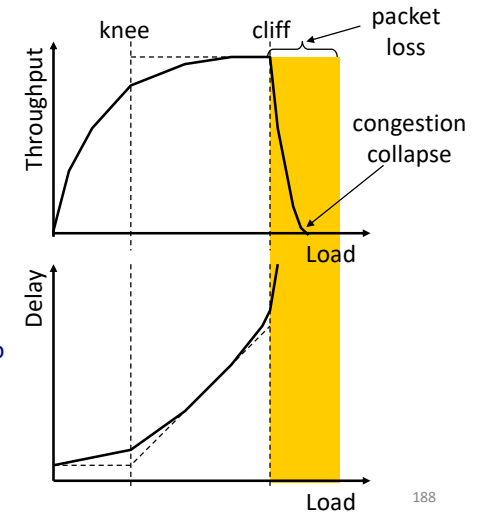


Congestion control is a resource allocation problem involving many flows, many links, and complicated global dynamics

187

View from a single flow

- Knee – point after which
 - Throughput increases slowly
 - Delay increases fast
- Cliff – point after which
 - Throughput starts to drop to zero (congestion collapse)
 - Delay approaches infinity



188

General Approaches

- (0) Send without care
 - Many packet drops

189

General Approaches

- (0) Send without care
- (1) Reservations
 - Pre-arrange bandwidth allocations
 - Requires negotiation before sending packets
 - Low utilization

190

General Approaches

- (0) Send without care
- (1) Reservations
- (2) Pricing
 - Don't drop packets for the high-bidders
 - Requires payment model

191

General Approaches

- (0) Send without care
- (1) Reservations
- (2) Pricing
- (3) Dynamic Adjustment
 - Hosts probe network; infer level of congestion; adjust
 - Network reports congestion level to hosts; hosts adjust
 - Combinations of the above
 - Simple to implement but suboptimal, messy dynamics

192

General Approaches

- (0) Send without care
- (1) Reservations
- (2) Pricing
- (3) Dynamic Adjustment

All three techniques have their place

- *Generality* of dynamic adjustment has proven powerful
- Doesn't presume business model, traffic characteristics, application requirements; does assume good citizenship

193

Who Takes Care of Congestion?

- Network? End hosts? Both?
- TCP's approach:
 - **End hosts** adjust sending rate
 - Based on **implicit feedback** from network
- Not the only approach
 - A consequence of history rather than planning

194

Some History: TCP in the 1980s

- Sending rate only limited by flow control
 - Packet drops → senders (repeatedly!) retransmit a full window's worth of packets
- Led to “congestion collapse” starting Oct. 1986
 - Throughput on the NSF network dropped from 32Kbits/s to 40bits/sec
- “Fixed” by Van Jacobson’s development of TCP’s congestion control (CC) algorithms

195

Jacobson’s Approach

- Extend TCP’s existing window-based protocol but adapt the window size in response to congestion
 - required no upgrades to routers or applications!
 - patch of a few lines of code to TCP implementations
- A pragmatic and effective solution
 - but many other approaches exist
- Extensively improved on since
 - topic now sees less activity in ISP contexts
 - but is making a comeback in datacenter environments

196

TCP’s Approach in a Nutshell

- TCP connection has window
 - Controls number of packets in flight
- Sending rate: $\sim \text{Window} / \text{RTT}$
- Vary window size to control sending rate

197

Windows, Buffers, and TCP



198

Windows, Buffers, and TCP

- TCP connection has a window
 - Controls number of packets in flight; filling a channel to improve throughput, and vary window size to control sending rate
- Buffers adapt mis-matched channels
 - Buffers smooth bursts
 - Adapt (re-time) arrivals for multiplexing

199

Windows, Buffers, and TCP

Buffers & TCP can make link utilization 100%

but

Buffers add delay, **variable** delay



200

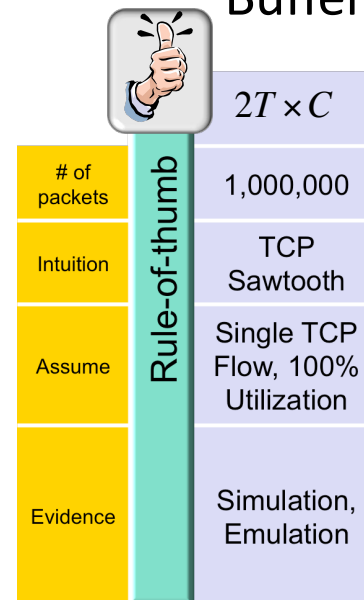
Sizing Buffers in Routers



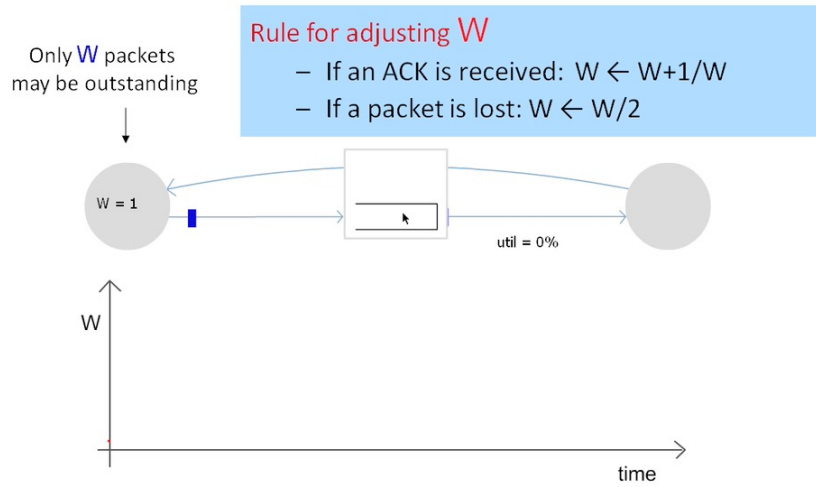
- Packet loss
 - Queue overload, and subsequent packet loss
- End-to-end delay
 - Transmission, propagation, and queueing delay
 - The only variable part is queueing delay
- Router architecture
 - Board space, power consumption, and cost
 - On chip buffers: higher density, higher capacity

201

Buffer Sizing Story

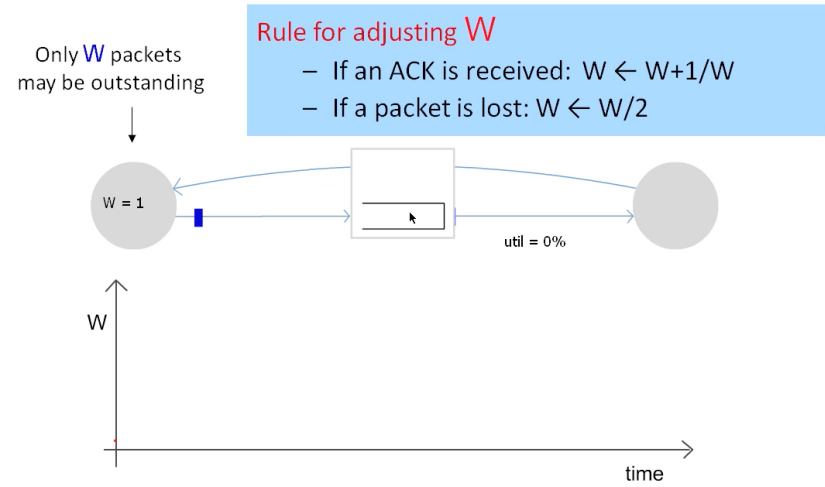


Continuous ARQ (TCP) adapting to congestion



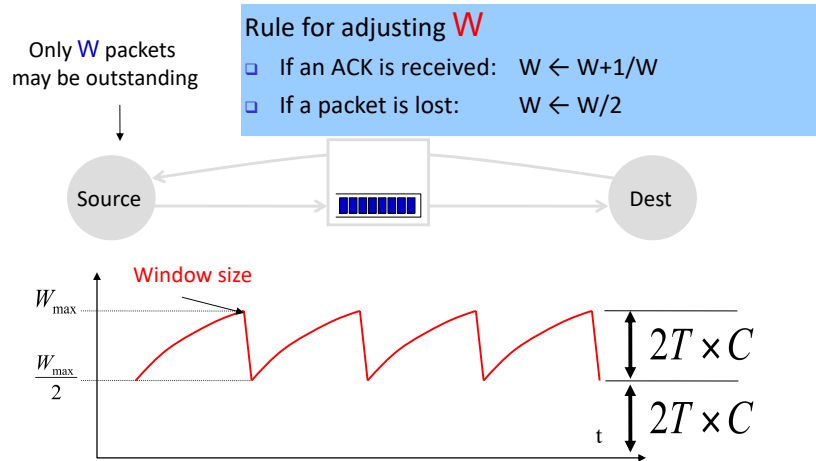
203

Continuous ARQ (TCP) adapting to congestion



204

Rule-of-thumb – Intuition



205

Buffers in Routers

So how large should the buffers be?

Buffer size matters

- Packet loss
 - Queue overload, and subsequent packet loss
- End-to-end delay
 - Transmission, propagation, and queuing delay
 - The only variable part is queuing delay

206



Buffers in Routers

So how large should the buffers be?

Buffer size matters

- Packet loss
 - Queue overload, and subsequent packet loss
- End-to-end delay
 - Transmission, propagation, and queuing delay
 - The only variable part is queuing delay
- Router architecture
 - Board space, power consumption, and cost
 - On chip buffers: higher density, higher capacity

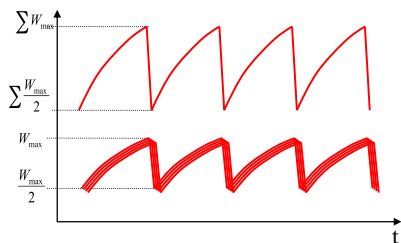
Buffer Sizing Story

		$2T \times C$		$\frac{2T \times C}{\sqrt{n}}$
# of packets	Rule-of-thumb	1,000,000	Small Buffers	10,000
Intuition		TCP Sawtooth		Sawtooth Smoothing
Assume		Single TCP Flow, 100% Utilization		Many Flows, 100% Utilization
Evidence		Simulation, Emulation		Simulations, Test-bed and Real Network Experiments

Small Buffers – Intuition

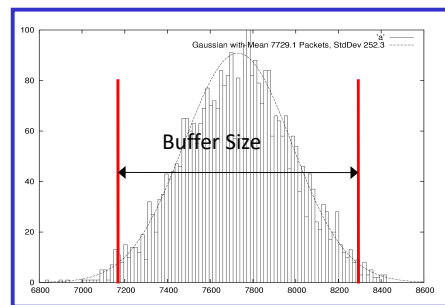
Synchronized Flows

- Aggregate window has same dynamics
- Therefore buffer occupancy has same dynamics
- Rule-of-thumb still holds.






Many TCP Flows

- Independent, desynchronized
- Central limit theorem says the aggregate becomes Gaussian
- Variance (buffer size) decreases as N increases



Buffer Sizing Story

		$2T \times C$		$\frac{2T \times C}{\sqrt{n}}$		$O(\log W)$
# of packets	Rule-of-thumb	1,000,000	Small Buffers	10,000	Large Buffers	20 - 50
Intuition		TCP Sawtooth		Sawtooth Smoothing	Non-bursty Arrivals	
Assume		Single TCP Flow, 100% Utilization		Many Synchronized TCP connections? Just TCP – what about other applications? Small BDP link? Large BDP link? How many devices? W of flows? How many flows?	Paced TCP, 85-90% Utilization	
Evidence		Simulation, Emulation		Simulations, Test-bed and Real Network Experiments	Simulations, Test-bed and Real Network Experiments	



TCP's Approach in a Nutshell

- TCP connection has window
 - Controls number of packets in flight
- Sending rate: $\sim \text{Window}/\text{RTT}$
- **Vary window size to control sending rate**

211

Note

- This lecture will talk about CWND in units of MSS
 - (Recall MSS: Maximum Segment Size, the amount of payload data in a TCP packet)
 - This is only for pedagogical purposes
- **In reality this is a LIE:** Real implementations maintain CWND in bytes

213

All These Windows...

- Congestion Window: **CWND**
 - How many bytes can be sent without overflowing routers
 - Computed by the sender using congestion control algorithm
- Flow control window: **AdvertisedWindow (RWND)**
 - How many bytes can be sent without overflowing receiver's buffers
 - Determined by the receiver and reported to the sender
- Sender-side window = **minimum{CWND, RWND}**
 - Assume for this material that $\text{RWND} \gg \text{CWND}$

212

Two Basic Questions

- How does the sender detect congestion?
- How does the sender adjust its sending rate?
 - **To address three issues**
 - Finding available bottleneck bandwidth
 - Adjusting to bandwidth variations
 - Sharing bandwidth

214

(Recall) Detecting Congestion

- Packet delays
 - Tricky: noisy signal (delay often varies considerably)
- Router tell end-hosts they're congested
- Packet loss
 - Fail-safe signal that TCP already has to detect
 - Complication: non-congestive loss (checksum errors)
- Two indicators of packet loss
 - No ACK after certain time interval: **timeout**
 - Multiple **duplicate ACKs**

215

Rate Adjustment

- Basic structure:
 - Upon receipt of ACK (of new data): increase rate
 - Upon detection of loss: decrease rate
- How we increase/decrease the rate depends on the phase of congestion control we're in:
 - Discovering available bottleneck bandwidth vs.
 - Adjusting to bandwidth variations

217

Not All Losses the Same

- Duplicate ACKs: isolated loss
 - Still getting ACKs
- Timeout: much more serious
 - Not enough packets in progress to trigger duplicate-acks, OR
 - Suffered several losses
- We will adjust rate differently for each case

216

Bandwidth Discovery with Slow Start

- Goal: estimate available bandwidth
 - start slow (for safety)
 - but ramp up quickly (for efficiency)
- Consider
 - RTT = 100ms, MSS=1000bytes
 - Window size to fill 1Mbps of BW = 12.5 packets
 - Window size to fill 1Gbps = 12,500 packets
 - Either is possible!

218

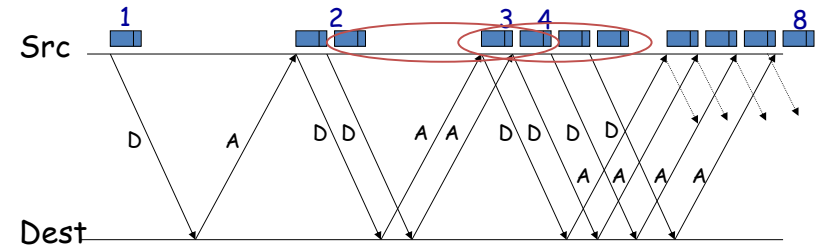
“Slow Start” Phase

- Sender starts at a slow rate but increases **exponentially** until first loss
- Start with a small congestion window
 - Initially, CWND = 1
 - So, initial sending rate is MSS/RTT
- Double the CWND for each RTT with no loss

219

Slow Start in Action

- For each RTT: double CWND
- Simpler implementation: for each ACK, CWND += 1



220

Adjusting to Varying Bandwidth

- Slow start gave an estimate of available bandwidth
- Now, want to track variations in this available bandwidth, oscillating around its current value
 - Repeated probing (rate increase) and backoff (rate decrease)
- TCP uses: “Additive Increase Multiplicative Decrease” (AIMD)
 - We’ll see why shortly...

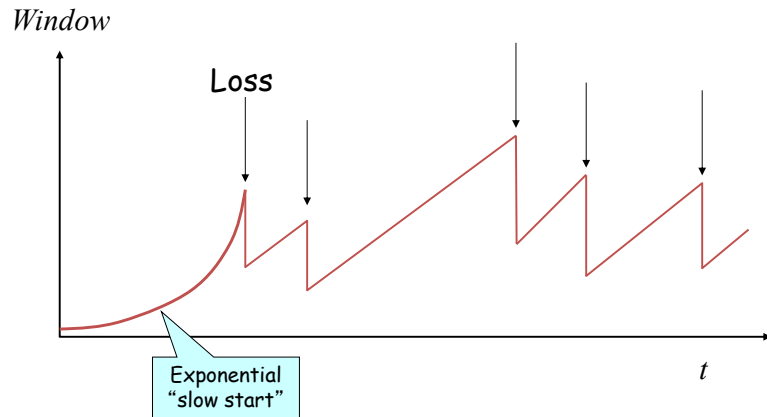
221

AIMD

- Additive increase
 - Window grows by one MSS for every RTT with no loss
 - For each successful RTT, $CWND = CWND + 1$
 - Simple implementation:
 - for each ACK, $CWND = CWND + 1/CWND$
- Multiplicative decrease
 - On loss of packet, divide congestion window in **half**
 - On loss, $CWND = CWND/2$

222

Leads to the TCP “Sawtooth”



223

- What does TCP do?
 - ARQ windowing, set-up, tear-down
- Flow Control in TCP
- Congestion Control in TCP
 - AIMD (slow-start, congestion avoidance)

225

Slow-Start vs. AIMD

- When does a sender stop Slow-Start and start Additive Increase?
- Introduce a “slow start threshold” (**ssthresh**)
 - Initialized to a large value
 - On timeout, $ssthresh = CWND/2$
- When $CWND = ssthresh$, sender switches from slow-start to AIMD-style increase

224
224

- What does TCP do?
 - ARQ windowing, set-up, tear-down
- Flow Control in TCP
- Congestion Control in TCP
 - AIMD (slow-start, congestion avoidance) and Fast-Recovery

226

One Final Phase: Fast Recovery

- The problem: congestion avoidance too slow in recovering from an isolated loss

227

The problem – A timeline

- ACK 101 (due to 102) $cwnd=10$ dupACK#1 (no xmit)
- ACK 101 (due to 103) $cwnd=10$ dupACK#2 (no xmit)
- ACK 101 (due to 104) $cwnd=10$ dupACK#3 (no xmit)
- RETRANSMIT 101 $ssthresh=5$ $cwnd=5$
- ACK 101 (due to 105) $cwnd=5 + 1/5$ (no xmit)
- ACK 101 (due to 106) $cwnd=5 + 2/5$ (no xmit)
- ACK 101 (due to 107) $cwnd=5 + 3/5$ (no xmit)
- ACK 101 (due to 108) $cwnd=5 + 4/5$ (no xmit)
- ACK 101 (due to 109) $cwnd=5 + 5/5$ (no xmit)
- ACK 101 (due to 110) $cwnd=6 + 1/5$ (no xmit)
- ACK 111 (due to 101) ← only now can we transmit new packets
- Plus no packets in flight so ACK “clocking” (to increase CWND) stalls for another RTT

229

Example (in units of MSS, not bytes)

- Consider a TCP connection with:
 - $CWND=10$ packets
 - Last ACK was for packet # 101
 - i.e., receiver expecting next packet to have seq. no. 101
- 10 packets [101, 102, 103,..., 110] are in flight
 - Packet 101 is dropped
 - What ACKs do they generate?
 - And how does the sender respond?

228

Solution: Fast Recovery

Idea: Grant the sender temporary “credit” for each dupACK so as to keep packets in flight

- If $dupACKcount = 3$
 - $ssthresh = cwnd/2$
 - $cwnd = ssthresh + 3$
- While in fast recovery
 - $cwnd = cwnd + 1$ for each additional duplicate ACK
- Exit fast recovery after receiving new ACK
 - set $cwnd = ssthresh$

230

Example

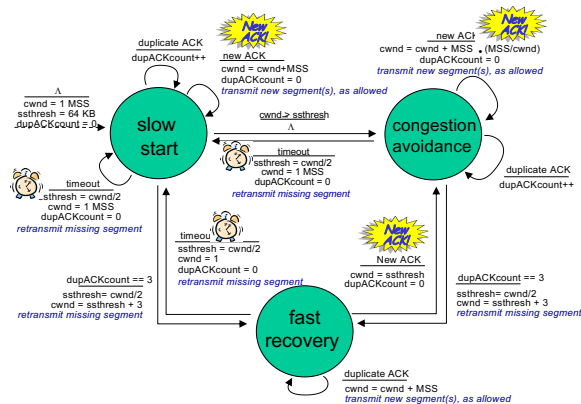
- Consider a TCP connection with:
 - CWND=10 packets
 - Last ACK was for packet # 101
 - i.e., receiver expecting next packet to have seq. no. 101
- 10 packets [101, 102, 103,..., 110] are in flight
 - Packet 101 is dropped

Timeline

- ACK 101 (due to 102) cwnd=10 dup#1
- ACK 101 (due to 103) cwnd=10 dup#2
- ACK 101 (due to 104) cwnd=10 dup#3
- REXMIT 101 ssthresh=5 cwnd= 8 (5+3)
- ACK 101 (due to 105) cwnd= 9 (no xmit)
- ACK 101 (due to 106) cwnd=10 (no xmit)
- ACK 101 (due to 107) cwnd=11 (xmit 111)
- ACK 101 (due to 108) cwnd=12 (xmit 112)
- ACK 101 (due to 109) cwnd=13 (xmit 113)
- ACK 101 (due to 110) cwnd=14 (xmit 114)
- ACK 111 (due to 101) cwnd = 5 (xmit 115) ← exiting fast recovery
- Packets 111-114 already in flight
- ACK 112 (due to 111) cwnd = 5 + 1/5 ← back in congestion avoidance

231

Summary: TCP congestion control

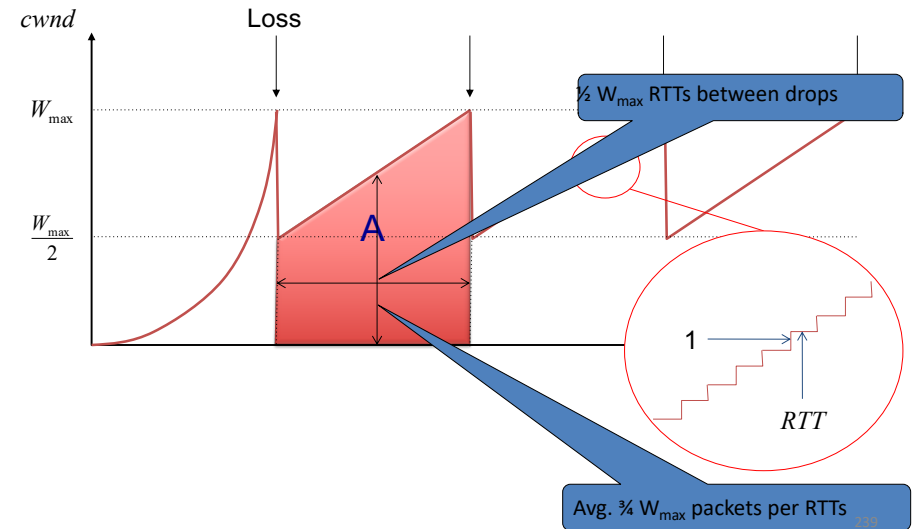


Variant	Feedback	Required changes	Benefits	Fairness
(New) Reno	Loss	—	—	Delay
Vegas	Delay	Sender	Less loss	Proportional
High Speed	Loss	Sender	High bandwidth	
BIC	Loss	Sender	High bandwidth	
CUBIC	Loss	Sender	High bandwidth	
C2TCP ^{[11][12]}	Loss/Delay	Sender	Ultra-low latency and high bandwidth	
NATCP ^[13]	Multi-bit signal	Sender	Near Optimal Performance	
Elastic-TCP	Loss/Delay	Sender	High bandwidth/short & long-distance	
Agile-TCP	Loss	Sender	High bandwidth/short-distance	
H-TCP	Loss	Sender	High bandwidth	
FAST	Delay	Sender	High bandwidth	Proportional
Compound TCP	Loss/Delay	Sender	High bandwidth	Proportional
Westwood	Loss/Delay	Sender	L	
Jersey	Loss/Delay	Sender	L	
BBR ^[14]	Delay	Sender	BLVC, Bufferbloat	
CLAMP	Multi-bit signal	Receiver, Router	V	Max-min
TFRC	Loss	Sender, Receiver	No Retransmission	Minimum delay
XCP	Multi-bit signal	Sender, Receiver, Router	BLFC	Max-min
VCP	2-bit signal	Sender, Receiver, Router	BLF	Proportional
MaxNet	Multi-bit signal	Sender, Receiver, Router	BLFSC	Max-min
JetMax	Multi-bit signal	Sender, Receiver, Router	High bandwidth	Max-min
RED	Loss	Router	Reduced delay	
ECN	Single-bit signal	Sender, Receiver, Router	Reduced loss	

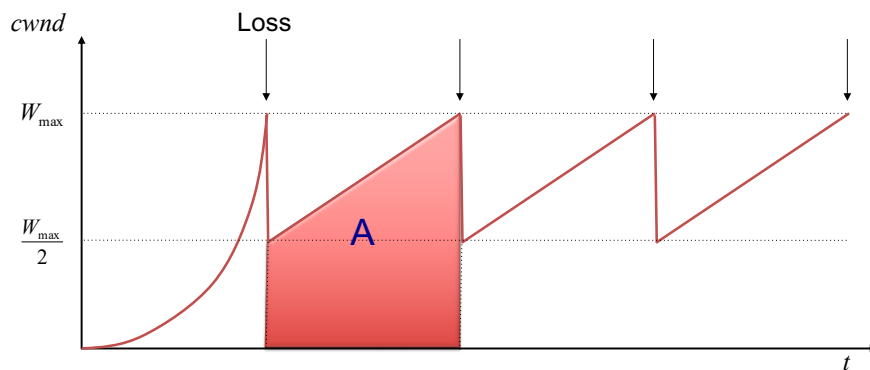
TCP Throughput Equation

238

A Simple Model for TCP Throughput



A Simple Model for TCP Throughput



Packet drop rate, $p = 1 / A$, where $A = \frac{3}{8} W_{max}^2$

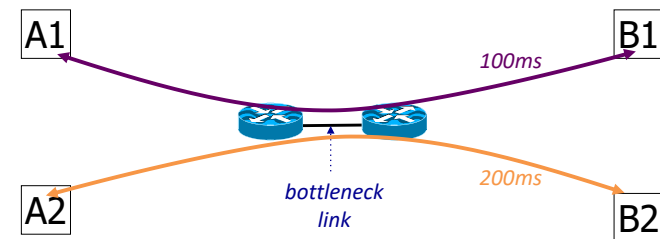
$$\text{Throughput, } B = \frac{A}{\left(\frac{W_{max}}{2}\right) RTT} = \sqrt{\frac{3}{2}} \frac{1}{RTT \sqrt{p}}$$

240

Implications (1): Different RTTs

$$\text{Throughput} = \sqrt{\frac{3}{2}} \frac{1}{RTT \sqrt{p}}$$

- Flows get throughput inversely proportional to RTT
- TCP unfair in the face of heterogeneous RTTs!



241

Implications (2): High Speed TCP

$$\text{Throughput} = \sqrt{\frac{3}{2}} \frac{1}{RTT \sqrt{p}}$$

- Assume $RTT = 100\text{ms}$, $MSS=1500\text{bytes}$
- What value of p is required to reach 100Gbps throughput
 - $\sim 2 \times 10^{-12}$
- How long between drops?
 - ~ 16.6 hours
- How much data has been sent in this time?
 - ~ 6 petabits
- These are not practical numbers!

242

Adapting TCP to High Speed

- Once past a threshold speed, increase CWND faster
 - A proposed standard [Floyd'03]: once speed is past some threshold, change equation to p^{-8} rather than p^{-5}
 - Let the additive constant in AIMD depend on CWND
- Other approaches?
 - Multiple simultaneous connections (*hacky* but works today)
 - Router-assisted approaches (will see shortly)

243

Implications (3): Rate-based CC

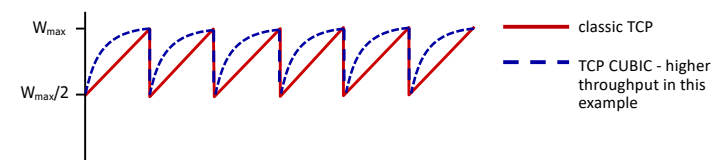
$$\text{Throughput} = \sqrt{\frac{3}{2}} \frac{1}{RTT \sqrt{p}}$$

- TCP throughput is “choppy”
 - repeated swings between $W/2$ to W
- Some apps would prefer sending at a steady rate
 - e.g., streaming apps
- A solution: “Equation-Based Congestion Control”
 - ditch TCP’s increase/decrease rules and just follow the equation
 - measure drop percentage p , and set rate accordingly
- Following the TCP equation ensures we’re “TCP friendly”
 - i.e., use no more than TCP does in similar setting

244

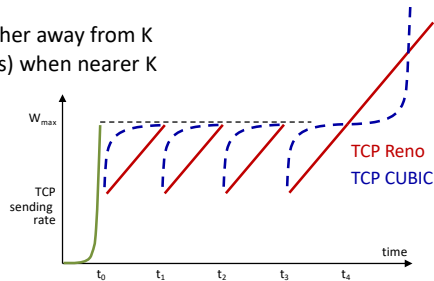
TCP CUBIC

- Is there a better way than AIMD to “probe” for usable bandwidth?
- Insight/intuition:
 - W_{\max} : sending rate at which congestion loss was detected
 - congestion state of bottleneck link probably (?) hasn’t changed much
 - after cutting rate/window in half on loss, initially ramp up to W_{\max} *faster*, but then approach W_{\max} more *slowly*



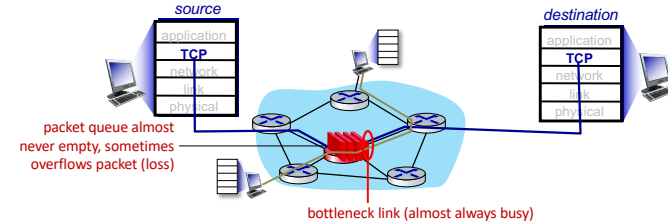
TCP CUBIC

- K: point in time when TCP window size will reach W_{max}
 - K itself is tuneable
- increase W as a function of the *cube* of the distance between current time and K
 - larger increases when further away from K
 - smaller increases (cautious) when nearer K
- TCP CUBIC default in Linux, most popular TCP for popular Web servers



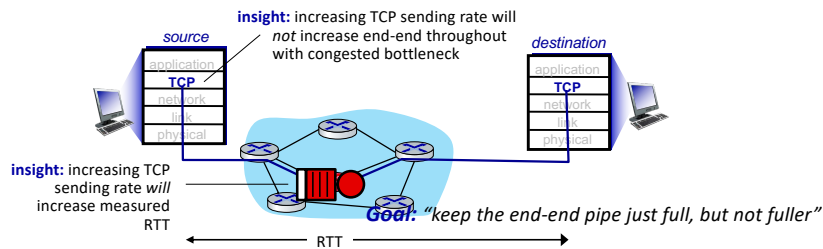
TCP and the congested “bottleneck link”

- TCP (classic, CUBIC) increase TCP’s sending rate until packet loss occurs at some router’s output: the *bottleneck link*



TCP and the congested “bottleneck link”

- TCP (classic, CUBIC) increase TCP’s sending rate until packet loss occurs at some router’s output: the *bottleneck link*
- understanding congestion: useful to focus on congested bottleneck link



Delay-based TCP Congestion Control

Keeping sender-to-receiver pipe “just full enough, but no fuller”: keep bottleneck link busy transmitting, but avoid high delays/buffering



Delay-based approach:

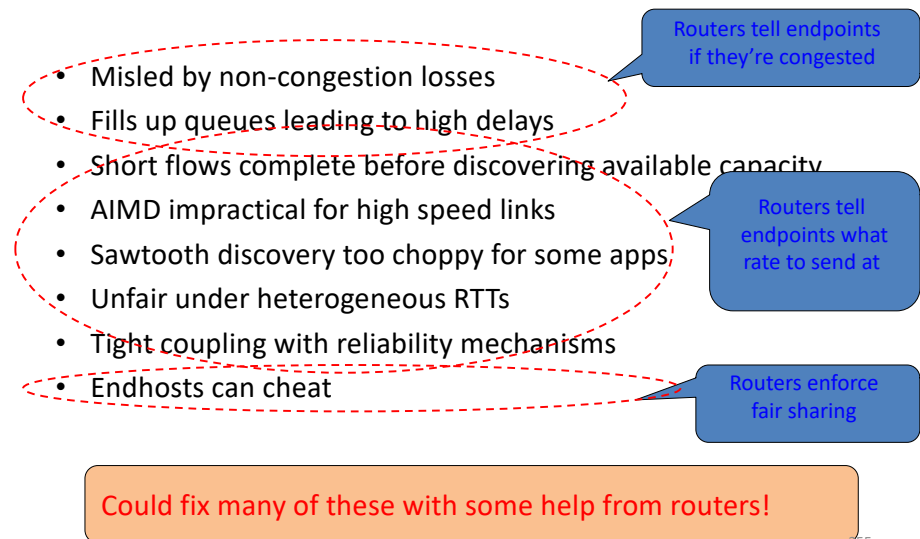
- RTT_{min} - minimum observed RTT (uncongested path)
- uncongested throughput with congestion window $cwnd$ is $cwnd/RTT_{min}$

if measured throughput “very close” to uncongested throughput
 increase $cwnd$ linearly /* since path not congested */
 else if measured throughput “far below” uncongested throughput
 decrease $cwnd$ linearly /* since path is congested */

Recap: TCP problems

Delay-based TCP Congestion Control

- congestion control without inducing/forcing loss
- maximizing throughput (“keeping the just pipe full...”) while keeping delay low (“...but not fuller”)
- a number of deployed TCPs take a delay-based approach
 - BBR deployed on Google’s (internal) backbone network



Router-Assisted Congestion Control

- Three tasks for CC:
 - Isolation/fairness
 - Adjustment*
 - Detecting congestion

How can routers ensure each flow gets its “fair share”?

* This may be *automatic* eg loss-response of TCP

Fairness: General Approach

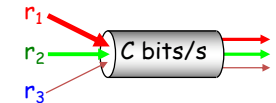
- Routers classify packets into “flows”
 - (For now) flows are packets between same source/destination
- Each flow has its own FIFO queue in router
- Router services flows in a fair fashion
 - When line becomes free, take packet from next flow in a fair order
- What does “fair” mean exactly?

258

Max-Min Fairness

- Given set of bandwidth demands r_i and total bandwidth C , max-min bandwidth allocations are:

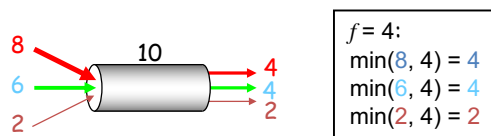
$$a_i = \min(f, r_i)$$
 where f is the unique value such that $\text{Sum}(a_i) = C$



259

Example

- $C = 10$; $r_1 = 8$, $r_2 = 6$, $r_3 = 2$; $N = 3$
- $C/3 = 3.33 \rightarrow$
 - Can service all of r_3
 - Remove r_3 from the accounting: $C = C - r_3 = 8$; $N = 2$
- $C/2 = 4 \rightarrow$
 - Can't service all of r_1 or r_2
 - So hold them to the remaining fair share: $f = 4$



260

Max-Min Fairness

- Given set of bandwidth demands r_i and total bandwidth C , max-min bandwidth allocations are:

$$a_i = \min(f, r_i)$$
- where f is the unique value such that $\text{Sum}(a_i) = C$
- Property:
 - If you don't get full demand, no one gets more than you
- This is what round-robin service gives if all packets are the same size

261

How do we deal with packets of different sizes?

- Mental model: Bit-by-bit round robin (“fluid flow”)
- Can you do this in practice?
- No, packets cannot be preempted
- But we can approximate it
 - This is what “fair queuing” routers do

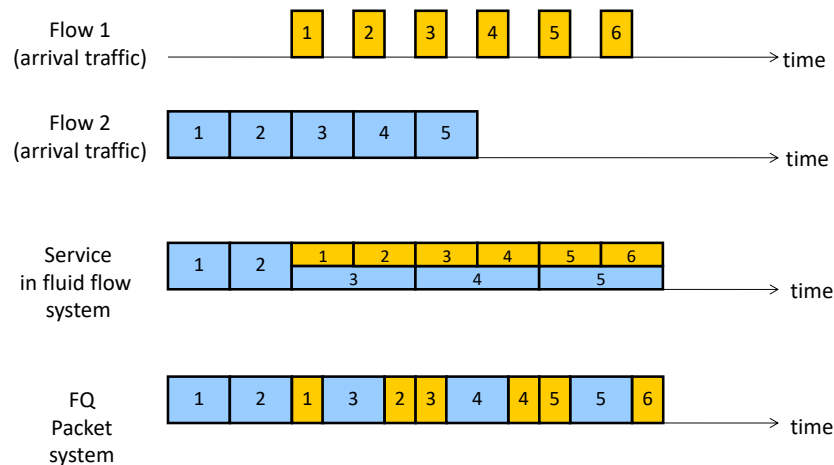
262

Fair Queuing (FQ)

- For each packet, compute the time at which the last bit of a packet would have left the router *if* flows are served bit-by-bit
- Then serve packets in the increasing order of their deadlines

263

Example



264

Fair Queuing (FQ)

- Think of it as an implementation of round-robin generalized to the case where not all packets are equal sized
- **Weighted** fair queuing (WFQ): assign different flows different shares
- Today, some form of WFQ implemented in almost all routers
 - Not the case in the 1980-90s, when CC was being developed
 - Mostly used to isolate traffic at larger granularities (e.g., per-prefix)

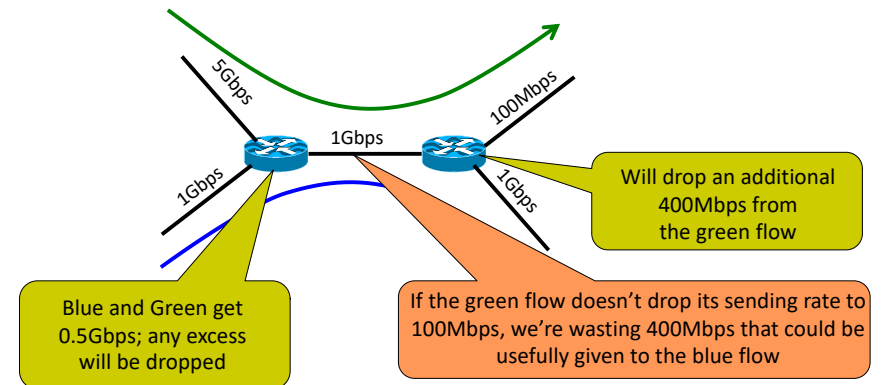
265

FQ vs. FIFO

- FQ advantages:
 - Isolation: cheating flows don't benefit
 - Bandwidth share does not depend on RTT
 - Flows can pick any rate adjustment scheme they want
- Disadvantages:
 - More complex than FIFO: per flow queue/state, additional per-packet book-keeping

FQ in the big picture

- FQ does not eliminate congestion → it just manages the congestion



FQ in the big picture

- FQ does not eliminate congestion → it just manages the congestion
 - robust to cheating, variations in RTT, details of delay, reordering, retransmission, *etc.*
- But congestion (and packet drops) still occurs
- And we still want end-hosts to discover/adapt to their fair share!
- What would the end-to-end argument say w.r.t. congestion control?

Fairness is a controversial goal

- What if you have 8 flows, and I have 4?
 - Why should you get twice the bandwidth
- What if your flow goes over 4 congested hops, and mine only goes over 1?
 - Why shouldn't you be penalized for using more scarce bandwidth?
- And what is a flow anyway?
 - TCP connection
 - Source-Destination pair?
 - Source?

Explicit Congestion Notification (ECN)

- Single bit in packet header; set by congested routers
 - If data packet has bit set, then ACK has ECN bit set
- Many options for when routers set the bit
 - tradeoff between (link) utilization and (packet) delay
- Congestion semantics can be exactly like that of drop
 - I.e., endhost reacts as though it saw a drop
- Advantages:
 - Don't confuse corruption with congestion; recovery w/ rate adjustment
 - Can serve as an early indicator of congestion to avoid delays
 - Easy (easier) to incrementally deploy
 - defined as extension to TCP/IP in RFC 3168 (uses diffserv bits in the IP header)

270

Transport Recap

A "big bag":

Multiplexing, reliability, error-detection, error-recovery, flow and congestion control,

- UDP:
 - Minimalist - multiplexing and error detection
- TCP:
 - somewhat hacky
 - but practical/deployable
 - good enough to have raised the bar for the deployment of new approaches
 - though the needs of datacenters change the status quos
- Beyond TCP (discussed in Topic 6):
 - QUIC / application-aware transport layers

277

Explicit congestion notification (ECN)

TCP deployments often implement *network-assisted* congestion control:

- two bits in IP header (ToS field) marked *by network router* to indicate congestion
 - *policy* to determine marking chosen by network operator
- congestion indication carried to destination
- destination sets ECE bit on ACK segment to notify sender of congestion
- involves both IP (IP header ECN bit marking) and TCP (TCP header C,E bit marking)

