



Transport Layer

Jennifer Rexford

COS 461: Computer Networks

Lectures: MW 10-10:50am in Architecture N101

<http://www.cs.princeton.edu/courses/archive/spr12/cos461/>

IP Protocol Stack: Key Abstractions

Application	Applications	
Transport	Reliable streams	Messages
Network	Best-effort <i>global</i> packet delivery	
Link	Best-effort <i>local</i> packet delivery	

- Transport layer is where we “pay the piper”
 - Provide applications with good abstractions
 - Without support or feedback from the network

2

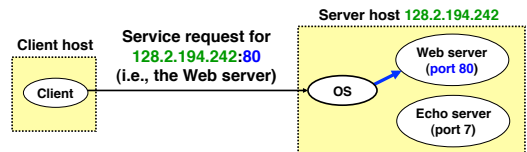
Transport Protocols

- Logical communication between processes
 - Sender divides a message into segments
 - Receiver reassembles segments into message
- Transport services
 - (De)multiplexing packets
 - Detecting corrupted data
 - Optionally: reliable delivery, flow control, ...

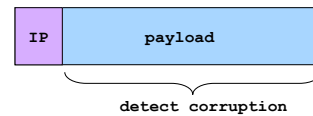
3

Two Basic Transport Features

- Demultiplexing: port numbers



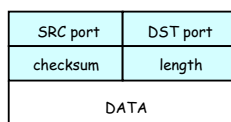
- Error detection: checksums



4

User Datagram Protocol (UDP)

- Datagram messaging service
 - Demultiplexing: port numbers
 - Detecting corruption: checksum
- Lightweight communication between processes
 - Send and receive messages
 - Avoid overhead of ordered, reliable delivery



5

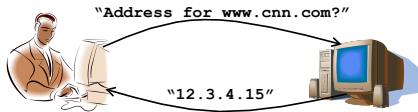
Advantages of UDP

- Fine-grain control
 - UDP sends as soon as the application writes
- No connection set-up delay
 - UDP sends without establishing a connection
- No connection state
 - No buffers, parameters, sequence #s, etc.
- Small header overhead
 - UDP header is only eight-bytes long

6

Popular Applications That Use UDP

- **Multimedia streaming**
 - Retransmitting packets is not always worthwhile
 - E.g., phone calls, video conferencing, gaming, IPTV
- **Simple query-response protocols**
 - Overhead of connection establishment is overkill
 - E.g., Domain Name System (DNS), DHCP, etc.



7

Transmission Control Protocol (TCP)

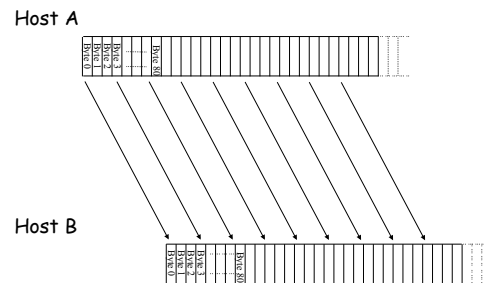
- **Stream-of-bytes service**
 - Sends and receives a stream of bytes
- **Connection oriented**
 - Explicit set-up and tear-down of TCP connection
- **Reliable, in-order delivery**
 - Corruption: checksums
 - Detect loss/reordering: sequence numbers
 - Reliable delivery: acknowledgments and retransmissions
- **Flow control**
 - Prevent overflow of the receiver's buffer space
- **Congestion control**
 - Adapt to network congestion for the greater good

8

Breaking a Stream of Bytes into TCP Segments

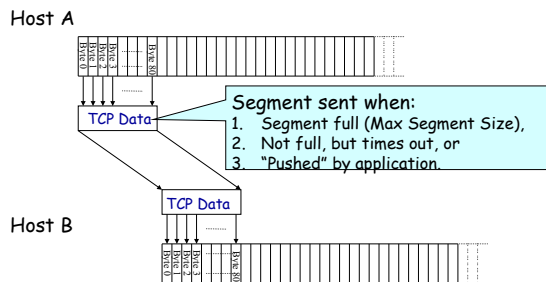
9

TCP "Stream of Bytes" Service



10

...Emulated Using TCP "Segments"

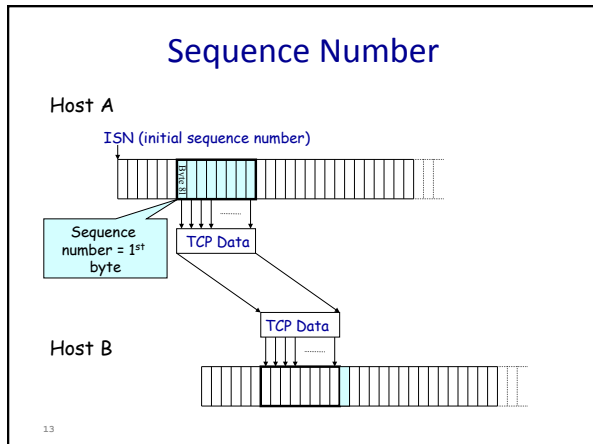


11

TCP Segment

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

12



- ### Initial Sequence Number (ISN)
- **Sequence number for the very first byte**
 - E.g., Why not a de facto ISN of 0?
 - **Practical issue: reuse of port numbers**
 - Port numbers must (eventually) get used again
 - ... and an old packet may still be in flight
 - ... and associated with the new connection
 - **So, TCP must change the ISN over time**
 - Set from a 32-bit clock that ticks every 4 microsec
 - ... which wraps around once every 4.55 hours!
- 14

Reliable Delivery on a Lossy Channel With Bit Errors

15

- ### Challenges of Reliable Data Transfer
- **Over a perfectly reliable channel**
 - Easy: sender sends, and receiver receives
 - **Over a channel with bit errors**
 - Receiver detects errors and requests retransmission
 - **Over a lossy channel with bit errors**
 - Some data are missing, and others corrupted
 - Receiver cannot always detect loss
 - **Over a channel that may reorder packets**
 - Receiver cannot distinguish loss from out-of-order
- 16

- ### An Analogy
- **Alice and Bob are talking**
 - What if Alice couldn't understand Bob?
 - Bob asks Alice to repeat what she said
 - **What if Bob hasn't heard Alice for a while?**
 - Is Alice just being quiet? Has she lost reception?
 - How long should Bob just keep on talking?
 - Maybe Alice should periodically say "uh huh"
 - ... or Bob should ask "Can you hear me now?" ☺
-
- 17

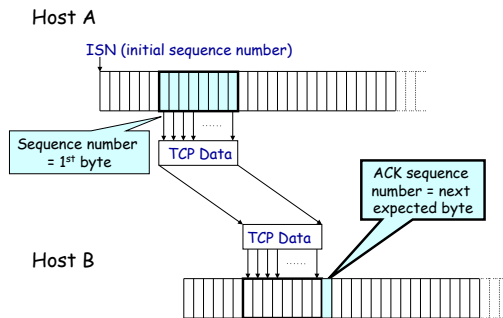
- ### Take-Aways from the Example
- **Acknowledgments from receiver**
 - Positive: "okay" or "uh huh" or "ACK"
 - Negative: "please repeat that" or "NACK"
 - **Retransmission by the sender**
 - After *not* receiving an "ACK"
 - After receiving a "NACK"
 - **Timeout by the sender ("stop and wait")**
 - Don't wait forever without some acknowledgment
- 18

TCP Support for Reliable Delivery

- **Detect bit errors: checksum**
 - Used to detect corrupted data at the receiver
 - ...leading the receiver to drop the packet
- **Detect missing data: sequence number**
 - Used to detect a gap in the stream of bytes
 - ... and for putting the data back in order
- **Recover from lost data: retransmission**
 - Sender retransmits lost or corrupted data
 - Two main ways to detect lost packets

19

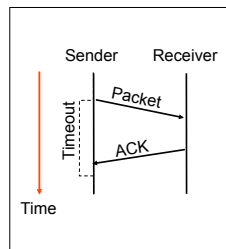
TCP Acknowledgments



20

Automatic Repeat reQuest (ARQ)

- **ACK and timeouts**
 - Receiver sends ACK when it receives packet
 - Sender waits for ACK and times out
- **Simplest ARQ protocol**
 - Stop and wait
 - Send a packet, stop and wait until ACK arrives



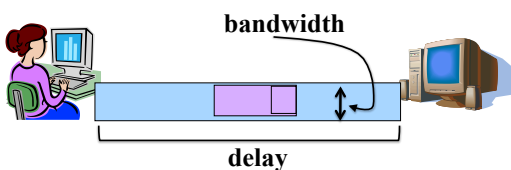
21

Flow Control: TCP Sliding Window

22

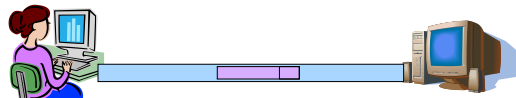
Motivation for Sliding Window

- **Stop-and-wait is inefficient**
 - Only one TCP segment is “in flight” at a time
 - Especially bad for high “delay-bandwidth product”



Numerical Example

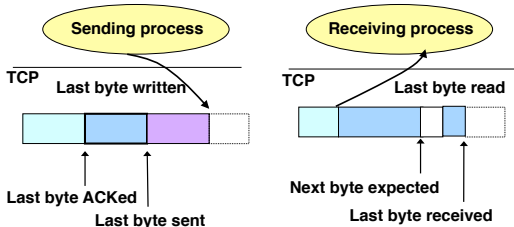
- **1.5 Mbps link with 45 msec round-trip time (RTT)**
 - Delay-bandwidth product is 67.5 Kbits (or 8 KBytes)
- **Sender can send at most one packet per RTT**
 - Assuming a segment size of 1 KB (8 Kbits)
 - 8 Kbits/segment at 45 msec/segment → 182 Kbps
 - That’s just *one-eighth* of the 1.5 Mbps link capacity



24

Sliding Window

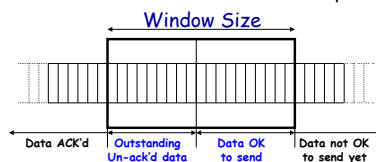
- Allow a larger amount of data “in flight”
 - Allow sender to get ahead of the receiver
 - ... though not too far ahead



25

Receiver Buffering

- Receive window size
 - Amount that can be sent without acknowledgment
 - Receiver must be able to store this amount of data
- Receiver tells the sender the window
 - Tells the sender the amount of free space left

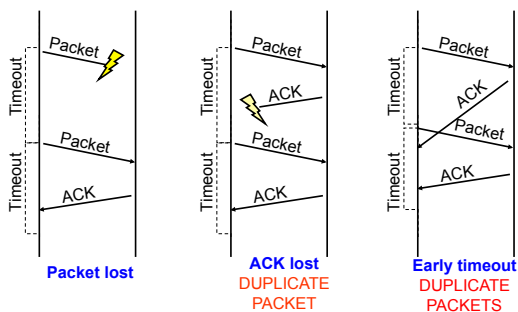


26

Optimizing Retransmissions

27

Reasons for Retransmission



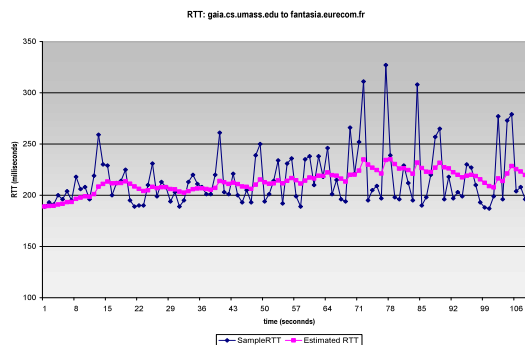
28

How Long Should Sender Wait?

- Sender sets a timeout to wait for an ACK
 - Too short: wasted retransmissions
 - Too long: excessive delays when packet lost
- TCP sets timeout as a function of the RTT
 - Expect ACK to arrive after an “round-trip time”
 - ... plus a fudge factor to account for queuing
- But, how does the sender know the RTT?
 - Running average of delay to receive an ACK

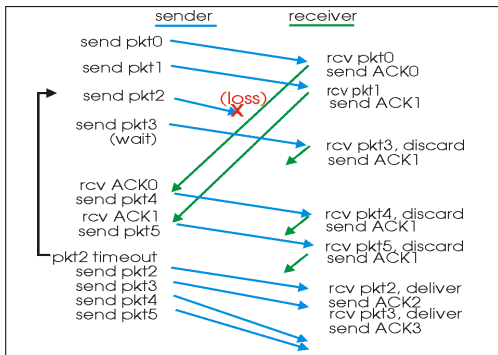
29

Example RTT Estimation



30

Still, Timeouts are Inefficient



31

Fast Retransmission

- When packet n is lost...
 - ... packets n+1, n+2, and so on may get through
- Exploit the ACKs of these packets
 - ACK says receiver is still awaiting nth packet
 - Duplicate ACKs suggest later packets arrived
 - Sender uses “duplicate ACKs” as a hint
- Fast retransmission
 - Retransmit after “triple duplicate ACK”

32

Effectiveness of Fast Retransmit

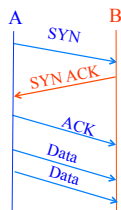
- When does Fast Retransmit work best?
 - High likelihood of many packets in flight
 - Long data transfers, large window size, ...
- Implications for Web traffic
 - Most Web transfers are short (e.g., 10 packets)
 - So, often there aren’t many packets in flight
 - Making fast retransmit is less likely to “kick in”
 - Forcing users to click “reload” more often... ☹️

33

Starting and Ending a Connection: TCP Handshakes

34

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) to the host B
 - Host B returns a SYN acknowledgment (**SYN ACK**)
 - Host A sends an **ACK** to acknowledge the SYN ACK

35

TCP Header

Flags: SYN
 FIN
 RST
 PSH
 URG
 ACK

Source port		Destination port	
Sequence number			
Acknowledgment			
HdrLen	0	Flags	Advertised window
Checksum		Urgent pointer	
Options (variable)			
Data			

36

Step 1: A's Initial SYN Packet

Flags: SYN
FIN
RST
PSH
URG
ACK

A's port		B's port	
A's Initial Sequence Number			
Acknowledgment			
20	0	Flags	Advertised window
Checksum		Urgent pointer	
Options (variable)			

A tells B it wants to open a connection...

37

Step 2: B's SYN-ACK Packet

Flags: SYN
FIN
RST
PSH
URG
ACK

B's port		A's port	
B's Initial Sequence Number			
A's ISN plus 1			
20	0	Flags	Advertised window
Checksum		Urgent pointer	
Options (variable)			

**B tells A it accepts, and is ready to hear the next byte...
... upon receiving this packet, A can start sending data**

38

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

Flags: SYN
FIN
RST
PSH
URG
ACK

A's port		B's port	
Sequence number			
B's ISN plus 1			
20	0	Flags	Advertised window
Checksum		Urgent pointer	
Options (variable)			

**A tells B it is okay to start sending
... upon receiving this packet, B can start sending data**

39

What if the SYN Packet Gets Lost?

- Suppose the SYN packet gets lost
 - Packet is lost inside the network, or
 - Server rejects the packet (e.g., listen queue is full)
- Eventually, no SYN-ACK arrives
 - Sender sets a timer and wait 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
 - Some TCPs use a default of 3 or 6 seconds

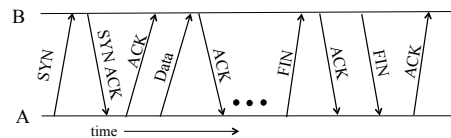
40

SYN Loss and Web Downloads

- User clicks on a hypertext link
 - Browser creates a socket and does a "connect"
 - The "connect" triggers the OS to transmit a SYN
- If the SYN is lost...
 - The 3-6 seconds of delay is very long
 - The impatient user may click "reload"
- User triggers an "abort" of the "connect"
 - Browser "connects" on a new socket
 - Essentially, forces a fast send of a new SYN!

41

Tearing Down the Connection



- Closing (each end of) the connection
 - Finish (FIN) to close and receive remaining bytes
 - And other host sends a FIN ACK to acknowledge
 - Reset (RST) to close and not receive remaining bytes

42

Sending/Receiving the FIN Packet

- **Sending a FIN: close()**
 - Process is done sending data via the socket
 - Process invokes “close()” to close the socket
 - Once TCP has sent all the outstanding bytes...
 - ... then TCP sends a FIN
- **Receiving a FIN: EOF**
 - Process is reading data from the socket
 - Eventually, the attempt to read returns an EOF

43

Conclusions

- **Transport protocols**
 - Multiplexing and demultiplexing
 - Checksum-based error detection
 - Sequence numbers
 - Retransmission
 - Window-based flow control
- **Precept on Friday**
 - Application-layer protocols: HTTP
 - HTTP proxy assignment

44