

CS162
Operating Systems and
Systems Programming
Lecture 15

Reliability, Transport Protocols

March 16, 2011
Ion Stoica
<http://inst.eecs.berkeley.edu/~cs162>

Goals for Today

- Finish e2e argument & fate sharing
- Transport: TCP/UDP
 - Reliability
 - Flow control

3/16

Ion Stoica CS162 ©UCB Spring 2011

Lec 15.2

Placing Network Functionality

- Hugely influential paper: “End-to-End Arguments in System Design” by Saltzer, Reed, and Clark ('84)
- “Sacred Text” of the Internet
 - Endless disputes about what it means
 - Everyone cites it as supporting their position

3/16

Ion Stoica CS162 ©UCB Spring 2011

Lec 15.3

Basic Observation

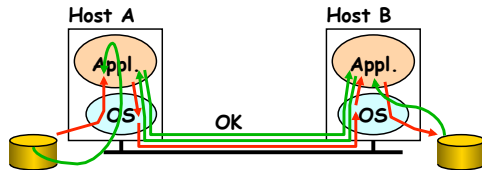
- Some types of network functionality can only be correctly implemented **end-to-end**
 - Reliability, security, etc
- Because of this, end hosts:
 - Can satisfy the requirement without network's help
 - Will/**must** do so, since can't **rely** on network's help
- Therefore **don't** go out of your way to implement them in the network

3/16

Ion Stoica CS162 ©UCB Spring 2011

Lec 15.4

Example: Reliable File Transfer



- Solution 1: make each step reliable, and then **concatenate** them
- Solution 2: end-to-end **check** and try again if necessary

3/16

Ion Stoica CS162 ©UCB Spring 2011

Lec 15.5

Discussion

- Solution 1 is **incomplete**
 - What happens if memory is corrupted?
 - Receiver has to do the check anyway!
- Solution 2 is **complete**
 - Full functionality can be entirely implemented at application layer with **no** need for reliability from lower layers
- *Is there any need to implement reliability at lower layers?*
 - Well, it could be **more efficient**

3/16

Ion Stoica CS162 ©UCB Spring 2011

Lec 15.6

Summary of End-to-End Principle

Implementing this functionality in the network:

- Doesn't reduce host implementation complexity
- Does increase network complexity
- Probably imposes delay and overhead on all applications, **even if they don't need functionality**
- However, implementing in network **can** enhance performance in some cases
 - E.g., very lossy link

3/16

Ion Stoica CS162 ©UCB Spring 2011

Lec 15.7

Conservative Interpretation of E2E

- Don't implement a function at the lower levels of the system unless it can be completely implemented at this level
- Unless you can relieve the burden from hosts, don't bother

3/16

Ion Stoica CS162 ©UCB Spring 2011

Lec 15.8

Moderate Interpretation

- Think twice before implementing functionality in the network
- If hosts can implement functionality correctly, implement it in a lower layer **only** as a performance enhancement
- But do so only if it **does not impose burden** on applications that do not require that functionality

3/16

Ion Stoica CS162 ©UCB Spring 2011

Lec 15.9

Related Notion of Fate-Sharing

- Idea: when storing **state** in a distributed system, keep it **co-located** with the entities that ultimately rely on the state
- Fate-sharing is a technique for dealing with **failure**
 - Only way that failure can cause loss of the critical state is if the entity that cares about it **also fails** ...
 - ... in which case **it doesn't matter**
- Often argues for keeping *network state* at end hosts rather than inside routers
 - In keeping with End-to-End principle
 - E.g., packet-switching rather than circuit-switching
 - E.g., NFS file handles, HTTP “cookies”

3/16

Ion Stoica CS162 ©UCB Spring 2011

Lec 15.10

Reliable Transfer

- Retransmit missing packets
 - Numbering of packets and ACKs
- Do this efficiently
 - Keep transmitting whenever possible
 - Detect missing ACKs and retransmit quickly
- Two schemes
 - Stop & Wait
 - Sliding Window (Go-back-n and Selective Repeat)

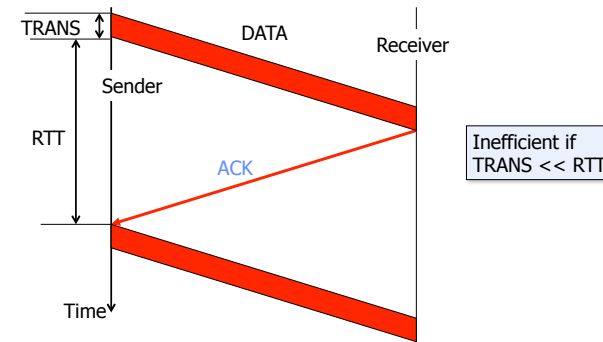
3/16

Ion Stoica CS162 ©UCB Spring 2011

Lec 15.11

Stop & Wait

- Send; wait for ack
- If timeout, retransmit; else repeat



3/16

Ion Stoica CS162 ©UCB Spring 2011

Lec 15.12

Sliding Window

- *window* = set of adjacent sequence numbers
- The size of the set is the *window size*
- Assume window size is n
- Let A be the last ack'd packet of sender without gap; then window of sender = $\{A+1, A+2, \dots, A+n\}$
- Sender can send packets in its window
- Let B be the last received packet without gap by receiver, then window of receiver = $\{B+1, \dots, B+n\}$
- Receiver can accept out of sequence, if in window

3/16

Ion Stoica CS162 ©UCB Spring 2011

Lec 15.13

Go-Back-n (GBN)

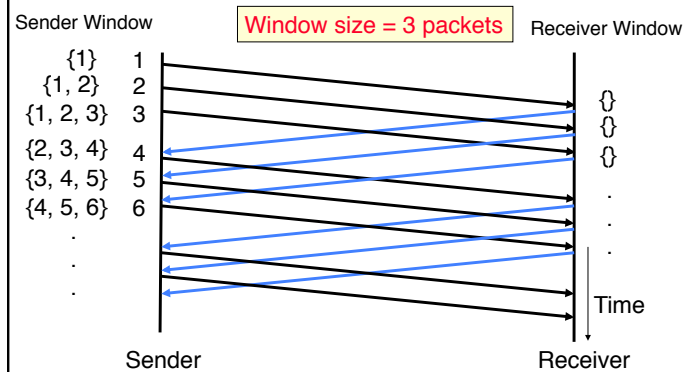
- Transmit up to n unacknowledged packets
- If timeout for $ACK(k)$, retransmit $k, k+1, \dots$

3/16

Ion Stoica CS162 ©UCB Spring 2011

Lec 15.14

GBN Example w/o Errors

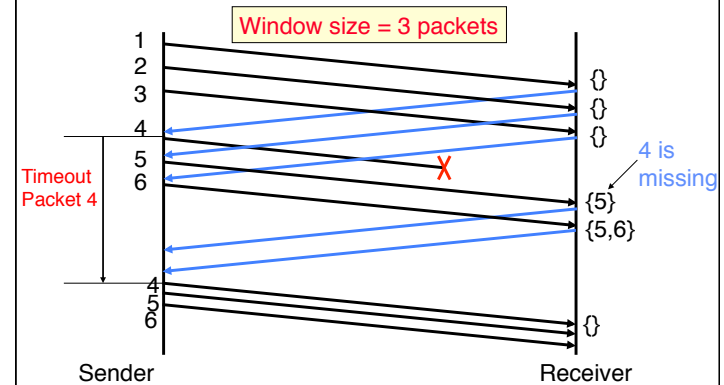


3/16

Ion Stoica CS162 ©UCB Spring 2011

Lec 15.15

GBN Example with Errors



3/16

Ion Stoica CS162 ©UCB Spring 2011

Lec 15.16

Selective Repeat (SR)

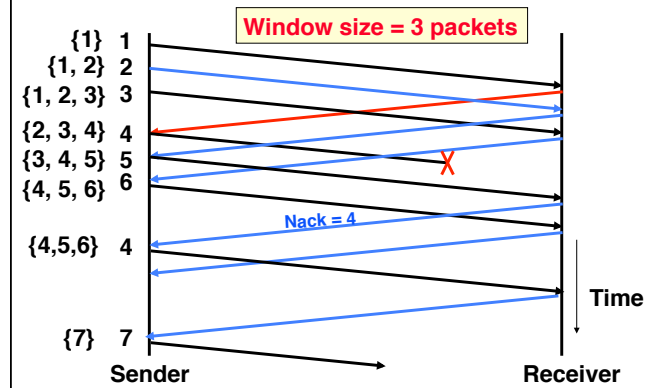
- Sender: transmit up to n unacknowledged packets; assume packet k is lost
- Receiver: indicate packet k is missing
- Sender: retransmit packet k

3/16

Ion Stoica CS162 ©UCB Spring 2011

Lec 15.17

SR Example with Errors



3/16

Ion Stoica CS162 ©UCB Spring 2011

Lec 15.18

Observations

- With sliding windows, it is possible to fully utilize a link, provided the window size 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

3/16

Ion Stoica CS162 ©UCB Spring 2011

Lec 15.19

Motivation for Transport Protocols

- IP provides a weak, but efficient service model (*best-effort*)
 - Packets can be delayed, dropped, reordered, duplicated
 - Packets have limited size (why?)
- IP packets are addressed to a host
 - How to decide which application gets which packets?
- How should hosts send packets into the network?
 - Too fast may overwhelm the network
 - Too slow is not efficient

3/16

Ion Stoica CS162 ©UCB Spring 2011

Lec 15.20

Transport Layer

- Provide a way to decide which packets go to which applications (*multiplexing/demultiplexing*)
- Can
 - Provide reliability, in order delivery, at most once delivery
 - Support messages of arbitrary length
 - Govern when hosts should send data → can implement congestion and flow control

3/16

Ion Stoica CS162 ©UCB Spring 2011

Lec 15.21

Congestion vs. Flow Control

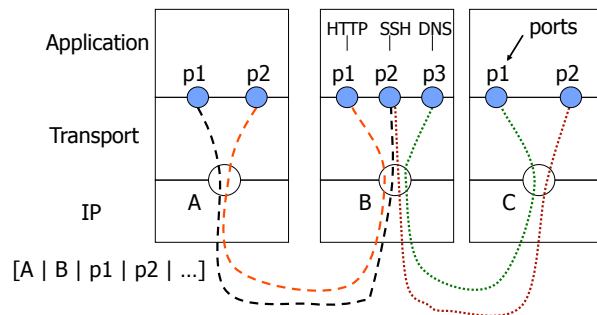
- **Flow Control** – avoid overflowing the receiver
- **Congestion Control** – avoid congesting the network
- What is network congestion?

3/16

Ion Stoica CS162 ©UCB Spring 2011

Lec 15.22

Transport Layer (cont'd)



UDP: Not reliable
TCP: Ordered, reliable, well-paced

3/16

Ion Stoica CS162 ©UCB Spring 2011

Lec 15.23

Ports

- Need to decide which application gets which packets
- Solution: map each socket to a *port*
- Client must know server's port
- Separate 16-bit port address space for UDP and TCP
 - (src_IP, src_port, dst_IP, dst_port) uniquely identifies TCP connection
- *Well known ports* (0-1023): everyone agrees which services run on these ports
 - e.g., ssh:22, http:80
 - On UNIX, must be root to gain access to these ports (why?)
- Ephemeral ports (most 1024-65535): given to clients
 - e.g. chat clients, p2p networks

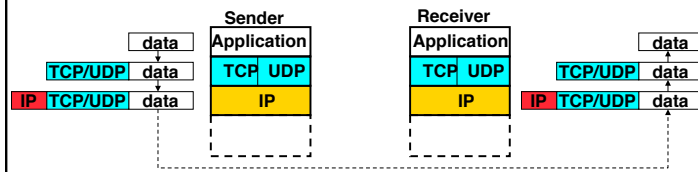
3/16

Ion Stoica CS162 ©UCB Spring 2011

Lec 15.24

Headers

- IP header → used for IP routing, fragmentation, error detection
- UDP header → used for multiplexing/demultiplexing, error detection
- TCP header → used for multiplexing/demultiplexing, flow and congestion control



3/16

Ion Stoica CS162 ©UCB Spring 2011

Lec 15.25

UDP: User (Unreliable) Data Protocol

- Minimalist transport protocol
- Same best-effort service model as IP
- Messages up to 64KB
- Provides multiplexing/demultiplexing to IP
- Does **not** provide flow and congestion control
- Application examples: video/audio streaming

3/16

Ion Stoica CS162 ©UCB Spring 2011

Lec 15.26

UDP Service & Header

- Service:
 - Send datagram from (IPa, Port1) to (IPb, Port2)
 - Service is unreliable, but error detection possible

• Header:

0	16	31
Source port	Destination port	
UDP length	UDP checksum	
Payload (variable)		

- UDP length is UDP packet length (including UDP header and payload, but not IP header)
- Optional UDP checksum is over UDP packet

3/16

Ion Stoica CS162 ©UCB Spring 2011

Lec 15.27

TCP: Transport Control Protocol

- Reliable, in-order, and at most once delivery
- Stream oriented: messages can be of arbitrary length
- Provides multiplexing/demultiplexing to IP
- Provides congestion control and avoidance
- Application examples: file transfer, chat

3/16

Ion Stoica CS162 ©UCB Spring 2011

Lec 15.28

TCP Service

- 1) Open connection: 3-way handshaking
- 2) Reliable byte stream transfer from (IPa, TCP_Port1) to (IPb, TCP_Port2)
 - Indication if connection fails: Reset
- 3) Close (tear-down) connection

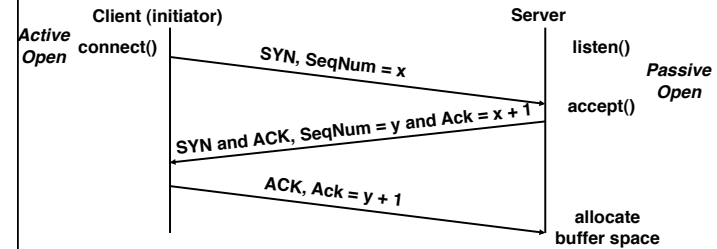
3/16

Ion Stoica CS162 ©UCB Spring 2011

Lec 15.29

Open Connection: 3-Way Handshaking

- Goal: agree on a set of parameters: the start sequence number for each side
 - Starting sequence numbers are random



3/16

Ion Stoica CS162 ©UCB Spring 2011

Lec 15.30

3-Way Handshaking (cont'd)

- Three-way handshake adds 1 RTT delay
- Why?
 - Congestion control: SYN (40 byte) acts as cheap probe
 - Protects against delayed packets from other connection (would confuse receiver)

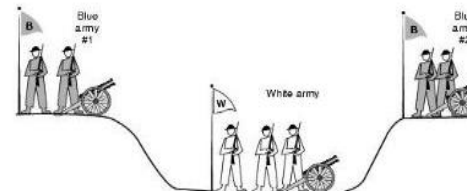
3/16

Ion Stoica CS162 ©UCB Spring 2011

Lec 15.31

Close Connection (Two-Army Problem)

- Goal: both sides agree to close the connection
- Two-army problem:
 - “Two blue armies need to simultaneously attack the white army to win; otherwise they will be defeated. The blue army can communicate only across the area controlled by the white army which can intercept the messengers.”



- What is the solution?

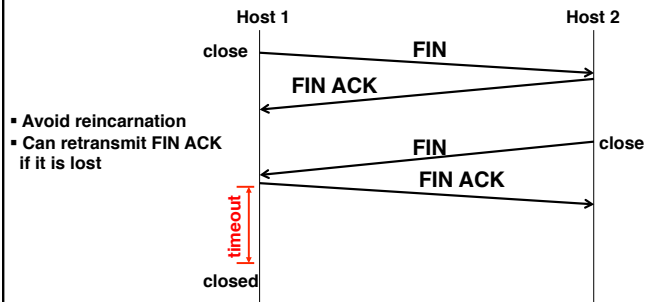
3/16

Ion Stoica CS162 ©UCB Spring 2011

Lec 15.32

Close Connection

- 4-way tear down connection



3/16

Ion Stoica CS162 ©UCB Spring 2011

Lec 15.33

- Avoid reincarnation
- Can retransmit FIN ACK if it is lost

TCP Flow Control

- Each **byte** has a sequence number
- Initial sequence numbers negotiated via SYN/SYN-ACK packets
- ACK contains the sequence number of the **next** byte expected by the receiver

3/16

Ion Stoica CS162 ©UCB Spring 2011

Lec 15.34

TCP Flow Control

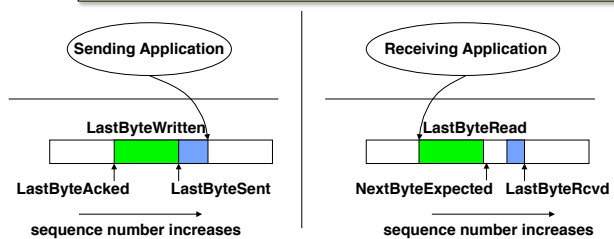
- Receiver window (MaxRcvBuf – maximum buffer size at receiver)

$$\text{AdvertisedWindow} = \text{MaxRcvBuffer} - (\text{LastByteRcvd} - \text{LastByteRead})$$

- Sender window (MaxSendBuf – maximum buffer size at sender)

$$\text{SenderWindow} = \text{AdvertisedWindow} - (\text{LastByteSent} - \text{LastByteAcked})$$

$$\text{MaxSendBuffer} \geq \text{LastByteWritten} - \text{LastByteAcked}$$



3/16

Ion Stoica CS162 ©UCB Spring 2011

Lec 15.35

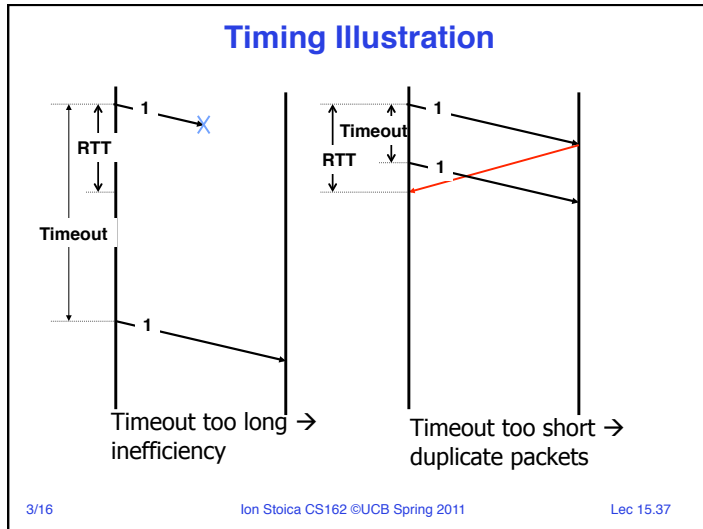
Retransmission Timeout

- If haven't received ack by timeout, retransmit packet after last acked packet
- How to set timeout?

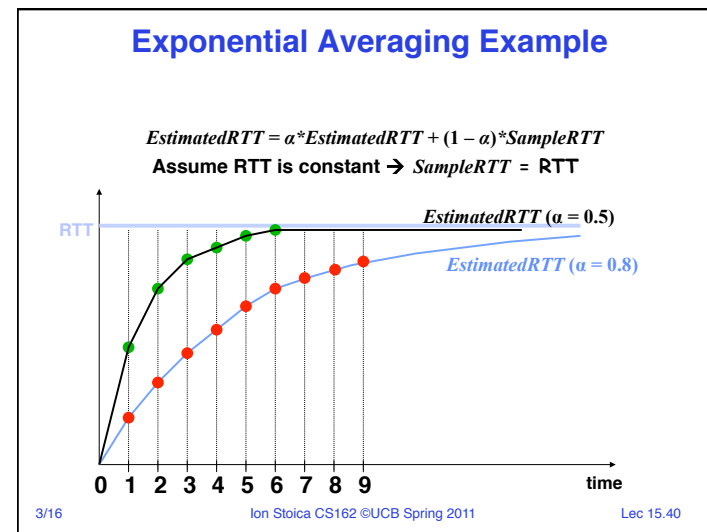
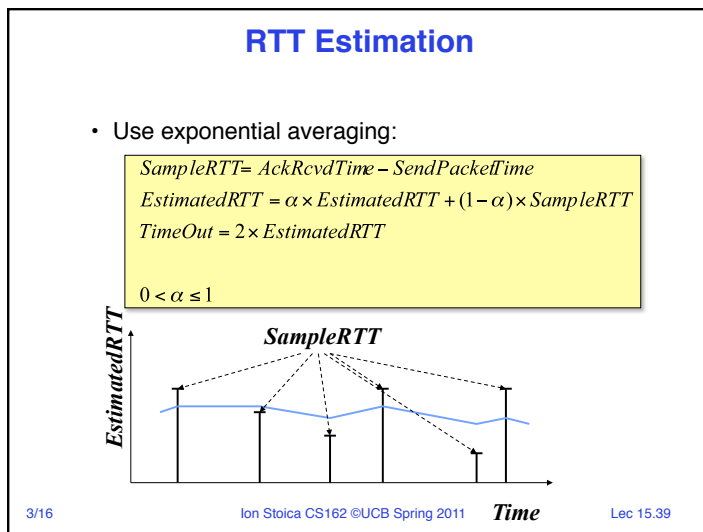
3/16

Ion Stoica CS162 ©UCB Spring 2011

Lec 15.36

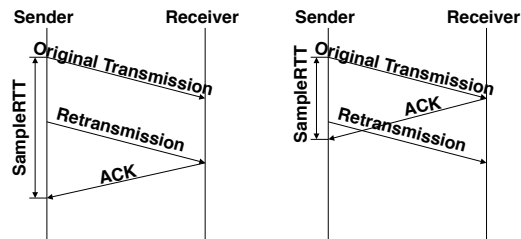


- ### Retransmission Timeout (cont'd)
- If haven't received ack by timeout, retransmit packet after last acked packet
 - How to set timeout?
 - Too long: connection has low throughput
 - Too short: retransmit packet that was just delayed
 - » Packet was probably delayed because of congestion
 - » Sending another packet too soon just makes congestion worse
 - Solution: make timeout proportional to RTT
- 3/16 Ion Stoica CS162 ©UCB Spring 2011 Lec 15.38



Problem

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



3/16

Ion Stoica CS162 ©UCB Spring 2011

41

Lec 15.41

Karn/Partridge Algorithm

- Measure *SampleRTT* only for original transmissions
- Exponential backoff → for each retransmission, double *EstimatedRTT*

3/16

Ion Stoica CS162 ©UCB Spring 2011

Lec 15.42

Jacobson/Karels Algorithm

- Problem: exponential average is not enough
 - One solution: use standard deviation (requires expensive square root computation)
 - Use mean deviation instead

$$\begin{aligned}
 \text{Difference} &= \text{SampleRTT} - \text{EstimatedRTT} \\
 \text{EstimatedRTT} &= \text{EstimatedRTT} + \delta \times \text{Difference} \\
 \text{Deviation} &= \text{Deviation} + \delta \times (|\text{Difference}| - \text{Deviation}) \\
 \text{TimeOut} &= \mu \times \text{EstimatedRTT} + \phi \times \text{Deviation} \\
 0 &< \delta \leq 1 \\
 \mu &= 1 \\
 \phi &= 4
 \end{aligned}$$

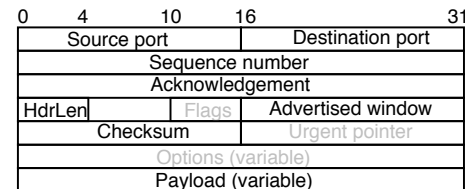
3/16

Ion Stoica CS162 ©UCB Spring 2011

Lec 15.43

TCP Header

- Sequence number, acknowledgement, and advertised window – used by sliding-window based flow control
- HdrLen: TCP header length in 4-byte words
- Checksum: checksum of TCP header + payload



3/16

Ion Stoica CS162 ©UCB Spring 2011

Lec 15.44

Summary

- Reliable transmission
 - S&W not efficient → Go-Back-n
 - What to ACK? (cumulative, ...)
- UDP: Multiplex, detect errors
- TCP: Reliable Byte Stream
 - 3-way handshaking
 - Flow control
 - Timer Value: based on measured RTT