

Lecture 6 - The Transport Layer

Computer Communication Networks
CS35201 - (001/ 002)
Fall 2022

Kent State University



Department of Computer Science
Betis Baheri
bbaheri@kent.edu

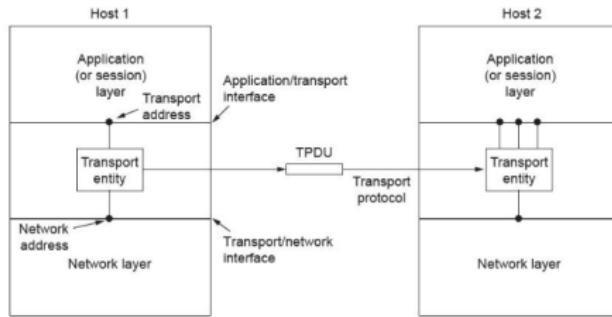
October 31, 2022

Outline I

- 1 Part I (Transport Layer)
- 2 Part II (UDP vs TCP)
- 3 Part III (TCP Flow and Congestion Control)
- 4 Part IV (TCP Congestion Control Approaches)
- 5 Part V (Packet Based Congestion Control)
- 6 References

Transport Layer I

- 1 Provides logical communication between application processes
 - ▶ Not hosts
- 2 Provides reliable connection-oriented services ⇒ TCP
- 3 Provides unreliable connection-less-oriented services ⇒ UDP
- 4 Provides Multiplexing/Demultiplexing of application services
 - ▶ Based on port# and Internet Protocol (IP) addresses
- 5 Provides parameters for specifying quality of services
- 6 Runs on end systems, but relies on network services ↓



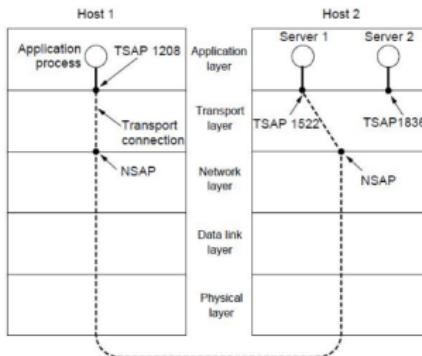
Elements of Transport Protocols

■ Functions

- 1 Addressing
- 2 Connection establishment
- 3 Connection release
- 4 Error control and flow control
- 5 Multiplexing
- 6 Crash recovery

■ Transport service access point

- ▶ Each layer has its own way of dealing with addresses
- ▶ In Internet, a transport service access point is an IP address with a port number



Transport Protocols

■ Point-to-Point communication ⇒ between processes

1 Connection-oriented ⇒ e.g. Transport Control Protocol (TCP)

- Reliable data transfer
- In-order delivery
- Flow Control (FC)
- Congestion Control

2 Connectionless-oriented ⇒ e.g. User Data Protocol (UDP)

- Unreliable
- Unordered
- No handshaking ⇒ independent segment handling

■ Application Needs

- ▶ Real-time
- ▶ Bandwidth guarantees
- ▶ Reliable multicast

Transport Architectural Requirements

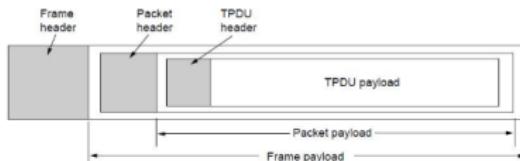
- Procedures to exchange data between devices
 - ▶ Can be complex
- High degree of cooperation between end systems

Example 6.1 (File Transfer)

- Transport needs:
 - 1 A data communication path ⇒ points of transfer
 - 2 An activation mechanism to activate the path
 - 3 A mechanism to determine the status of the destination
 - 4 A mechanism for the destination file management to store the file
 - 5 A format conversion mechanism
 - 6 A few other details such Congestion Control (CC) and FC

Transport Layer Operations

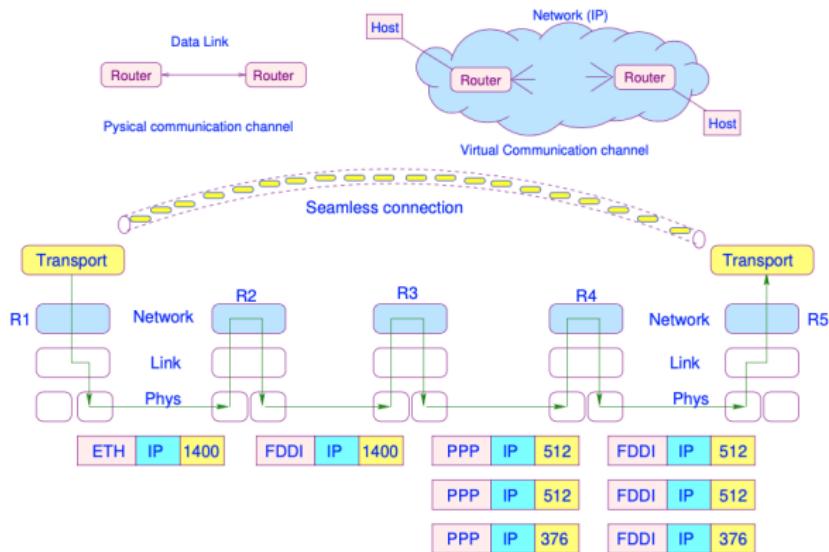
- Messages are encapsulated as Transport Protocol Data Units (TPDUs) to the network layer



- ▶ The network layer is in the hands of carriers
- ▶ Transport layer is in the hand of host (computer)
 - Has no say on what the carrier actually offers
- A real hard part is establishing and releasing connections
 - 1 Symmetric: one side sends a disconnect request, and waits for the other to acknowledge that the connection is closed
 - 2 Asymmetric: one side just closes the connection, no Ack
 - Simple, but unacceptable ⇒ data may be lostHow?
- Transport protocols strongly resemble those in the data link layer
 - ▶ Error Control (EC)
 - ▶ FC/CC
 - ▶ Reliability ⇒ sequencing
- Messages are encapsulated as TPDUs to the network layer

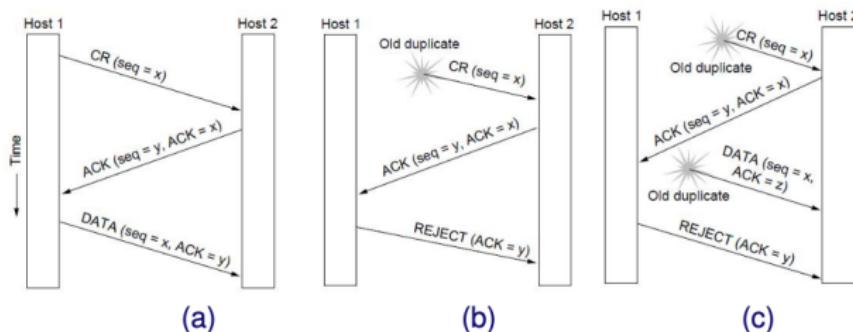
Major differences with L2

- Explicit addressing
- Establishing, maintaining, and releasing connections
 - ▶ TCP, not UDP
- Many connections (applications) require different solutions
 - ▶ Rates, losses, etc.



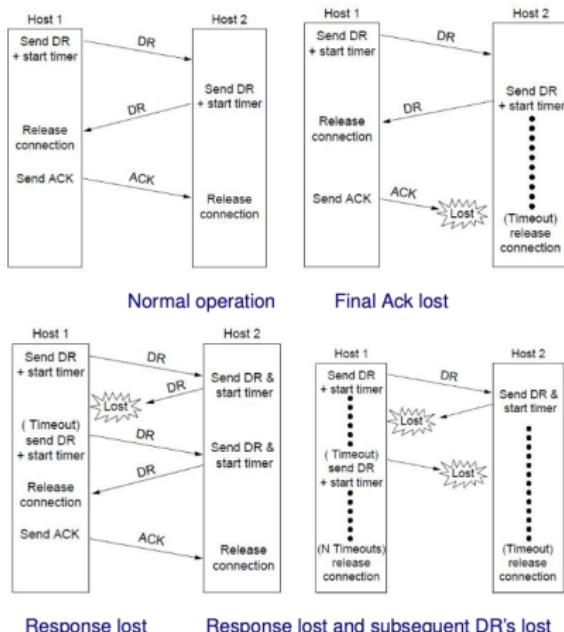
Connection Establishment (3-Way Handshaking)

- 1 Host 1 chooses a sequence number, x , and sends a CONNECTION REQUEST TPDU containing it to host 2.
- 2 Host 2 replies with an ACK TPDU acknowledging x and announcing its own initial sequence number, y .
- 3 Finally, host 1 acknowledges host 2's choice of an initial sequence number in the first data TPDU that it sends



- (a) Normal Operation
 (b) Old duplicate CONNECTION REQUEST appearing out of nowhere
 (c) Duplicate CONNECTION REQUEST and duplicate ACK

Connection Release



Common Transport Protocols

- 1 UDP provides just integrity and demux
- 2 TCP adds:
 - ▶ Connection-oriented
 - ▶ Reliable and ordered delivery
 - ▶ end-to-end
 - ▶ Byte-stream
 - ▶ Full duplex
 - ▶ FCCC

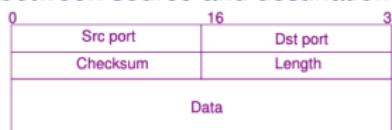
Outline I

- 1 Part I (Transport Layer)
- 2 Part II (UDP vs TCP)
- 3 Part III (TCP Flow and Congestion Control)
- 4 Part IV (TCP Congestion Control Approaches)
- 5 Part V (Packet Based Congestion Control)
- 6 References

User Datagram Protocol (UDP)

1 Bare bone Internet transport protocol \Rightarrow IP + Port

- ▶ Endpoints are identified by ports
- ▶ Best effort service between source and destination



2 Connection-less

- ▶ No connection setup/tear down between Src/Dst

see [/etc/services](#) on Unix

- ▶ No handshaking
- ▶ Checksum is optional

- UDP segments can be lost \Rightarrow Unreliable \Rightarrow loss
- UDP can be delivered out of order

- ▶ A segment is a sequence of 16-bit integer
 - Checksum= addition (1's complement) of 16-bit integers

Why?

- ▶ Each UDP segment handled independently

3 UDP is trade-off between reliability and speed

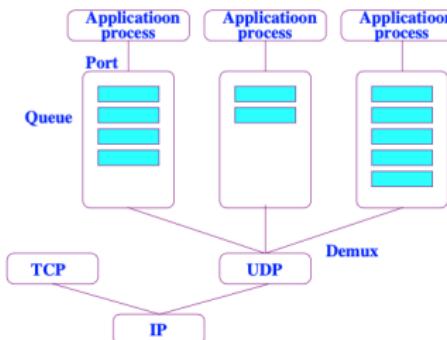
- ▶ Not necessarily a good idea
- ▶ Small segment header

Why?

4 State-less

- ▶ No state information \Rightarrow less overhead \Rightarrow faster than TCP
- ▶ No congestion or flow control
- ▶ No flow control, no Ack, no error recovery

Simple UDP Multiplexing/Demultiplexing



- UDP is often used for streaming multimedia applications
 - ▶ Loss tolerant
 - ▶ Rate sensitive
- Other UDP usage
 - ▶ Domain Name System (DNS)
 - ▶ Simple Network Management Protocol (SNMP)
- Reliability is obtained by application layer
 - ▶ Application specification error recovery

Why?

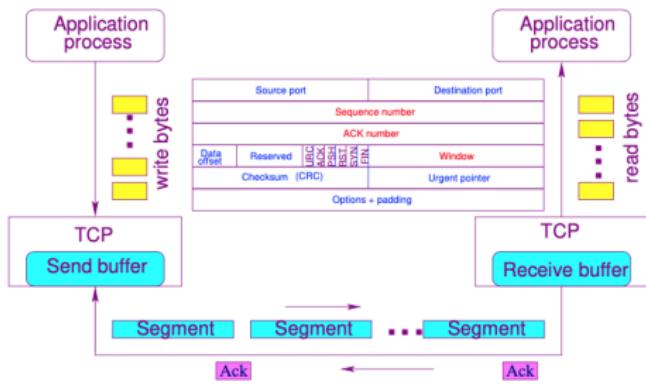
UDP Applications

Implications:

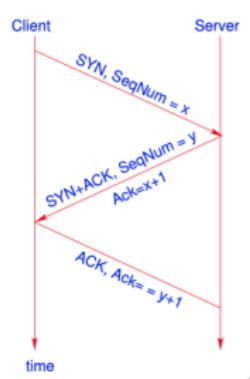
- Applications
 - ▶ Inward data collection
 - ▶ Outward data dissemination
 - ▶ Request/Reply
 - ▶ Real-time data
 - ▶ Used by SNMP, DNS, Trivial File Transfer Protocol (TFTP) etc.
- Application message encapsulated and sent to IP
 - ▶ Can result in fragmentation How?
- Implications
 - ▶ Responsiveness
 - TCP responds to congestion, UDP does not
 - ▶ Source Quench
 - UDP sources ignore source quench due to congestion
 - ▶ ARP Flooding (Poisoning)
 - When UDP datagram is fragmented, each fragment may generate an Address Resolution Protocol (ARP) request How?
 - Security concern
 - ▶ Datagram Truncation
 - A destination may truncate a large UDP datagram
 - TCP does not have this problem

Transport Control Protocol (TCP)

Operation



3-Way Handshake



- A reliable transport layer protocol that preserves
 - ▶ Zero loss of bytes
 - Link Automatic Repeat Request (ARQ)
 - ▶ Zero duplication of bytes
 - Acknowledgment (ACK) (**sliding window**)
- Octets of the data stream are numbered sequentially

TCP Header

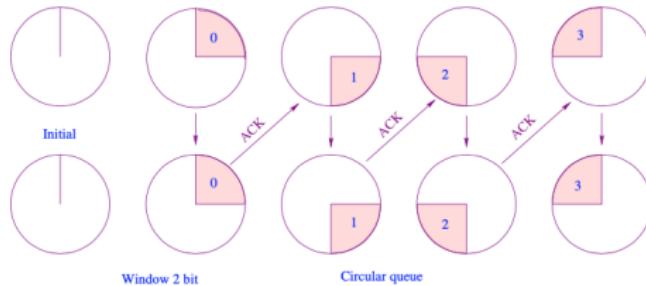
- Windows: \Rightarrow Operate at the octet level
 - ▶ Flow Window (W_F) \Rightarrow controlled by the rate of Acknowledgments
 - Flow control: keep sender from overrunning the receiver *How?*
 - ★ Receiver controls \Rightarrow Hold ACKs
 - ▶ Congestion Window (W_C) \Rightarrow Sliding Window
 - The number of bytes that the receiver is prepared to accept *How?*
 - ★ Advertised window (W_A) \Rightarrow $\min\{W_F, W_C\}$ $W_A \leq W_C \leq W_F = W$
 - Congestion control: keep sender from overrunning network
- Connection-oriented \Rightarrow Full duplex
 - ▶ Sending process writes some number of bytes
 - ▶ TCP breaks into segments and sends via IP
 - ★ Further fragmentation possible by Local Area Networks (LANs)
- Checksum
 - ▶ Pseudo header + TCP header + Data
- Each connection identified with 4-tuple:

$$<SrcPort, SrcIPAddr, DstPort, DstIPAddr>$$
- Sliding window + flow control:

$$Ack, SeqNum, AdverWindow$$

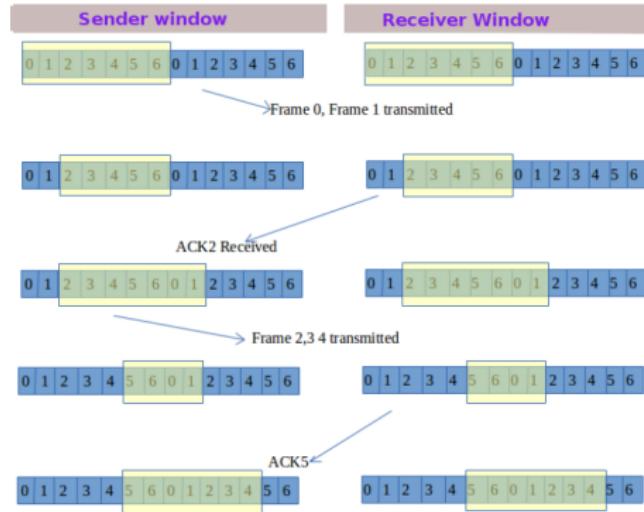
TCP Window I

- The sender/receiver each has two windows:
 - ▶ Sending window
 - ▶ Receiving window
- It is a duplex mechanism
 - ▶ One can send/receive at the same time
- Windows don't have to be the same size
 - ▶ Window size = 4, but uses Stop-and-Wait protocol



TCP Window II

Recall Window from Data Link Layer



- Sending ACK2 – the receiver is expecting segment 2

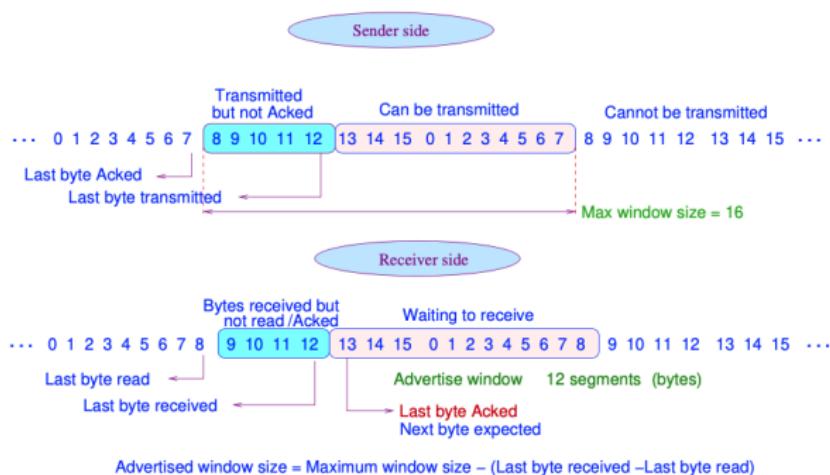
TCP Window III

■ Sender keeps three pointers

- ▶ Last byte Acked
- ▶ Last byte Sent(transmitted)
- ▶ Max Window

■ Receiver keeps three pointers

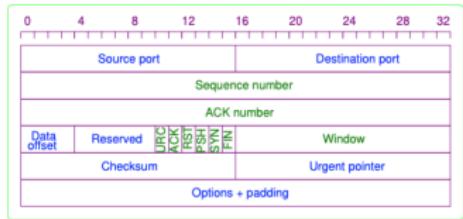
- ▶ Last byte Read
- ▶ Last byte Received
- ▶ Next Byte Expected



TCP Segment Format

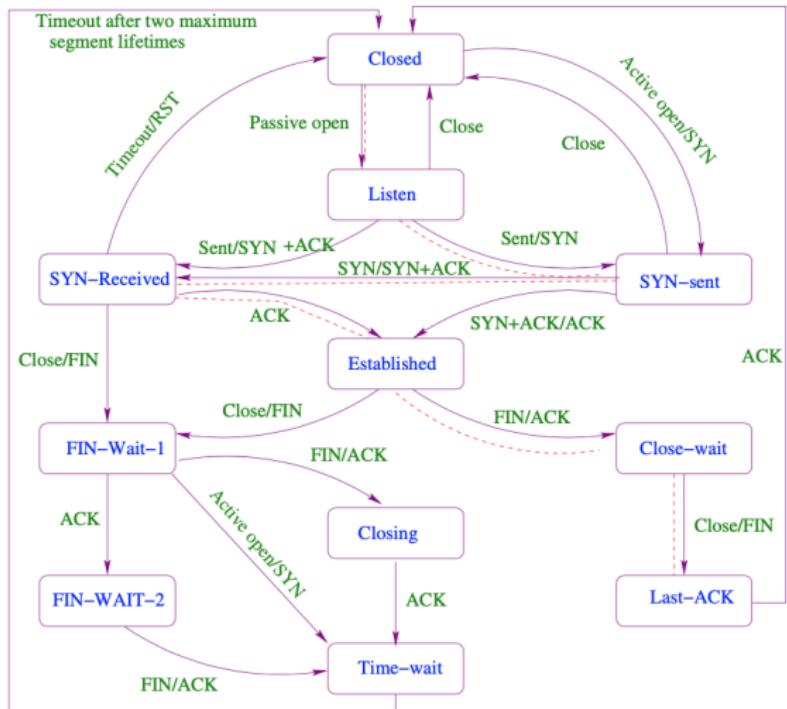
- The sender/receiver each has two windows:
- The windows move to the right not necessarily by 1 Why?
- Cumulative Acknowledgments \Rightarrow source of bursty
 - ▶ Can work with Go-back-N retransmission protocol
 - If a byte is missing, the receiver can ACK it to be re-transmitted
 - ▶ Can not identify which block of data is lost
 - Duplicate ACK can identify the lost block
- Initially
 - ▶ Sender buffer size: MaxSendBuffer
 - ▶ Receive buffer size: MaxRcvBuffer

Segment Format



- SYN: Connection establishing
- FIN: Connection terminating
- RESET: Receiver is confused
- PUSH: Sender invokes push
- URG: segment contains urgent data
- ACK: Acknowledgment

TCP Finite State Diagram



TCP Significance

- 1 Many applications run on TCP/IP suite
- 2 TCP/IP protocol engine becomes critical
 - Many believe boosting hardware ⇒ better network performance
 - TCP has total control of applications
 - TCP is a complex protocol
 - Depends on many network elements
 - Unless TCP is optimized, hardware alone cannot do much
- 3 Emergence of new technologies and environment
 - Proliferation of different networking technologies
 - Wireless, satellite, optical, etc.
 - TCP algorithms for one environment may not work best with another
 - A huge body of literature
- 4 IP convergence: Many non TCP/IP industries converging to TCP/IP
 - Cellular
- 5 Understanding TCP/IP fundamental performance is essential

End-to-End Issues

- Based on the sliding window (Data Link Control (DLC)), but very different
- Potentially connects many different hosts
 - ▶ Explicit connection establishment/termination
- Potentially different Round Trip Time (RTT)
 - ▶ Need adaptive timeout mechanism
- Potentially long delay in network
 - ▶ Need to be prepared for arrival of very old packets
- Potentially different capacity at destination
 - ▶ Buffering accommodation at the receiver
- Potentially different network capacity ⇒ congestion
 - ▶ Need to be prepared for network congestion

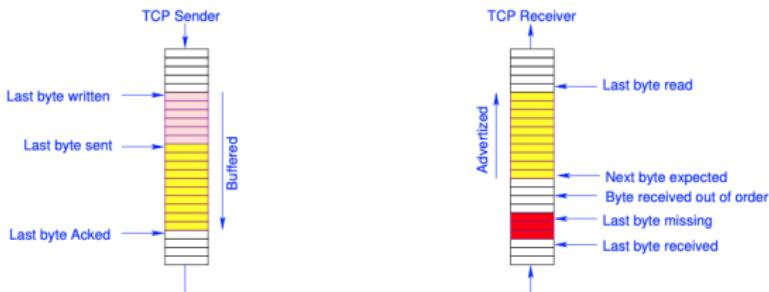
Outline I

- 1 Part I (Transport Layer)
- 2 Part II (UDP vs TCP)
- 3 Part III (TCP Flow and Congestion Control)
- 4 Part IV (TCP Congestion Control Approaches)
- 5 Part V (Packet Based Congestion Control)
- 6 References

TCP Flow and Congestion Control

- Receiver tells Sender how much it can receive

- ▶ 16 bit *Advertise Window*



How?

- Receiving side: \Rightarrow in control of flow by:

- 1 Holding off Ack \Rightarrow not for too long

- Could cause retransmission by sender if held too long
- Cumulative Acknowledgments cause traffic burstiness
- ★ Adjusting Cumulative Acknowledgments to reduce burstiness may cause buffer overflow

Why?

- 2 Slow the sender by persisting *Advertise Window*=0

- Blocked sender if $(LastByteReceived - LastByteRead) + y > MaxSendBufferSize$

How?

- Sending side

- 1 Adjusts its Time Out (TO) by the rate of Acknowledgments arriving

- It reduces unnecessary re-transmissions

RTT and Sequence Number Wrap-around

- Sequence Number Wrap Around: 32-bit *SequenceNum*

- Time to Wrap-around

- ▶ 1 Mbps (2^{20} bps) $\Rightarrow \frac{2^{32}}{2^{20}/2^3 \text{ B/s}} = 2^{15} \text{ s} \approx 9.1 \text{ hours}$
- ▶ 10 Mbps $\Rightarrow \frac{2^{32}}{10 \times 2^{20}/2^3 \text{ B/s}} = 2^{15} \text{ s} \approx 0.91 \text{ hour}$
- ▶ 1 Gbps $\Rightarrow \frac{2^{32}}{2^{30}/2^3 \text{ B/s}} = 2^5 \text{ s} = 32 \text{ sec.}$

- Wrap-around puts a limit how much we can fill the pipe

- ▶ Each byte has a sequence number
- ▶ Cumulative \Rightarrow send several bytes \Rightarrow how many?

- Bytes in Transit: 16-bit *AdvertisedWindow*

- ▶ You can't ask for more 2^{16} bytes \Rightarrow how many?
- ▶ This will put a limit on how much you can fill the pipe
- ▶ keeping the pipe full, assume RTT= 100 ms cross-country delay in US
 - You can't ask for more 2^{16} bytes \Rightarrow then how many?
 - For 1 Mbps $\text{Delay} \times \text{Bandwidth} = 2^{20}/8 \text{ B/s} \times 0.1 \text{ s} = 13,107 \text{ B}$
 - For 1 Gbps $\text{Delay} \times \text{Bandwidth} = 2^{30}/8 \text{ B/s} \times 0.1 \text{ s} = 13,421,772 \text{ B}$
- ▶ In *AdvertisedWindow* receiver can ask for maximum 2^{16} segments
 \Rightarrow 65,536 bytes
- ▶ How do we solve this problem (not being able to fill the pipe)
 - Send larger segments (more than one byte)
 - 2^{16} segments, each 8 bytes $\Rightarrow 8 \times 2^{16}$ segments = 6,524,288 bytes

TCP Congestion Control

- Congestion: when traffic load > network capacity
 - ▶ Traffic aggregate load exceeds the capacity
- Effects
 - ▶ Packet loss (or multiple packet losses)
 - Retransmission
 - Wasted resources due to packet loss
 - ▶ Reduced throughput \Rightarrow low throughput \Rightarrow low link utilization
 - ▶ Long delays \Rightarrow high queuing delay
 - ▶ Congestion Collapse

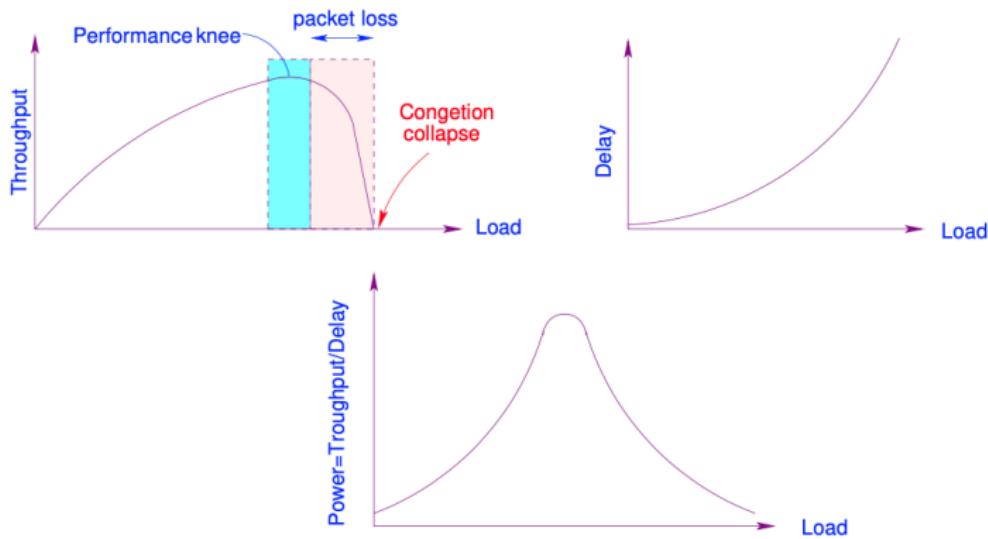
Definition 6.4.3.1 (Congestion collapse)

Congestion Collapse occurs when the network gets to a point in which it only forwards retransmissions while most of the retransmitted packets are delivered only half-way before being discarded again due to congestion

- Congestion collapse due to
 - ▶ Unnecessarily retransmitted packets
 - ▶ Un-delivered or unusable packets

Congestion Collapse

■ Saturation point



What to Do with Congestion

- Congestion Control remains a critical issue given
 - ▶ Network growing size
 - ▶ Demand bandwidth and speed
 - ▶ Increasingly integrated services
 - ▶ Aggregate traffic
- Despite all research efforts, there is no perfect solution
- Current solutions
 - ▶ Mainly intuition driven
 - ▶ Non-linear simple schemes
 - ▶ Exhibits poor performance ⇒ not strongly responsive to
 - Oscillation, chaotic behavior, fairness
 - ▶ Ineffective, particularly for non-elastic (delay/loss sensitive) traffic
 - Multimedia
 - ▶ Not scalable

Two Strategies at Saturation Point

- 1 Discard any incoming packet if no buffer available
- 2 Saturated routers exercises flow control over neighbors
 - ▶ Called back-pressure
 - ▶ May cause congestion to propagate throughout network
 - It passes congestion to other nodes
 - ▶ Oscillation

What Makes CC a Difficult Problem?

- 1 A very complex system
- 2 Deals with a very large distributed system
- 3 Diverse in term of
 - ▶ Traffic carried
 - ▶ Application requirements
 - ▶ QoS requirements
- 4 Large feed-back delays \Rightarrow large bandwidth \times delay
- 5 Unpredictable user behavior in terms of
 - ▶ Timing
 - ▶ Space \Rightarrow locality
 - ▶ Size \Rightarrow bursts
- 6 Lack of appropriate dynamic models
 - ▶ Mathematically a difficult problem
- 7 Fairness provisioning without computational complexity

Solutions

- CC \Rightarrow congestion recovery
 - ▶ Reactive approach after congestion is detected
- Congestion Avoidance (CA)
 - ▶ Proactive approach before congestion occurs
- CC involves statically or dynamically limiting demand-capacity mismatch
- Static solutions such as additional buffer, links, and processors are not effective

Why?

War Between Mice and Elephants

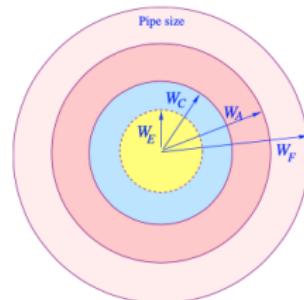
- TCP traffic is still dominant ⇒ 90%
 - ▶ But UDP is growing in a faster pace
- 50-70% of TCP are short-lived connections
 - ▶ A large number connections are short in terms of traffic they carry
 - ▶ A small fraction of connections carry a large portion of traffic
 - ⇒ elephants
 - ▶ A large fraction of connections carry a small portion of traffic ⇒ mice

Outline I

- 1 Part I (Transport Layer)
- 2 Part II (UDP vs TCP)
- 3 Part III (TCP Flow and Congestion Control)
- 4 Part IV (TCP Congestion Control Approaches)
- 5 Part V (Packet Based Congestion Control)
- 6 References

TCP Congestion Approaches

- End-to-End (E2E) congestion control
- Implicit feed-back \Rightarrow time-out, duplicate ACKs, etc.
- Uses window based flow control



- 1 Advertise Windo Window(W_A)
 - ▶ Advertized by the reciever in the header of TCP segment
 - 2 Flow Window(W_F) $\Rightarrow \text{Min } \{ \text{Pipe}_{size}, W_A \}$
 - ▶ Determined by the source
 - 3 Congestion Window ($W_C \leq W_F$)
 - ▶ Number of segments which won't congestion
 - 4 Effective Window ($W_E \leq W_C$)
 - ▶ Number of segments actually sent beased on all information
- Many variations of TCP
 - ▶ TCP Tahoe, TCP Reno, TCP Vegas, etc.
 - ▶ TCP uses Slow start + CA

How?

How TCP Handles Congestion?

- A TCP host inject a segment(s) into the networks without reservation
- It reacts (adjust rate) by observing events
 - ▶ Time-out
 - ▶ Rate of ACKs arrivals
 - ▶ Duplicate ACKs
 - ▶ Triple duplicate ACKs \Rightarrow 4 ACKs
- It determines how much capacity is available
- Rate of ACKs arrival determine how much a source can transmit
 - ▶ TCP is self-clocking
 - Sender is clocked (controlled) by the receiver through ACK rate

How?

Additive Increase/Multiplicative Decrease (AIMD) I

- Additive Increase, Multiplicative Decrease (AIMD) is a feedback mechanism used in TCP
- AIMD combines linear growth of the congestion window with an exponential reduction when a congestion takes place
- Multiple flows using AIMD congestion control will eventually converge to use equal amounts of a contended link
- Let $W_c(t)$ be the sending rate during time slot t , then

$$W_c(t+1) = \begin{cases} W_c(t) + a & \text{If congestion is not detected} \\ W_c(t) \times d & \text{If congestion is detected} \end{cases}$$

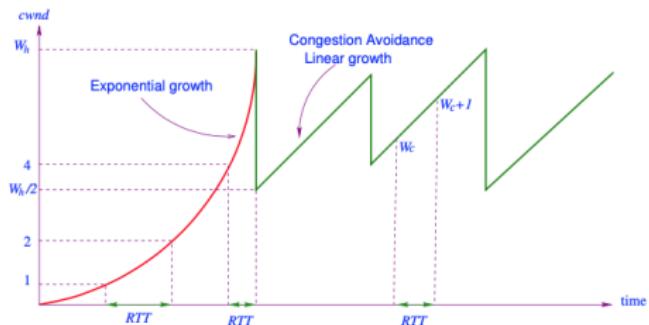
$a > 0, \quad 0 < d < 1 \quad (1)$

- It has been shown that AIMD is necessary for stability of TCP CC
- Accurate time-out is important to
 - ▶ Prevent unnecessary retransmission
 - ▶ Prevent under-utilization
- Time-out is a function of average and standard deviation of RTT

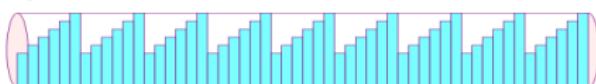
Slow Start I

■ Two phases

- 1 Congestion detection \Rightarrow exponential growth of W_C
 - W_C grows exponentially \Rightarrow effectively $W_C = 2 \times W_C$ every RTT
- 2 Congestion Avoidance \Rightarrow linear growth of W_C
 - W_C grows linearly $W_C = W_C + 1$

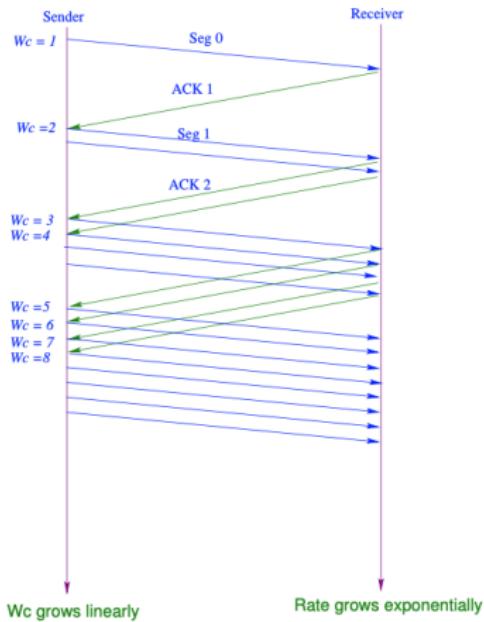


■ Does it do the job?



- Cold Start (beginning of a connection), $W_C = 1 \downarrow$
- Dead/broken connection results in advertised window goes to zero \downarrow

Slow Start II



Slow Start III

Slow Start Algorithm

- Slow Start has two phases

- 1 Exponential growth phase

- Start with $W_C = 1$
- For each ACK, $W_C := W_C + 1$ segment
 - ★ It increases W_C by the number of segments acknowledged
 - ★ W_C is effectively doubled per RTT
 - ★ Note on cumulative ACKs
- Until it reaches the threshold W_h , i.e., $W_C \geq W_h$ or an ACK is lost after a time out

- ★ Then $W_h = W_C$; $W_C = W_C/2$,
- ★ Enters Congestion Avoidance phase \Rightarrow linear growth \Rightarrow step 2
- ★ Problem: When a segment is lost TCP assumes that it is due to congestion
→ Not always true Why?

- 2 Congestion avoidance phase (linear growth)

- For each successful ACK $W_C := W_C + 1/W_C$ How?
- For each RTT $W_C := W_C + 1$ Until a segment is lost

- Slow Start is quite aggressive than congestion avoidance phase

- ▶ Could drive the network into saturation How?

Slow Start IV

- ▶ However, slow start is slower than sending the full advertised window of segments
- Current W_c can be used as **congestion threshold W_h (ssthresh)**
⇒ implicit information

Algorithm 1: Slow Start

```

input :  $W_c = 1, W_h = \infty$ 
1 for every ACK do
2   | if  $W_c < W_h$  then
3   |   |  $W_c \leftarrow W_c + 1;$  // 1 max segment size (MSS) per ACK
4   |   | // Additive increase
5   |   | else
6   |   |   |  $W_c \leftarrow W_c + 1/W_c;$  // Congestion avoidance
7   |   | end
8   | end
9 end
10 for every time-out do
11   |  $W_h \leftarrow W_c/2;$  // Multiplication decrease
12   |  $W_c \leftarrow 1;$  // Slow Start
13 end

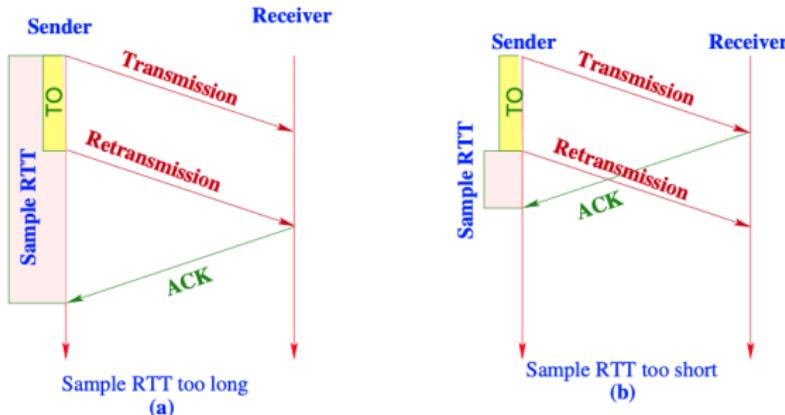
```

Adaptive Retransmission

- Measure SampleRTT for each segment/ACK pair
 - ▶ Needed to set the TO
- RTT is never accurate due to dynamic load
 - ▶ We estimate its average, but not arithmetic average \Rightarrow aged samples
- Compute weighted average of RTT
 - ▶ $EstimatedRTT = \alpha \times EstimatedRTT + \beta \times SampleRTT$
 - $\alpha + \beta = 1$
 - $0.8 \leq \alpha \leq 0.9$
 - $0.1 \leq \beta \leq 0.2$
- Set timeout based on EstimatedRTT
 - ▶ TimeOut = $2 \times EstimatedRTT$

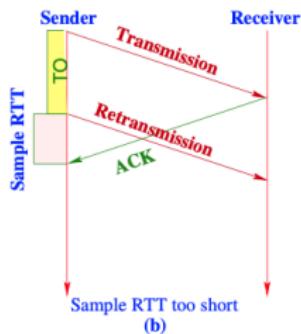
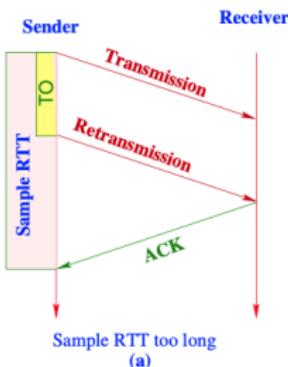
RTT Flaw

- Does ACK really acknowledges a transmission?
- How many retransmissions had taken place before ACK arrived?



- ▶ Wrong samples
 - ▶ Non-representative
- How do we fix this?

Karn/Partridge Algorithm



- Which Sample?
 - ▶ In (a) sample should be for the second attempt
 - ▶ In (b) sample should be for the first attempt

Why?
Why?
- Karn/Partridge suggest
 - ▶ Do not sample RTT when retransmitting
 - ▶ Double TO after each retransmission
 - Similar to back-off algorithm
- Karn/Partridge algorithm was introduced (87) when the Internet was not suffering the current congestion
- Jacobson/Karels came up with a new calculation for average RTT

Jacobson/Karels Algorithm

$$\text{Difference} = \text{SampleRTT} - \text{EstimatedRTT}$$

$$\text{EstimatedRTT} = \text{EstimatedRTT} + \delta \times \text{Difference}$$

$$\text{Deviation} = \text{Deviation} + \delta | \text{Difference} - \text{Deviation} | \quad 0 < \delta < 1$$

- Consider variance when setting timeout value

$$\text{TimeOut} = \mu \times \text{EstimatedRTT} + \phi \times \text{Deviation} \quad \mu = 1 \quad \text{and} \quad \phi = 4$$

- Notes: accurate timeout is important in congestion control

TCP Performance Metrics

- 1 Round trip delay
 - ▶ Access delay + Propagation Delay + Transmission Delay + Queuing Delay
 - ▶ One way delay \Rightarrow Asymmetric?
 - ▶ Maximum tolerable delay
- 2 Delay variation (or jitter)
- 3 Packet loss rate
- 4 Bandwidth
- 5 Throughput variation
 - ▶ Variability in the received bandwidth over time
- 6 File transfer time
- 7 Fairness
- 8 Resource consumption \Rightarrow CPU cycles, memory, battery, etc.

TCP Limitations

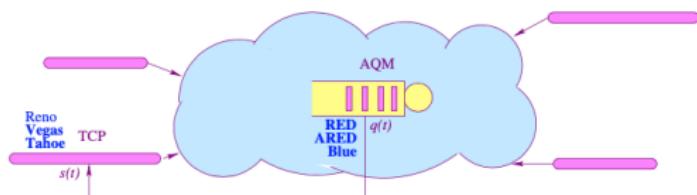
- 1 Limitation of existing transport protocols
 - ▶ Depends on the application domain
 - Some need reliability and ordered delivery
 - Some are sensitive to delay
- 2 TCP was designed with certain assumptions
 - ▶ Segment loss is due to congestion \Rightarrow may not be true
 - ▶ Could be link/Cyclic Redundancy Check (CRC) error
- 3 TCP Performance depends on $delay \times bandwidth$
 - ▶ Could lead to inefficient link utilization
- 4 Standard TCP does not scale well in large $delay \times bandwidth$
 - ▶ Optical and grid applications

Outline I

- 1 Part I (Transport Layer)
- 2 Part II (UDP vs TCP)
- 3 Part III (TCP Flow and Congestion Control)
- 4 Part IV (TCP Congestion Control Approaches)
- 5 Part V (Packet Based Congestion Control)
- 6 References

CC Techniques

- So far, we controlled congestion with **flow** (W_F, W_C, W_A)
 - ▶ Not scalable
 - ▶ Does not distinguish between
 - Packet drop as a result of congestion
 - Packet drop as a result CRC error
- Other measures of network control performance
 - ▶ Delay, loss, and throughput are not enough
 - ▶ Must look at the robustness, complexity, etc.
- Packet level \Rightarrow doesn't look into flows
 - ▶ Packet are too microscopic
 - ▶ Coarser grained models needed
 - ▶ Resurgence of fluid models
- Algorithms
 - ▶ Window Flow Control
 - TCP: Tahoe, Reno, Vegas
 - Receiver window flow control
 - Self Clocking , Vegas
 - ▶ Active Queue Management (AQM) Congestion Control
 - Random Early Detection (RED), Adaptive RED (ARED), etc.



Internet Congestion Control

- Current CCs result in unsatisfactory performance
- Many CC algorithms have been proposed based on AQM
 - ▶ Executed by routers
- Network algorithm \Rightarrow What to feedback?
 - ▶ Implicit/Explicit congestion notification
- Source algorithms \Rightarrow How to react?
 - ▶ Adjusts data rate accordingly

Four Basic Components of CC

- 1 Packet scheduler
 - 2 Buffer management
 - 3 Feedback \Rightarrow best possible explicit feedback
 - 4 Source algorithms \Rightarrow end-system, FC
- Major problems with AQM
- ▶ Parameter setting
 - ▶ The insensitivity to the input traffic load variation.
 - ▶ The mismatch between macroscopic and microscopic behavior of queue length dynamics

Active Queue Management (AQM)

- AQM Provide CC information by marking probabilistically
- Issues
 - ▶ How to measure congestion?
 - ▶ How to embed congestion measures?
 - ▶ How to feed-back congestion information?
- Traditional voice model based on Poisson process fails
 - ▶ Aggregated traffic smooths out
- Internet traffic; File Transfer Protocol (FTP), Web, video are bursty and self-similar
- To cope with the bursty traffic, intelligent CC are based on
 - ▶ FIFO queue management
 - ▶ Per-flow queue management
 - ▶ Scheduling
- Traditional (tail drop) has two drawbacks
 - ▶ Lock-out \Rightarrow a few flows monopolize the queue
 - As a result of tail drop
 - ▶ Full queue may persist for a long time

Early drop may fix the problem \Rightarrow AQM

Random Early Detection (RED) I

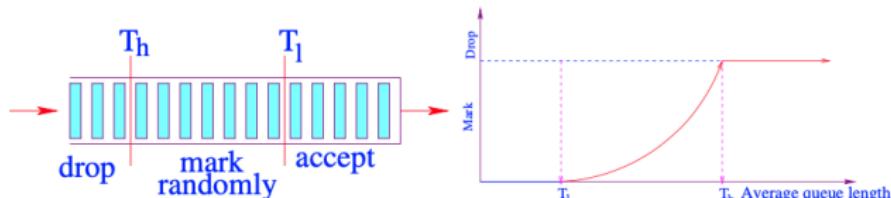
- TCP (Tahoe, Reno, Vegas) detect congestion after a buffer at an intermediate router is full
- RED provides a warning to sources ⇒
 - ▶ It detects **incipient** congestion
- Congestion measure ⇒ Average queue length

$$Q(t+1) = [Q(t) + \lambda(t) - \mu(t)]^+$$

- RED objectives
 - ▶ Congestion avoidance ⇒ proactive rather than reactive
 - Detect the onset of congestion, rather than react to it
 - Maintain the optimal region of high throughput and low delay
 - ▶ Avoid bias against bursty traffic
 - Probabilistic drop/mark
 - ▶ Preventing global synchronization or traffic oscillation
 - ▶ Bounded average queue length
 - ▶ Accommodating an equitable distribution of packet loss
 - Fairness
 - ▶ Providing a lower delay or a lower jitter (delay variation)

Random Early Detection (RED) II

- *p*-linear probability function



- Uses Exponentially Weighted Moving Average (EWMA) *Why?*
- Probabilistically drop packets
- Probabilistically mark packets
 - ▶ Marking require Explicit Congestion Notification (ECN) bit (RFC 2481)

$$\bar{Q} = \begin{cases} (1 - W_q) \times \bar{Q} + W_q \times L_q & \text{if } L_q \neq 0 \\ (1 - W_q)^m \times \bar{Q} & \text{otherwise,} \end{cases} \quad (2)$$

- ▶ W_q determines how rapidly \bar{Q} w.r.t L_q ,
- ▶ m accounts for the periods when the queue has been empty
- It uses two thresholds T_h and T_l to control the growth of the queue

Random Early Detection (RED) III

If	$\bar{Q} < T_\ell$	accept the packet	
If	$\bar{Q} \geq T_h$	mark/drop the packet	(3)
If	$T_\ell \leq \bar{Q} < T_h$	mark/drop probabilistically with, P_a	

- \bar{Q} is calculated at the packet arrival times
 - ▶ Rather than at a fixed time interval *Why?*
- m estimates the number of packets that could have been transmitted during an idle period.
- Within the region $[T_\ell, T_h]$, the probability of packet discard depends on the proximity of \bar{Q} to T_h and it is bounded by P_{max}
- When $T_\ell \leq \bar{Q} < T_h$, P_b increases linearly from zero to P_{max}

$$P_b = P_{max} \times \frac{(\bar{Q} - T_\ell)}{(T_h - T_\ell)} \quad (4)$$

Outline I

- 1 Part I (Transport Layer)
- 2 Part II (UDP vs TCP)
- 3 Part III (TCP Flow and Congestion Control)
- 4 Part IV (TCP Congestion Control Approaches)
- 5 Part V (Packet Based Congestion Control)
- 6 References

References

[Tanenbaum and Wetherall, 2011] Tanenbaum, A. S. and Wetherall, D. J. (2011). Computer Networks: 5th Edition. Prentice Hall PTR.