

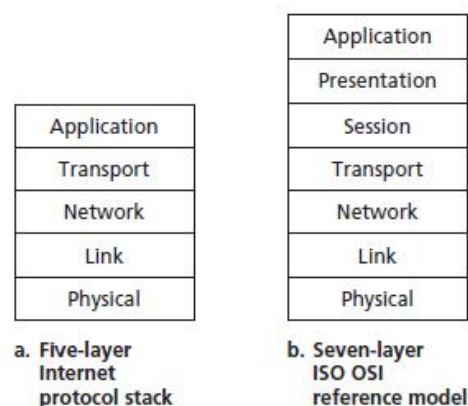
1. Introduction to computer network

1.1 概念

略

1.2 重点

- Protocol Layering
 - application-layer **message**
 - transport-layer **segment**
 - network-layer **datagram**
 - link-layer **frame**
 - **header fields** and a *payload field*
 - Pros and cons
 - Pro: provides a common abstraction for various network technologies
 - Con: Duplicate functionality
 - Con: Performance
 - Con: Header gets big
 - Layer violation
 - Which one does what?



- Application Layer: where network applications and their application-layer protocols reside.
- Transport Layer:
 - TCP provides a connection-oriented service to its applications, **guaranteed** delivery (reliable), flow control, segmentation, and congestion control.
 - The UDP protocol provides a connectionless service to its applications. This is a no-frills service that provides no reliability, no flow control, and no congestion control.
- Network layer:
 - Best effort delivery from one host to another(global)
- Link layer:
 - Best effort delivery from one host to another(local)
- Packet switching: pros and cons

Circuit Switching Vs Packet Switching

Circuit Switching	Packet Switching
Physical path between source and destination	No physical path
All packets use same path	Packets travel independently
Reserve the entire bandwidth in advance	Does not reserve
Bandwidth Wastage	No Bandwidth wastage
No store and forward transmission	Supports store and forward transmission

3

- Delay: analysis and calculation
 - $d_{\text{end-end}} = \sum (d_{\text{proc}} + d_{\text{trans}} + d_{\text{prop}})$
 - Down to the bottom, data are just electronic signal! It travels through some media.
 - Understand the whole process:
 - packet arrive at a switch (input buffer)
 - switch processes it d_{proc}
 - switch dispatch it to a port, entering an output queue d_{queue}
 - packet is transmitted to the media (cable/fibre/air) $d_{\text{trans}}(\text{us-ms}): L/R$
 - the signal travels through media $d_{\text{prop}}(\text{ms}): d/s$

1.3 习题

1. Suppose users share a 2 Mbps link. Also suppose each user transmits continuously at 1 Mbps when transmitting, but each user transmits only 20 percent of the time.

- a. When circuit switching is used, how many users can be supported?
- b. For the remainder of this problem, suppose packet switching is used. Why will there be essentially no queuing delay before the link if two or fewer users transmit at the same time? Why will there be a queuing delay if three users transmit at the same time?
- c. Find the probability that a given user is transmitting.
- d. Suppose now there are three users. Find the probability that at any given time, all three users are transmitting simultaneously. Find the fraction of time during which the queue grows.

Answer: a) 2 users can be supported because each user requires half of the link bandwidth.

b) Since each user requires 1Mbps when transmitting, if two or fewer users transmit simultaneously, a maximum of 2Mbps will be required. Since the available bandwidth of the shared link is 2Mbps, there will be no queuing delay before the link. Whereas, if three users transmit simultaneously, the bandwidth required will be 3Mbps which is more than the available bandwidth of the shared link. In this case, there will be queuing delay before the link.

c) Probability that a given user is transmitting = 0.2

d) Probability that all three users are transmitting simultaneously = $\binom{3}{3} p^3 (1-p)^{3-3}$

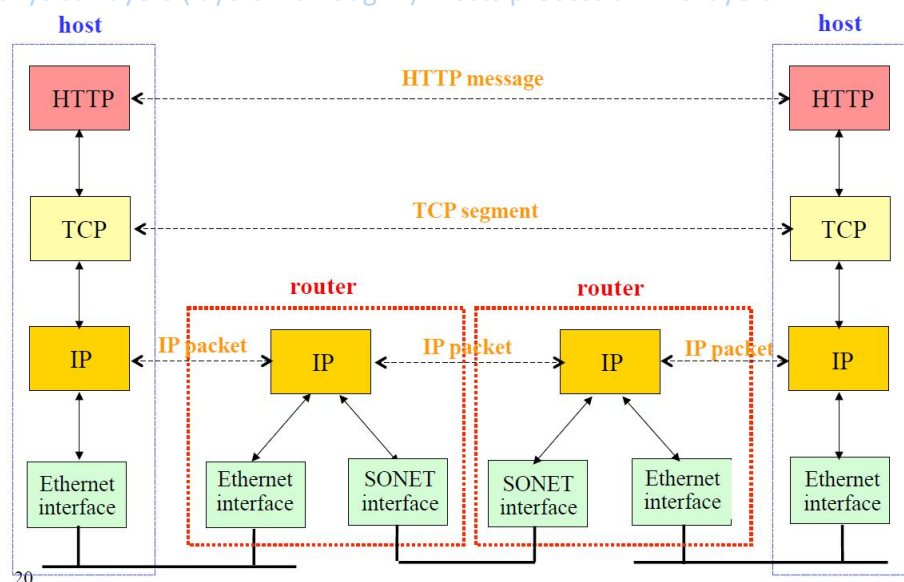
$= (0.2)^3 = 0.008$. Since the queue grows when all the users are transmitting, the fraction of time during which the queue grows (which is equal to the probability that all three users are transmitting simultaneously) is 0.008.

2. List five tasks that a layer can perform. Is it possible that one (or more) of these tasks could be performed by two (or more) layers?

Answer: Five generic tasks are error control, flow control, segmentation and reassembly, multiplexing, and connection setup. Yes, these tasks can be duplicated at different layers. For example, error control is often provided at more than one layer.

6. Which layers in the Internet protocol stack does a router process? Which layers does a link-layer switch process? Which layers does a host process?

Answer: Routers process network, link and physical layers (layers 1 through 3). (This is a little bit of a white lie, as modern routers sometimes act as firewalls or caching components, and process Transport layer as well.) Link layer switches process link and physical layers (layers 1 through 2). Hosts process all five layers.



7. Which of the following services does the Internet network layer provide for the Internet transport layer? (possible multiple answers)

- a. in-order delivery of data segments between processes
- b. best effort delivery of data segments between communicating hosts*
- c. multiplexing and demultiplexing of transport layer segments*
- d. congestion control

2. Application Layer

2.1 概念

- Application architecture: client/server architecture, P2P architecture
 - **Server:**
 - Exports well-defined request/response interface

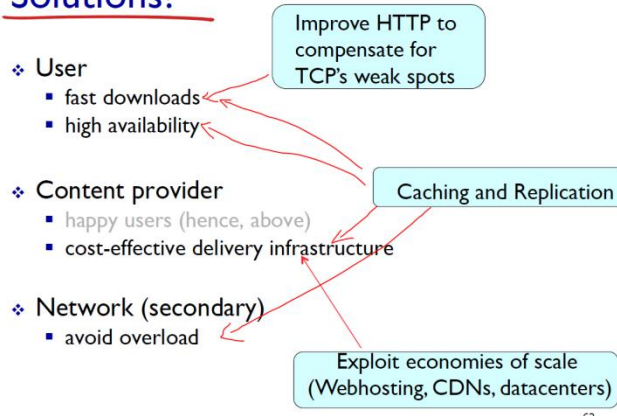
- long-lived process that waits for requests
 - Upon receiving request, carries it out
- **Client:**
 - Short-lived, user side, initiate communication
- **P2P:**
 - no always-on server
 - arbitrary end systems (peers) directly communicate
 - Symmetric responsibility
- Purpose of computer network: allow Inter-Process Communication (IPC)
- *Socket: interface between the application process and the transport-layer protocol. Defined by IP address + port number*
- 4 dimension for choosing transport-layer protocol (important):
 - *Reliability, throughput, timing and security*
 - loss-tolerant applications
 - bandwidth-sensitive applications
 - elastic applications
- **HTTP(important)**
 - Web page: often consist of a base HTML file and several referenced objects
 - URL: hostname and path
 - HTTP is over TCP, port 80
 - Stateless
 - Non-Persistent and Persistent Connections
 - RTT
 - Request: *request line*, header lines, entity body
 - *request line*: method, URL, HTTP version
 - header lines: tell server more info about browser and what I want
 - Entity body: after one blank line. Can be empty for GET
 - Response: status line, header lines, entity body
 - status : 200 OK, 301 moved, 404 Not found
 - *Last-Modified, conditional GET*
 - *Set-cookie: used for user-server interaction*
 - *Caching, institutional cache, CDN*
 - All Text: simple(readable) but not efficient
 - HTTP 1.1: more request method(delete etc.), persistent connection with pipelining
- **HTTPS**
 - Provide security by authentication and bidirectional encryption
- **FTP**
 - TCP, port 20/21
 - Control connection, data connection; out-of-band, in-band
 - NAT
- **SMTP**
 - Used for mail

- *Similarity and difference to HTTP:*
 - SMTP pushes file to server
 - HTTP pulls data from server
 - SMTP require 7-bit ASCII encoding
 - HTTP send images/document independently, SMTP combine them together
- **DNS** domain name system (important):
 - Distributed, hierarchical
 - Application layer
 - Resolve domain name to IP address
 - Usually UDP
- **DASH**: Dynamic, adaptive streaming over HTTP
 - Server divide video into chunks and encode it using different rates
 - Server provide manifest file for different chunk/rate
 - Client request suitable chance/rate based on periodic measurement of C/S bandwidth
- **P2P**(important)
 - Bit torrent: Tracker and torrent
 - rarest first: what and why?
 - tit-for-tat: what and why?
 - optimistically unchoke: what and why?
 - Free-riding: Can someone who never send, receive whole file?
 - DHT (Distributed hash table), p2p-transport p13
 - Hash table: key-value pair
 - Peer has ID, store key-value in peer that has closest ID to key.
 - Peer churn: store 2 successor
 - Short cuts

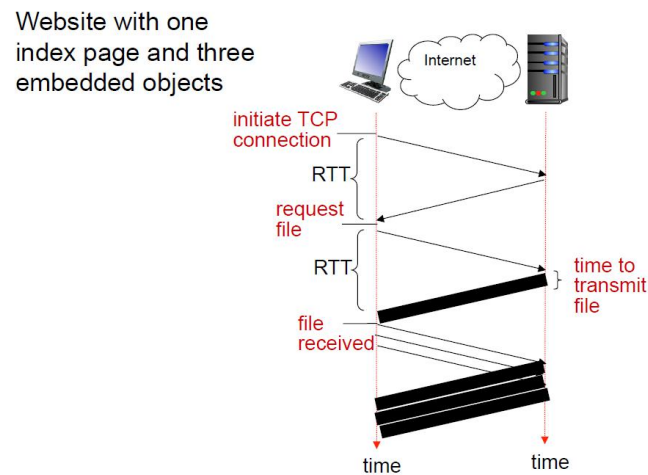
2.2 重点

- Few method to improve HTTP performance:

Solutions?



- *Persistence connection + pipelining: introduced in HTTP1.1*



- Web caches (proxy): satisfy client request without involving origin server
 - acts as both client and server
 - typically cache is installed by ISP (university, company, residential ISP)
 - reduces response time and traffic
 - PPT 74-77: cheaper and faster
- Conditional GET
 - utilising Last-modified and If-modified-since header fields
 - reduces traffic
- Replication and CDN
 - Replicate popular Web site across many machines
 - CDN: Caching and replication as a service
 - Basically, stores content in a server close to user
 - distributed, application-level infrastructure
 - pull and push: predict access rate
 - different with cache: application level, dynamic, smarter(push)
- HTTP: How many socket for n client?
 - $N+1$
- DNS: Goals
 - No naming conflicts
 - Scalable
 - Distributed, autonomous administration
 - Highly available

- Lookups should be fast
- DNS: Three intertwined hierarchies
 - Hierarchical namespace
 - Hierarchically administered
 - (Distributed) hierarchy of servers
 - Root servers
 - Top-level domain(TLD) servers
 - Authoritative DNS servers
 - Example: .edu, berkeley.edu, eecs.berkeley.edu
- DNS resolution:
 - Recursive or iterative query: example
 - Cache and TTL
 - Types of DNS record:
 - Type A: (relay1.bar.foo.com, 145.37.93.126, A)
 - Type NS: (foo.com, dns.foo.com, NS)
 - Type CNAME: (foo.com, relay1.bar.foo.com, CNAME)
 - Type MX: (foo.com, mail.bar.foo.com, MX)
 - PTR
- Insert DNS records
 - Register NS and A record to TLD
 - (networkutopia.com, dns1.networkutopia.com, NS)
 - (dns1.networkutopia.com, 212.212.212.1, A)
 - Insert A(www.networkutopia.com, xxx.xxx.xxx, A) and MX records to authoritative DNS server of network Utopia

Visiting webpage of Network Utopia

1. Alice in Sydney wants to visit web page of Network Utopia
2. Alice's browser sends DNS query to local DNS server
3. Local DNS server contacts TLD server for com (assuming TLD server address is cached in local server to avoid going to root)
4. TLD server sends reply to local server with two RRs
 1. (networkutopia.com, dns1.networkutopia.com, NS)
 2. (dns1.networkutopia.com, 212.212.212.1, A)
5. Local DNS server sends DNS query to 212.212.212.1 asking for Type A record for www.networkutopia.com
6. The response provides RR that contains say 212.212.71.4
7. Alice's browser now initiates HTTP request to 212.212.71.4 to get the web pages for Network Utopia

2.3 习题

Tutorial I 2, 3, II 3,4,5

1. How does SMTP mark the end of a message body? How about HTTP? Can HTTP use the same method as SMTP to mark the end of a message body?

Answer: SMTP uses a line containing only a period to mark the end of a message body.

HTTP uses "Content-Length header field" to indicate the length of a message body.

No, HTTP cannot use the method used by SMTP, because HTTP message could be binary data, whereas in SMTP, the message body must be in 7-bit ASCII format.

3. Transport Layer

3.1 概念

- Application layer is our boss, network layer is ours to command!
- Because...
 - Network layer only find path, doesn't guarantee delivery/order/path
- Therefore...
 - Transport layer need to be able to do it! Or not...
- Transport layer provide logical connection between **processes**; network layer provide communication between hosts
 - Segmentation and reassembling
 - Multiplexing and de-multiplexing: multiple sockets(processes/port)
 - Source port / **destination port**
 - Source IP / **destination IP**
 - Non-persistent HTTP will have **different** socket for each request(connection)
 - Welcoming socket/connection socket
- Connectionless demux/connection oriented demux
 - Connectionless: UDP. Defined by 2-tuple (dest IP + dest port)
 - Connection-oriented: Defined by 4-tuple (src IP and src port is also used)
- UDP
 - Segment header: 64 bits. Source port(return address), dest port, length(in byte, including header), checksum, payload
 - Checksum: 1's complements
 - Detect single (or odd number) bit error. Cannot detect 2 (or even) bit error!
 - Calculation: add together, wraparound carry bit, then take complement.
 - Calculated over header and data
 - Error correction unnecessary (periodic messages)
 - Application: DNS, DHCP, SNMP, Gaming, video/audio
- TCP (important)
 - Handshake/acknowledgement
 - Why these concerns?

❖ Concerns

- Message corruption
- Message duplication
- Message loss
- Message reordering
- Performance

❖ Our toolbox



- Checksums
- Timeouts
- Acks and Nacks
- Sequence numbering
- Pipelining

- Feature: transport_part2 p21
 - Reliable, duplex, connection-oriented, flow controlled
- segment header: 160+ bits
 - source port, dest port, seq(counting byte), ack#(counting byte), flags, checksum, receive window(expected byte#)
- MSS(maximum segment size):
 - $1460 = \text{MTU}(\text{maximum transmission unit, limit of IP datagram, } 1500) - \text{IP header}(20) - \text{TCP header}(20)$
- Sequence number: $\text{ISN} + k$
 - ISN: initial sequence number
 - K: kth byte of stream
- ACK sequence number
 - = next expected byte
 - = seq# + length
 - Cumulative ACK
- Buffers out-of-order packet
- Single timer
- Fast retransmit: use duplicate(3) acks to trigger early retransmission
- Rwnd: limit amount of in-flight data to rwnd value
- Connection management
- Syn-flooding and syn cookie

3.2 重点

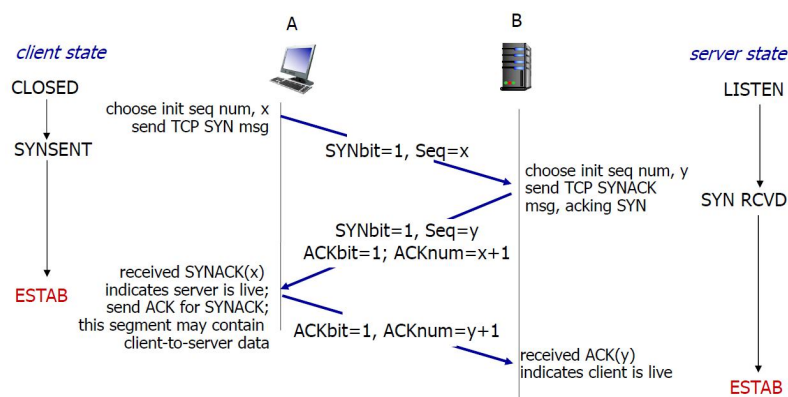
- Why 4 tuple for a TCP socket?

Everybody can connect to google.com:80 (even same computer can open different window), but how does Google's server know what to send to everyone? Therefore src IP and src port are needed.

UDP: also need src IP and src port, but user application receives everything, need to handle it. TCP socket doesn't receive these noise.

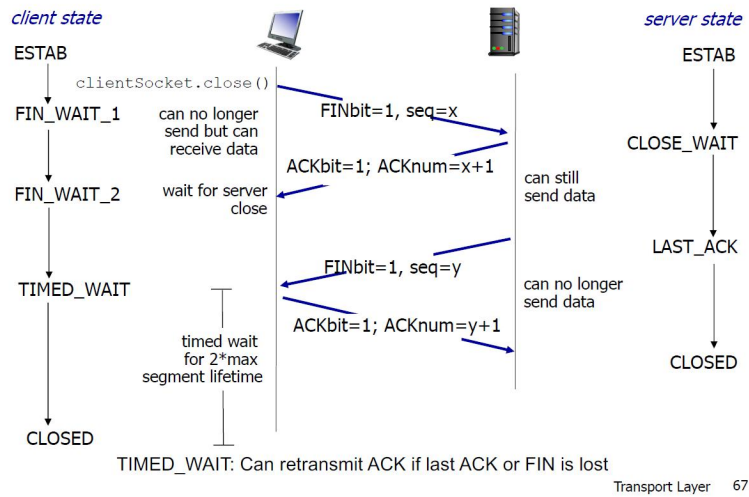
- Develop a reliable transport protocol(rdt) step by step: p2p_transport p54+
 - 1.0: assume network layer is reliable, transport layer doesn't have to do anything!
 - 2.0: deals with error bit. (partial) solution:
 - Error detection (with checksum)
 - Feedback: ack/nak
 - 2.1: deals with corruption of ack/nak, and possible duplication by:
 - Sequence number: sender add a number to pkt, receiver keeps track of which to expect
 - Stop and wait
 - Q: why only (0,1) sufficient for seq#?

- 2.2 Nak free!
 - Ack last packet received
- 3.0: deals with loss:
 - Timeout is introduced
- Pipelined(sliding window) 3.0:
 - Deals with utilization: how?
 - Go-Back-N:
 - timer for oldest in-flight pkt
 - ack(n) means all packets up to n are received (cumulative ACK)
 - on timeout(n), retransmit packet n and all higher seq# in window
 - for receiver: out-of order pkt: re-ack
 - Selective repeat
 - Timer for each pkt
 - for receiver: Buffer out-of-order pkt, Ack individually
 - window size: $\leq (\text{seq\# range})/2$
 - Study the state machine
- TCP timeout estimation
 - $\text{EstimatedRTT} = (1 - \alpha) * \text{EstimatedRTT} + \alpha * \text{SampleRTT}$
 - $\text{DevRTT} = (1 - \beta) * \text{DevRTT} + \beta * |\text{SampleRTT} - \text{EstimatedRTT}|$
(typically, $\beta = 0.25$)
 - $\text{TimeoutInterval} = \text{EstimatedRTT} + 4 * \text{DevRTT}$
 - Retransmission is excluded(p44)
- TCP connection management
 - 3-way handshake



Syn loss: timeout(3s)

4-way termination



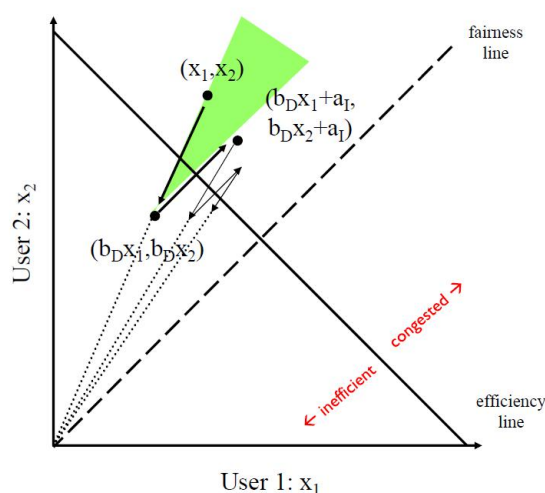
3.3 练习

R8, R9, R14, R15, P14, P25

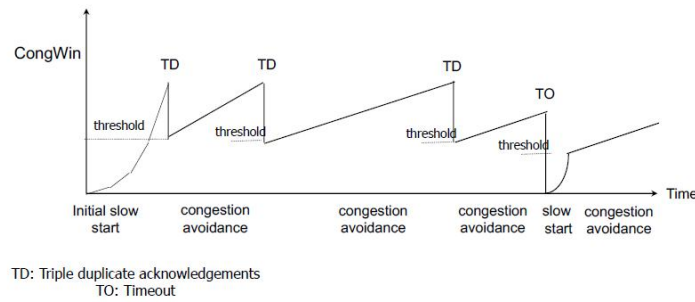
4. Transport Layer – Congestion control（期中考试之后）

4.1 概念

- Congestion Collapse
- Load-Throughput-Delay graph
 - Why is the Cliff? Packet loss and congestion collapse
 - Solution: senders limit sending rate
- 2 approaches:
 - end2end(congestion inferred from loss & delay)
 - network-assisted(routers provide feedback)
- TCP's approach: controls number of packets in flight
 - $\text{rate} \approx \frac{cwnd}{RTT} \text{ byte/sec}$
- CWND: Congestion window, calculated by sender using algorithms below(SS, AIMD etc.)
- RWND: Advertised window, feedback from receiver
- Sender window = $\text{LastByteSent} - \text{LastByteAcked} = \min(\text{CWND}, \text{RWND})$
- MSS: Maximum Segment Size (TCP payload size), only for pedagogical purposes
- Dup ACK and Timeout: isolated loss vs. severe loss
- Rate adjustment:
 - 2 action: increase(upon receipt of ack)/decrease(upon loss)
 - 2 motive/phase: discovering bandwidth/adjusting to variation
- SS: Slow Start
 - Goal: discover/estimate bandwidth
 - Start slow/ramp up quickly
 - Double cwnd every **RTT**: exponential
- AIMD: Additive-Increase, Multiplicative Decrease
 - Adjust to varying bandwidth
 - Increase CWND by 1 every **RTT** (or $\text{cwnd} = \text{cwnd} + 1/\text{cwnd}$ on every ACK)
 - Cut CWND in half after loss
 - When to stop slow start: slow start threshold
 - $\text{ssthresh} = \text{CWND}/2$ on every TD/TO
 - When $\text{CWND} = \text{ssthresh}$, switch from SS to AIMD
- Fairness(Important): AIAD vs AIMD



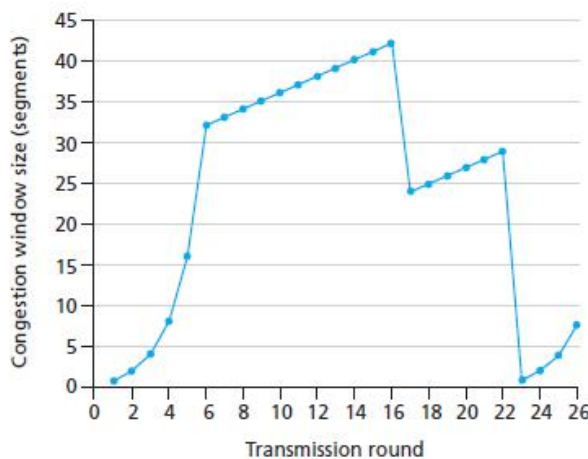
- Reno vs Tahoe: Tahoe set CWND to 1 on any loss, reno only cut in half on TD (Below is Reno)



4.2 习题

P40. Consider Figure 3.58. Assuming TCP Reno is the protocol experiencing the behavior shown above, answer the following questions. In all cases, you should provide a short discussion justifying your answer.

- Identify the intervals of time when TCP slow start is operating.
- Identify the intervals of time when TCP congestion avoidance is operating.
- After the 16th transmission round, is segment loss detected by a triple duplicate ACK or by a timeout?
- After the 22nd transmission round, is segment loss detected by a triple duplicate ACK or by a timeout?
- What is the initial value of $ssthresh$ at the first transmission round?
- What is the value of $ssthresh$ at the 18th transmission round?
- What is the value of $ssthresh$ at the 24th transmission round?
- During what transmission round is the 70th segment sent?
- Assuming a packet loss is detected after the 26th round by the receipt of a triple duplicate ACK, what will be the values of the congestion window size and of $ssthresh$?



Answer:

- TCP slowstart is operating in the intervals [1,6] and [23,26]
- TCP congestion avoidance is operating in the intervals [6,16] and [17,22]
- After the 16th transmission round, packet loss is recognized by a triple duplicate ACK. If there was a timeout, the congestion window size would have dropped to 1.

- d) After the 22nd transmission round, segment loss is detected due to timeout, and hence the congestion window size is set to 1.
- e) The threshold is initially 32, since it is at this window size that slow start stops and congestion avoidance begins.
- f) The threshold is set to half the value of the congestion window when packet loss is detected. When loss is detected during transmission round 16, the congestion windows size is 42. Hence the threshold is 21 during the 18th transmission round.
- g) The threshold is set to half the value of the congestion window when packet loss is detected. When loss is detected during transmission round 22, the congestion windows size is 29. Hence the threshold is 14 (taking lower floor of 14.5) during the 24th transmission round.
- h) During the 1st transmission round, packet 1 is sent; packet 2-3 are sent in the 2nd transmission round; packets 4-7 are sent in the 3rd transmission round; packets 8-15 are sent in the 4th transmission round; packets 16-31 are sent in the 5th transmission round; packets 32-63 are sent in the 6th transmission round; packets 64 – 96 are sent in the 7th transmission round. Thus packet 70 is sent in the 7th transmission round.
- i) The threshold will be set to half the current value of the congestion window (8) when the loss occurred and congestion window will be set to the new threshold value + 3 MSS . Thus the new values of the threshold and window will be 4 and 7 respectively.
- j) Threshold is 21, and congestion window size is 1.
- k) Round 17, 1 packet; round 18, 2 packets; round 19, 4 packets; round 20, 8 packets; round 21, 16 packets; round 22, 21 packets. So, the total number is 52.

P46. Consider that only a single TCP (Reno) connection uses one 10Mbps link which does not buffer any data. Suppose that this link is the only congested link between the sending and receiving hosts. Assume that the TCP sender has a huge file to send to the receiver, and the receiver's receive buffer is much larger than the congestion window. We also make the following assumptions: each TCP segment size is 1,500 bytes; the two-way propagation delay of this connection is 150 msec; and this TCP connection is always in congestion avoidance phase, that is, ignore slow start.

- What is the maximum window size (in segments) that this TCP connection can achieve?
- What is the average window size (in segments) and average throughput (in bps) of this TCP connection?
- How long would it take for this TCP connection to reach its maximum window again after recovering from a packet loss?

Answer:

- Let W denote the max window size measured in segments. Then, $W \cdot \text{MSS} / \text{RTT} = 10\text{Mbps}$, as packets will be dropped if the maximum sending rate exceeds link capacity. Thus, we have $W \cdot 1500 \cdot 8 / 0.15 = 10 \cdot 10^6$, then W is about 125 segments.
- As congestion window size varies from $W/2$ to W , then the average window size is $0.75W = 94$ (ceiling of 93.75) segments. Average throughput is $94 \cdot 1500 \cdot 8 / 0.15 = 7.52\text{Mbps}$.
- $94/2 \cdot 0.15 = 7.05$ seconds, as the number of RTTs (that this TCP connections needs in order to increase its window size from $W/2$ to W) is given by $W/2$. Recall the window size increases by one in each RTT.

5. Network Layer

5.1 概念

- Common IP protocol, supports any data-link/transport layer protocol
- Handled in every host/router
- Basic function: forwarding and routing
 - forwarding: move packets from router's input to appropriate output
 - routing: determine route from source to dest(can be done in an ad-hoc way or centralised way)
- Data plane vs Control Plane
 - Traditional or Per-router control plane
 - SDN control plane
- Network Layer service models
 - Bandwidth, loss, order, timing and congestion feedback
 - IP guarantees none: Best effort delivery of packet(globally)
- Router:
 - Input ports(queuing, look up and forwarding done here, decentralised switching)
 - Routing processor(control plane)
 - High-speed switching fabric
 - Output ports
- Input port
 - Queuing: Head-of-the-Line (HOL) blocking: queued datagram at front of queue prevents others in queue from moving forward
 - Destination based forwarding
 - Longest prefix matching
 - TCAM
- Switching fabric
 - Via memory
 - Via bus
 - Via interconnection network (crossbar)
- Output ports
 - Buffering and scheduling
 - Buffer size proportional to link capacity
 - Scheduling: FIFO, priority based, **Round Robin**(RR)