

Transport Services and Protocols

- ↳ Provides logical comm. b/w app processes running on diff hosts
 - ↳ during comm the location of processes isn't necessary → no physical comm.
- ↳ so we believe the devices are connected intermediately and have direct comm path



Sender side

- ↳ breaks app msg into segments
- ↳ assign header to segments
- ↳ passes to network layer

receiver side

- ↳ reassembles segments into msgs
- ↳ passes to app layer

Application Layer

- ↳ user browser requests to send a msg

Transport Layer

- ↳ always implemented on host
- ↳ breaks msg into segments
- ↳ each segment is given a header



Network Layer

- ↳ puts network layer header on segments
- ↳ the segments are routed through network

Destination Side

TRANSPORT LAYER

- ↳ segments are reassembled into msgs
- ↳ then passed to application layer

Transport Layer

↳ logical comm. b/w processes

household analogy:

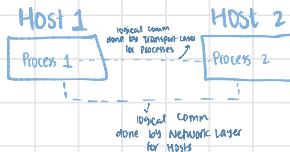
12 kids in Ann's house sending letters to 12 kids in Bill's house:

- hosts = houses
- processes = kids
- app messages = letters in envelopes
- transport protocol = Ann and Bill who demux to in-house siblings
- network-layer protocol = postal service

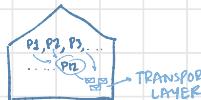
VS

Network Layer

↳ logical comm. b/w hosts



Ann's Host



Bills host



Transport Layer Protocols

- ↳ TCP
- ↳ UDP

Limitations

- ↳ delay guarantees
- ↳ bandwidth guarantees

Transmission Control Protocol (TCP)

- ↳ reliable services
- ↳ in order delivery → data received in order
- ↳ congestion control → hints sender to a specific rate to reduce traffic and blocks
- ↳ flow control → data is sent at a rate which receiver can process
- ↳ connection setup → handshaking
- ↳ connection oriented transport

User Datagram Protocol (UDP)

- ↳ unreliable → no guarantee whether you receive what or not
- ↳ unordered delivery
- ↳ best effort protocol
- ↳ connection less transport

PDoS

- ↳ faster as less delay → due to fewer packet losses

PROS

- ↳ faster
- ↳ more efficient

in transmission

Send next ACK

in retransmission

Send cumulative ACK

in retransmission

Send Prev ACK

in transmission

Send current ACK

A

Multiplexing at sender

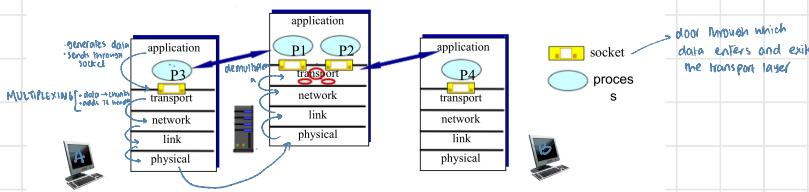
- ↳ handle data from multiple sockets
- ↳ add transport layer header
- ↳ send to network layer

B

DeMultiplexing

- ↳ use header to deliver received segments to correct socket

→ distribute among multiple processes



- ↳ each host can have multiple sockets and run multiple processes

HOW DEMULTIPLEXING WORKS

- ↳ host receives IP datagram

- ↳ each datagram has

- ↳ source IP address

] to identify host

- ↳ destination IP address

- ↳ each datagram carries

- ↳ 1 transport layer segment

- ↳ each segment has

- ↳ source Port no.

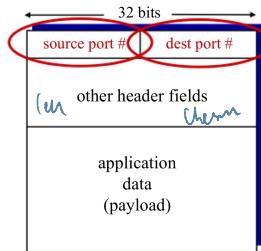
- ↳ destination Port no.

- ↳ host uses

- IP address + Port no.

↳ identification of each process

- to direct segment to correct socket



TCP/UDP segment format

Connectionless demultiplexing

how UDP does demultiplexing

- recall: created socket has host-local port #:

```
DatagramSocket mySocket1 = new DatagramSocket(1234);
```

port

- when host receives UDP segment:

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

- recall: when creating datagram to send into UDP socket, must specify

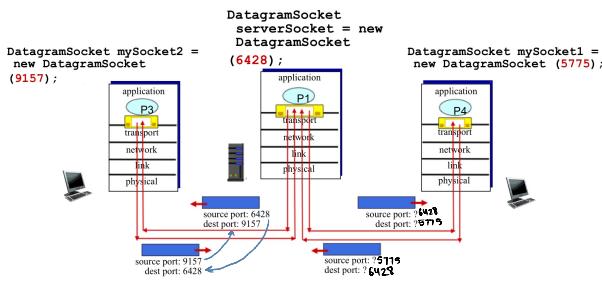
- destination IP address
- destination port #

YOU NEED
For identification

↳ same destination IP

↳ diff source IP

↳ Same socket



Transport Layer 3-19

know source
by looking at
TCP segment in header
containing sourceIP, source port

Connection oriented Demultiplexing

→ HOW TCP does demultiplexing

- TCP socket identified by 4-tuple:

- source IP address
- source port number
- dest IP address
- dest port number

- server host may support many simultaneous TCP sockets:

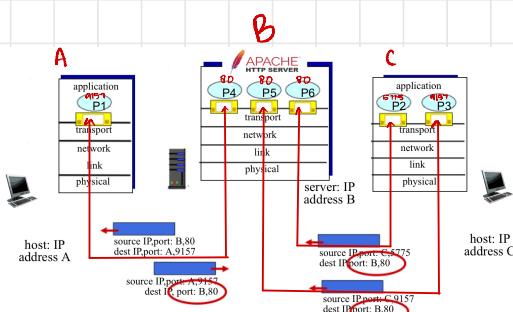
- each socket identified by its own 4-tuple
- web servers have different sockets for each connecting client
- non-persistent HTTP will have different socket for each request

- demux: receiver uses all four values to direct segment to appropriate socket

↳ same destination IP

↳ diff source IP

↳ diff socket, same port



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 happen at *all* layers

Connectionless Transport: UDP

User Datagram Protocol (UDP)

↳ frills protocol → segments may be lost
↳ bare bones

↳ unordered delivery

↳ best effort service → send and hope for the best

↳ not reliable → no priority queue, you receive data or not

↳ UDP use:

↳ streaming multimedia apps → online gaming, video conferencing

↳ DNS → loss tolerant, rate sensitive

↳ SNMP

↳ UDP is a connectionless protocol

↳ no handshaking b/w UDP sender and receiver

↳ each UDP segment handled independently

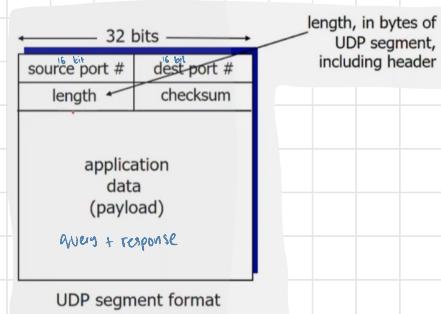
↳ How to make UDP reliable

↳ add reliability at application layer

↳ add congestion control at application layer

↓

CON: more work for application developer → add reliability, debug error



PROS

↳ no delay, faster → as no connection establishment which adds RTT delay

↳ no congestion control → can be as fast as desired

↳ simple: no connection state at sender/receiver

↳ small header size → 20 bytes → good for sending huge amount of data

→ TCP header is 40 bytes

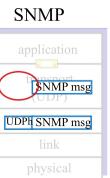
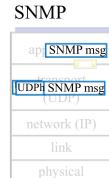
↳ helps with reliability → checksum

↳ can function when network service is compromised

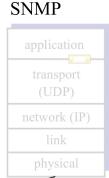
UDP: Transport Layer Actions



- UDP sender actions:
- is passed an application-layer message
 - determines UDP segment header fields values
 - creates UDP segment
 - passes segment to IP



- UDP receiver actions:
- receives segment from IP
 - checks UDP checksum header value
 - extracts application-layer message
 - demultiplexes message up to application via socket



UDP checksum

↳ a method to detect errors in transmitted segment

sender:

- treat contents of UDP segment (including UDP header fields and IP addresses) as sequence of 16-bit integers
- **checksum:** addition (one's complement sum) of segment content
- checksum value put into UDP checksum field
- compute checksum of received segment
- check if computed checksum equals checksum field value:
 - Not equal - error detected
 - Equal - no error detected. *But maybe errors nonetheless? More later*

example: add two 16-bit integers

$$\begin{array}{r}
 1'1'1'0'0'1'1'0'0'1'1'0'0'1'1'0 \\
 + 1'1'0'1'0'1'0'1'0'1'0'1'0'1'0'1'0 \\
 \hline
 \text{wraparound} \quad \text{odd if} \quad 0'1'1'0'1'1'1'0'1'1'1'0'1'1'1'0'1'1 \\
 \text{sum} \quad \text{here is complemented} \quad + 1 \\
 \text{checksum} \quad \hline
 1'0'1'1'1'0'1'1'1'0'1'1'1'1'1'0'0 \\
 + 0'1'0'0'0'1'0'0'0'1'0'0'0'0'1'1 \\
 \hline
 \text{if all } 1's : \text{no change} \\
 \text{if any } 0' : \text{change} \rightarrow \text{error}
 \end{array}$$

Even though numbers have changed (bit flips), **no** change in checksum!

$$\begin{array}{r}
 0'1'1'1'0'0'1'0'1'0'1'1'0'0'1'1 \\
 + 1'0'1'1'0'0'1'1'1'0'1'0'1'0'0'0 \\
 \hline
 0'0'0'1'0'0'1'1'0'0'1'0'1'1'0'1'1 \\
 \text{wraparound} \quad + 1 \rightarrow \text{wraparound}
 \end{array}$$

$$\begin{array}{r}
 0'0'1'0'0'1'1'0'0'1'0'1'1'0'0 \\
 + 1'0'1'1'1'0'1'1'0'0'1'1'0'1'0'1 \\
 \hline
 1'1'1'0'0'0'0'1'1'0'0'1'0'0'0'1 \\
 \text{sum} \\
 0'0'0'1'1'1'1'0'0'1'1'0'1'1'1'0 \\
 \text{1's complement} \rightarrow \text{checksum}
 \end{array}$$

hence no change

Question 2:

[10 Marks]

Checksum is an error-detecting code used in many Internet standard protocols, including IP, TCP, and UDP. You have to generate Checksum for the following data blocks, which are transmitted using UDP.

N1	N2	N3
0111001010110011	1011001110101000	1011101100110101

Q2)
$$\begin{array}{r}
 0'1'110'0'1'01'0\ 1100011 \\
 + 101100111010101000 \\
 \hline
 00100110010110'1'1
 \end{array} \rightarrow N1$$

$$\begin{array}{r}
 0'1'110'0'1'01'0\ 1100011 \\
 + 1011101100110101 \\
 \hline
 1110000110010001
 \end{array} \rightarrow N2$$

$$\begin{array}{r}
 0'1'110'0'1'01'0\ 1100011 \\
 + 1011101100110101 \\
 \hline
 0001111001101110
 \end{array} \rightarrow N3$$

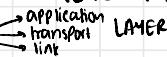
$$\begin{array}{r}
 0'1'110'0'1'01'0\ 1100011 \\
 + 1011101100110101 \\
 \hline
 1110000110010001 \\
 0001111001101110
 \end{array} \rightarrow \text{sum}$$

$$\begin{array}{r}
 0'1'110'0'1'01'0\ 1100011 \\
 + 1011101100110101 \\
 \hline
 0001111001101110
 \end{array} \rightarrow 1's \text{ compliment}$$

↓

CHECKSUM

Principles of reliable data transfer protocols (rdt)

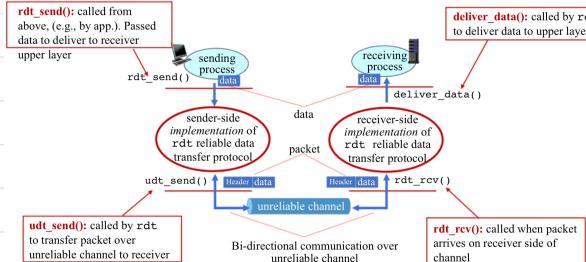
↳ important in  LAYER

a) provided service



reliable service *abstraction*

b) service implementation



AL generate data → data

call rdt_send() → sends data to TL

divide in chunks, add header → packet

call udt_send() → transfers packet from unreliable channel

call rdt_receive() → receives packet

call deliver_data() → delivers data to upper layer

Timing seq diagrams

Reliable Data Transfer (rdt)

↳ we'll:

↳ incrementally develop sender, receiver side rdt

↳ unidirectional data transfer

↳ bidirectional control info transfer

↳ USE FSM to specify sender, receiver

rdt Protocols

1. rdt 1.0

2. rdt 2.0

3. rdt 2.1

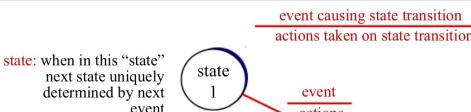
4. rdt 2.2

5. rdt 3.0

stop and wait
protocols

data info: video

control info: OK / NOT OK msg on video being received

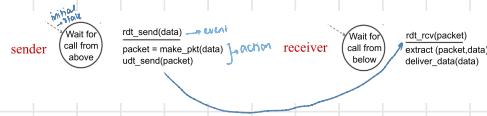


Trans

Rdt 1.0: reliable transfer over reliable channel

- ↳ reliable transfer
- ↳ reliable channel
- ↳ no bit errors
- ↳ no loss of packets

- ↳ separate FSMs for sender, receiver
- ↳ sends data into underlying channel
- ↳ receives data from underlying channel



Rdt 2.0: channel with bit errors

- ↳ underlying channel may flip bits in packet $101 \rightarrow 010$

↳ checksum to detect errors

↳ HOW TO RECOVER FROM ERRORS

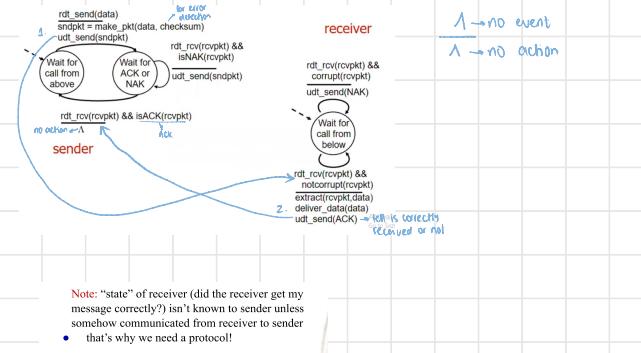
1. Acknowledgements (ACKs): receiver explicitly tells sender that packet received OK $\rightarrow 0$
 2. Negative Acknowledgements (NAKs): receiver explicitly tells sender that packet had ERRORS $\rightarrow 1$
- ↳ sender retransmits packet on receipt of NAK

PROS

- ↳ error detection
- ↳ feedback: control msg from r to s

↳ FSM specifications

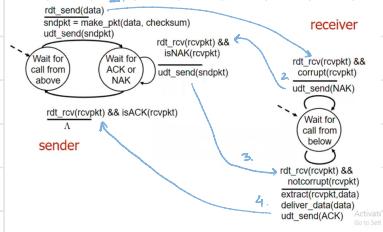
Operation with no errors



Sending: 2 states

Receiver: 1 state

Corrupted Packet scenario



Stop and wait protocol

- ↳ sender sends one packet then waits for receiver response
- ↳ has low response time \rightarrow as waiting

FLAW

↳ What happens if ACK/NAK corrupted \rightarrow ACK/NAK flipped

↳ sender doesn't know what happened to receiver

↳ can't just transmit a duplicate

SOLUTION 1

Handling duplicates

↳ sender retransmits current pkt if ACK/NAK corrupted

↳ sender adds sequence number to each pkt \rightarrow Packet

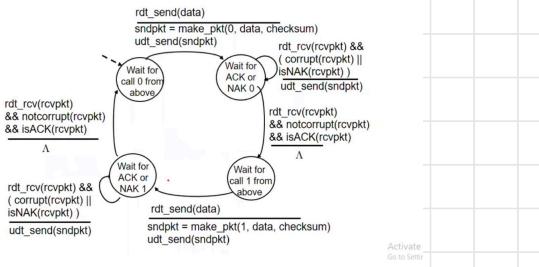
↳ If receiver receives pkt of same seq.no.

↳ it will discard one and deliver the other

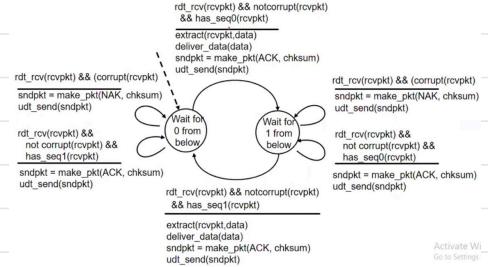


rdt 2.1:

rdt 2.1: sender, handles garbled ACK/NAKs



rdt 2.1: receiver, handles garbled ACK/NAKs



↳ discussion

Sender

↳ seq no. added to PKT \rightarrow 2 seq nos: (0,1) will suffice

↳ check if ACK/NAK corrupted

↳ twice as many states

↳ as state must remember whether

'expected' PKT should have

seq no. of 0 or 1

receiver

↳ checks if received packet is duplicate

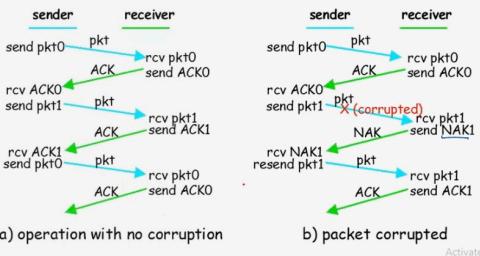
↳ state indicates whether 0 or 1 is expected PKT seq. no.

↳ receiver can't know if its last ACK/NAK received OK at sender

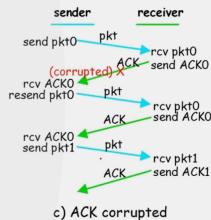
Sending: 4 states

Receiver: 2 state

rdt 2.1 in action



rdt 2.1 in action (cont)



SOLUTION 2

rdt 2.2: a NAK free protocol

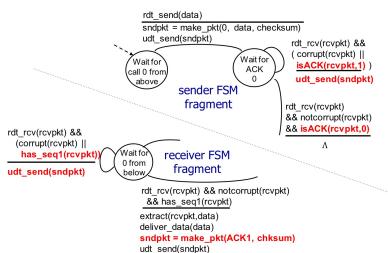
↳ same functionality as rdt 2.1, using ACKs only

↳ instead of NAK, receiver sends ACK for last pkt received OK

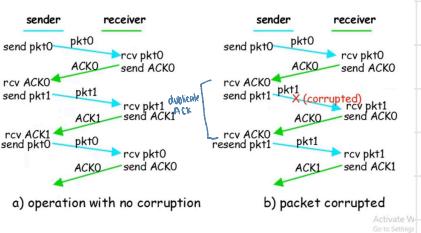
↳ receiver must explicitly include seq_no. of pkt being ACKed

↳ duplicate ACK at sender results in same actions as NAK:

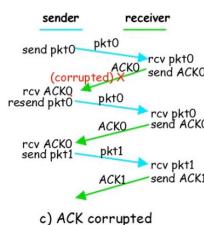
↳ retransmits current PKT



rdt 2.2 in action



rdt 2.2 in action (cont)



if PKT corrupted → send ^R last ACK

if ACK corrupted → resend ACK1

rdt3.0 : channels with errors and loss



↳ New channel assumption

↳ underlying channel can also lose packets (data, ACKs)

↳ checksum, seq no., ACKs, retransmission will be of help but not quite enough

↳ How do humans handle lost sender to rec words in conversation

↓ Approach

use timer

↳ Sender waits reasonable amount of time for ACK

↳ retransmits if no ACK received in this time

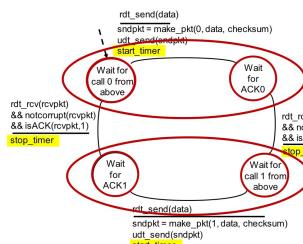
↳ If Pkt (or ACK) just delayed (not lost):

↳ retransmission will be duplicate → but seq no. already handles this 2.1, 2.2

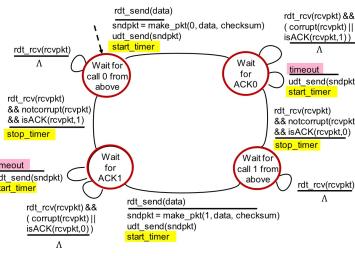
↳ receiver must specify seq no of pkt being ACKed

↳ use countdown timer to interrupt after "reasonable" amount of time

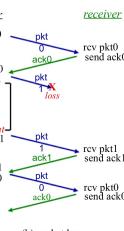
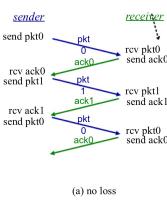
t3.0 sender



s.u sener

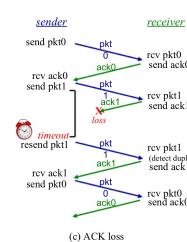


rdt3.0 in action

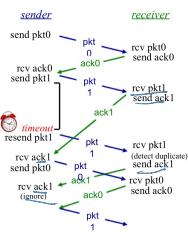


Transport Layer

rdt3.0 in action



Transport Layer



Transport Layer 3-4

if Pkt corrupted → send new ACK
if ACK corrupted → resnd ACK-1

if Pkt loss → after timeout resnd Pkt 1
if ACK loss → after timeout resnd ACK-1
if delayed ACK → delete dupl ACK
ignore dupl ACK

↳ Performance of fdt3.0

↳ correct but slow performance → as stop and wait protocol

↳ has capacity to transfer more data
but protocol limits it down

conversions

$$1 \text{ byte} = 8 \text{ bit}$$

$$1 \text{ Gb} = 10^9$$

$$1 \text{ micros} = 10^{-6}$$

$$1 \text{ milli s} = 10^{-3}$$

Q) link: 1 Gbps, Prob delay: 15ms, Packet: 8000 bit

$$D_{\text{tran}} = \frac{L}{R} = \frac{8000 \text{ bits}}{10^9 \text{ bits/sec}} \cdot 8 \text{ micro s} = 0.008 \text{ ms}$$

↓
time to transmit
packet into channel

$$d_{\text{trans}} = \frac{L}{R} \quad \begin{matrix} \rightarrow \text{packet length} \\ \text{bits} \end{matrix}$$

↓
link transmission rate bps

$$d_{\text{PROP}} = \frac{d}{S} \quad \begin{matrix} \rightarrow \text{length of} \\ \text{physical link} \end{matrix}$$

↓
propagation speed
27000 m/s

$$RTT = 2 \times \text{Prob delay} = 2(15) = 30 \text{ msec}$$

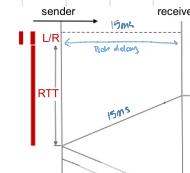
↳
Round trip time
Waiting from when
sender sent first
request received
receiver sent first

↳ U sender: utilization - fraction of time sender busy sending

$$U_{\text{sender}} = \frac{L/R}{RTT + L/R} = \frac{0.008}{30 + 0.008} = \frac{0.008}{30.008} = 0.00027 \text{ ms}$$

• 1KB pkt every 30ms

$$\text{Throughput} = \frac{1 \text{ KB}}{30 \text{ ms}} = 33 \text{ KB/s}$$



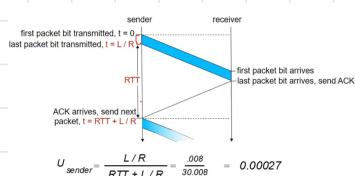
$$U_{\text{sender}} = \frac{\text{Window size}}{RTT + \text{Packet size}}$$

↓
RTT: 2 × Prob delay

$$\frac{10,000}{5 \times 10^9}$$

SD

↳ STOP AND WAIT operation → low performance



Part c) A point-to-point satellite transmission link of 5 Gbps and 30ms propagation delay connecting two computers uses a stop-and-wait ARQ strategy. All packets consist of 10,000 bits and average round trip time RTT is 50ms.

- i) Calculate the link efficiency (utilization) when only one packet is sent.
- ii) Calculate the new link efficiency (utilization) when pipelining is applied, now sender send three packets back-to-back.

Part c)

i) Answer:

$$U_{\text{sender}} = \frac{L/R}{RTT+L/R} = \frac{10000/5\text{Gbps}}{50\text{ms}+10000/5\text{Gbps}} = 0.0000399$$

ii) Answer:

$$U_{\text{sender}} = \frac{3L/R}{RTT+L/R} = \frac{3 \times 10000/5\text{Gbps}}{50\text{ms}+10000/5\text{Gbps}} = 0.0001199$$

- b) How the channel utilization may be increased?

Solution:-

The problem is the low utilization of the system which could be solved by pipelining protocol

Question 4: CLO-01

[5x2=10 points]

Write the difference, in maximum 30 words each (answer will not be marked after 30 words)

- a) (In Forwarding table computing/learning) process at switches versus routers

Solution:-

routers: compute tables using routing algorithms, IP addresses

switches: learn forwarding table using flooding, learning, MAC addresses

- b) intra-AS routing versus intra-domain routing

Solution:-

No difference.

- c) Time-to-Live field of IPv4 Header versus Hop-Limit field of IPv6 Header

Solution:-

No difference, same function, just different naming convention.

- d) Flow control mechanism in link layer versus flow control in transport layer

Solution:-

On the link layer, pacing between adjacent sending and receiving nodes on a link is handled, while on transport layer, overflow at the final destination node is avoided.

- e) Connectionless link versus Unreliable link

Solution:-

Connection-less: no handshaking between sending and receiving NICs

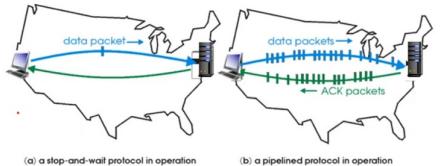
Unreliable: receiving NIC doesn't send ACKs or NAKs to sending NIC

Solves low performance

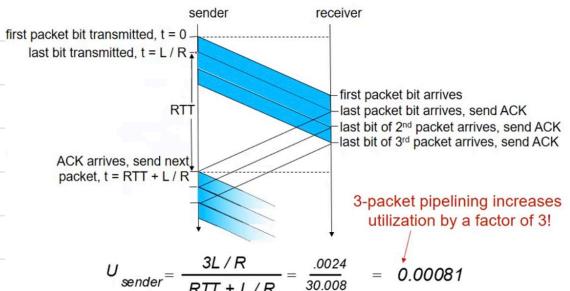
good approach

↳ Pipelined Protocols Operation

- ↳ Pipelining: sender allows multiple, 'in-flight', yet to be acknowledged packets
 - ↳ range of seq no. must be increased → k bit



Pipelining: increased utilization



Generic forms of Pipeline Protocols

1. GO-Back-N
 2. Selective Repeat

Pipelined protocols: overview

Go-back-N:

- ❖ sender can have up to N unacked packets in pipeline
 - ❖ receiver only sends **cumulative ack**
 - doesn't ack packet if there's a gap
 - ❖ sender has timer for oldest unacked packet
 - when timer expires, retransmit *all* unacked packets

Selective Repeat:

- ❖ sender can have up to N unacked packets in pipeline
 - ❖ rcvr sends *individual ack* for each packet
 - ❖ sender maintains timer for each unacked packet
 - when timer expires, retransmit only that unacked packet

1. Go-Back-N: sender

retransmitt PKT from
first out of order

↳ sender can have upto N unacked \rightarrow sender window

PKS in Pipeline

↳ k-bit seq. no in Pkt header



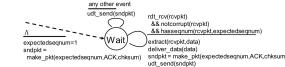
Cumulative ACK:

↳ ACK(n): ACKs all packets up to, including seq no 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. no. packets in window



lost
↑
01/3
↳ will only receive ACK of 0,1 even though S was also sent
↳ retransmit 2,3

Go Back-N: receiver (GBN)

LACK ONLY: always send ACK for correctly-received packet so far

with highest in-order seq no

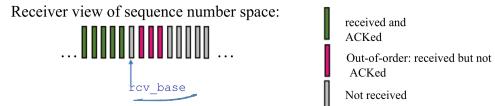
↳ may generate duplicate ACKs

↳ need only remember expected sequence

↳ On receipt of out-of-order packet: → 3 received but 2 were so not in order hence discarded

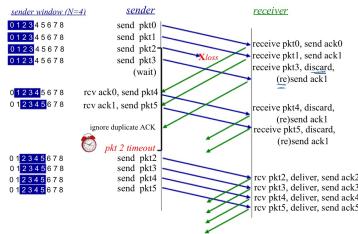
↳ can discard (don't buffer) or buffer → an implementation decision

↳ re-ACK Pkt with highest in-order seq. no



sends ack of last
acked value

Go-Back-N in action



CON

↳ Cumulative ack

↳ will have to reflect from first
out of order

↳ retransmission high

↳ more bandwidth used

R
resend last ACK
S
ignore duplicate ACK

2. Selective Repeat (SR)

↳ Sender can have up to N unacked → sender window

Pkts in Pipeline

↳ receiver individually acknowledge all correctly received packets

↳ buffers packets → DS needed for eventual in-order delivery to upper layer

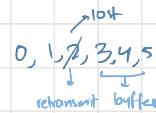
↳ sender times out / retransmits individually for unACKed packets

↳ sender maintains times for each unACKed PKT

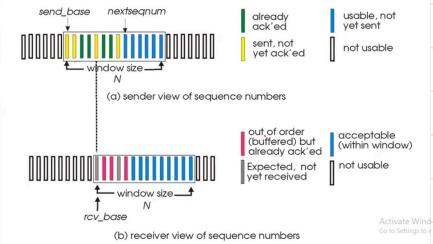
↳ sender window

↳ N consecutive seq. no.'s

↳ limits seq. no.'s of sent, unACKed packets



Selective repeat: sender, receiver windows



Selective repeat: sender and receiver

sender

data from above:

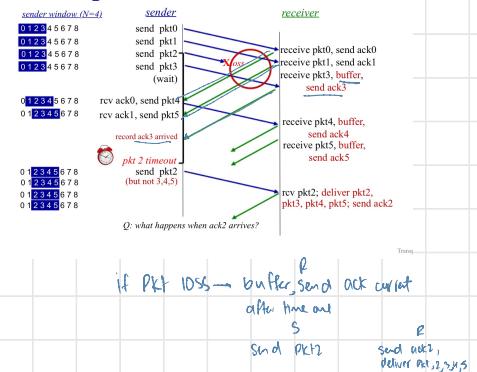
- if next available seq # in window, send packet timeout(t_r):
- resend packet n, restart timer ACK(n) in [sendbase, sendbase+N]: mark packet n as received
- if n smallest unACKed packet, advance window base to next unACKed seq #

receiver

- packet n in [rcvbase, revbase+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 [revbase-N, revbase-1]
 - ACK(n)
 - otherwise:
 - ignore

Transport Layer

Selective Repeat in action

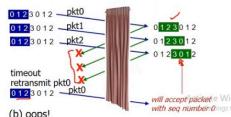
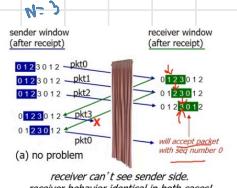


Selective repeat: dilemma

example:

- seq #'s: 0, 1, 2, 3
- window size=3
- receiver sees no difference in two scenarios!
- duplicate data accepted as new in (b)

Q: what relationship between seq # size and window size to avoid problem in (b)?



↳ RS b/w seq. no and window size

↳ seq. no and window size should have

a diff of > 1

SOLUTION

a) window = 3 → diff of 1
Seq. no : 4

b) window = 3 → diff > 1 → then the problem is avoided
Seq. no : 5

TCP: Overview

- ↳ point to point → host to host comm only
 - ↳ one sender, one receiver
 - ↳ reliable, in-order byte stream:
 - ↳ no "msg boundaries"
 - ↳ full duplex data
 - ↳ bidirectional data flow
 - ↳ MSS: max segment size

Sequence no.

- ↳ first byte in segment data

outgoing segment from sender

source port	dest port
sequence number	
acknowledgement number	
	rwnd
checksum	urg pointer

Acknowledgements

- ↳ seq. no. of next byte expected from other side

Cumulative ACK

outgoing segment from receiver

source port	dest port
sequence number	
acknowledgement number	
A	rwnd
checksum	urg pointer

Cumulative ACKs

Pipelining

- ↳ TCP congestion and flow control
 - set window size

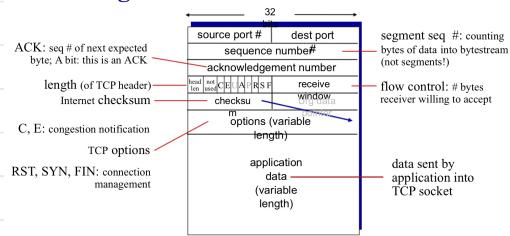
Connection-oriented

- ↳ handshaking initializes sender/receiver state before data exchange

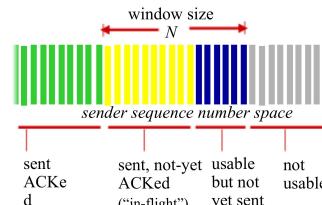
Flow controlled

- ↳ sender will not overwhelm receiver

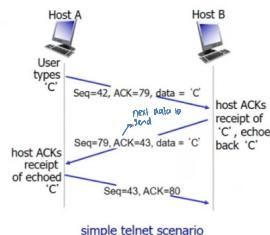
TCP segment structure



Transport Layer 3-78



TCP seq. numbers, ACKs



↓
use拥塞避免技术
GBN, SR

Q4) UDP header.

Source Port: 0019
 Destination Port: D36A
 User Datagram Length: 001C
 IP: 202.28.33.21

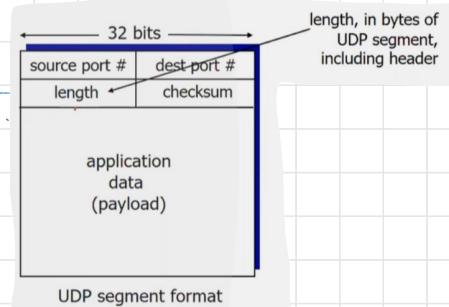
a) Source Port number: 25

b) socket address of the sender end: 202.28.33.21

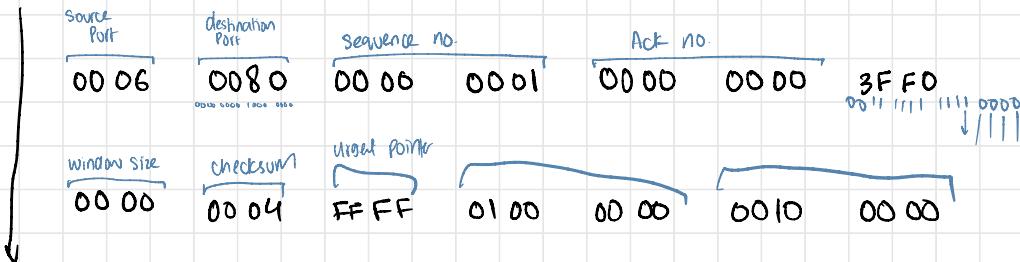
c) User datagram length: 28

d) length of data: $28 - 8 = 20$ bytes

e) packet direction: server to client $\rightarrow 25 \rightarrow$ SMTP



8) TCP header



1) Source Port number 6

$$x = 4 \quad y = 6 \quad z = 8$$

21k 4688

2) Destination Port 128

3) Sequence Number 1

4) Ack Number 0

5) TCP header length 3

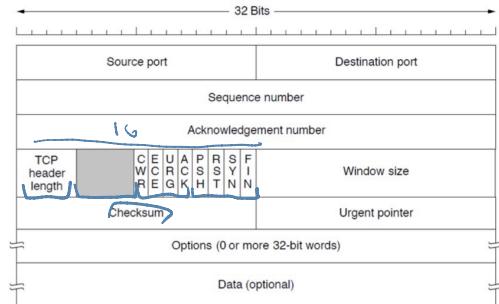
6) Ack bit 1

7) FIN bit 0

8) SYN bit 0

9) Window size 0

10) checksum 4



lec 19

TCP round trip time, timeout

Q) How to set TCP timeout value?

↳ Time out value shouldn't be

↳ too short → premature timeout → unnecessary retransmissions

↳ too long → slow reaction to segment loss

↳ longer than RTT

↳ but RTT varies

↳ Round trip time
Time from when
sender sent RTT
- receive received
- receiver send ACK

Q) How to estimate RTT?

SAMPLE RTT

↳ A variable used to measure time

from segment transmission until ACK receipt

↳ ignore retransmission

↳ It will vary → depends on network congestion

↳ So calculate estimatedRTT using
SAMPLERTT to make it smoother

Avg of multiple SampleRTT

$$\text{EstimatedRTT} = (1 - \alpha) * \text{EstimatedRTT} + \alpha * \text{sampleRTT}$$

$$\text{DevRTT} = (1 - \beta) * \text{DevRTT} + \beta * |\text{sampleRTT} - \text{EstimatedRTT}|$$

typical value 0.125

↳ deviation of sampleRTT from EstimatedRTT

↳ typically 0.25

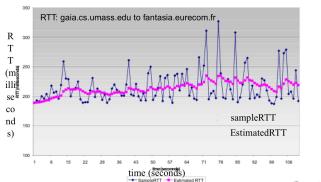
↳ safety margin

+ large variation in EstimatedRTT → want a larger safety margin

→ Put new values

TimeoutInterval: EstimatedRTT + 4 * DevRTT

- exponential weighted moving average (EWMA)
- influence of past sample decreases exponentially fast
- typical value: $\alpha = 0.125$



Q) ERTT = 30ms

SRTT = 15ms

$\alpha = 0.125$

ERTT: $(1 - 0.125) + 30 + 0.125(15)$

TCP reliable data transfer

- TCP creates rdt service on top of IP's unreliable service

- pipelined segments
- cumulative acks ✓
- single retransmission timer ✓

- retransmissions triggered by:
 - timeout events
 - duplicate acks

let's initially consider simplified TCP sender:

- ignore duplicate acks
- ignore flow control, congestion control

TCP Sender

- ↳ event: data received from application
 - ↳ create segment with seq no.
 - ↳ start timer if not already running
 - ↳ assume timer as oldest ACKed segment
 - ↳ retransmit initial Timer duration

↳ event · timeout

↳ retransmit segment that caused timeout

↳ restart timer

↳ event: ACK received

bif ACK acknowledges previously unAcked segment

↳ Update what is known to be ACKed

↳ start timer if there are still unAcked segments

TCP Receiver: ACK generation [RFC 5681]

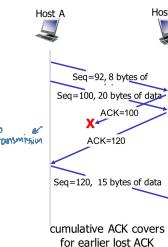
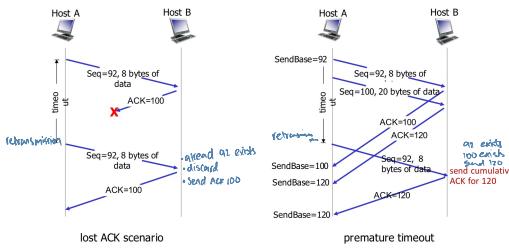
Event at receiver	TCP receiver action

↳ when retransmission

↳ Timeout is doubled

send ACK of next value

TCP: retransmission scenarios



TCP fast retransmit

↳ If sender receives 3 additional ACKs, for same data

↳ resend unAcked segment with smallest seq no.

It's most likely that unACKed segment got lost,
so don't wait for timeout



PROS

↳ saves time, as time out period relatively long

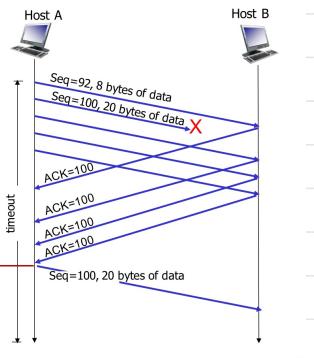
↳ Pkt not received

Sent 3 additional duplicate ACK

Let them retransmit before timeout



 Receipt of three duplicate ACKs indicates 3 segments received after a missing segment – lost segment is likely. So retransmit!



TCP Flow Control

- ↳ receiver controls sender, so sender won't overflow
receivers buffer by transmitting too much, too fast

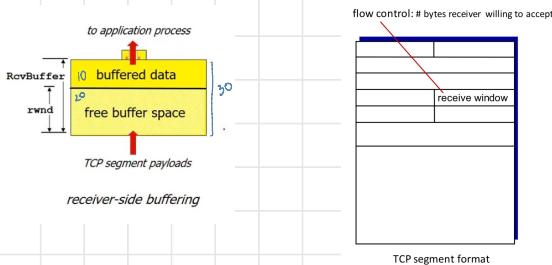
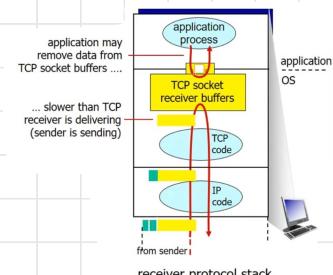
- ↳ TCP receiver advises free buffer space
in rwnd field in TCP header

↳ RCVBuffer size set via socket options

→ default
TCP buffer
→ auto adjusted
by main OS

↳ Sender limits amount of unACKed data to received rwnd → last byte sent - last byte ACK <= receiving window

↳ guarantees receive buffer will not overflow



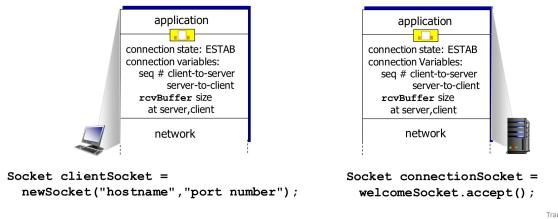
- ↳ what happens if network layer delivers data faster than application layer removes data from socket buffers

TCP Connection Management

↳ before exchanging data sender/receiver handshake

↳ agree to establish connection

↳ agree on connection parameters e.g. starting seq. no.



Agreeing to establish connection

↳ 2-way handshake

↳ 3-way handshake

2-way handshake

CONS

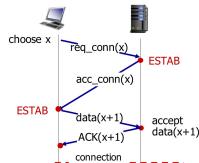
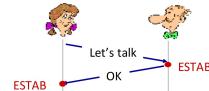
↳ variable delays

↳ retransmitted messages e.g. req_conn(x) due to msg loss

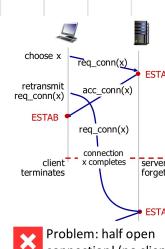
↳ message reordering

↳ can't see other side

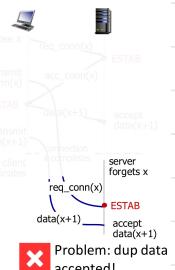
2-way handshake:



No problem!



Problem: half open connection! (no client)

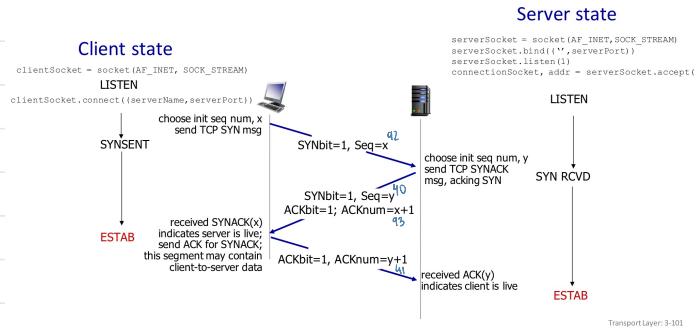


Problem: dup data accepted!

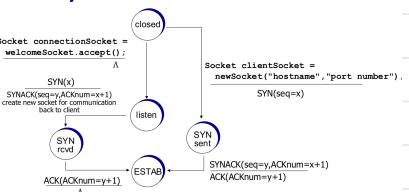
3-way handshake

↳ chooses ISN → client initial Seq no

↳ Sends TCP SYN msg → to open connection



TCP 3-way handshake FSM



CLOSING TCP CONNECTION

↳ client, server each close their side of connection

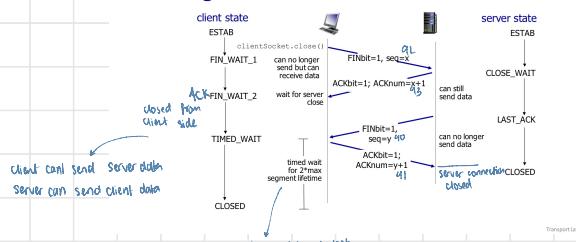
↳ send TCP segment with FIN bit = 1 → to close connection

↳ respond to received FIN with ACK

↳ on receiving FIN, ACK can be combined with own FIN

↳ simultaneous FIN exchanges can be handled

Closing a TCP connection



Question 4: Because of the connection-oriented nature of TCP, a connection setup phase is required at the beginning of each session, as well as a connection tear-down phase at the end of the session. Enumerate the events below in the order they occur as host A opens a TCP connection to host B, transmits data and then closes the connection. Write a 1 next to the event that occurs first and continue like that until all occurring events are enumerated (the first event has been enumerated for you). You may assume that no segments are lost. Also indicate at which host the event happens. Please note that there might be events listed below that are not a part of the above data transfer and hence should not be enumerated. [10 points]

Event	Host	Order
Send an ACK segment	A	8
Do the rest of the data exchange	A, B	4
Close the connection	A	11
Send an ACK segment	B	6
Send a FIN segment	A	5
Send a SYN segment	A	1
Send a FIN segment	B	7
Send a RST segment		
Send a SYN-ACK segment	B	2
Enter the TIME-WAIT state	A	9
Send an ACK+DATA segment	A	3
Close the connection	B	10

Question 5: The Transmission Control Protocol uses a method called connection control to regulate the traffic.



PRINCIPLES OF CONGESTION CONTROL

Congestion

- ↳ too many senders, sending too much data
- ↳ too fast for network to handle

flow control

- ↳ 1 sender, too fast for 1 receiver

Manifestations

- ↳ long delays → queuing in router buffers
- ↳ packet loss → buffer overflow at routers

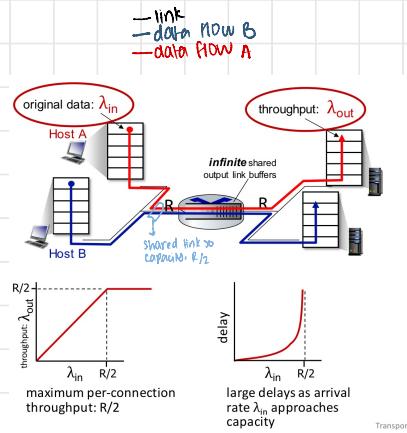
CAUSES / COSTS OF CONGESTION

Scenario 1: Simple scenario

- ↳ 2 senders
- ↳ 2 receivers
- ↳ 1 router, infinite buffers, no overflow
- ↳ input, output link capacity: R
- ↳ 2 flows
- ↳ no retransmissions needed → all data is being sent

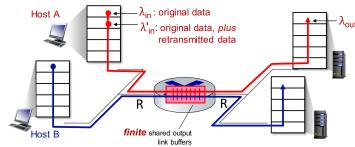
Q) what happens as arrival rate

λ_{in} approaches $R/2$?



Scenario 2

- ↳ 1 router, finite buffers → can overflow
- ↳ sender retransmits lost/timed-out packet → PRR may be lost
- ↳ application layer input = application layer output
 $\lambda_{in} = \lambda_{out}$



- ↳ transport layer input includes retransmissions

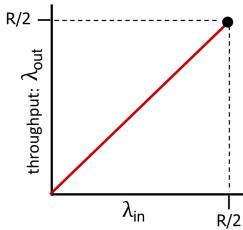
$$\lambda'_{in} \geq \lambda_{in}$$

↑
bg data +
retransmited
data
↓
bg data

↳ 1 hop

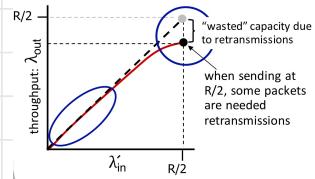
↳ Idealization: Perfect knowledge

- ↳ sender only sends when router buffers available



↳ Idealization: Perfect knowledge some dropped at router

- ↳ packets can be lost due to full buffers
- ↳ sender knows when packet has been dropped
- ↳ only resends if packet known to be lost → timeout



↳ realistic scenario: Un-needed duplicates

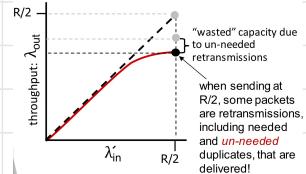
- ↳ packets can be lost/dropped at router due to full buffers

↳ this req. retransmission

- ↳ but senders time can time out prematurely

↙ sending 2 copies, both of which are delivered

- when timeout
 - * send another copy
 - * for gets 2 copies



↳ costs of congestion

- ↳ more retransmission for given receiver throughput → more retransmissions

↳ un-needed retransmission → link carries multiple copies of a packet

↳ decreasing throughput

Scenario 3

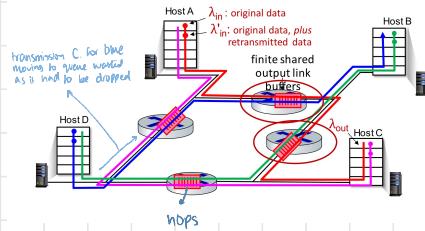
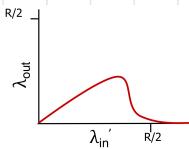
- ↳ 4 senders
- ↳ multi hop paths
- ↳ timeout/retransmit

Q) what happens if λ_{in} and λ_{out} increase?

- ↳ as red λ_{in} increases, λ_{in} increases → fed capacity ↑
↳ red takes more space in $R/2$
- ↳ all arriving blue PKTs at upper avenue are dropped → blue ↓
- ↳ blue throughput → 0

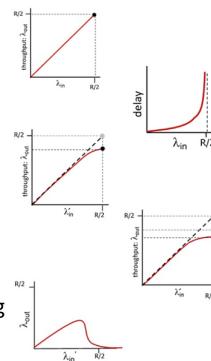
↳ another cost of congestion

- ↳ when Pkt 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



Transport Layer: 3-11

Throughput: rate at which bits are being sent from sender to receiver

Question # 3: [CLO-3]**[7.5 + 7.5 + 5 = 20 Points]**

Part a) Consider there is only one router, with finite buffer capacity. The input and output link capacity of router is R. Two flows are attached to the router, each flow sharing half of the router capacity and each flow have:

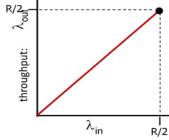
- i) The perfect knowledge of the router buffer capacity.
- ii) Have partial knowledge, and knows only if the packet is dropped at the router (retransmission needed for lost packet).
- iii) Have no knowledge about router buffer, nor about any dropped packet (guess only by acknowledgement packets).

Plot the graphs for each of the above case for only one flow, showing transport-layer input on x-axis and maximum achievable receiver throughput at y-axis.

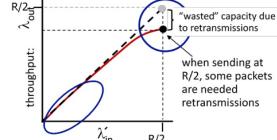
Question # 3: [CLO-2]**[7.5 + 7.5 + 5 = 20 Points]**

Part a) Answer:

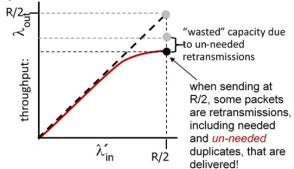
i)



ii)



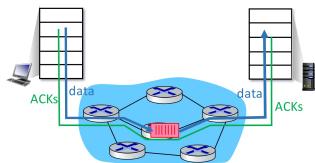
iii)



APPROACHES TOWARDS CONGESTION CONTROL

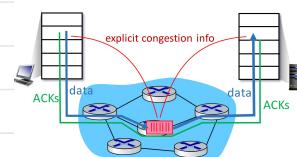
End-end congestion control

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



Network Assisted congestion control

- ↳ routers provide direct feedback to end systems
sending/receiving hosts
- ↳ with flow passing through the congested router
- ↳ may indicate congestion by
 1. single bit indicating congestion
 - ↳ TCP ECN, ATM, DECBIT protocols
 2. explicitly set sending rate for sender



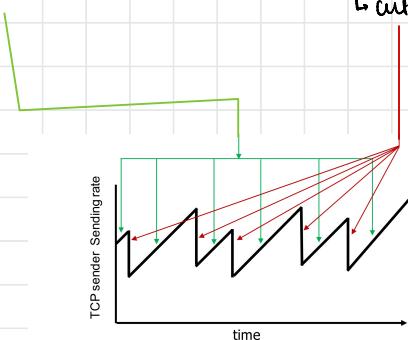
TCP CONGESTION CONTROL : AIMD

Approach.

- ↳ senders can increase sending rate
congestion
until packet loss occurs,
- ↳ then decrease sending rate on loss event

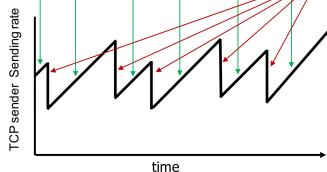
Active Increase

- ↳ increase sending rate by 1 MSS every RTT, until loss detected



Multiplicative Decrease

- ↳ sending rate is
round → congestion window
- ↳ cut in half on loss detected by triple duplicate ACK
- ↳ cut to 1 MSS when loss detected by timeout
more important
TCP Read
TCP Timer



Additive Increase
Multiplicative Decrease

AIMD

↳ a distributed, asynchronous algorithm

↳ that

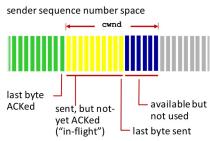
↳ optimises congested flow rates network wide

↳ has desirable stability properties

↳ by probing for bandwidth

TCP congestion control: Details

TCP congestion control: details



TCP sending behavior:

- roughly: send $cwnd$ bytes, wait RTT for ACKs, then send more bytes

$$\text{TCP rate} \approx \frac{cwnd}{RTT} \text{ bytes/sec}$$

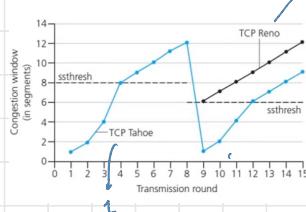
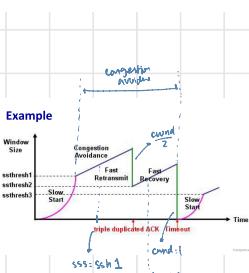
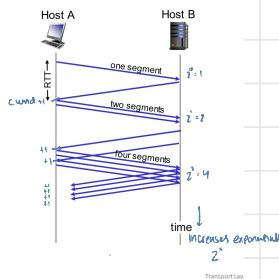
- TCP sender limits transmission: $\text{LastByteSent} - \text{LastByteAcked} \leq cwnd$
- $cwnd$ is dynamically adjusted in response to observed network congestion (implementing TCP congestion control)

Transport Layer: S-121

TCP Slow Start

TCP slow start

- when connection begins, increase rate exponentially until first loss event:
 - initially $cwnd = 1$ MSS
 - double $cwnd$ every RTT
 - done by incrementing $cwnd$ for every ACK received
- **summary:** initial rate is slow, but ramps up exponentially fast



$$ssthresh = \frac{1}{2} cwnd \text{ before loss event}$$

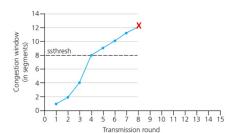
TCP: from slow start to congestion avoidance

Q: when should the exponential increase switch to linear?

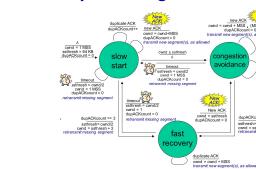
A: when $cwnd$ gets to 1/2 of its value before timeout.

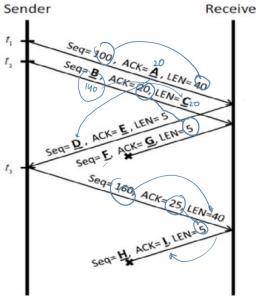
Implementation:

- variable $ssthresh$
- on loss event, $ssthresh$ is set to 1/2 of $cwnd$ just before loss event



Summary: TCP congestion control



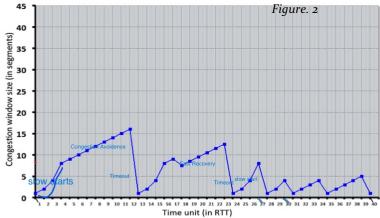


A	20
B	100+40 = 140
C	20
D	20
E	100+40 = 140
F	100+5 = 105
G	100+10 = 160
H	30
I	100+40 = 140

Solution:

$$B = 140; D = 20; E = 140; F = 25; H = 30; I = 200$$

Question 5: Consider the figure. 2, which plots the evolution of TCP's congestion window at the beginning of each time unit (where the unit of time is equal to the RTT). In the abstract model for this problem, TCP sends a "flight" of packets of size $cwnd$ at the beginning of each time unit. The result of sending this flight of packets is that either (i) all packets are ACKed at the end of the time unit, (ii) there is a timeout for the first packet, or (iii) there is a triple duplicate ACK for the first packet. In this problem, you are asked to reconstruct the sequence of events (ACKs, losses) that results in the evolution of TCP's $cwnd$ shown below.



(a) Give the time unit(s) in (RTT) at which TCP is in slow start.

Solution:

The times where TCP is in slow start are: 1,2,3,13,14,15,24,25,26,28,29,31,35,40.

(b) Give the time unit(s) in (RTT) at which TCP is in congestion avoidance.

Solution:

The times where TCP is in congestion avoidance are:

4,5,6,7,8,9,10,11,12,16,17,19,20,21,22,23,27,30,32,33,34,36,37,38,39

(c) Give the time unit(s) in (RTT) at which TCP is in fast recovery.

Solution:

The times where TCP is in fast recovery are: 18.

(d) Give the time unit(s) in (RTT) at which packets are lost via timeout.

Solution:

The times where TCP has a loss by timeout are: 12,23,27,30,34,39.

(e) Give the time unit(s) in (RTT) at which packets are lost via triple ACK.

Solution:

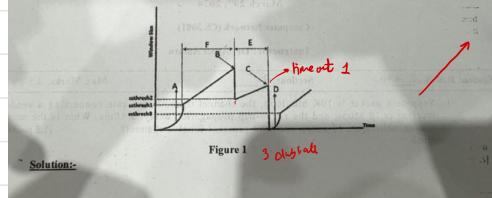
The times where TCP has a loss by triple duplicate ACK are: 17.

(f) Give the time unit(s) in (RTT) at which the value of $ssthresh$ changes (hint: if it changes between t=3 and t=4, use t=4 in your answer)

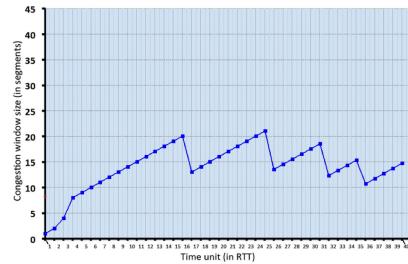
Solution:

The times where the $ssthresh$ changes are: 18,24,28,31,40.

2. Examine the TCP throughput graph depicted in Figure 1, with the sender's TCP window size represented on the y-axis. Identify the labeled points from A to I and address the following questions: [0.5*3=1.5 points]
- Identify the region where TCP slow-start is operating. **A,D**
 - Identify the region where TCP congestion-avoidance is operating. **F**
 - Identify the operation at point B and point C.



Solution:-



Answer the following questions.

The initial value of $cwnd$ is 1 and the initial value of $ssthresh$ (Little + sign) is 8.

1. Give the times at which TCP is in slow start. **1,2,3**

2. Give the times at which TCP is in congestion avoidance.

4,5,6,7,8,9,10,11,12,13,14,15,16,18,19,20,21,22,23,24,25,27,28,29,30,31,33,34,35,37,38,39,40

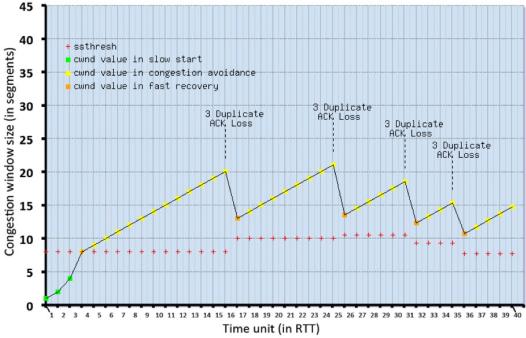
3. Give the times at which TCP is in fast recovery. **17,26,32,36**

4. Give the times at which packets are lost via triple ACK. **16,25,31,35**

5. Give the times at which the value of $ssthresh$ changes (if it changes between t=3 and t=4, use t=4 in your answer) **17,26,32,36**

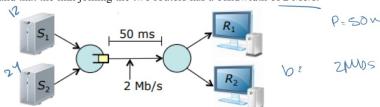
The complete solution is shown in the figure below:

- For intervals of time when TCP is in slow start, the plotted value of cwnd is shown as a green square
- For intervals of time when TCP is in congestion avoidance, the plotted value of cwnd is shown as a yellow square
- For intervals of time when TCP is in fast recovery, the plotted value of cwnd is shown as an orange square
- The values for ssthresh are shown following a change as a red plus sign
- A flight of packets experiencing a loss has the loss type (which determines the next value of cwnd) labeled above



Question 3 [10 Points]

The following diagram shows two TCP senders at left and the corresponding receivers at right. Both senders use TCP Reno. Assume that the MSS is 1 KB, that the one-way propagation delay for both connections is 50 ms and that the link joining the two routers has a bandwidth of 2 Mb/s.

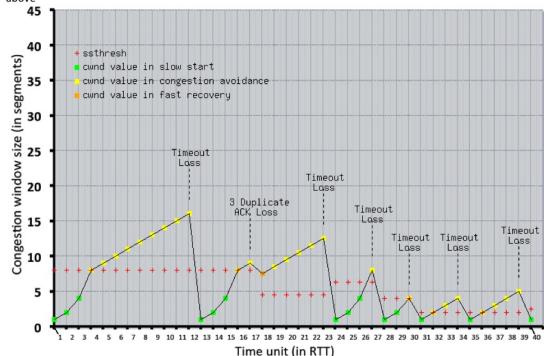


(a) Assume that the link buffer overflows whenever $cwnd1 + cwnd2 \geq 36$ KB and that at time 0, $cwnd1=12$ KB and $cwnd2=24$ KB. Approximately, what are the values of $cwnd1$ and $cwnd2$ one RTT later? Assume that all packet losses are detected by a triple duplicate ack.

Solution	[6 points]
Cwnd = 36 KB when buffer overflows and congestion is detected. Therefore, ssthreshold would be updated and taken half.	
New_ssthreshold1 = $12/2$ KB = 6 KB New_cwnd1 = New_ssthreshold1 + 3MSS New_cwnd1 = 9 KB	
New_ssthreshold2 = $24/2$ KB = 12 KB New_cwnd2 = New_ssthreshold2 + 3MSS New_cwnd2 = 15 KB	

The complete solution is shown in the figure below:

- For intervals of time when TCP is in slow start, the plotted value of cwnd is shown as a green square
- For intervals of time when TCP is in congestion avoidance, the plotted value of cwnd is shown as a yellow square
- For intervals of time when TCP is in fast recovery, the plotted value of cwnd is shown as an orange square
- The values for ssthresh are shown following a change as a red plus sign
- A flight of packets experiencing a loss has the loss type (which determines the next value of cwnd) labeled above



a) **Slow Start:** identify the intervals of time when TCP slow start is operating.

1-6, 7-11, 23-26

b) **Congestion Avoidance:** identify the intervals of time when TCP congestion avoidance is operating.

11-15, 16-22, 26-32
3 duplicate

c) **Fast Retransmission:** identify the intervals of time when TCP fast retransmission is used.

15-16

d) **Fast Recovery:** identify the intervals of time when TCP fast recovery is operating.
15-16

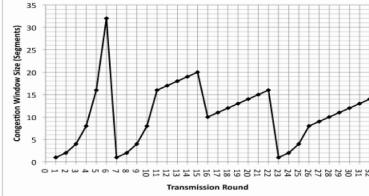
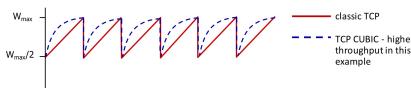


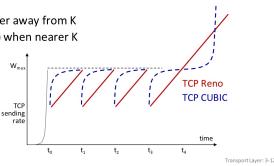
Figure 2

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 to to W_{max} **faster**, but then approach W_{max} more **slowly**



- 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



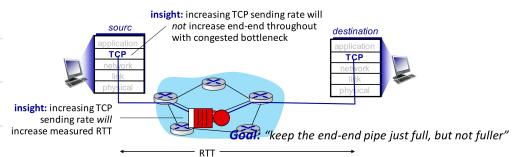
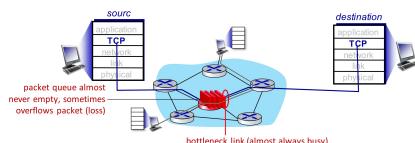
TCP and the congested "Bottleneck link"

↳ TCP increase TCP's sending rate

until Pkt loss occurs at some routers output → the bottleneck link

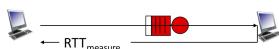
↳ understanding congestion

↳ 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



$$\text{measured throughput} = \frac{\# \text{bytes sent in last RTT interval}}{\text{RTT measure}}$$

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 */

- congestion control without inducing/forcing loss
- maximizing throughout ("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

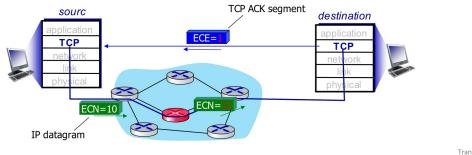
Throughput

↳ rate at which bits are being sent from sender to receiver

Explicit Congestion Notification (ECN)

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

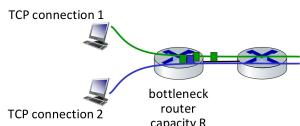
- two bits in IP header (ToS field) marked by *network router* to indicate congestion
 - *policy* to determine marking chosen by network operator
- congestion indication carried to destination
- destination sets ECN bit on ACK segment to notify sender of congestion
- involves both IP (IP header ECN bit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



Transport Layer: 3-132

TCP Fairness

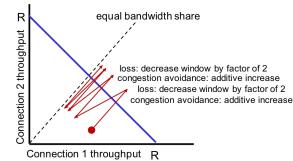
Fairness goal: if K TCP sessions share same bottleneck link of bandwidth R , each should have average rate of R/K



Q: is TCP Fair?

Example: two competing TCP sessions:

- additive increase gives slope of 1, as throughput increases
- multiplicative decrease decreases throughput proportionally



Transport Layer: 3-134

Is TCP fair?
A: Yes, under idealized assumptions:
 • same RTT
 • fixed number of sessions
 only in congestion avoidance

Fairness: must all network apps be “fair”?

Fairness and UDP

- multimedia apps often do not use TCP
 - do not want rate throttled by congestion control
- instead use UDP:
 - send audio/video at constant rate, tolerate packet loss
 - there is no “Internet police” policing use of congestion control

Fairness, parallel TCP connections

- application can open *multiple* parallel connections between two hosts
- web browsers do this, e.g., link of rate R with 9 existing connections:
 - new app asks for 1 TCP, gets rate $R/10$
 - new app asks for 11 TCPs, gets $R/2$

Transport Layer: 3-135

Evolving Transport Layer Functionality

- TCP, UDP: principal transport protocols for 40 years
- different “flavors” of TCP developed, for specific scenarios:

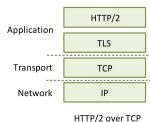
Scenario	Challenges
Long, fat pipes (large data transfers)	Many packets “in flight”; loss shuts down pipeline
Wireless networks	Loss due to noisy wireless links, mobility; TCP treats this as congestion loss
Long-delay links	Extremely long RTTs
Data center networks	Latency sensitive
Background traffic flows	Low priority, “background” TCP flows

- moving transport-layer functions to application layer, on top of UDP
- HTTP/3: QUIC

Transport Layer 3

QUIC: Quick UDP Internet Connections

- application-layer protocol, on top of UDP
- increase performance of HTTP
- deployed on many Google servers, apps (Chrome, mobile YouTube app)

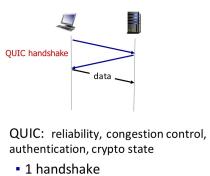
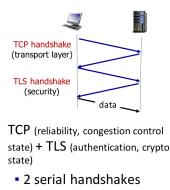


adopts approaches we've studied in this chapter for connection establishment, error control, congestion control

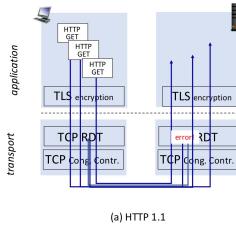
- **error and congestion control:** “Readers familiar with TCP’s loss detection and congestion control will find algorithms here that parallel well-known TCP ones.” (from QUIC specification)
- **connection establishment:** reliability, congestion control, authentication, encryption, state established in one RTT
- multiple application-level “streams” multiplexed over single QUIC connection
 - separate reliable data transfer, security
 - common congestion control

Transport Layer

QUIC: Connection establishment



QUIC: streams: parallelism, no HOL blocking



Transport

Chapter 3: summary

- principles behind transport layer services:
 - multiplexing, demultiplexing
 - reliable data transfer
 - flow control
 - congestion control
- instantiation, implementation in the Internet
 - UDP
 - TCP

Up next:

- leaving the network “edge” (application, transport layers)
- into the network “core”
- two network-layer chapters:
 - data plane
 - control plane

Transport Layer 3-142

ishmaq hafeez
notes
represent

NETWORK LAYER : DATA PLANE

CMP 4

u) Network Layer Services and Protocol

- ↳ transport segment from sending to receiving host

↳ sender

- ↳ encapsulates segments into datagrams

- ↳ passes to link layer

↳ receiver

- ↳ delivers segments to transport layer protocol

↳ router

- ↳ examines header files in all IP datagrams

- ↳ moves datagrams from input to output ports to transfer

TWO KEY NETWORK layer functions

1. FORWARDING

- ↳ move packets from routers input link

- ↳ to appropriate router output link

2. ROUTING

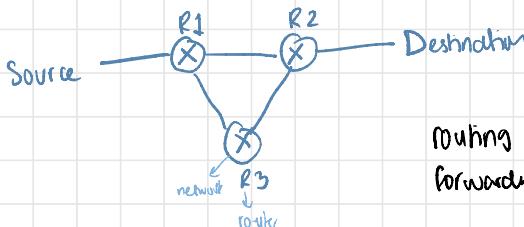
- ↳ determine route taken by packets
from source to destination

- ↳ routing algorithms

- ↳ has forwarding table

analogy: taking a trip

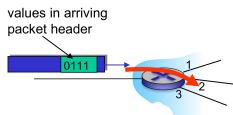
- *forwarding*: process of getting through single interchange
- *routing*: process of planning trip from source to destination



routing chooses which Path
forwarding forwards the packet

DATA PLANE

- ↳ local, per-route function
- ↳ determines how data arrives on router input port and is forwarded to router output port



CONTROL PLANE

- ↳ network wide logic
- ↳ determines how datagram is routed among routers along end-end paths from source to destination host

2 control plane approaches

1. Traditional Routing algos

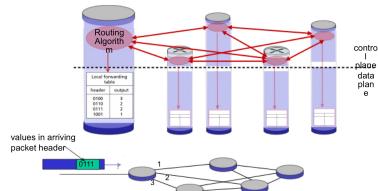
- ↳ implemented in routers

2. Software defined networking (SDN)

- ↳ implemented in remote servers

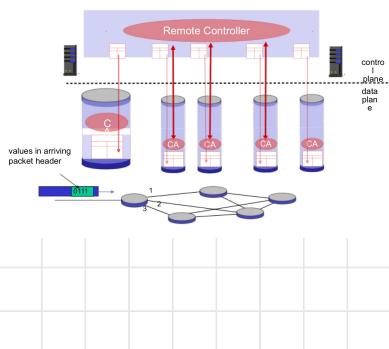
Per Router Control Plane

- ↳ individual routing algo components in each and every router interact in the control plane



SDN Control Plane

- ↳ remote controller computes, installs forwarding tables in routers



Question

Question 3: SDN provides a programmable control plane, however, it network infrastructure needs changes as we studied in chapter # 5 of our textbook. Now, answer the following questions. [10 points]

- a) How the data plane is different while using SDN? Explain.

Two fold: i) data plane is based on SDN controlled switches and ii) uses flow table to forward traffic.

Figure. 1

- a) Show in steps, how data packets entering one interface of a router are forwarded to an output interface of the same router.

Solution:

Step 1: IP address is extracted from the IP packet header.

Step 2: IP prefix matched with forwarding table and corresponding output interface is read.

Step 3: IP packets is sent to the output interface queue for further processing.

- b) How routing helps forwarding? Explain.

Solution:

The data packets using forwarding tables to traverse down the shortest paths which are calculated by routing algorithms.

4.1.2 Network Service Model

for individual datagrams

- ↳ guaranteed delivery
- ↳ with <40 msec delay

for a flow of datagrams

- ↳ in order datagram delivery
- ↳ guaranteed min bandwidth to flow
- ↳ restrictions on changes in inter-packet spacing

best effort service model

- ↳ NO guarantees on
 - ↳ successful datagram delivery to destination
 - ↳ timing/order of delivery
 - ↳ bandwidth avail to end-end flow

e.g. Internet

Network-layer service model

Network Architecture	Service Model	Quality of Service (QoS) Guarantees ?			
		Bandwidth	Loss	Order	Timing
Internet	best effort	none	no	no	no
ATM	Constant Bit Rate	Constant rate	yes	yes	yes
ATM	Available Bit Rate	Guaranteed min	no	yes	no
Internet	Intserv Guaranteed (RFC 1633)	yes	yes	yes	yes
Internet	Diffserv (RFC 2475)	possible	possibly	possibly	no

PROS best effort service

1. simplicity of mechanism
 - ↳ allows internet to be widely deployed adopted
2. sufficient provisioning of bandwidth
 - ↳ allows performance of real time applications
 - e.g. interactive voice, video
3. replicated application layer distributed services
 - ↳ connects close to clients networks
 - allowing services to be provided from multiple locations
4. congestion control of "elastic" services helps

4.2) Destination Based Forwarding

forwarding table	
Destination Address Range	Link Interface
11001000 00010111 00010000 00000000 through	n
11001000 00010111 00010000 00000100 through	3
11001000 00010111 00010000 00000111	
11001000 00010111 00011000 11111111	
11001000 00010111 00011001 00000000 through	2
11001000 00010111 00011111 11111111	
otherwise	3

Q: but what happens if ranges don't divide up so nicely?

Longest Prefix Matching

longest prefix match

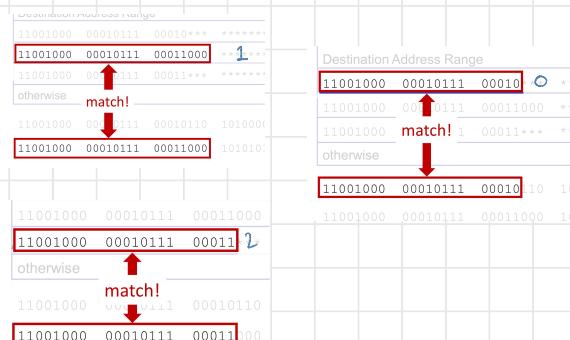
when looking for forwarding table entry for given destination address, use *longest* address prefix that matches destination address.

Destination Address Range	Link interface
11001000 00010111 00010*** *****	0
11001000 00010111 00011000 *****	1
11001000 00010111 00011*** *****	2
otherwise	3

examples:

11001000 00010111 00010110 10100001 which interface?

11001000 00010111 00011000 10101010 which interface? 1



4.2.5 Packet Scheduling

↳ decides which packet to send next on link

↳ FCFS

↳ Priority

↳ Round Robin

↳ Weighted fair queuing

1. FCFS

↳ first come first serve

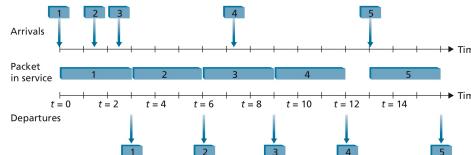
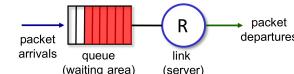


Figure 4.12 • The FIFO queue in operation

Abstraction: queue



→ every one allocated
Same time

2. Priority Scheduling

↳ arriving traffic classified,
queued by class

↳ any header fields can be used for classification

↳ sends packet from highest priority queue
that has buffered packets

↳ FCFS with Priority class

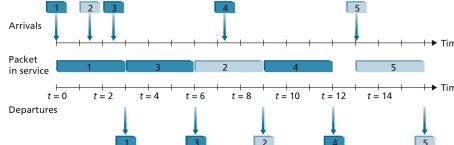
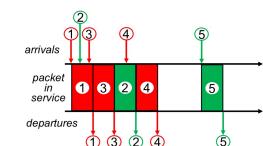
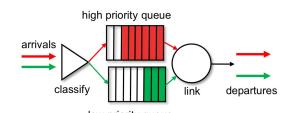


Figure 4.14 • The priority queue in operation



→ every one allocated
Same time

→ 3 had more priority than 2
and 3 and 2 were both arrived

3. Round Robin Scheduling (RR)

↳ arriving traffic classified.

queued by class

↳ any header fields can be used for classification

↳ server cyclically scans class queues, sending 1 complete packet from each class in turn

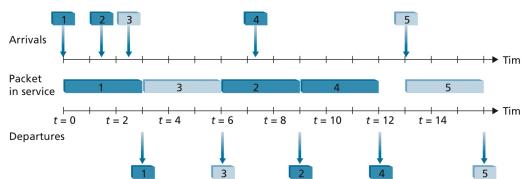
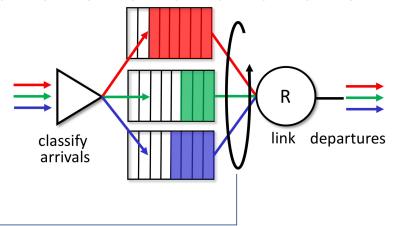


Figure 4.15 • The two-class robin queue in operation

CLASS A: 1, 2, 4

CLASS B: 3, 5

Picks alternate classes

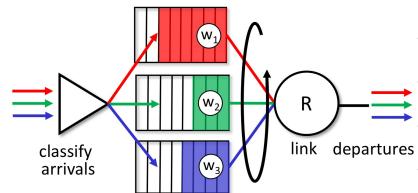
WEIGHTED FAIR QUEUING (WFQ)

↳ generalized Round Robin

↳ each class i , has weight w_i

and gets weight amount of service in each cycle

↳ min bandwidth guarantee → per traffic class



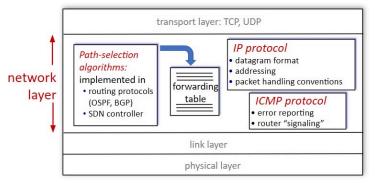
$$\frac{w_i}{\sum_j w_j}$$

u.3

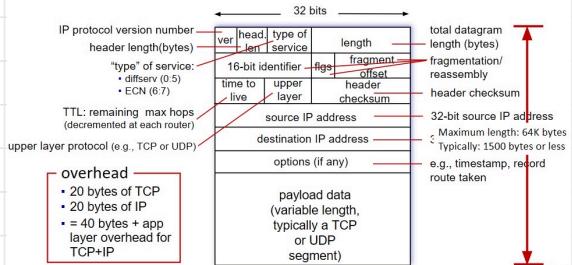
Network Layer: Internet

↳ host functions

↳ router network layer functions



IP Datagram format



↳ needs 3 mins to transfer data S → D

1. IP Protocol

↳ what format of datagram

↳ how to handle Pkt

2. Path selection Algo / Routing Protocol

↳ get shortest Path

3. ICMP Protocol

↳ tells errors

∴ UDP header = 8 bytes

∴ TCP header = 20 bytes

∴ Network Layer header OR IP = 20 bytes

overhead: TCP + IP + application layer data
20 + 20

4.1

IP ADDRESSING

IP address - IPv4

↳ 32 bit identifier

↳ associated with each host/router interface

interface

↳ connection b/w host/router and physical link



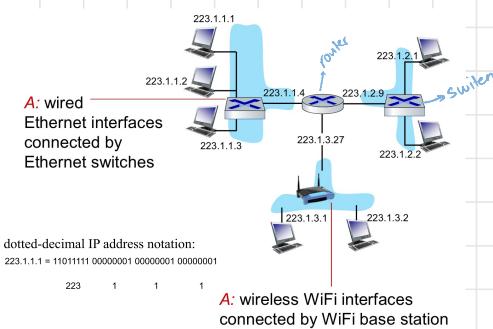
↳ routers have multiple interfaces

↳ host have 1 or 2 interfaces

↳ wired Ethernet

↳ wireless

↳ every interface has an IP address



Subnet

↳ device interfaces

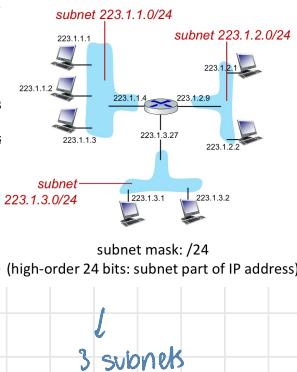
that physically reach each other

w/o passing through an intervening router

Subnets

Recipe for defining subnets:

- detach each interface from its host or router, creating “islands” of isolated networks
- each isolated network is called a **subnet**



CIDR

→ Classless Inter Domain Routing

↳ Subnet portion of address of arbitrary length

↳ same address format

a . b . c . d /x → no of bits in subnet portion of address



IP address structure

↳ subnet part - high order e.g. 223.1.1.

↳ devices in same subnet

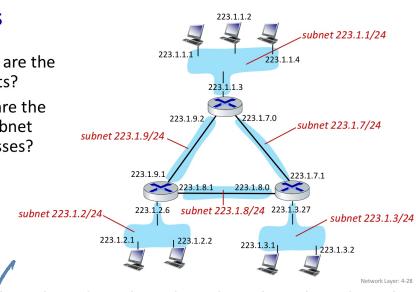
have common high order bits

↳ host part - low order e.g. 223.1.1.1

↳ remaining lower bits

Subnets

- where are the subnets?
- what are the /24 subnet addresses?



6 subnets

IPv6: motivation

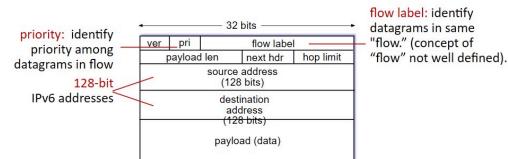
Initial motivation

- ↳ 32 bit IPv4 address space would be completely allocated

Additional motivation

- ↳ speed processing: 40 byte fixed length header
- ↳ enable diff network layer treatment of flows

IPv6 datagram format



IPv6 doesn't have compared to IPv4

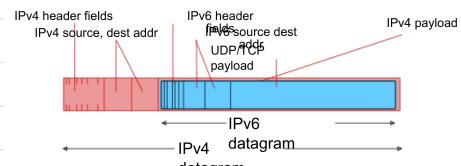
- ↳ no checksum → to speed processing at routers
- ↳ no fragmentation/reassembly
- ↳ no options → out as upperlayer, next header protocol at routers

Transition from IPv4 to IPv6

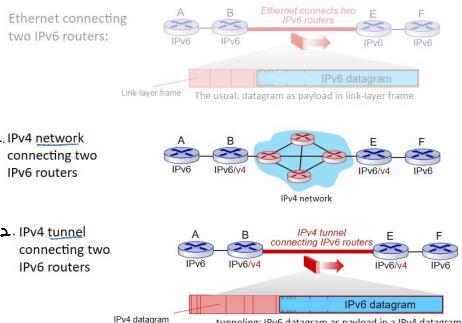
- ↳ not all routes can be upgraded simultaneously
- ↳ no 'flag' days

Tunneling

- ↳ allows network to operate with mixed IPv4 and IPv6 routes
- ↳ IPv6 datagram carried as payload in IPv4 datagram among IPv6 routes *Packet within a packet*



Tunneling and encapsulation



Tunneling

