

CS-470

Chapter 1: introduction

our goal:

- get “feel” and terminology
- more depth, detail later in course
- approach:
 - use Internet as example

overview:

- what’s the Internet?
- what’s a protocol?
- network edge; hosts, access net, physical media
- network core: packet/circuit switching, Internet structure
- performance: loss, delay, throughput
- security
- protocol layers, service models
- history

Introduction 1-1

1

Chapter 1: roadmap

1.1 what is the Internet?

1.2 network edge

- end systems, access networks, links

1.3 network core

- packet switching, circuit switching, network structure

1.4 delay, loss, throughput in networks

1.5 protocol layers, service models

1.6 networks under attack: security

1.7 history

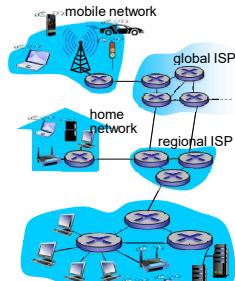
Introduction 1-2

2

What’s the Internet: “nuts and bolts” view



- billions of connected computing devices:
 - **hosts** = end systems
 - running **network apps**
- **communication links**
 - fiber, copper, radio, satellite
 - transmission rate: **bandwidth**
- **packet switches**: forward packets (chunks of data)
 - **routers** and **switches**



Introduction 1-3

3

“Fun” Internet-connected devices



IP picture frame
<http://www.ceiva.com/>



Web-enabled toaster +
weather forecaster



Tweet-a-watt:
monitor energy use



Slingbox: watch,
control cable TV remotely



sensorized,
bed mattress



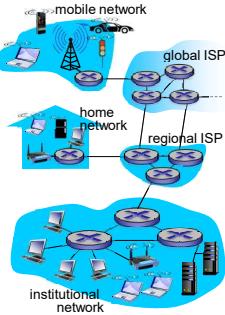
Internet phones

Introduction 1-4

4

What’s the Internet: “nuts and bolts” view

- **Internet: “network of networks”**
 - Interconnected ISPs
- **protocols** control sending, receiving of messages
 - e.g., TCP, IP, HTTP, Skype, 802.11
- **Internet standards**
 - RFC: Request for comments
 - IETF: Internet Engineering Task Force



Introduction 1-5

5

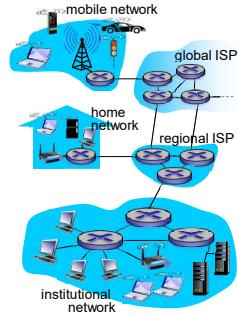
What’s the Internet: a service view

- **infrastructure that provides services to applications:**

- Web, VoIP, email, games, e-commerce, social nets, ...

- **provides programming interface to apps**

- hooks that allow sending and receiving app programs to “connect” to Internet
- provides service options, analogous to postal service



Introduction 1-6

6

What's a protocol?

human protocols:

- "what's the time?"
- "I have a question"
- introductions

- ... specific messages sent
- ... specific actions taken when messages received, or other events

network protocols:

- machines rather than humans
- all communication activity in Internet governed by protocols

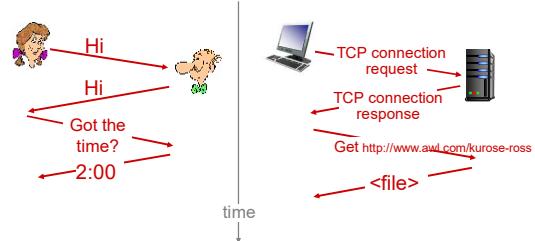
protocols define **format, order of messages sent and received among network entities, and actions taken on message transmission, receipt**

Introduction 1-7

7

What's a protocol?

a human protocol and a computer network protocol:



Q: other human protocols?

Introduction 1-8

8

Chapter 1: roadmap

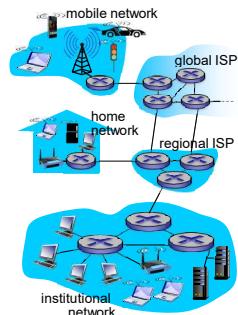
- 1.1 what is the Internet?
- 1.2 network edge
 - end systems, access networks, links
- 1.3 network core
 - packet switching, circuit switching, network structure
- 1.4 delay, loss, throughput in networks
- 1.5 protocol layers, service models
- 1.6 networks under attack: security
- 1.7 history

Introduction 1-9

9

A closer look at network structure:

- **network edge:**
 - hosts: clients and servers
 - servers often in data centers
- **access networks, physical media:** wired, wireless communication links
- **network core:**
 - interconnected routers
 - network of networks



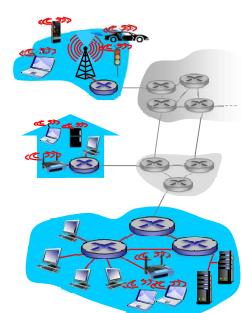
Introduction 1-10

10

Access networks and physical media

Q: How to connect end systems to edge router?

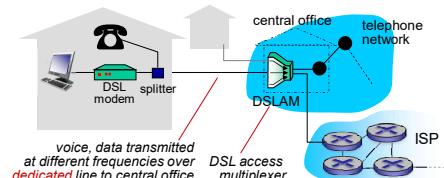
- residential access nets
- institutional access networks (school, company)
- mobile access networks



Introduction 1-11

11

Access network: digital subscriber line (DSL)

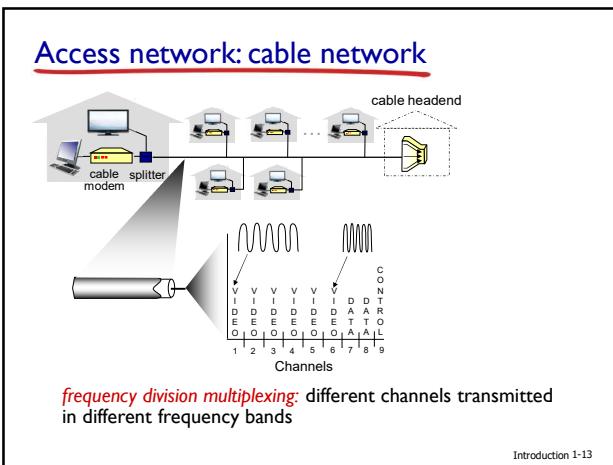


- use **existing** telephone line to central office DSLAM
 - data over DSL phone line goes to Internet
 - voice over DSL phone line goes to telephone net
- < 2.5 Mbps upstream transmission rate (typically < 1 Mbps)
- < 24 Mbps downstream transmission rate (typically < 10 Mbps)

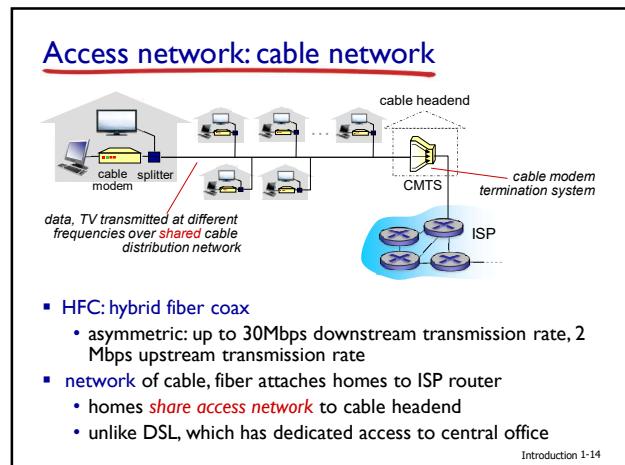
<http://www.ebay.com/gds/What-is-the-Difference-Between-a-DSL-and-a-Cable-Modem-/1000000177629327/g.html>

Introduction 1-12

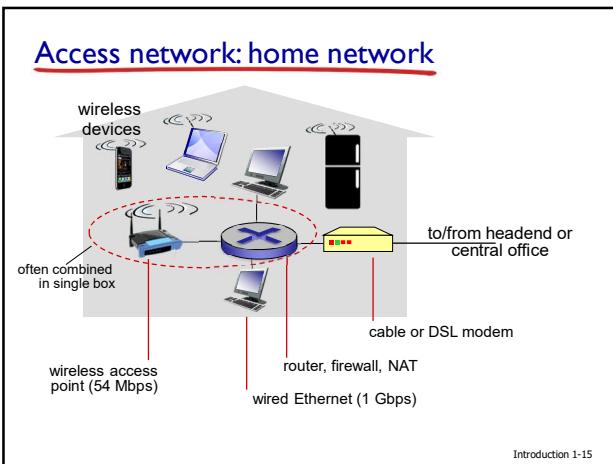
12



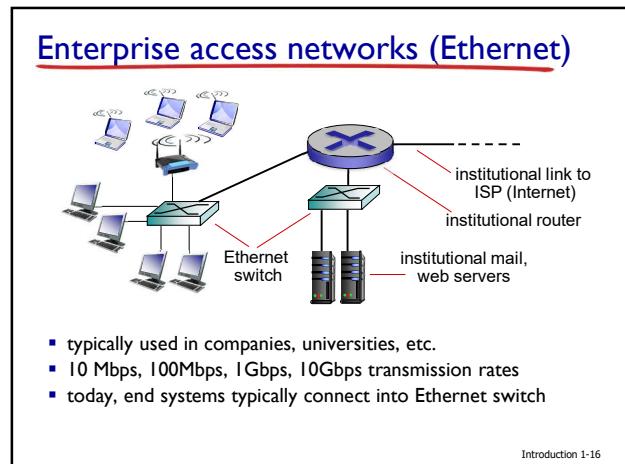
13



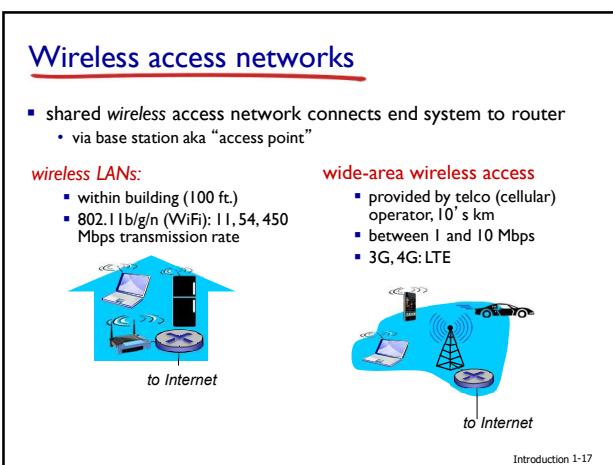
14



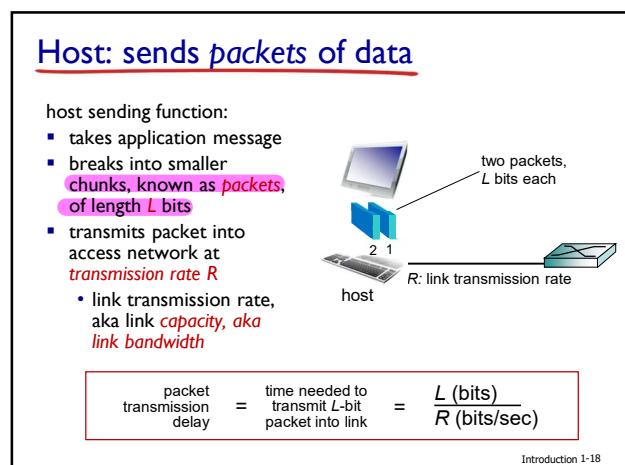
15



16



17



18

Physical media

- **bit:** propagates between transmitter/receiver pairs
- **physical link:** what lies between transmitter & receiver
- **guided media:**
 - signals propagate in solid media: copper, fiber, coax
- **unguided media:**
 - signals propagate freely, e.g., radio

twisted pair (TP)

- two insulated copper wires
 - Category 5: 100 Mbps, 1 Gbps Ethernet
 - Category 6: 10Gbps



Introduction 1-19

19

Physical media: coax, fiber

coaxial cable:

- two concentric copper conductors
- bidirectional
- broadband:
 - multiple channels on cable
 - HFC



fiber optic cable:

- glass fiber carrying light pulses, each pulse a bit
- high-speed operation:
 - high-speed point-to-point transmission (e.g., 10' s-100' s Gbps transmission rate)
- low error rate:
 - repeaters spaced far apart
 - immune to electromagnetic noise



Introduction 1-20

20

Physical media: radio

- signal carried in electromagnetic spectrum
- no physical "wire"
- bidirectional
- propagation environment effects:
 - reflection
 - obstruction by objects
 - interference

radio link types:

- **terrestrial microwave**
 - e.g. up to 45 Mbps channels
- **LAN** (e.g., WiFi)
 - 54 Mbps
- **wide-area** (e.g., cellular)
 - 4G cellular: ~ 10 Mbps
- **satellite**
 - Kbps to 45Mbps channel (or multiple smaller channels)
 - 270 msec end-end delay
 - geosynchronous versus low altitude

Introduction 1-21

21

Chapter 1: roadmap

1.1 what is the Internet?

1.2 network edge

- end systems, access networks, links

1.3 network core

- packet switching, circuit switching, network structure

1.4 delay, loss, throughput in networks

1.5 protocol layers, service models

1.6 networks under attack: security

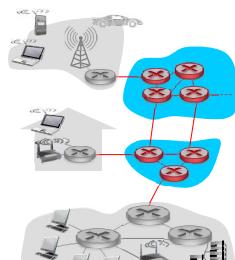
1.7 history

Introduction 1-22

22 Router is layer 3 device: network layer

The network core

- mesh of interconnected routers
- **packet-switching:** hosts break application-layer messages into **packets**
 - forward packets from one router to the next, across links on path from source to destination
 - each packet transmitted at full link capacity



Introduction 1-23

23

Packet-switching: store-and-forward

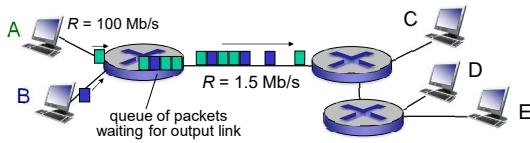
- L bits per packet
- source R bps destination
-
- takes L/R seconds to transmit (push out) L -bit packet into link at R bps
 - **store and forward:** entire packet must arrive at router before it can be transmitted on next link
 - end-end delay = $2L/R$ (assuming zero propagation delay)
- one-hop numerical example:
- $L = 7.5$ Mbits
 - $R = 1.5$ Mbps
 - one-hop transmission delay = 5 sec
- } more on delay shortly ...

Introduction 1-24

24

Routing forwarding

Packet Switching: queueing delay, loss



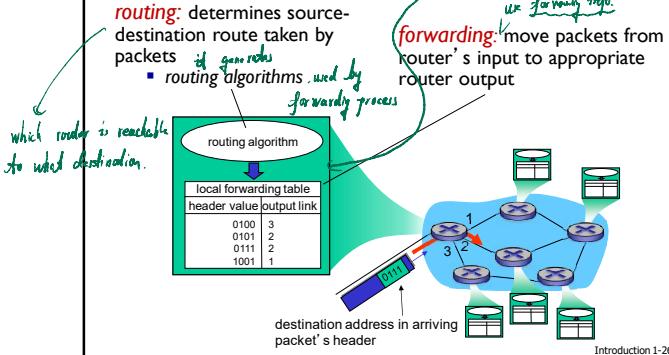
queueing and loss:

- if arrival rate (in bits) to link exceeds transmission rate of link for a period of time:
 - packets will queue, wait to be transmitted on link
 - packets can be dropped (lost) if memory (buffer) fills up

Introduction 1-25

25

Two key network-core functions

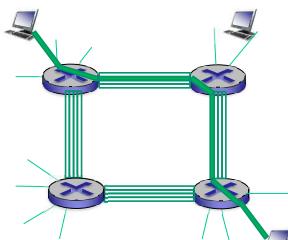


26 *routing algorithms: does not do routing . only produce table.
Forwarding is not intelligent operation*

Alternative core: circuit switching

end-end resources allocated to, reserved for “call” between source & dest:

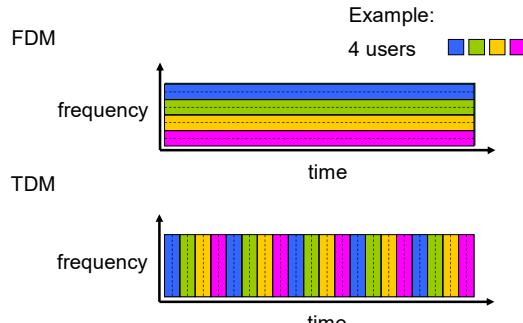
- in diagram, each link has four circuits.
 - call gets 2nd circuit in top link and 1st circuit in right link.
- dedicated resources: no sharing
 - circuit-like (guaranteed) performance
- circuit segment idle if not used by call (*no sharing*)
- commonly used in traditional telephone networks



Introduction 1-27

27

Circuit switching: FDM versus TDM



Introduction 1-28

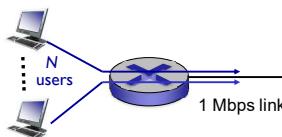
28

Packet switching versus circuit switching

packet switching allows more users to use network!

example:

- 1 Mb/s link
- each user:
 - 100 kb/s when “active”
 - active 10% of time
- circuit-switching:**
 - 10 users
- packet switching:**
 - with 35 users, probability > 10 active at same time is less than .0004 *



Introduction 1-29

* Check out the online interactive exercises for more examples: http://gaia.cs.umass.edu/kurose_ross/interactive/

29

Packet switching versus circuit switching

is packet switching a “slam dunk winner?”

- great for bursty data
 - resource sharing
 - simpler, no call setup
- excessive congestion possible:** packet delay and loss
 - protocols needed for reliable data transfer, congestion control
- Q: How to provide circuit-like behavior?**
 - bandwidth guarantees needed for audio/video apps
 - still an unsolved problem (chapter 7)

Q: human analogies of reserved resources (circuit switching) versus on-demand allocation (packet-switching)?

Introduction 1-30

30

Internet structure: network of networks

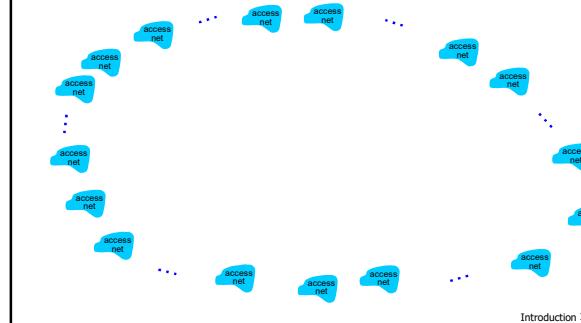
- End systems connect to Internet via **access ISPs** (Internet Service Providers)
 - residential, company and university ISPs
- Access ISPs in turn must be interconnected.
 - so that any two hosts can send packets to each other
- Resulting network of networks is very complex
 - evolution was driven by **economics** and **national policies**
- Let's take a stepwise approach to describe current Internet structure

Introduction 1-31

31

Internet structure: network of networks

Question: given millions of access ISPs, how to connect them together?

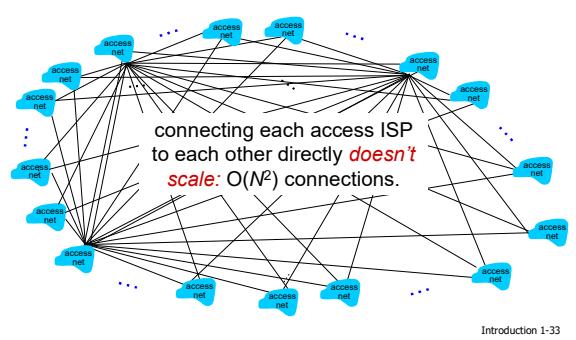


Introduction 1-32

32

Internet structure: network of networks

Option: connect each access ISP to every other access ISP?

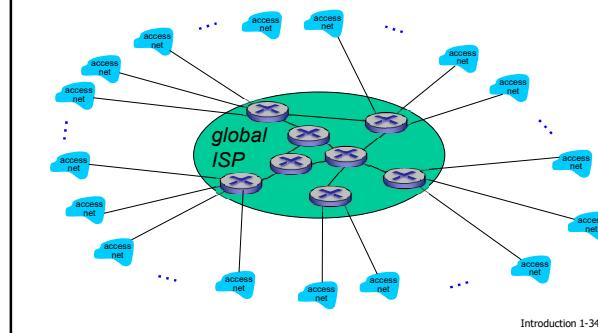


Introduction 1-33

33

Internet structure: network of networks

Option: connect each access ISP to one global transit ISP?
Customer and **provider** ISPs have economic agreement.

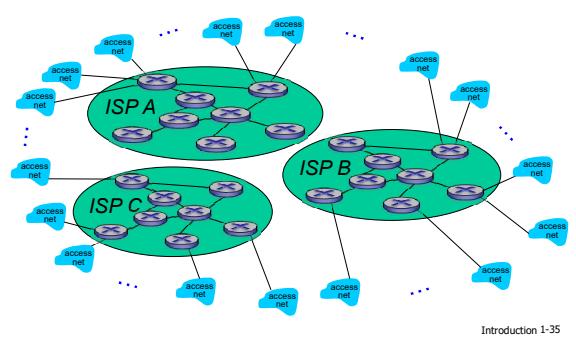


Introduction 1-34

34

Internet structure: network of networks

But if one global ISP is viable business, there will be competitors
....

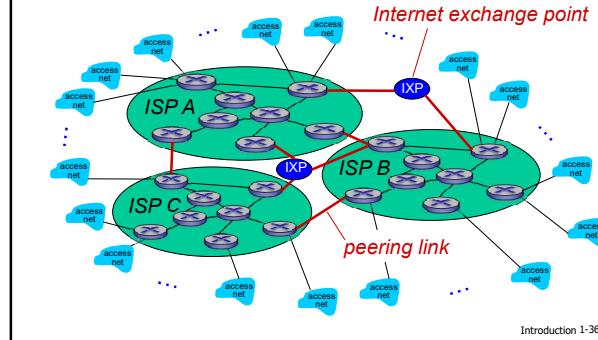


Introduction 1-35

35

Internet structure: network of networks

But if one global ISP is viable business, there will be competitors
.... which must be interconnected

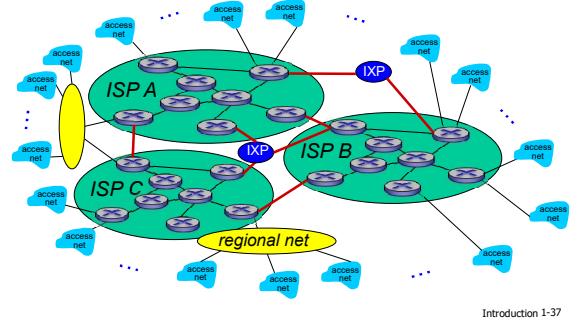


Introduction 1-36

36

Internet structure: network of networks

... and regional networks may arise to connect access nets to ISPs

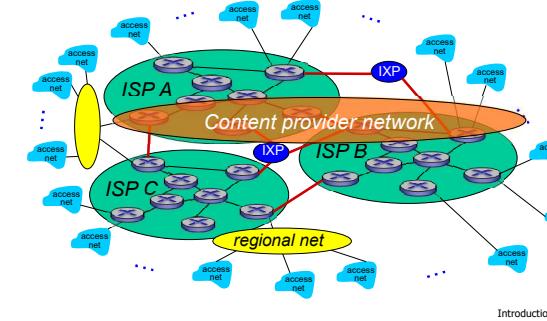


Introduction 1-37

37

Internet structure: network of networks

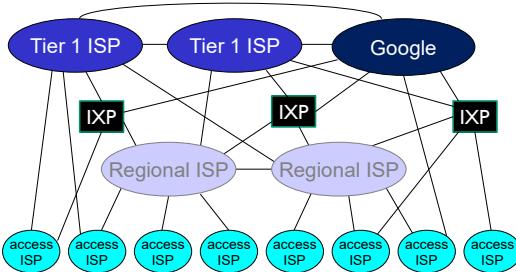
... and content provider networks (e.g., Google, Microsoft, Akamai) may run their own network, to bring services, content close to end users



Introduction 1-38

38

Internet structure: network of networks

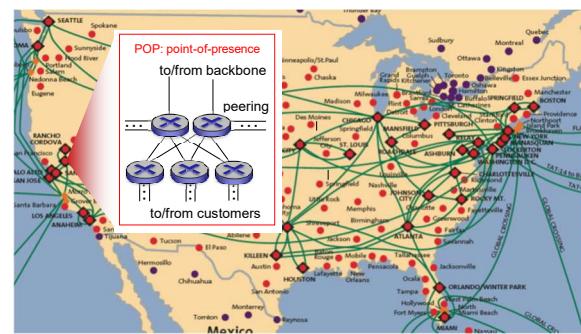


- at center: small # of well-connected large networks
 - "tier-1" commercial ISPs (e.g., Level 3, Sprint, AT&T, NTT), national & international coverage
 - content provider network (e.g., Google): private network that connects it data centers to Internet, often bypassing tier-1, regional ISPs

Introduction 1-39

39

Tier-1 ISP: e.g., Sprint



Introduction 1-40

40

Chapter 1: roadmap

- 1.1 what is the Internet?
- 1.2 network edge
 - end systems, access networks, links
- 1.3 network core
 - packet switching, circuit switching, network structure
- 1.4 delay, loss, throughput in networks
- 1.5 protocol layers, service models
- 1.6 networks under attack: security
- 1.7 history

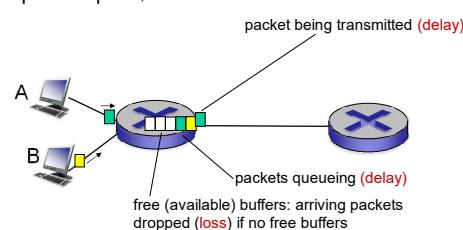
Introduction 1-41

41

How do loss and delay occur?

packets queue in router buffers

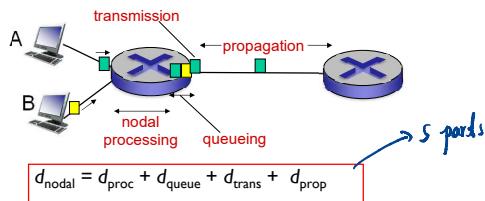
- packet arrival rate to link (temporarily) exceeds output link capacity
- packets queue, wait for turn



Introduction 1-42

42

Four sources of packet delay

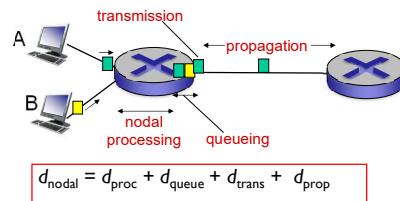


- d_{proc} : nodal processing
 - check bit errors
 - determine output link
 - typically < msec

- d_{queue} : queueing delay
 - time waiting at output link for transmission
 - depends on congestion level of router

Introduction 1-43

Four sources of packet delay



d_{trans} : transmission delay:

- L : packet length (bits)
- R : link bandwidth (bps)
- $d_{\text{trans}} = L/R$

d_{prop} : propagation delay:

- d : length of physical link
- s : propagation speed ($\sim 2 \times 10^8$ m/sec)
- $d_{\text{prop}} = d/s$

* Check out the online interactive exercises for more examples: http://gaia.cs.umass.edu/kurose_ross/interactive/
* Check out the Java applet for an interactive animation on trans vs. prop delay

Introduction 1-44

43

44

Caravan analogy



- cars "propagate" at 100 km/hr
- toll booth takes 12 sec to service car (bit transmission time)
- car ~ bit; caravan ~ packet
- Q: How long until caravan is lined up before 2nd toll booth?

- time to "push" entire caravan through toll booth onto highway = $12 \times 10 = 120$ sec
- time for last car to propagate from 1st to 2nd toll booth: $100\text{km}/(100\text{km/hr}) = 1$ hr
- A: 62 minutes

Introduction 1-45

Introduction 1-46

45

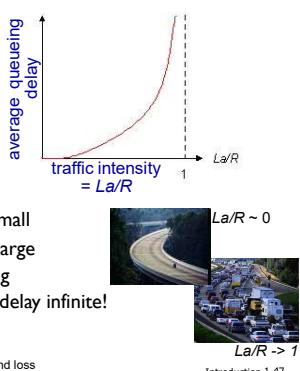
Caravan analogy (more)

-
- suppose cars now "propagate" at 1000 km/hr
 - and suppose toll booth now takes one min to service a car
 - Q: Will cars arrive to 2nd booth before all cars serviced at first booth?
 - A: Yes! after 7 min, first car arrives at second booth; three cars still at first booth

参考后面的 notes

Queueing delay (revisited)

- R : link bandwidth (bps)
- L : packet length (bits)
- a : average packet arrival rate



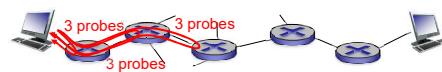
Introduction 1-47

* Check online interactive animation on queuing and loss

47

"Real" Internet delays and routes

- what do "real" Internet delay & loss look like?
- traceroute** program: provides delay measurement from source to router along end-end Internet path towards destination. For all i :
 - sends three packets that will reach router i on path towards destination
 - router i will return packets to sender
 - sender times interval between transmission and reply.



<https://gsuite.tools/traceroute>

<https://centralops.net/co/>

Introduction 1-48

48

"Real" Internet delays, routes

traceroute: gaia.cs.umass.edu to www.eurecom.fr

```

1 cs-gw (128.119.240.254) 1 ms 1 ms 2 ms
2 border1-rt-fa5-1-0.gw.umass.edu (128.119.3.130) 6 ms 5 ms 5 ms
3 cht-vbns.gw.umass.edu (128.119.3.130) 6 ms 5 ms 5 ms
4 in1-at1-0-19.wor.vbns.net (204.147.132.129) 16 ms 11 ms 13 ms
5 in1-so7-0-0-wae.vbns.net (204.147.136.136) 21 ms 18 ms 18 ms
6 abilene-vbns.abilene.uchicago.edu (198.32.11.9) 22 ms 18 ms 22 ms
7 nycm-wash.abilene.uchicago.edu (198.32.8.46) 22 ms 22 ms 22 ms
8 62.40.103.253 (62.40.103.253) 104 ms 109 ms 106 ms
9 de2-1.de1.de1.geant.net (62.40.96.50) 113 ms 121 ms 114 ms
10 de1-fr1.fr.geant.net (62.40.96.50) 113 ms 121 ms 114 ms
11 renater-gw.fr1.fr.geant.net (62.40.103.54) 112 ms 114 ms 112 ms
12 nfr2.cs.renater.fr (194.220.98.109) 13 ms 14 ms 10 ms
13 nfr2.cs.renater.fr (194.220.98.110) 123 ms 125 ms 124 ms
14 r3t2-nice.cs.renater.fr (195.220.98.110) 126 ms 126 ms 124 ms
15 eurecom-valbonne.r3t2.fr.net (193.48.50.54) 135 ms 128 ms 133 ms
16 194.214.211.25 (194.214.211.25) 126 ms 128 ms 126 ms
17 ***
18 *** * means no response (probe lost, router not replying)
19 fantasia.eurecom.fr (193.55.113.142) 132 ms 128 ms 136 ms

```

* means no response (probe lost, router not replying)

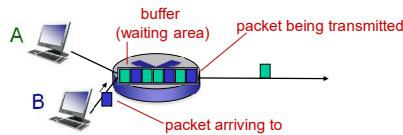
* Do some traceroutes from exotic countries at www.traceroute.org

Introduction 1-49

49

Packet loss

- queue (aka buffer) preceding link in buffer has finite capacity
- packet arriving to full queue dropped (aka lost)
- lost packet may be retransmitted by previous node, by source end system, or not at all



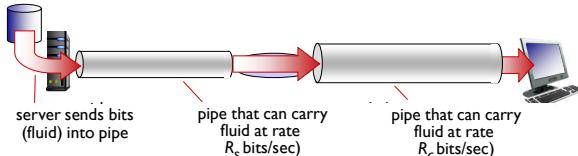
* Check out the Java applet for an interactive animation on queuing and loss

Introduction 1-50

50

Throughput

- throughput:** rate (bits/time unit) at which bits transferred between sender/receiver
 - instantaneous:** rate at given point in time
 - average:** rate over longer period of time

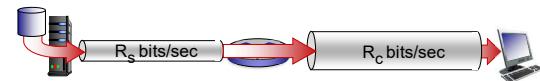


Introduction 1-51

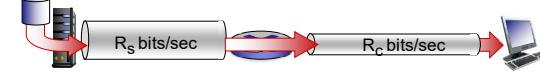
51

Throughput (more)

- $R_s < R_c$ What is average end-end throughput?



- $R_s > R_c$ What is average end-end throughput?



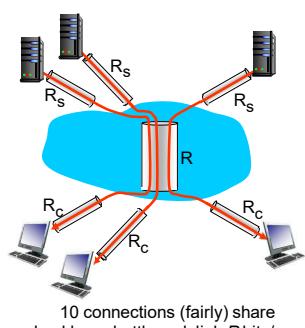
bottleneck link
link on end-end path that constrains end-end throughput

Introduction 1-52

52

Throughput: Internet scenario

- per-connection end-end throughput: $\min(R_c, R_s)R/10$)
- in practice: R_c or R_s is often bottleneck



10 connections (fairly) share backbone bottleneck link R bits/sec

Introduction 1-53

* Check out the online interactive exercises for more examples: http://gaia.cs.umass.edu/kurose_ross/interactive/

53

Chapter 1: roadmap

- 1.1 what is the Internet?
- 1.2 network edge
 - end systems, access networks, links
- 1.3 network core
 - packet switching, circuit switching, network structure
- 1.4 delay, loss, throughput in networks
- 1.5 protocol layers, service models
- 1.6 networks under attack: security
- 1.7 history

Introduction 1-54

54

Protocol “layers”

Networks are complex,
with many “pieces”:

- hosts
- routers
- links of various media
- applications
- protocols
- hardware, software

Question:

is there any hope of organizing structure of network?

.... or at least our discussion of networks?

Introduction 1-55

55

Organization of air travel

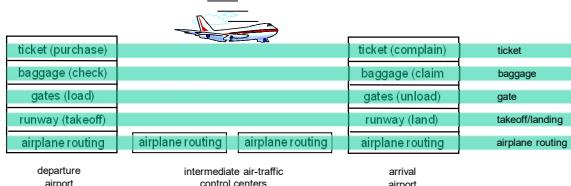


- a series of steps

Introduction 1-56

56

Layering of airline functionality



- layers:** each layer implements a service
- via its own internal-layer actions
 - relying on services provided by layer below

Introduction 1-57

57

Why layering?

dealing with complex systems:

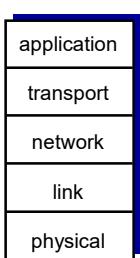
- explicit structure allows identification, relationship of complex system’s pieces
 - layered *reference model* for discussion
- modularization eases maintenance, updating of system
 - change of implementation of layer’s service transparent to rest of system
 - e.g., change in gate procedure doesn’t affect rest of system
- layering considered harmful?

Introduction 1-58

58

Internet protocol stack

- **application:** supporting network applications
 - FTP, SMTP, HTTP
- **transport:** process-process data transfer
 - TCP, UDP
- **network:** routing of datagrams from source to destination
 - IP, routing protocols
- **link:** data transfer between neighboring network elements
 - Ethernet, 802.111 (WiFi), PPP
- **physical:** bits “on the wire”

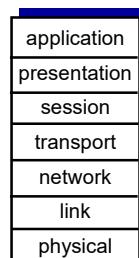


Introduction 1-59

59

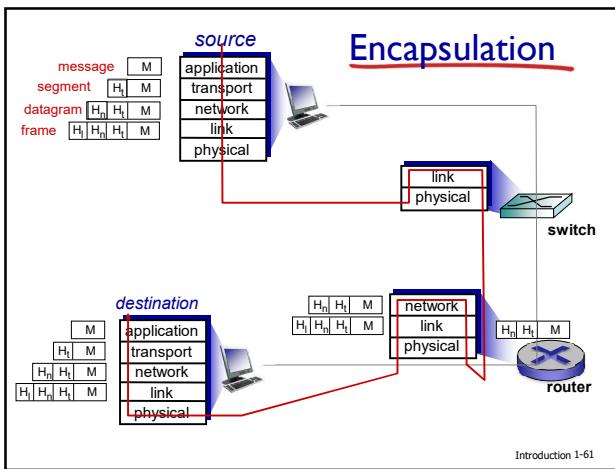
ISO/OSI reference model

- **presentation:** allow applications to interpret meaning of data, e.g., encryption, compression, machine-specific conventions
- **session:** synchronization, checkpointing, recovery of data exchange
- Internet stack “missing” these layers!
 - these services, *if needed*, must be implemented in application
 - needed?



Introduction 1-60

60



61

Chapter 1: roadmap

- 1.1 what is the Internet?
- 1.2 network edge
 - end systems, access networks, links
- 1.3 network core
 - packet switching, circuit switching, network structure
- 1.4 delay, loss, throughput in networks
- 1.5 protocol layers, service models
- 1.6 networks under attack: security
- 1.7 history

Introduction 1-62

62

Network security

- field of network security:
 - how bad guys can attack computer networks
 - how we can defend networks against attacks
 - how to design architectures that are immune to attacks
- Internet not originally designed with (much) security in mind
 - original vision: “a group of mutually trusting users attached to a transparent network” ☺
 - Internet protocol designers playing “catch-up”
 - security considerations in all layers!

Introduction 1-63

63

Bad guys: put malware into hosts via Internet

- malware can get in host from:
 - virus: self-replicating infection by receiving/executing object (e.g., e-mail attachment)
 - worm: self-replicating infection by passively receiving object that gets itself executed
- spyware malware can record keystrokes, web sites visited, upload info to collection site
- infected host can be enrolled in botnet, used for spam, DDoS attacks

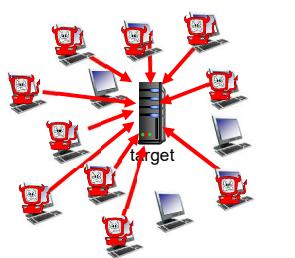
Introduction 1-64

64

Bad guys: attack server, network infrastructure

Denial of Service (DoS): attackers make resources (server, bandwidth) unavailable to legitimate traffic by overwhelming resource with bogus traffic

1. select target
2. break into hosts around the network (see botnet)
3. send packets to target from compromised hosts

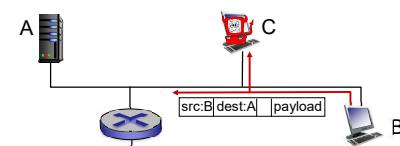


65

Bad guys can sniff packets

packet “sniffing”:

- broadcast media (shared Ethernet, wireless)
- promiscuous network interface reads/records all packets (e.g., including passwords!) passing by

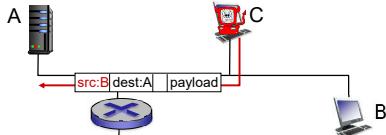


- wireshark software used for end-of-chapter labs is a (free) packet-sniffer

66

Bad guys can use fake addresses

IP spoofing: send packet with false source address



... lots more on security (throughout, Chapter 8)

Introduction 1-67

67

Chapter 1: roadmap

1.1 what is the Internet?

1.2 network edge

- end systems, access networks, links

1.3 network core

- packet switching, circuit switching, network structure

1.4 delay, loss, throughput in networks

1.5 protocol layers, service models

1.6 networks under attack: security

1.7 history

Introduction 1-68

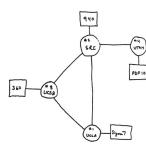
68

Internet history

1961-1972: Early packet-switching principles

- **1961:** Kleinrock - queueing theory shows effectiveness of packet-switching
- **1964:** Baran - packet-switching in military nets
- **1967:** ARPAnet conceived by Advanced Research Projects Agency
- **1969:** first ARPAnet node operational

- **1972:**
 - ARPAnet public demo
 - NCP (Network Control Protocol) first host-host protocol
 - first e-mail program
 - ARPAnet has 15 nodes



Introduction 1-69

69

Internet history

1972-1980: Internetworking, new and proprietary nets

- **1970:** ALOHAnet satellite network in Hawaii
- **1974:** Cerf and Kahn - architecture for interconnecting networks
- **1976:** Ethernet at Xerox PARC
- **late 70's:** proprietary architectures: DECnet, SNA, XNA
- **late 70's:** switching fixed length packets (ATM precursor)
- **1979:** ARPAnet has 200 nodes

Cerf and Kahn's internetworking principles:

- minimalism, autonomy - no internal changes required to interconnect networks
- best effort service model
- stateless routers
- decentralized control

define today's Internet architecture

Introduction 1-70

70

Internet history

1980-1990: new protocols, a proliferation of networks

- **1983:** deployment of TCP/IP
- **1982:** smtp e-mail protocol defined
- **1983:** DNS defined for name-to-IP-address translation
- **1985:** ftp protocol defined
- **1988:** TCP congestion control
- **new national networks:** CSNet, BITnet, NSFnet, Minitel
- **1983:** 100,000 hosts connected to confederation of networks

Introduction 1-71

71

Internet history

1990, 2000's: commercialization, the Web, new apps

- **early 1990's:** ARPAnet decommissioned
- **1991:** NSF lifts restrictions on commercial use of NSFnet (decommissioned, 1995)
- **early 1990s:** Web
 - hypertext [Bush 1945, Nelson 1960's]
 - HTML, HTTP: Berners-Lee
 - 1994: Mosaic, later Netscape
 - late 1990's: commercialization of the Web
- **late 1990's – 2000's:**
 - more killer apps: instant messaging, P2P file sharing
 - network security to forefront
 - est. 50 million host, 100 million+ users
 - backbone links running at Gbps

Introduction 1-72

72

12

Internet history

2005-present

- ~5B devices attached to Internet (2016)
 - smartphones and tablets
- aggressive deployment of broadband access
- increasing ubiquity of high-speed wireless access
- emergence of online social networks:
 - Facebook: ~ one billion users
- service providers (Google, Microsoft) create their own networks
 - bypass Internet, providing “instantaneous” access to search, video content, email, etc.
- e-commerce, universities, enterprises running their services in “cloud” (e.g., Amazon EC2)

Introduction 1-73

73

Introduction: summary

covered a “ton” of material!

- Internet overview
- what’s a protocol?
- network edge, core, access network
 - packet-switching versus circuit-switching
- Internet structure
- performance: loss, delay, throughput
- layering, service models
- security
- history

you now have:

- context, overview, “feel” of networking
- more depth, detail to follow!

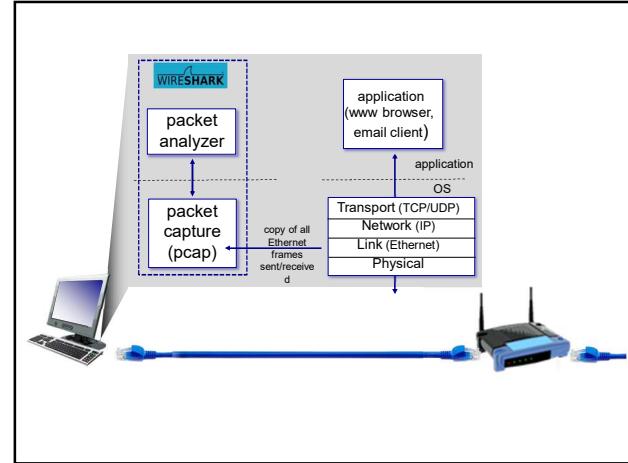
Introduction 1-74

74

Chapter I Additional Slides

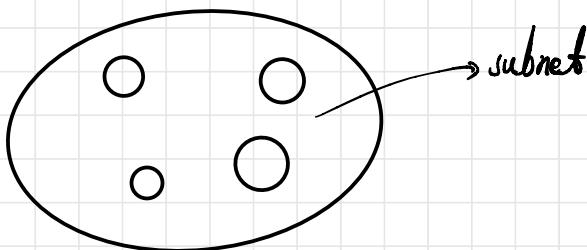
Introduction 1-75

75



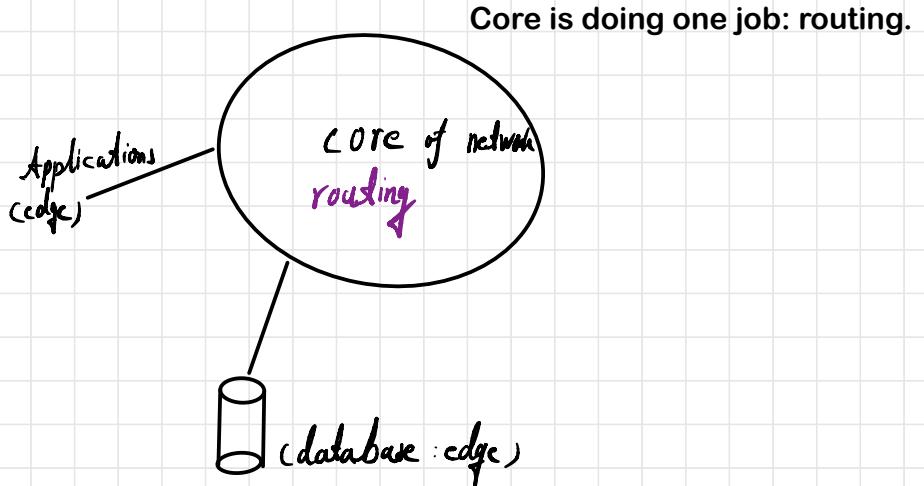
76

1/11/2022



- Network backbone: IETF(Internet Engineering Task Force)

- Major designer of internet: CISCO



Flag day: The day to reboot network.(never gonna happen).

- Modem: Mojolater and DEMojolater.

- Protocol:

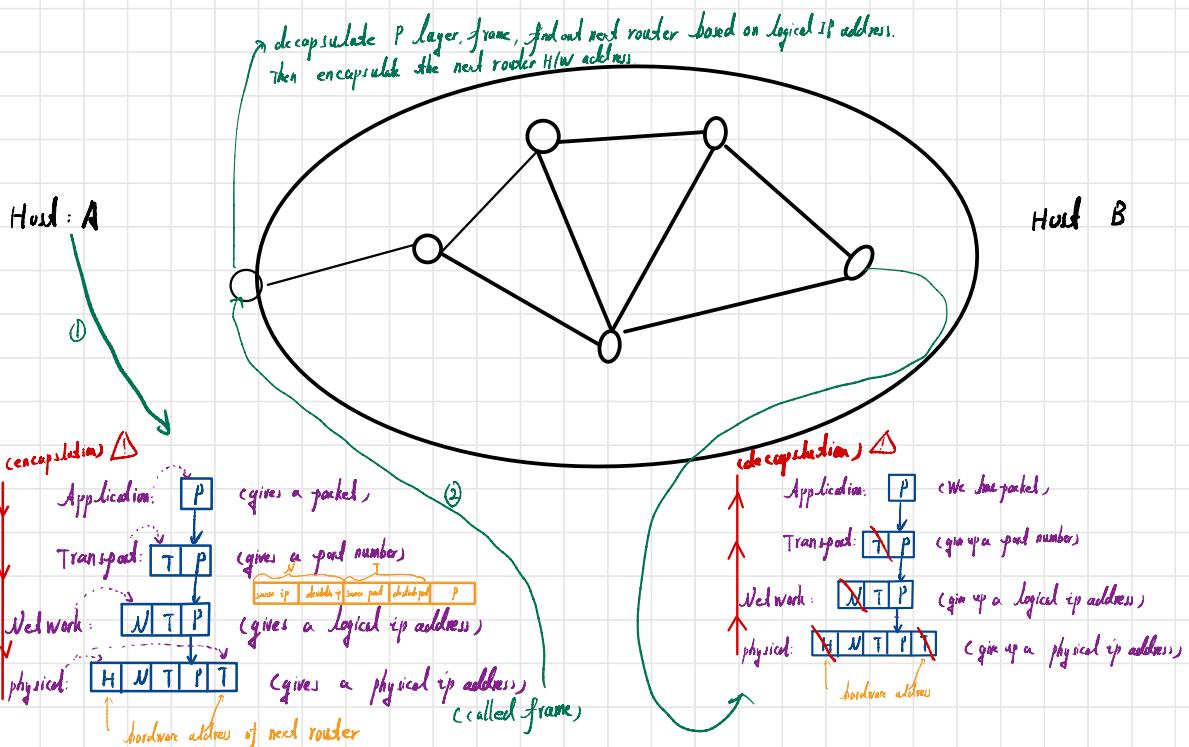
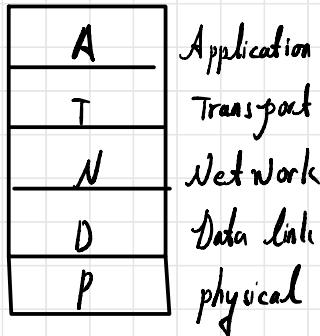
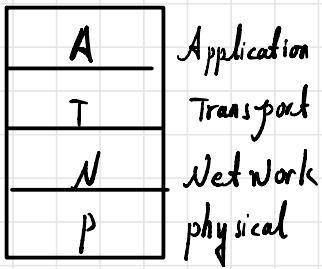
protocols define format, order of messages sent and received among network entities, and actions taken on message transmission, receipt

Introduction 1-7

Packet: chunk of data.

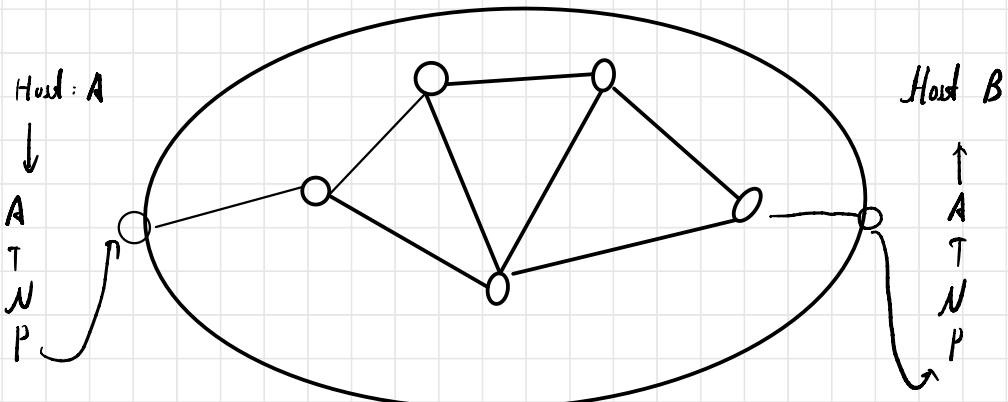
Transmission rate: how fast your pc dump data on wires.

Propagation rate: how fast the single transfer on wire.



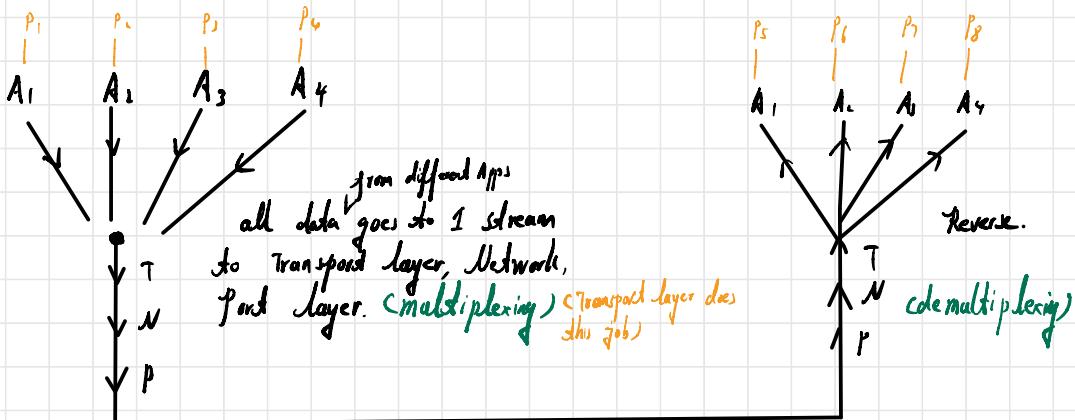
{ dynamic IP address (logical address)
physical address (Mac address)

next router = next hop



It is possible to skip Transport layer, but it turns out pc can only run one application, no multiple programs.

Mud & Demux: multiplexing and demultiplexing.



- Q: Does the transport need to know the destination port ahead?

Yes, it must know it before you can communicate that machine.

?

- Q: Could A1 and A5 use same port number like P1?

Yes, since the port number is unique on one computer.

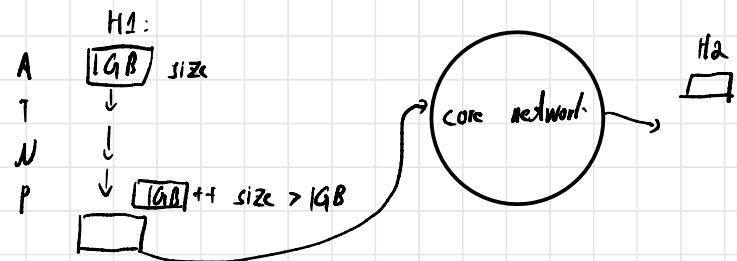
But you can not let two different app runs on same pc and use same port number, that will be disaster.

- Loop back IP address: 127.0.0.1

When we send this data to this IP address, it comes back to you for testing purpose.

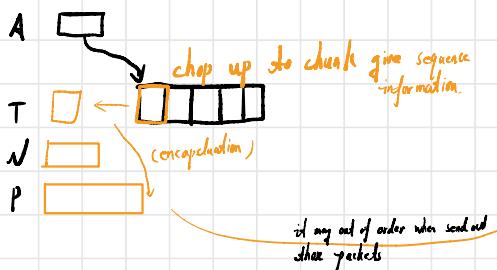
- Three switching methods.

- Message switching: No limited for data size which sent out to network.



Bad things:
 ① If the connection in the network disconnect then the whole data will be discarded.
 ② Fair share of internet will be interrupted.

- Packet switching: There is limit for data size which sent out to network. And chop it up to chunk.



- Circuit switching: guarantee the resources is from h1 to h2, but low source reservation.



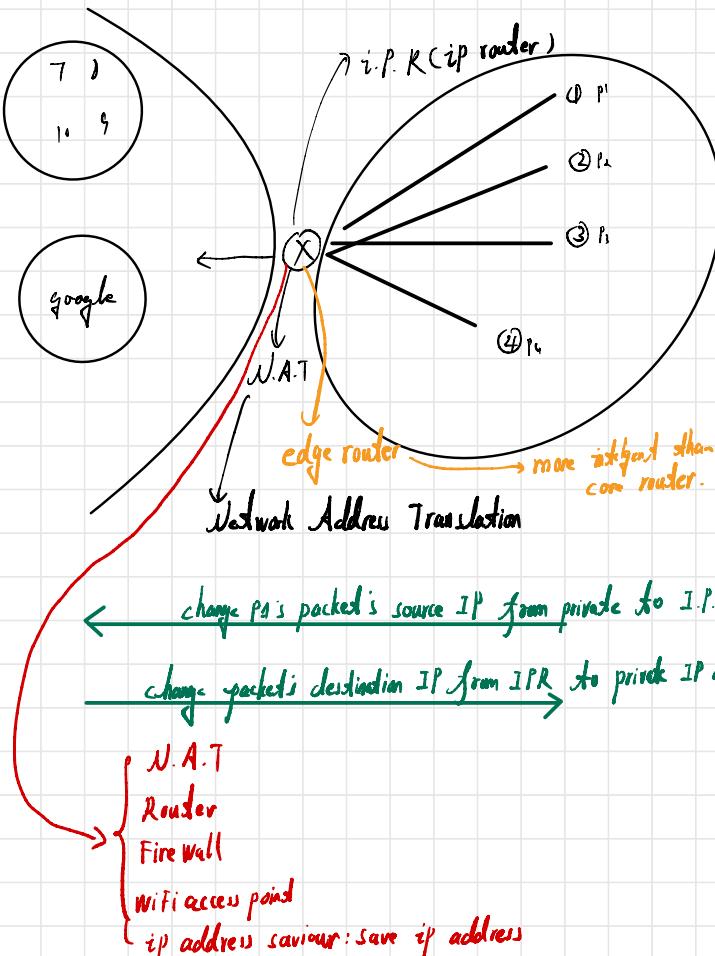
Switch: dump router, never change the connect automatically. (static)

eg: Truman WiFi accessing point, WiFi Extender

Router: intelligent, and dynamic automatically change the links. Make decisions where to forward.

eg: Truman network outgoing router, home wifi router, NAT

Home network:



private ip: 10 / 192

public ip: Globally unique

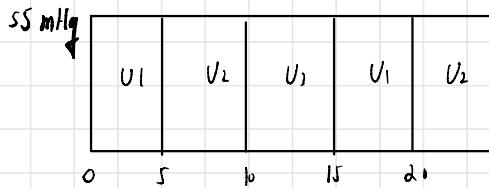
edge router may not always be NAT

eg: Truman: Gate V, NAT v
Home: Gate V, NAT v

TDM(Time division multiplexing)

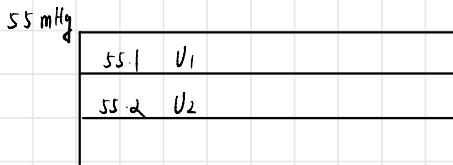
Collision will happen if multiple users use same frequency at same time.

Same Frequency different time

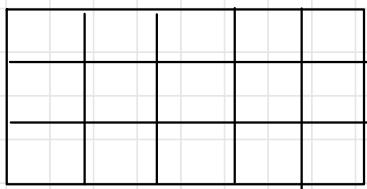


Timeline : interval 5ms

FDM(Frequency devision multiplexing)



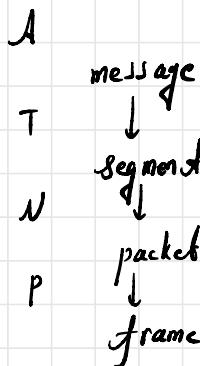
If more users then combine fdm and tdm.



Some notes are on slides.

Segmentation and reassembly

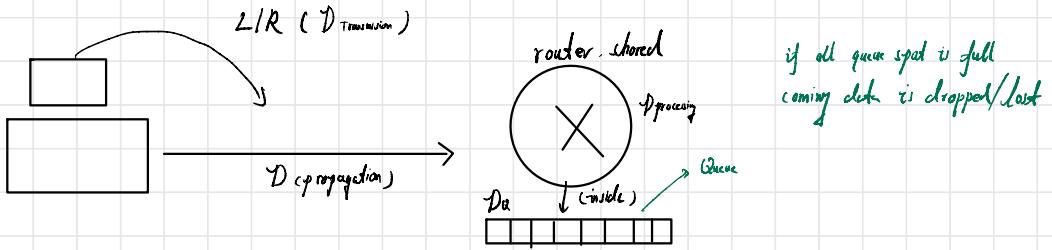
{ segmentation: Transport layer segment message into packets
reassembly: reassemble all packets by using sequence information



- **store and forward:** entire packet must arrive at router before it can be transmitted on next link

The router need to wait whole packet to arrive even they can read the ip information. Why?

Because the packet may have error and router will know until whole packet has arrived. Thats why we need store and forward.



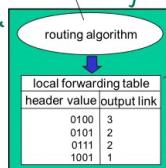
$$\text{total delay: } D_T + D_{prop} + D_q + D_{proc} = D_{\text{total}}$$

Routing / Forwarding

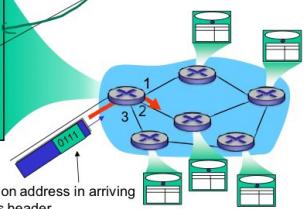
Routing forwarding

Two key network-core functions

routing: determines source-destination route taken by packets of multiple nodes
 - routing algorithms used by forwarding process
 which router is reachable to what destination.



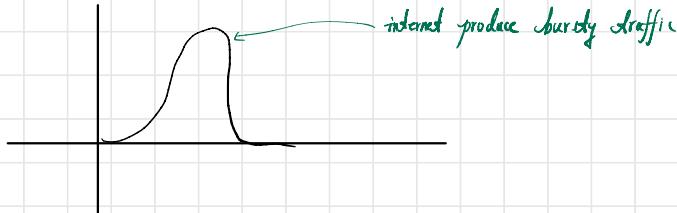
forwarding: move packets from router's input to appropriate router output
we forward info.



Introduction 1-26

26 routing Algorithms does not do routing - only produce table.
Forwarding is not intelligent operation

Internet traffic = bursty traffic



Packet switching is good for public internet, like bursty data, and it could be shared for many users.

Circuit switching: is good for dedicated link communication like government or special comms.

R: Link bandwidth(processing capacity)

L: Packet length.

a: packet arrival rate.

B_{ps} = bytes per second
 b_{ps} = bits per second

congestion: 滞塞

$$\text{Traffic intensity} = L \cdot a / R$$

$$\frac{La}{R} \frac{0}{\bigcirc} \approx 0, \quad \frac{0}{0} \approx 1 \cdot \frac{\bigcirc}{0} > 1$$

\downarrow full \downarrow not lost data \downarrow begin to lose data

Conclusion: No matter how to decrease the a or L or increase the R, there is no way to decline the congestion.

- That's why we have congestion control and install it on Transport layer.

Because network layer is used to prepare routing, physical layer doesn't make any decisions. Install on application will require more effort from programmers. Only transport is left.

- Congestion control: Slow down the network speed that won't be noticed by users but enough for router to process datas.

TCP is more sensitive, and UDP doesn't care about congestion.

Max-hop-count: max number of intermediate router your packet should visit till destination in order to make sure the packet will not stuck in loop of routers.

If the count is to 0, then router drop the packet.

ICMP: internet control message protocol: used to determine if the packet reached to the destination in timely manner. And not every router has it.

Review trace route

Throughout:

C₁, C₂, C₃, C₄, C₅

The effective internet speed is affected by the smallest intermediate link rate.

瓶颈效应 (bottleneck)

ISO/OSI model: 7 layers could be found on slides.

Why do we have layers?

Why layering?

dealing with complex systems:

- explicit structure allows identification, relationship of complex system's pieces } separate the network into specific areas
 - layered *reference model* for discussion
- modularization eases maintenance, updating of system
 - change of implementation of layer's service transparent to rest of system
 - e.g., change in gate procedure doesn't affect rest of system
- layering considered harmful?

Job of Application layer: syntax semantic meaning of transport message.

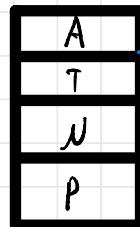
Transport layer: Deliver packet from one application to application.

Network: host to host delivery.

Physical layer: node to node delivery.

- TLS: Transport Layer Security.

Used to be SSL but now is TLS.



TLS is here, encrypt the whole message, not single packet since this is more economical and faster

Chapter 2

Chapter 2: outline

- | | |
|---|--|
| 2.1 principles of network applications | 2.5 P2P applications |
| 2.2 Web and HTTP | 2.6 video streaming and content distribution networks |
| 2.3 electronic mail | |
| • SMTP, POP3, IMAP | |
| 2.4 DNS | 2.7 socket programming with UDP and TCP |

Application Layer 2-1

1

Chapter 2: application layer

our goals:

- conceptual, implementation aspects of network application protocols
 - transport-layer service models
 - client-server paradigm
 - peer-to-peer paradigm
 - content distribution networks
- learn about protocols by examining popular application-level protocols
 - HTTP
 - FTP
 - SMTP / POP3 / IMAP
 - DNS
- creating network applications
 - socket API

Application Layer 2-2

2

Some network apps

- e-mail
- web
- text messaging
- remote login
- P2P file sharing
- multi-user network games
- streaming stored video (YouTube, Hulu, Netflix)
- voice over IP (e.g., Skype)
- real-time video conferencing
- social networking
- search
- ...
- ...

Application Layer 2-3

3

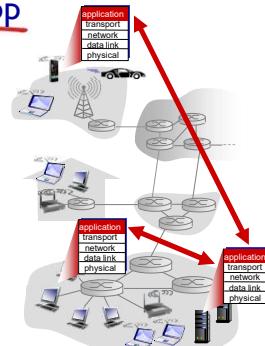
Creating a network app

write programs that:

- run on (different) end systems
- communicate over network
- e.g., web server software communicates with browser software

no need to write software for network-core devices

- network-core devices do not run user applications
- applications on end systems allows for rapid app development, propagation



Application Layer 2-4

4

Application architectures

possible structure of applications:

- client-server
- peer-to-peer (P2P)

Application Layer 2-5

5

Client-server architecture

server:

- always-on host
- permanent IP address
- data centers for scaling

clients:

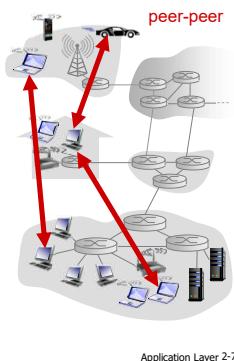
- communicate with server
- may be intermittently connected
- may have dynamic IP addresses
- do not communicate directly with each other

Application Layer 2-6

6

P2P architecture

- no always-on server
- arbitrary end systems directly communicate
- peers request service from other peers, provide service in return to other peers
 - self scalability* – new peers bring new service capacity, as well as new service demands
- peers are intermittently connected and change IP addresses
 - complex management



Application Layer 2-7

7

Processes communicating

- process:** program running within a host
- within same host, two processes communicate using **inter-process communication** (defined by OS)
 - processes in different hosts communicate by exchanging **messages**

clients, servers
client process: process that initiates communication
server process: process that waits to be contacted

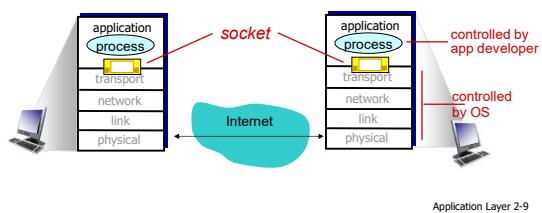
- aside: applications with P2P architectures have client processes & server processes

Application Layer 2-8

8

Sockets

- process sends/receives messages to/from its **socket**
- socket analogous to door
 - sending process shoves message out door
 - sending process relies on transport infrastructure on other side of door to deliver message to socket at receiving process



Application Layer 2-9

9

Data is generated by App layer

App-layer protocol defines

- types of messages exchanged,
 - e.g., request, response
- message syntax:**
 - what fields in messages & how fields are delineated
- message semantics**
 - meaning of information in fields
- rules** for when and how processes send & respond to messages

- open protocols:**
- defined in RFCs
 - allows for interoperability
 - e.g., HTTP, SMTP
- proprietary protocols:**
- e.g., Skype

需要速度，但是可能掉包，说的是fps游戏

Application Layer 2-11

11

10

就是说需要完整的data, 但是速度不快: gmail

T offer services like this

What transport service does an app need?

data integrity

- some apps (e.g., file transfer, web transactions) require 100% reliable data transfer
- other apps (e.g., audio) can tolerate some loss

timing

- some apps (e.g., Internet telephony, interactive games) require low delay to be "effective"

throughput

- some apps (e.g., multimedia) require minimum amount of throughput to be "effective"
- other apps ("elastic apps") make use of whatever throughput they get

吞吐量

- encryption, data integrity, ...

Application Layer 2-12

12

Transport service requirements: common apps

application	data loss	throughput	time sensitive
file transfer	no loss	elastic	no
e-mail	no loss	elastic	no
Web documents	no loss	elastic	no
real-time audio/video	loss-tolerant	audio: 5kbps-1Mbps video:10kbps-5Mbps	yes, 100's msec
stored audio/video	loss-tolerant	same as above	
interactive games	loss-tolerant	few kbps up	yes, few secs
text messaging	no loss	elastic	yes, 100's msec yes and no

Application Layer 2-13

13

TCP and UDP protocol:

Internet transport protocols services

TCP service:

- **reliable transport** between sending and receiving process
- **flow control**: sender won't overwhelm receiver
- **congestion control**: throttle sender when network overloaded
- **does not provide**: timing, minimum throughput guarantee, security
- **connection-oriented**: setup required between client and server processes

UDP service:

- **unreliable data transfer** between sending and receiving process
- **does not provide**: reliability, flow control, congestion control, timing, throughput guarantee, security, or connection setup,

Q: why bother? Why is there a UDP?

Application Layer 2-14

14

Internet apps: application, transport protocols

application	application layer protocol	underlying transport protocol
e-mail	SMTP [RFC 2821]	TCP
remote terminal access	Telnet [RFC 854]	TCP
Web	HTTP [RFC 2616]	TCP
file transfer	FTP [RFC 959]	TCP
streaming multimedia	HTTP (e.g., YouTube), RTP [RFC 1889]	TCP or UDP
Internet telephony	SIP, RTP, proprietary (e.g., Skype)	TCP or UDP

Application Layer 2-15

15

Securing TCP

TCP & UDP

- no encryption
- cleartext passwds sent into socket traverse Internet in cleartext
- SSL
- provides encrypted TCP connection
- data integrity
- end-point authentication

SSL is at app layer

- apps use SSL libraries, that "talk" to TCP
- **SSL socket API**
- cleartext passwords sent into socket traverse Internet encrypted
- see Chapter 8

Application Layer 2-16

16

Chapter 2: outline

- | | |
|--|---|
| 2.1 principles of network applications | 2.5 P2P applications |
| 2.2 Web and HTTP | 2.6 video streaming and content distribution networks |
| 2.3 electronic mail | 2.7 socket programming with UDP and TCP |
| • SMTP, POP3, IMAP | |
| 2.4 DNS | |

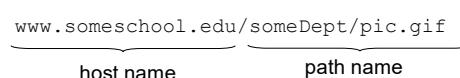
Application Layer 2-17

17

Web and HTTP

First, a review...

- **web page** consists of **objects**
- object can be HTML file, JPEG image, Java applet, audio file,...
- web page consists of **base HTML-file** which includes **several referenced objects**
- each object is addressable by a **URL**, e.g.,



Application Layer 2-18

18

HTTP overview

HTTP: hypertext transfer protocol

- Web's application layer protocol
- client/server model
 - **client:** browser that requests, receives, (using HTTP protocol) and "displays" Web objects
 - **server:** Web server sends (using HTTP protocol) objects in response to requests



Application Layer 2-19

19

HTTP overview (continued)

uses TCP:

- client initiates TCP connection (creates socket) to server, port 80
- server accepts TCP connection from client
- HTTP messages (application-layer protocol messages) exchanged between browser (HTTP client) and Web server (HTTP server)
- TCP connection closed

HTTP is "stateless"

- server maintains no information about past client requests

protocols that maintain "state" are complex!

- past history (state) must be maintained
- if server/client crashes, their views of "state" may be inconsistent, must be reconciled

Application Layer 2-20

20

HTTP connections

non-persistent HTTP

- at most one object sent over TCP connection
 - connection then closed
- downloading multiple objects required multiple connections

persistent HTTP

- multiple objects can be sent over single TCP connection between client, server

Application Layer 2-21

21

Non-persistent HTTP

suppose user enters URL:

`www.someSchool.edu/someDepartment/home.index` (contains text, references to 10 jpeg images)

- 1a. HTTP client initiates TCP connection to HTTP server (process) at `www.someSchool.edu` on port 80
- 1b. HTTP server at host `www.someSchool.edu` waiting for TCP connection at port 80. "accepts" connection, notifying client
2. HTTP client sends HTTP *request message* (containing URL) into TCP connection socket. Message indicates that client wants object `someDepartment/home.index`
3. HTTP server receives request message, forms *response message* containing requested object, and sends message into its socket

Application Layer 2-22

22

Non-persistent HTTP (cont.)

4. HTTP server closes TCP connection.
5. HTTP client receives response message containing html file, displays html. Parsing html file, finds 10 referenced jpeg objects

6. Steps 1-5 repeated for each 10 jpeg objects

Application Layer 2-23

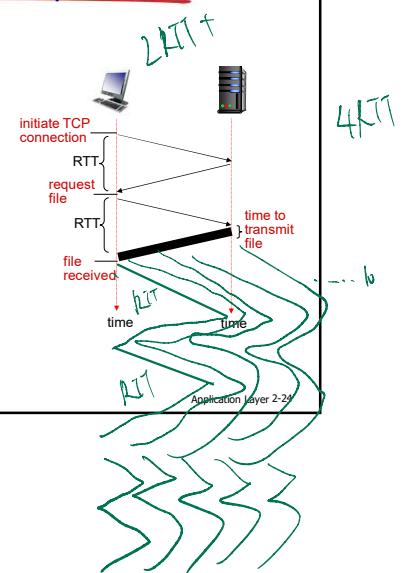
23

Non-persistent HTTP: response time

RTT (definition): time for a small packet to travel from client to server and back

HTTP response time:

- one RTT to initiate TCP connection
- one RTT for HTTP request and first few bytes of HTTP response to return
- file transmission time
- non-persistent HTTP response time = $2\text{RTT} + \text{file transmission time}$



24

Persistent HTTP

non-persistent HTTP issues:

- requires 2 RTTs per object
- OS overhead for each TCP connection
- browsers often open parallel TCP connections to fetch referenced objects

persistent HTTP:

- server leaves connection open after sending response
- subsequent HTTP messages between same client/server sent over open connection
- client sends requests as soon as it encounters a referenced object
- as little as one RTT for all the referenced objects

Application Layer 2-25

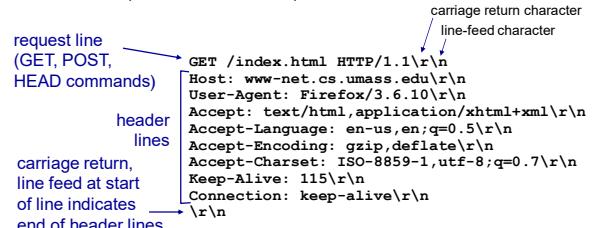
25

HTTP request message

- two types of HTTP messages: *request, response*

HTTP request message:

- ASCII (human-readable format)

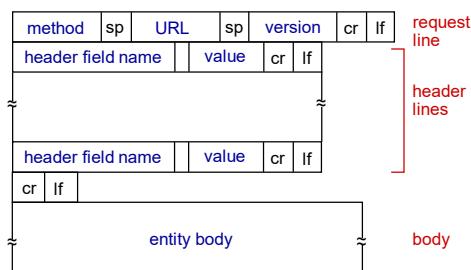


* Check out the online interactive exercises for more examples: http://gaia.cs.umass.edu/kurose_ross/interactive/

Application Layer 2-26

26

HTTP request message: general format



Application Layer 2-27

27

Uploading form input

POST method:

- web page often includes form input
- input is uploaded to server in entity body

URL method:

- uses GET method
- input is uploaded in URL field of request line:

www.somesite.com/animalsearch?monkey&banana

Application Layer 2-28

28

Method types

HTTP/1.0:

- GET
- POST
- HEAD
 - asks server to leave requested object out of response

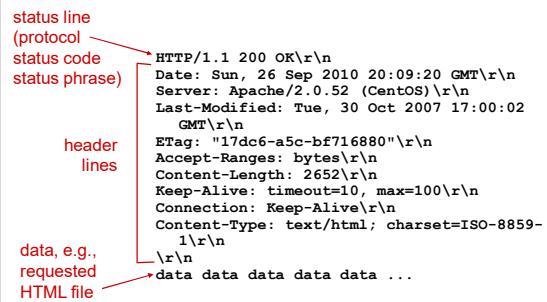
HTTP/1.1:

- GET, POST, HEAD
- PUT
 - uploads file in entity body to path specified in URL field
- DELETE
 - deletes file specified in the URL field

Application Layer 2-29

29

HTTP response message



* Check out the online interactive exercises for more examples: http://gaia.cs.umass.edu/kurose_ross/interactive/

Application Layer 2-30

30

HTTP response status codes

- status code appears in 1st line in server-to-client response message.
- some sample codes:
 - 200 OK**
 - request succeeded, requested object later in this msg
 - 301 Moved Permanently**
 - requested object moved, new location specified later in this msg (Location:)
 - 400 Bad Request**
 - request msg not understood by server
 - 404 Not Found**
 - requested document not found on this server
 - 505 HTTP Version Not Supported**

Application Layer 2-31

31

Trying out HTTP (client side) for yourself

1. Telnet to your favorite Web server:

```
telnet gaia.cs.umass.edu 80
```

opens TCP connection to port 80 (default HTTP server port) at gaia.cs.umass.edu. anything typed in will be sent to port 80 at gaia.cs.umass.edu

2. type in a GET HTTP request:

```
GET /kurose_ross/interactive/index.php HTTP/1.1
Host: gaia.cs.umass.edu
```

by typing this in (hit carriage return twice), you send this minimal (but complete) GET request to HTTP server

3. look at response message sent by HTTP server!
(or use Wireshark to look at captured HTTP request/response)

Application Layer 2-32

32

User-server state: cookies

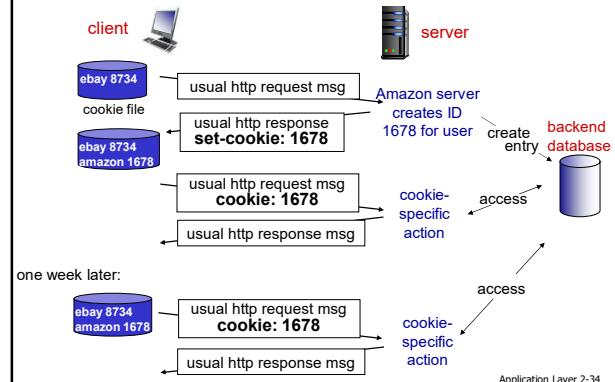
many Web sites use cookies
four components:

- 1) cookie header line of HTTP response message
 - 2) cookie header line in next HTTP request message
 - 3) cookie file kept on user's host, managed by user's browser
 - 4) back-end database at Web site
- example:
- Susan always access Internet from PC
 - visits specific e-commerce site for first time
 - when initial HTTP requests arrives at site, site creates:
 - unique ID
 - entry in backend database for ID

Application Layer 2-33

33

Cookies: keeping "state" (cont.)



Application Layer 2-34

34

Cookies (continued)

what cookies can be used for:

- authorization
- shopping carts
- recommendations
- user session state (Web e-mail)

cookies and privacy: aside

- cookies permit sites to learn a lot about you
- you may supply name and e-mail to sites

how to keep "state":

- protocol endpoints: maintain state at sender/receiver over multiple transactions
- cookies: http messages carry state

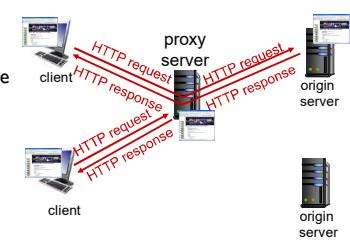
Application Layer 2-35

35

Web caches (proxy server)

goal: satisfy client request without involving origin server

- user sets browser: Web accesses via cache
- browser sends all HTTP requests to cache
 - object in cache: cache returns object
 - else cache requests object from origin server, then returns object to client



Application Layer 2-36

36

why Web cache

More about Web caching

- cache acts as both client and server
 - server for original requesting client
 - client to origin server
- typically cache is installed by ISP (university, company, residential ISP)
- why Web caching?
 - reduce response time for client request
 - reduce traffic on an institution's access link
 - Internet dense with caches: enables "poor" content providers to effectively deliver content (so too does P2P file sharing)

Application Layer 2-37

37

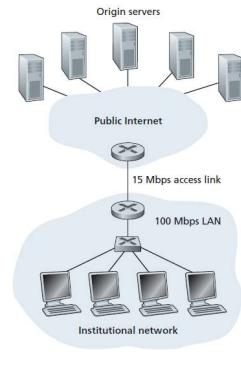
Caching example:

assumptions:

- avg object size: 1M bits
- avg request rate from browsers to origin servers: 15/sec
- RTT from institutional router to any origin server: 2 sec
- LAN rate: 100 Mbps
- access link rate: 15 Mbps

consequences:

- LAN utilization: 15% *problem!*
- access link utilization = 100% *100%*
- total delay = Internet delay + access delay + LAN delay
 $= 2 \text{ sec} + \text{minutes} + \text{usecs}$



Application Layer 2-38

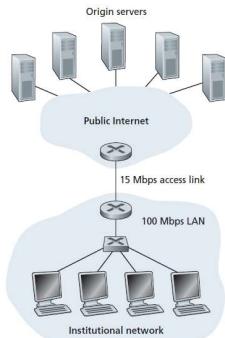
38

minutes : for long

Caching example:

consequences:

- LAN utilization: 15% (Traffic Intensity)
 $(15 \text{ requests/sec}) * (1 \text{ Mbits/request}) / (100 \text{ Mbps}) = 0.15$
- access link utilization = 100% (Traffic Intensity)
 $(15 \text{ requests/sec}) * (1 \text{ Mbits/request}) / (15 \text{ Mbps}) = 1$
- total delay = Internet delay + access delay + LAN delay
 $= 2 \text{ sec} + \text{minutes} + \text{usecs}$



Application Layer 2-39

39

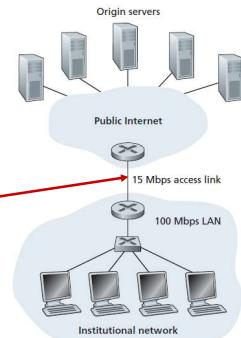
Caching example:

assumptions:

- avg object size: 1M bits
- avg request rate from browsers to origin servers: 15/sec
- RTT from institutional router to any origin server: 2 sec
- LAN rate: 100 Mbps
- access link rate: 100 Mbps

consequences:

- LAN utilization: 15%
- access link utilization = 15%
- total delay = Internet delay + access delay + LAN delay
 $= 2 \text{ sec} + \text{usecs} + \text{usecs}$



Application Layer 2-40

40

microseconds

Caching example: install local cache

assumptions:

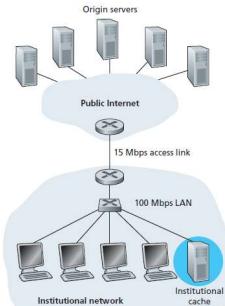
- avg object size: 1M bits
- avg request rate from browsers to origin servers: 15/sec
- RTT from institutional router to any origin server: 2 sec
- access link rate: 15 Mbps

consequences:

- LAN utilization: 15%
- access link utilization = 100%
- total delay = Internet delay + LAN delay
 $= 2 \text{ sec} + \text{usecs}$

How to compute link utilization, delay?

Cost: web cache (cheap!)



Application Layer 2-41

41

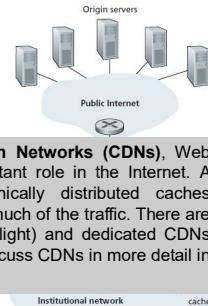
Caching example: install local cache

Calculating access link utilization, delay with cache:

- suppose cache hit rate is 0.4
 - 40% requests satisfied at cache, 60% requests satisfied at origin
- access link utilization:

Through the use of **Content Distribution Networks (CDNs)**, Web caches are increasingly playing an important role in the Internet. A CDN company installs many geographically distributed caches throughout the Internet, thereby localizing much of the traffic. There are shared CDNs (such as Akamai and Limelight) and dedicated CDNs (such as Google and Microsoft). We will discuss CDNs in more detail in Chapter 7.

less than 15 Mbps link (and cheaper too!)



Application Layer 2-42

42

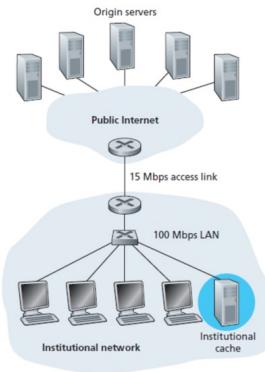
Why CDN's are getting so popular?

Caching example: install local cache

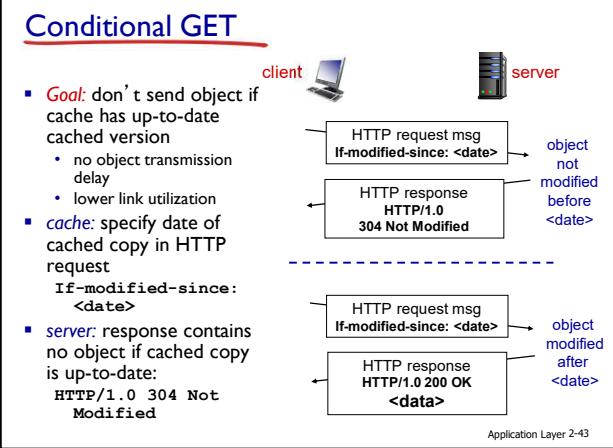
Calculating access link utilization, delay with cache:

- suppose cache hit rate is 0.4
 - 40% requests satisfied at cache, 60% requests satisfied at origin
- access link utilization:
 - 60% of requests use access link
 - Traffic Intensity at access link: 0.6
 - Utilization: 60%
- total delay:
 - $= 0.6 * (\text{delay from origin servers}) + 0.4 * (\text{delay when satisfied at cache})$
 - $= 0.6 (2.01) + 0.4 (\sim \text{msecs}) = \sim 1.2 \text{ secs}$
 - less than with 154 Mbps link (and cheaper too!)

Why CDN's are getting so popular?



Application Layer 2-42



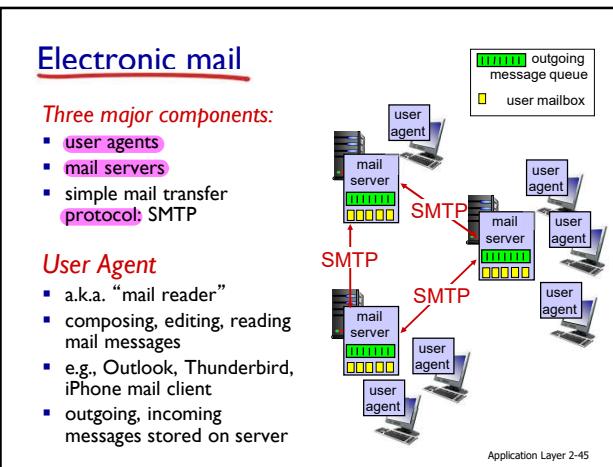
43

Chapter 2: outline

- | | |
|--|---|
| 2.1 principles of network applications | 2.5 P2P applications |
| 2.2 Web and HTTP | 2.6 video streaming and content distribution networks |
| 2.3 electronic mail | 2.7 socket programming with UDP and TCP |
| 2.4 DNS | |

Application Layer 2-44

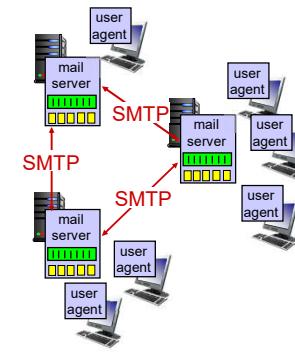
44



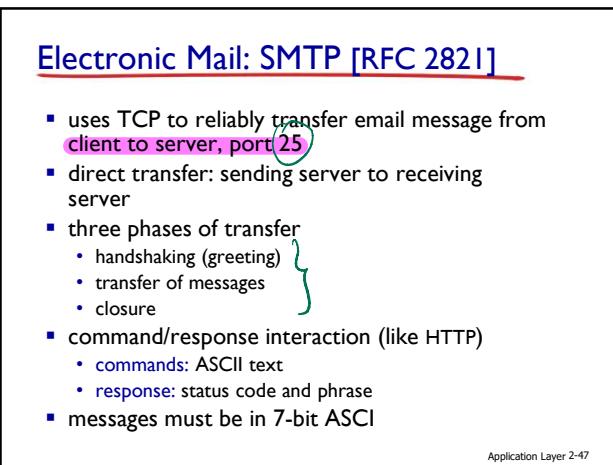
45

Electronic mail: mail servers

- mail servers:**
- mailbox contains incoming messages for user
 - message queue of outgoing (to be sent) mail messages
 - SMTP protocol between mail servers to send email messages
 - client: sending mail server
 - "server": receiving mail server



46

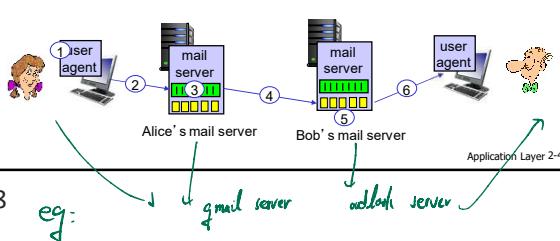


Application Layer 2-47

47

Scenario: Alice sends message to Bob

- Alice uses UA to compose message "to" `bob@some school.edu`
- Alice's UA sends message to her mail server; message placed in message queue
- client side of SMTP opens TCP connection with Bob's mail server
- 4) SMTP client sends Alice's message over the TCP connection
- Bob's mail server places the message in Bob's mailbox
- Bob invokes his user agent to read message



48

8

Sample SMTP interaction

```
S: 220 hamburger.edu
C: HELO crepes.fr
S: 250 Hello crepes.fr, pleased to meet you
C: MAIL FROM: <alice@crepes.fr>
S: 250 alice@crepes.fr... Sender ok
C: RCPT TO: <bob@hamburger.edu>
S: 250 bob@hamburger.edu ... Recipient ok
C: DATA
S: 354 Enter mail, end with "." on a line by itself
C: Do you like ketchup?
C: How about pickles?
C: .
S: 250 Message accepted for delivery
C: QUIT
S: 221 hamburger.edu closing connection
```

Application Layer 2-49

49

Try SMTP interaction for yourself:

- telnet servername 25
- see 220 reply from server
- enter HELO, MAIL FROM, RCPT TO, DATA, QUIT commands

above lets you send email without using email client (reader)

Application Layer 2-50

50

SMTP: final words

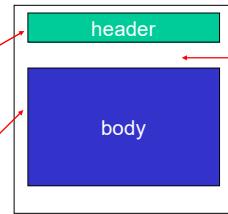
- SMTP uses persistent connections
 - SMTP requires message (header & body) to be in 7-bit ASCII
 - SMTP server uses CRLF, CRLF to determine end of message
- comparison with HTTP:*
- HTTP: pull
 - SMTP: push
 - both have ASCII command/response interaction, status codes
 - HTTP: each object encapsulated in its own response message
 - SMTP: multiple objects sent in multipart message

Application Layer 2-51

51

Mail message format

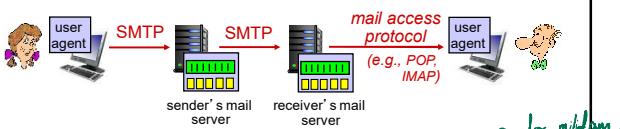
- SMTP: protocol for exchanging email messages
RFC 822: standard for text message format:
▪ header lines, e.g.,
 - To:
 - From:
 - Subject:*different from SMTP MAIL FROM, RCPT TO: commands!*
▪ Body: the “message”
 - ASCII characters only



Application Layer 2-52

52

Mail access protocols



- **SMTP:** delivery/storage to receiver's server
- mail access protocol: retrieval from server
 - **POP:** Post Office Protocol [RFC 1939]: authorization, download
 - **IMAP:** Internet Mail Access Protocol [RFC 1730]: more features, including manipulation of stored messages on server
 - **HTTP:** gmail, Hotmail, Yahoo! Mail, etc.

Application Layer 2-53

53

POP3 protocol

authorization phase

- client commands:
 - user: declare username
 - pass: password
- server responses
 - +OK
 - -ERR

transaction phase, client:

- list: list message numbers
- retr: retrieve message by number
- dele: delete
- quit

```
S: +OK POP3 server ready
C: user bob
S: +OK
C: pass hungry
S: +OK user successfully logged on
C: list
S: 1 498
S: 2 912
S: .
C: retr 1
S: <message 1 contents>
S: .
C: dele 1
C: retr 2
S: <message 1 contents>
S: .
C: dele 2
C: quit
S: +OK POP3 server signing off
```

Application Layer 2-54

54

2.4:

Chapter 2: outline

- | | |
|--|---|
| 2.1 principles of network applications | 2.5 P2P applications |
| 2.2 Web and HTTP | 2.6 video streaming and content distribution networks |
| 2.3 electronic mail | 2.7 socket programming with UDP and TCP |
| 2.4 DNS | |

Application Layer 2-56

POP3 (more) and IMAP

more about POP3

- previous example uses POP3 “download and delete” mode
 - Bob cannot re-read e-mail if he changes client
- POP3 “download-and-keep”: copies of messages on different clients
- POP3 is stateless across sessions

IMAP

- keeps all messages in one place: at server
- allows user to organize messages in folders
- keeps user state across sessions:
 - names of folders and mappings between message IDs and folder name

Application Layer 2-55

55

56

Based on UDP protocol because it's fast

DNS: domain name system

people: many identifiers:

- SSN, name, passport #

Internet hosts, routers:

- IP address (32 bit) - used for addressing datagrams
- “name”, e.g., www.yahoo.com - used by humans

Q: how to map between IP address and name, and vice versa ?

Domain Name System:

- *distributed database* implemented in hierarchy of many *name servers*
- *application-layer protocol:* hosts, name servers communicate to *resolve* names (address/name translation)
 - note: core Internet function, implemented as application-layer protocol
 - complexity at network’s “edge”

Application Layer 2-57

57

DNS: services, structure

DNS services

- hostname to IP address translation
- host aliasing
 - canonical, alias names
- mail server aliasing
- load distribution
 - replicated Web servers: many IP addresses correspond to one name

why not centralize DNS?

- single point of failure
- traffic volume
- distant centralized database
- maintenance

A: doesn't scale!

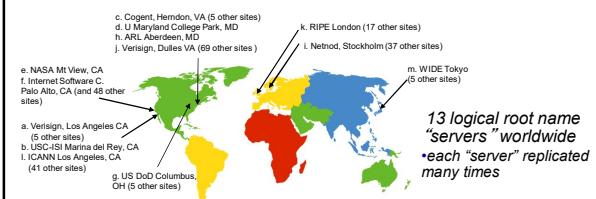
Application Layer 2-58

58

→ knows all translation

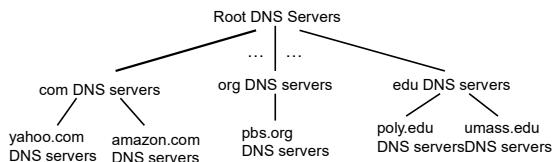
DNS: root name servers

- contacted by local name server that can not resolve name
- root name server:
 - contacts authoritative name server if name mapping not known
 - gets mapping
 - returns mapping to local name server



Application Layer 2-60

DNS: a distributed, hierarchical database



client wants IP for www.amazon.com; 1st approximation:

- client queries root server to find com DNS server
- client queries .com DNS server to get amazon.com DNS server
- client queries amazon.com DNS server to get IP address for www.amazon.com

Application Layer 2-59

59

60

TLD, authoritative servers

top-level domain (TLD) servers:

- responsible for com, org, net, edu, aero, jobs, museums, and all top-level country domains, e.g. uk, fr, ca, jp
- Network Solutions maintains servers for .com TLD
- Educause for .edu TLD

authoritative DNS servers:

- organization's own DNS server(s), providing authoritative hostname to IP mappings for organization's named hosts
- can be maintained by organization or service provider

Application Layer 2-61

61

Local DNS name server

- does not strictly belong to hierarchy
- each ISP (residential ISP, company, university) has one
 - also called "default name server"
- when host makes DNS query, query is sent to its local DNS server
 - has local cache of recent name-to-address translation pairs (but may be out of date!)
 - acts as proxy, forwards query into hierarchy

Application Layer 2-62

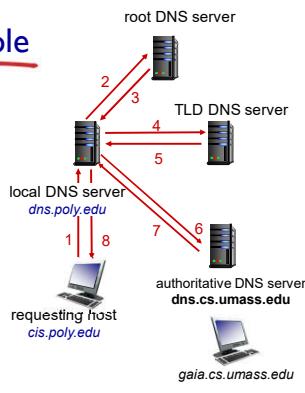
62

DNS name resolution example

- host at cis.poly.edu wants IP address for gaia.cs.umass.edu

iterated query:

- contacted server replies with name of server to contact
- "I don't know this name, but ask this server"



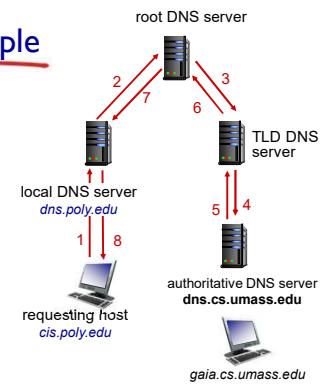
Application Layer 2-63

63

DNS name resolution example

recursive query:

- puts burden of name resolution on contacted name server
- heavy load at upper levels of hierarchy?



Application Layer 2-64

64

DNS: caching, updating records

- once (any) name server learns mapping, it *caches* mapping
 - cache entries timeout (disappear) after some time (TTL)
 - TLD servers typically cached in local name servers
 - thus root name servers not often visited
- cached entries may be *out-of-date* (best effort name-to-address translation!)
 - if name host changes IP address, may not be known Internet-wide until all TTLs expire
- update/notify mechanisms proposed IETF standard
 - RFC 2136

Application Layer 2-65

65

DNS records

DNS: distributed database storing resource records (RR)

RR format: (name, value, type, ttl)

type=A

- name is hostname
- value is IP address

type=NS

- name is domain (e.g., foo.com)
- value is hostname of authoritative name server for this domain

type=CNAME

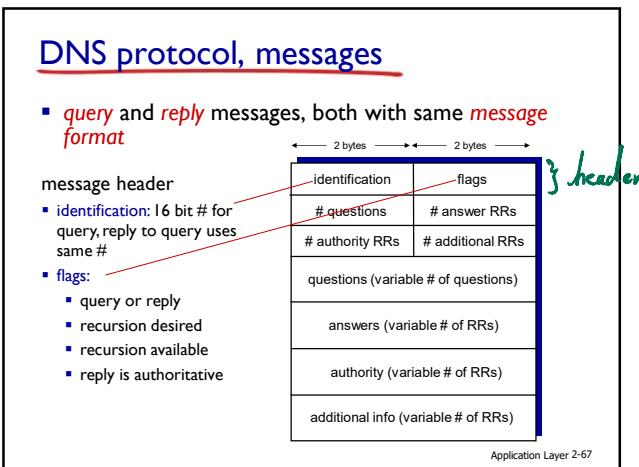
- name is alias name for some "canonical" (the real) name
- www.ibm.com is really severtheast.backup2.ibm.com
- value is canonical name

type=MX

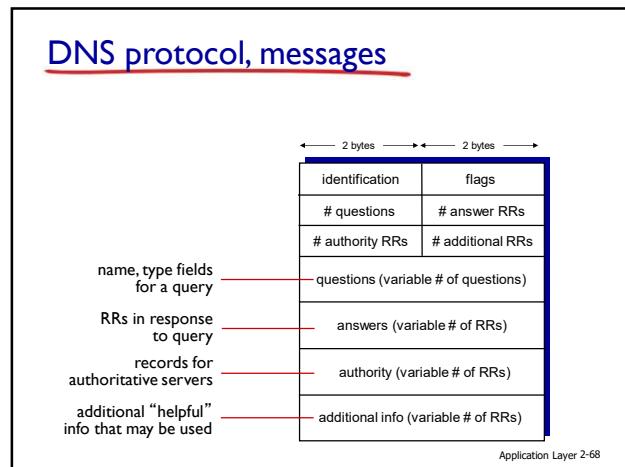
- value is name of mailserver associated with name

Application Layer 2-66

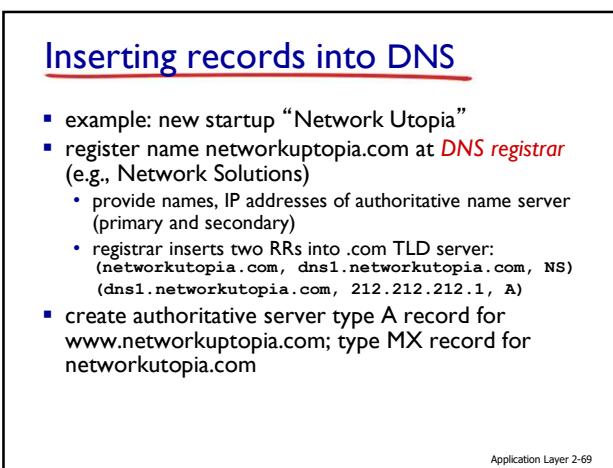
66



67

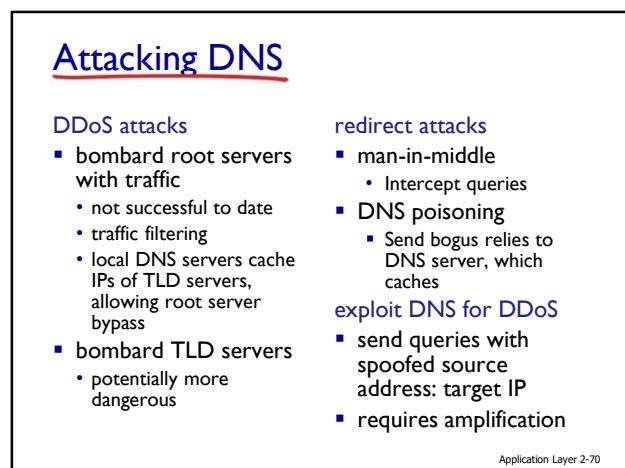


68

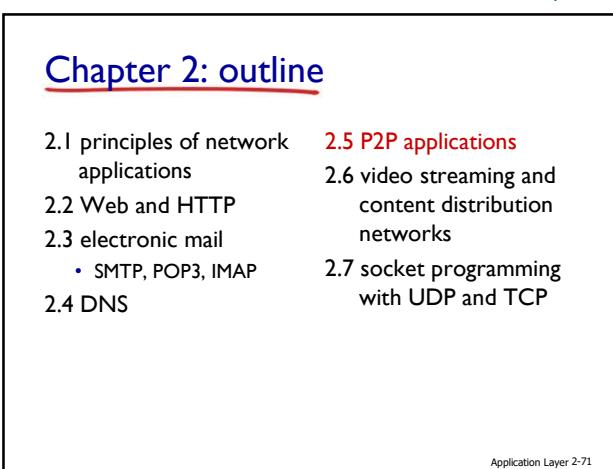


69

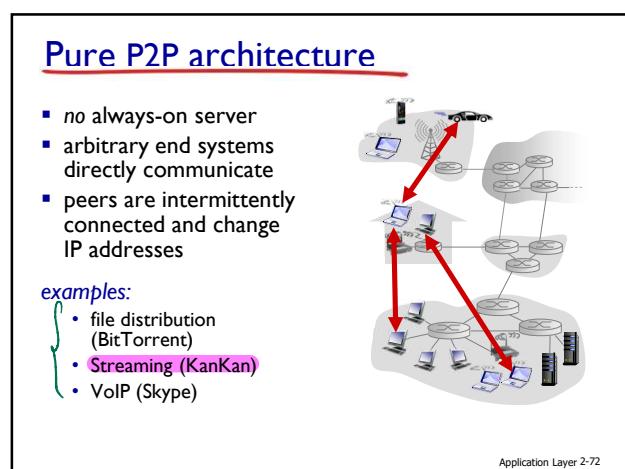
P2P:



70 example of P2P: Gmail, uTorrent, BitTorrent, skype

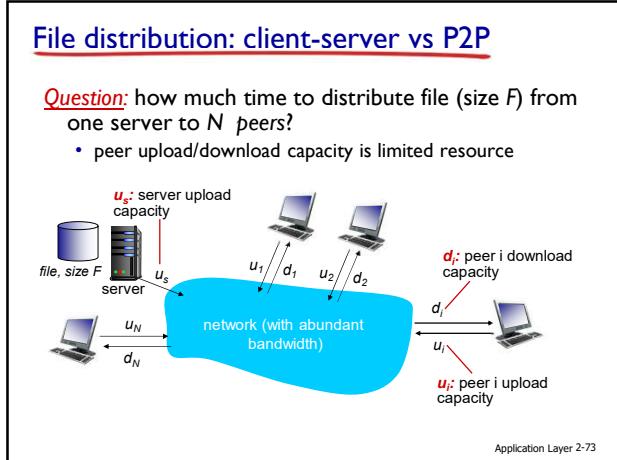


71



72 Why P2P: No third person read this information.
And only this protocol use upload speed

File distribution:



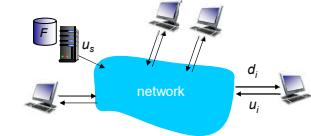
73

↗ file upload speed: u_s

File distribution time: client-server

- server transmission:** must sequentially send (upload) N file copies:

- time to send one copy: F/u_s
- time to send N copies: NF/u_s



- client:** each client must download file copy

- d_{\min} = min client download rate
- min client download time F/d_{\min}

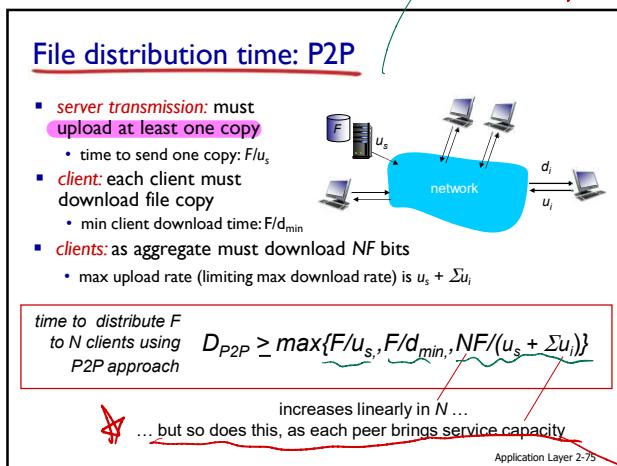
$$\text{time to distribute } F \text{ to } N \text{ clients using client-server approach} \quad D_{c-s} \geq \max\{NF/u_s, F/d_{\min}\}$$

increases linearly in N

Application Layer 2-74

74

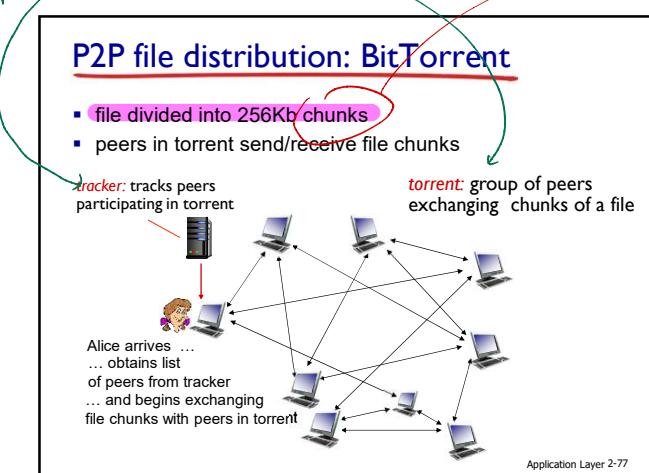
all of us upload file to everyone
The radically, P2P is faster than server-client.



75

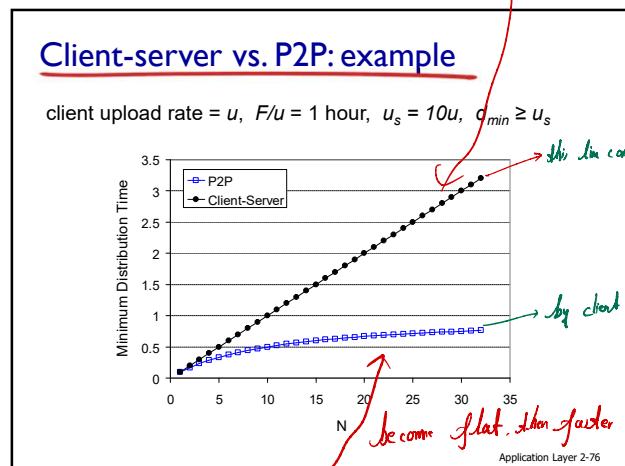
Tracker: users who have file
Seeds: chunks users have

Torrent: or seeds



77

We need a server for tracking purpose
not data purpose



76

P2P file distribution: BitTorrent

- peer joining torrent:**

- has no chunks, but will accumulate them over time from other peers

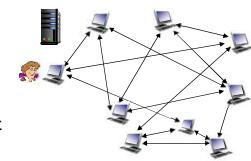
- registers with tracker to get list of peers, connects to subset of peers ("neighbors")

- while downloading, peer uploads chunks to other peers

- peer may change peers with whom it exchanges chunks

- churn:** peers may come and go

- once peer has entire file, it may (selfishly) leave or (altruistically) remain in torrent



Application Layer 2-78

78

BitTorrent: requesting, sending file chunks

requesting chunks:

- at any given time, different peers have different subsets of file chunks
- periodically, Alice asks each peer for list of chunks that they have
- Alice requests missing chunks from peers, rarest first

sending chunks: tit-for-tat

- Alice sends chunks to those four peers currently sending her chunks at highest rate
 - other peers are choked by Alice (do not receive chunks from her)
 - re-evaluate top 4 every 10 secs
- every 30 secs: randomly select another peer, starts sending chunks
 - "optimistically unchoke" this peer
 - newly chosen peer may join top 4

Application Layer 2-79

79

2.6 CDNs

Chapter 2: outline

- 2.1 principles of network applications
- 2.2 Web and HTTP
- 2.3 electronic mail
 - SMTP, POP3, IMAP
- 2.4 DNS

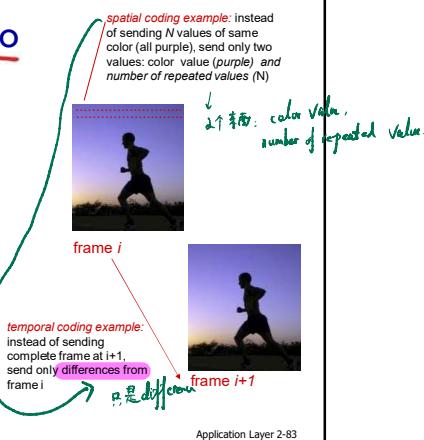
- 2.5 P2P applications
- 2.6 video streaming and content distribution networks (CDNs)**
- 2.7 socket programming with UDP and TCP

Application Layer 2-81

81

Multimedia: video

- video:** sequence of images displayed at constant rate
 - e.g., 24 images/sec
- digital image:** array of pixels
 - each pixel represented by bits
- coding: use redundancy **within** and **between** images to decrease # bits used to encode image
 - spatial (within image)
 - temporal (from one image to next)

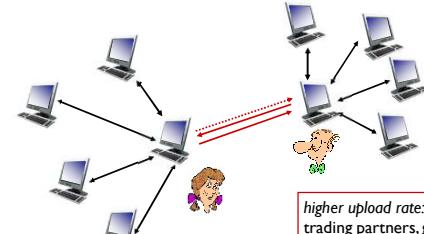


Application Layer 2-83

83

BitTorrent: tit-for-tat

- (1) Alice "optimistically unchoke" Bob
- (2) Alice becomes one of Bob's top-four providers; Bob reciprocates
- (3) Bob becomes one of Alice's top-four providers



Application Layer 2-80

80

Video Streaming and CDNs: context

- video traffic: major consumer of Internet bandwidth
 - Netflix, YouTube: 37%, 16% of downstream residential ISP traffic
 - ~1B YouTube users, ~75M Netflix users
- challenge: scale - how to reach ~1B users?
 - single mega-video server won't work (why?)
- challenge: heterogeneity
 - different users have different capabilities (e.g., wired versus mobile; bandwidth rich versus bandwidth poor)
- solution:** distributed, application-level infrastructure

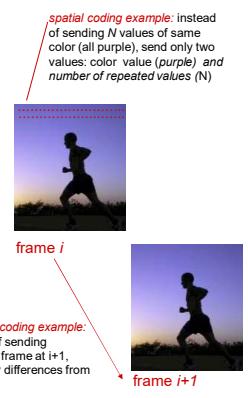


Application Layer 2-82

82

Multimedia: video

- CBR: (constant bit rate):** video encoding rate fixed
- VBR: (variable bit rate):** video encoding rate changes as amount of spatial, temporal coding changes
- examples:
 - MPEG 1 (CD-ROM) 1.5 Mbps
 - MPEG2 (DVD) 3-6 Mbps
 - MPEG4 (often used in Internet, < 1 Mbps)



Application Layer 2-84

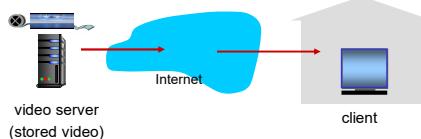
84

Audio and video is separated

Since you can not record video with audio

Streaming stored video:

simple scenario:



Application Layer 2-85

85

Streaming multimedia: DASH

▪ **DASH: Dynamic, Adaptive Streaming over HTTP**

▪ **server:**

- divides video file into multiple chunks
- each chunk stored, encoded at different rates
- **manifest file:** provides URLs for different chunks

▪ **client:**

- periodically measures server-to-client bandwidth
- consulting manifest, requests one chunk at a time
 - chooses maximum coding rate sustainable given current bandwidth
- can choose different coding rates at different points in time (depending on available bandwidth at time)

Application Layer 2-86

86

one video could be stored multiple version at different rates for different users, who has different bandwidth, or resolution

Streaming multimedia: DASH

▪ **DASH: Dynamic, Adaptive Streaming over HTTP**

▪ “intelligence” at client: client determines

- **when** to request chunk (so that **buffer starvation**, or overflow does not occur)
 - **what encoding rate** to request (**higher quality when more bandwidth available**)
 - **where** to request chunk (can request from URL server that is “close” to client or has high available bandwidth)
- the complexity is pushed onto edge, so server doesn't have to deal with it.*
- so fucking smart to choose which URL of chunk*

Application Layer 2-87

87

Content distribution networks

▪ **challenge:** how to stream content (selected from millions of videos) to hundreds of thousands of simultaneous users?

▪ **option 1:** single, large “mega-server” *不行*

- single point of failure
- point of network congestion
- long path to distant clients
- multiple copies of video sent over outgoing link

....quite simply: this solution **doesn't scale**

Application Layer 2-88

88

Content distribution networks

▪ **challenge:** how to stream content (selected from millions of videos) to hundreds of thousands of

Akamai: 100,000+ servers in 1000+ clusters

in 1000+ networks in 70+ countries serving

trillions of requests a day.

- **enter deep:** push CDN servers deep into many access networks
 - close to users
 - used by Akamai, 1700 locations

- **bring home:** smaller number (10's) of larger clusters in POPs (point of presence) near (but not within) access networks
 - used by Limelight

Application Layer 2-89

89

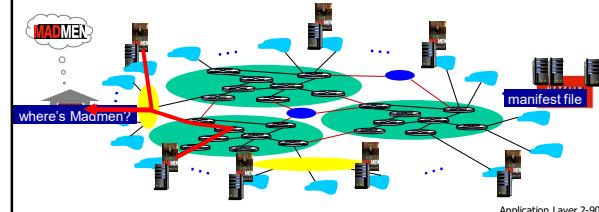
Content Distribution Networks (CDNs)

▪ **CDN:** stores copies of content at CDN nodes

- e.g. Netflix stores copies of MadMen

▪ **subscriber requests content from CDN**

- directed to nearby copy, retrieves content
- may choose different copy if network path congested



Application Layer 2-90

90

Content distribution networks

- **challenge:** how to stream content (selected from millions of videos) to hundreds of thousands of simultaneous users?

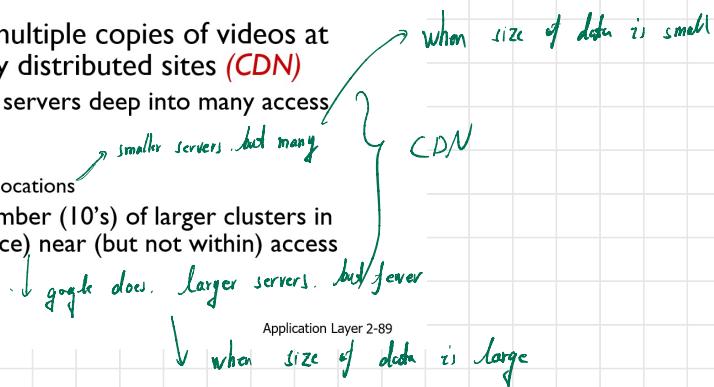
- **option 2:** store/serve multiple copies of videos at multiple geographically distributed sites (**CDN**)

- **enter deep:** push CDN servers deep into many access networks

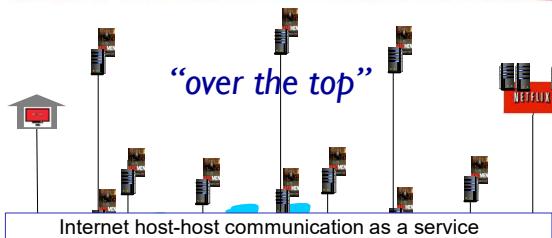
- close to users
 - used by Akamai, 1700 locations

- **bring home:** smaller number (10's) of larger clusters in POPs (point of presence) near (but not within) access networks

- used by Limelight



Content Distribution Networks (CDNs)



OTT challenges: coping with a congested Internet

- from which CDN node to retrieve content?
- viewer behavior in presence of congestion?
- what content to place in which CDN node?

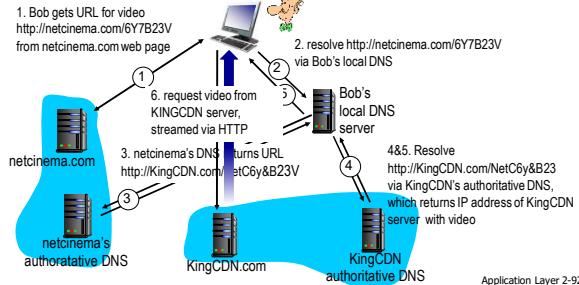
more .. in chapter 7

91

CDN content access: a closer look

Bob (client) requests video <http://netcinema.com/6Y7B23V>

- video stored in CDN at <http://KingCDN.com/NetC6y&B23V>



92

2.7:

Chapter 2: outline

2.1 principles of network applications

2.2 Web and HTTP

2.3 electronic mail

- SMTP, POP3, IMAP

2.4 DNS

2.5 P2P applications

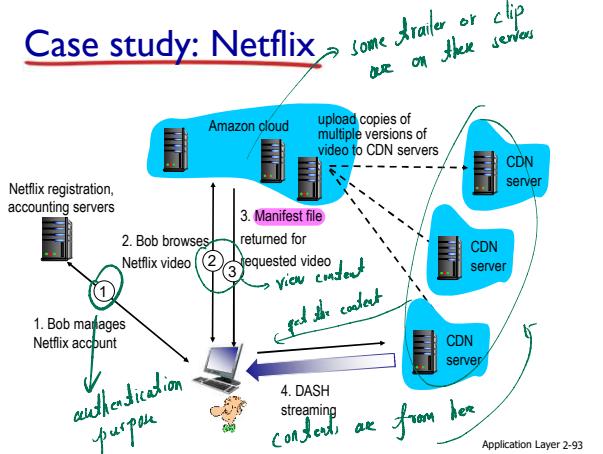
2.6 video streaming and content distribution networks

2.7 socket programming with UDP and TCP

Application Layer 2-94

94

Case study: Netflix

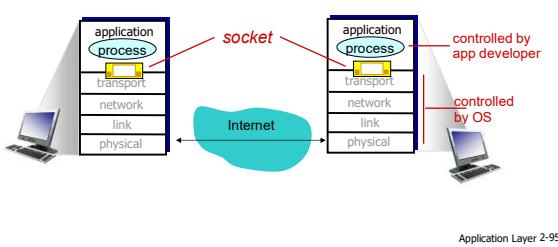


93

Socket programming

goal: learn how to build client/server applications that communicate using sockets

socket: door between application process and end-end-transport protocol



95

Socket programming

Two socket types for two transport services:

- **UDP:** unreliable datagram
- **TCP:** reliable, byte stream-oriented

Application Example:

1. client reads a line of characters (data) from its keyboard and sends data to server
2. server receives the data and converts characters to uppercase
3. server sends modified data to client
4. client receives modified data and displays line on its screen

Application Layer 2-96

96

16

Socket programming with UDP

- UDP:** no “connection” between client & server
- no handshaking before sending data
 - sender explicitly attaches IP destination address and port # to each packet
 - receiver extracts sender IP address and port# from received packet

UDP: transmitted data may be lost or received out-of-order

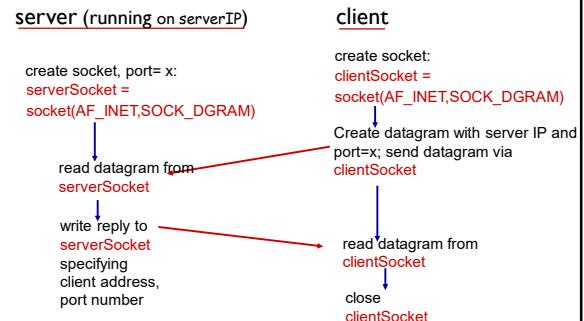
Application viewpoint:

- UDP provides *unreliable* transfer of groups of bytes (“datagrams”) between client and server

Application Layer 2-97

97

Client/server socket interaction: UDP



Application 2-98

98

Example app: UDP client

```

Python UDPClient
include Python's socket library
from socket import *
serverName = 'hostname'
serverPort = 12000
create UDP socket for server
clientSocket = socket(AF_INET,
                     SOCK_DGRAM)
get user keyboard input
message = raw_input('Input lowercase sentence:')
Attach server name, port to message; send into socket
clientSocket.sendto(message.encode(),
                    (serverName, serverPort))
read reply characters from socket into string
modifiedMessage, serverAddress =
clientSocket.recvfrom(2048)
print out received string and close socket
print modifiedMessage.decode()
clientSocket.close()
Application Layer 2-99

```

99

Example app: UDP server

```

Python UDPSServer
from socket import *
serverPort = 12000
create UDP socket
serverSocket = socket(AF_INET, SOCK_DGRAM)
bind socket to local port number 12000
serverSocket.bind(("0.0.0.0", serverPort))
print ("The server is ready to receive")
loop forever
while True:
Read from UDP socket into message, getting client's address (client IP and port)
message, clientAddress = serverSocket.recvfrom(2048)
modifiedMessage = message.decode().upper()
send upper case string back to this client
serverSocket.sendto(modifiedMessage.encode(),
                   clientAddress)
Application Layer 2-100

```

100

Socket programming with TCP

client must contact server

- server process must first be running
- server must have created socket (door) that welcomes client’s contact

client contacts server by:

- Creating TCP socket, specifying IP address, port number of server process
- when client creates socket:** client TCP establishes connection to server TCP

- when contacted by client, **server TCP creates new socket** for server process to communicate with that particular client
 - allows server to talk with multiple clients
 - source port numbers used to distinguish clients (more in Chap 3)

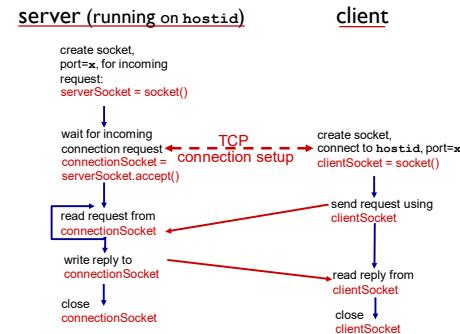
application viewpoint:

TCP provides reliable, in-order byte-stream transfer (“pipe”) between client and server

Application Layer 2-101

101

Client/server socket interaction: TCP



Application Layer 2-102

102

Example app:TCP client

Python TCPClient

```

from socket import *
serverName = 'servername'
serverPort = 12000
create TCP socket for
server, remote port 12000
clientSocket = socket(AF_INET,SOCK_STREAM)
clientSocket.connect((serverName,serverPort))
sentence = raw_input('Input lowercase sentence:')
clientSocket.send(sentence.encode())
modifiedSentence = clientSocket.recv(1024)
print ('From Server:', modifiedSentence.decode())
clientSocket.close()
No need to attach server
name, port
  
```

Application Layer 2-103

103

Example app:TCP server

Python TCPServer

```

from socket import *
serverPort = 12000
create TCP welcoming
socket
serverSocket = socket(AF_INET,SOCK_STREAM)
serverSocket.bind(('',serverPort))
serverSocket.listen(1)
print 'The server is ready to receive'
loop forever
server waits on accept()
for incoming requests, new
socket created on return
connectionSocket, addr = serverSocket.accept()
read bytes from socket (but
not address as in UDP)
sentence = connectionSocket.recv(1024).decode()
capitalizedSentence = sentence.upper()
close connection to this
client (but not welcoming
socket)
connectionSocket.send(capitalizedSentence.
encode())
connectionSocket.close()
  
```

Application Layer 2-104

104

Chapter 2: summary

our study of network apps now complete!

- application architectures
 - client-server
 - P2P
- application service requirements:
 - reliability, bandwidth, delay
- Internet transport service model
 - connection-oriented, reliable: TCP
 - unreliable, datagrams: UDP
- specific protocols:
 - HTTP
 - SMTP, POP, IMAP
 - DNS
 - P2P: BitTorrent
 - video streaming, CDNs
 - socket programming: TCP, UDP sockets

Application Layer 2-105

105

Chapter 2: summary

most importantly: learned about protocols!

- typical request/reply message exchange:
 - client requests info or service
 - server responds with data, status code
 - message formats:
 - **headers:** fields giving info about data
 - **data:** info(payload) being communicated
- important themes:**
- control vs. messages
 - in-band, out-of-band
 - centralized vs. decentralized
 - stateless vs. stateful
 - reliable vs. unreliable message transfer
 - “complexity at network edge”

Application Layer 2-106

106

Two protocol services offered from T to A:

TCP and UDP: check the details on slides.

TCP and UDP protocol:

Internet transport protocols services

TCP service:

- **reliable transport** between sending and receiving process
- **Based on throttles**
- **flow control**: sender won't overwhelm receiver
- **congestion control**: throttle sender when network overloaded
- **does not provide**: timing, minimum throughput guarantee, security
- **connection-oriented**: setup required between client and server processes

UDP service:

- **unreliable data transfer** between sending and receiving process
- **does not provide**: reliability, flow control, congestion control, timing, throughput guarantee, security, or connection setup,

Q: why bother? Why is there a UDP?

Application Layer 2-14

14

Congestion: could be only control by sender, which need to be sensed by TCP since it would be worse if we add feedback feature.

Stream delivery offered by TCP: If send 1, 2, 3, 4 then receiver will receive 1,2,3,4 not 4,3,2,1 by sequence number feature.

UDP: Does not care about anything.

If UDP dropped packet, then it will do nothing.

It has no chop-up feature for example: send 1GB packet it wouldn't chop up, has no segmentation and reassemble feature.

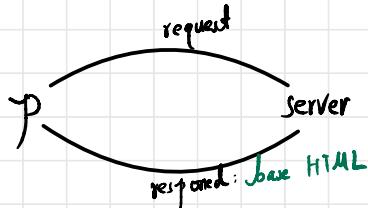
Simple, Fast

TCP: Not fast, but reliable, send the packet again when dropped.
Has segmentation and reassembly.

Section 2:HTML and Web

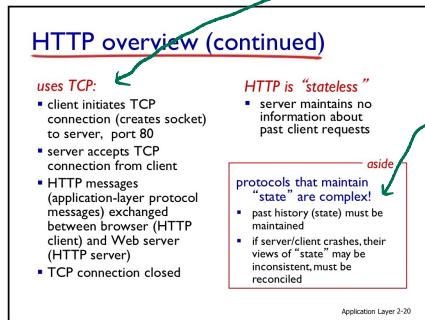
Browser only knows one language: HTML

When 1st time P request Server.



Base html contains objects: JS, Vedio....links, then P keep going request server for more objects.

TCP is stateful protocol: can save data



Why HTTP should be stateless: Too much information.

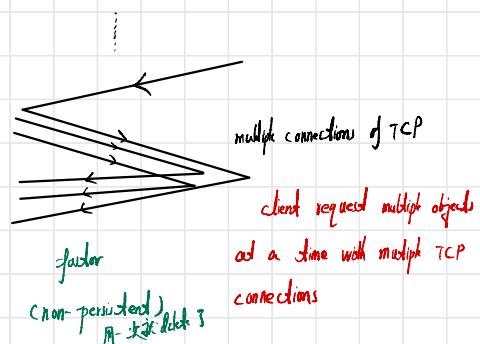
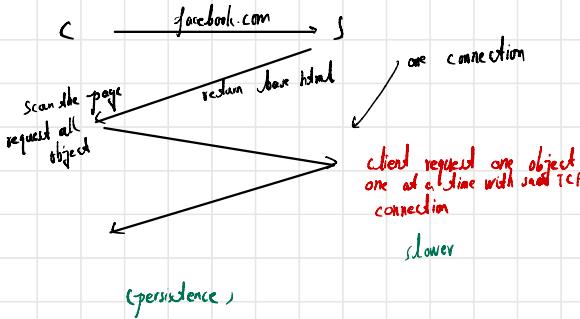
continue

TCP协议对应于传输层，而HTTP协议对应于应用层，从本质上来说，二者没有可比性。Http协议是建立在TCP协议基础之上的，当浏览器需要从服务器获取网页数据的时候，会发出一次Http请求。Http会通过TCP建立起一个到服务器的连接通道，当本次请求需要的数据完毕后，Http会立即将TCP连接断开，这个过程是很短的。所以Http连接是一种短连接，是一种无状态的连接。所谓的无状态，是指浏览器每次向服务器发起请求的时候，不是通过一个连接，而是每次都建立一个新的连接。如果是一个连接的话，服务器进程中就能保持住这个连接并且在内存中记住一些信息状态。而每次请求结束后，连接就关闭，相关的内容就释放了，所以记不住任何状态，成为无状态连接。

20

HTTP versions: 1.0, 1.1, 2.0, 3.0
most used
spdy
made by google
Quick
and do stand use you

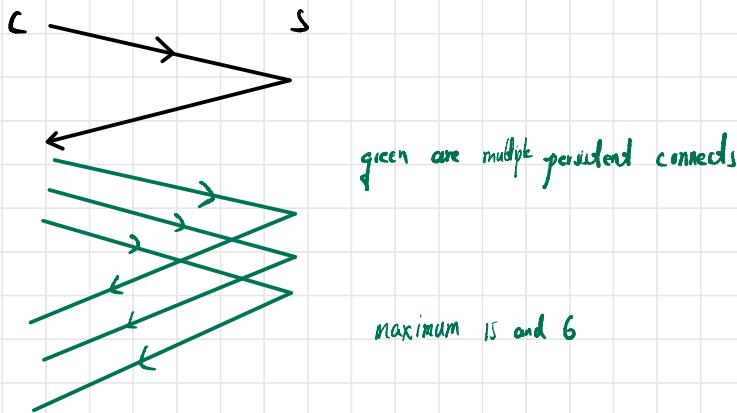
When browser scan the base page it has two ways to request objects from servers. Persistence and Non-persistence.



No one is better than other, since non-persistent need to maintain multiple TCP even it is fast.

⚠ Heavy, complex, need more sources

Hybrid mode:



Web Proxy

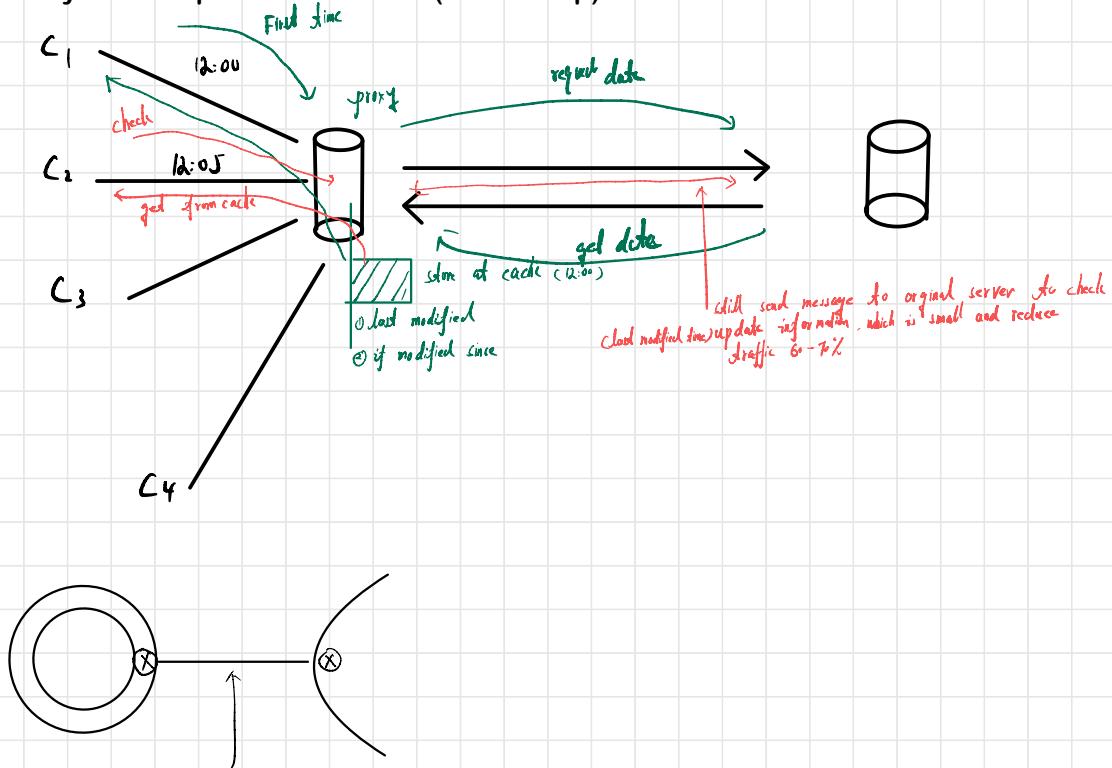
Gate way router is a multiplexers and Demultiplexer.

On site cache server: server records what request user made and respond users got, whenever the receive the same request then it doesn't go to original server, instead give you the data in cache server.

Why need cache server: faster, higher response time, reduce redundant of traffic network.

Session log site: where cache server is useless because everyone's data is different.

TTL: Every website provides time info (time stamp) of cache.



Access link: the connection between your network and to the public network

Local network is faster than public

- { Size is small
- Connection is good

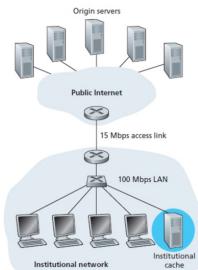
Why is not sustainable: More students, more traffic, more congestion, pay more
It is like black hole

Install local cache

Caching example: install local cache

Calculating access link utilization, delay with cache:

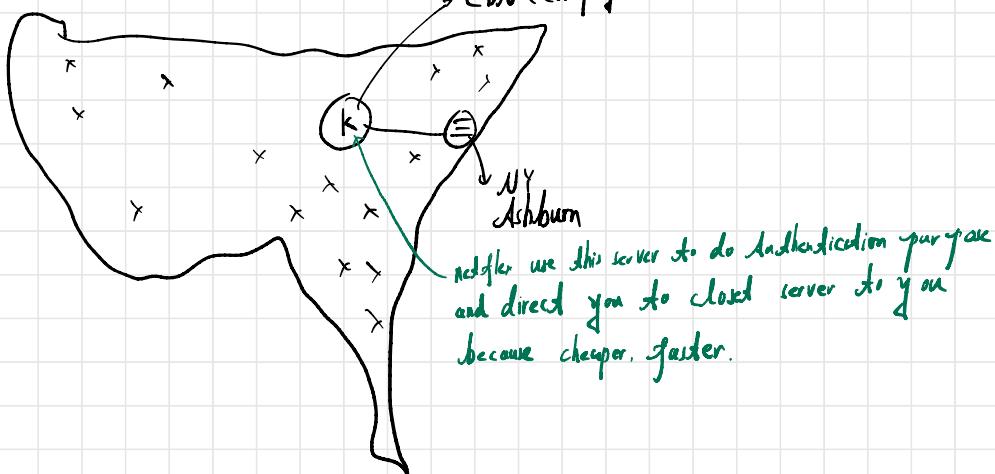
- suppose cache hit rate is 0.4
 - 40% requests satisfied at cache, 60% requests satisfied at origin
- access link utilization:
 - 60% of requests use access link
 - Traffic Intensity at access link: 0.6
 - Utilization: 60%
- total delay:
 - = 0.6 * (delay from origin servers) + 0.4 * (delay when satisfied at cache)
 - = 0.6 (2.01) + 0.4 (1 msec) = ~ 1.2 secs
 - less than with 154 Mbps link (and cheaper tool!)



Application Layer 2-42

Why CDN's are getting so popular?

CDN: content delivery network



DNS: Convert domain name to IP address.

IP address: unique identifier of a machine on any network.

02/03/2022

{ open public DNS:
Google DNS

Why DNS shouldn't be centralized: Fast response , preventing one down the whole system down.

DNS: is a edge server. It is a application layer protocol, which implements DNS. It is a core internet functionality but implemented as application protocol.

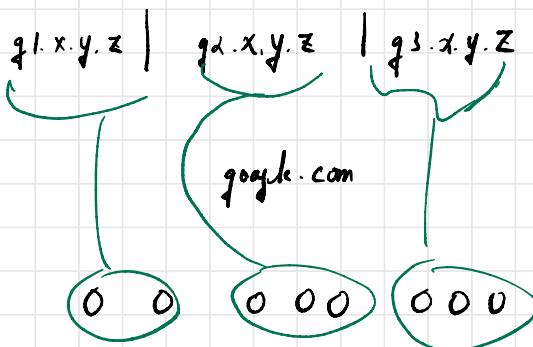
Edge layers: Application and transport.

Alias and Canonical name: If we want IP address then we need to know canonical name.

eg: wiki.org (easy to remember) d1.71.wiki.org (actual name) And convert cname to IP

DNS also convert alias to canonical name

Interesting phenomenon:



同一个domain name 他出现不同的 ip address

TLD:

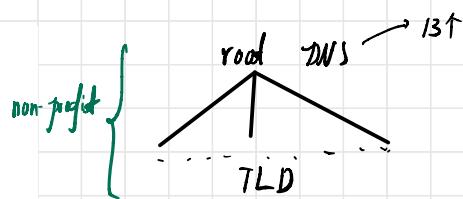
TLD, authoritative servers

top-level domain (TLD) servers:

- responsible for com, org, net, edu, aero, jobs, museums, and all top-level country domains, e.g.: uk, fr, ca, jp
- Network Solutions maintains servers for .com TLD
- Educause for .edu TLD

authoritative DNS servers:

- organization's own DNS server(s), providing authoritative hostname to IP mappings for organization's named hosts
- can be maintained by organization or service provider



Application Layer 2-61

61

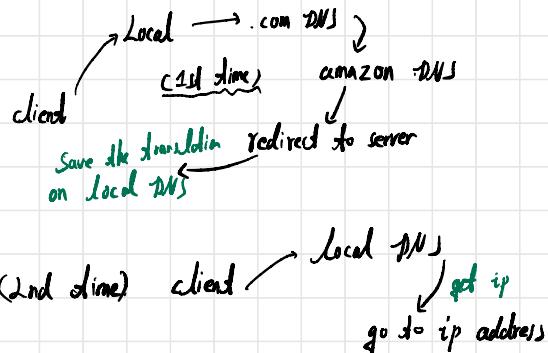
Local DNS:

Local DNS name server

- does not strictly belong to hierarchy
- each ISP (residential ISP, company, university) has one
 - also called "default name server"
- when host makes DNS query, query is sent to its local DNS server
 - has local cache of recent name-to-address translation pairs (but may be out of date!)
 - acts as proxy, forwards query into hierarchy

Application Layer 2-62

Cloud.amazon.com



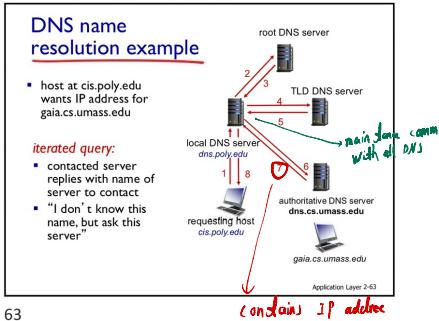
当然如果突然ip变了，不再是原来的了，那这个local的存的ip就用不了了，所以每个translation都有个 TTL. (Time To Live)
当过了TTL，它会自己discard。

除了那些非常常用的网站，不如不会自己refresh。要不然就等再cache一遍。

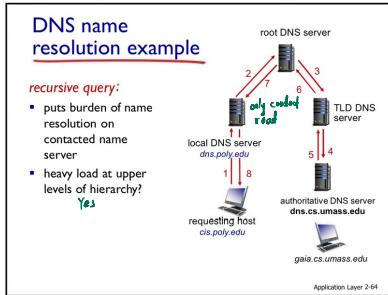
Authoritative dns balance the loads

DNS can convert ip to domain name as well.

Two types of resolutions:



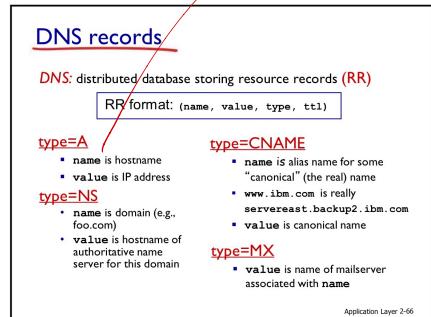
63



64

If the server is down, the authoritative dns will know first, and the local dns won't know until it refreshes or TTL expires.

DNS Records



Authoritative server gives Type A record and CName

TLD will give type NS and A

Root gives type A and NS

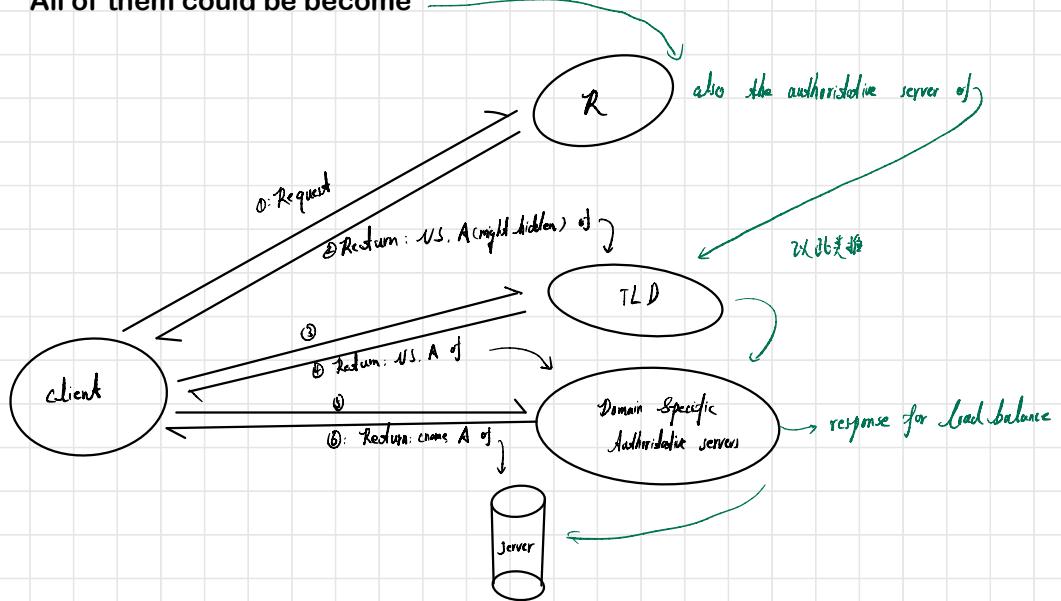
02/08/2022

Iterated and recursive, which is one is better.

Both has pro and cons, but Iterative is slightly better since no pressure on top servers.

Could Root, TLD, and authoritative server become authoritative servers?

All of them could be become



Build own website: no need for root to know. TLD is enough.

DNS use UDP to transport: Why?

The DNS has to be the fastest, and TCP will slow down DNS transmission.

DNS poisoning. External learning. Since DNS use UDP has no secure measurement.

Why we still have UDP and TCP?

UDT has all features of UCP but more secured and safe.

DNS will be the crucial content on EXAN and quiz

E-mail:

The protocol email uses: SMTP, POP3, IMAP.

Simple Mail Transport Protocol:

HTTP is a pool protocol and SMTP is a punch protocol

POP3: used for military.

02/10/2022

Most of notes will be on slides

Why P2P: No third person read this information.

And only this protocol use upload speed

P2P is not favored by those companies: Since they can not make money like Netflix

File distribution: on slides

02/15/2022

都在 PPt 上

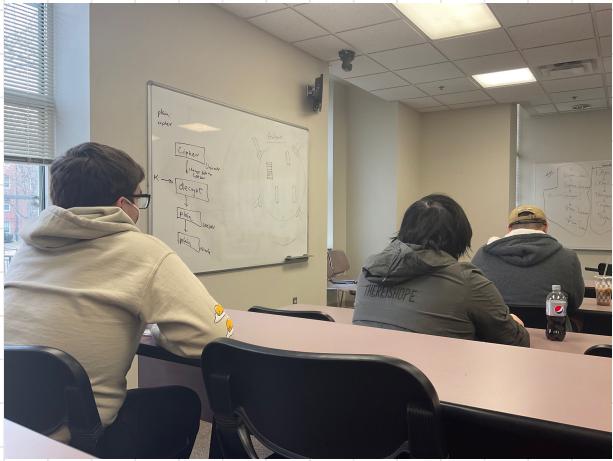
(~08:09)

Bob example

(~12:21) 2.6 end

preference 發音

(~19:09) DNS message on slides



Network OS: NOS

HAC: High Availability Cluster (-23:48)

(30:45~) project

(46: 43~)

UDP sever class

socket = new Datagram (9876)

incoming Pack .get ... ()

Quiz 2 summary:

What server will the url goes first?

Local dns is a cache which is not a server. Thus the correct answer is root dns.

Chapter 3

Chapter 3: Transport Layer

our goals:

- understand principles behind transport layer services:
 - multiplexing, demultiplexing
 - reliable data transfer
 - flow control
 - congestion control
- learn about Internet transport layer protocols:
 - UDP: connectionless transport
 - TCP: connection-oriented reliable transport
 - TCP congestion control

Transport Layer 3-1

1

Chapter 3 outline

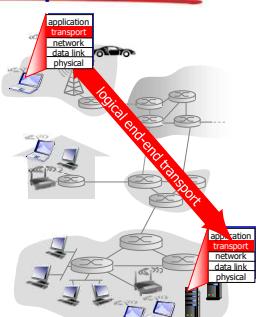
- 3.1 transport-layer services
- 3.2 multiplexing and demultiplexing
- 3.3 connectionless transport: UDP
- 3.4 principles of reliable data transfer
- 3.5 connection-oriented transport: TCP
 - segment structure
 - reliable data transfer
 - flow control
 - connection management
- 3.6 principles of congestion control
- 3.7 TCP congestion control

Transport Layer 3-2

2

Transport services and protocols

- provide *logical communication* between app processes running on different hosts
- transport protocols run in end systems
 - send side: breaks app messages into *segments*, passes to network layer
 - rcv side: reassembles segments into messages, passes to app layer
- more than one transport protocol available to apps
 - Internet: TCP and UDP



Transport Layer 3-3

3

Transport vs. network layer

- **network layer:** logical communication between hosts
- **transport layer:** logical communication between processes
 - relies on, enhances, network layer services

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

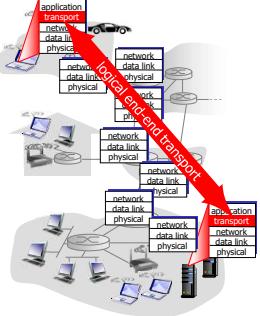
Transport Layer 3-4

4

3.2:

Internet transport-layer protocols

- reliable, in-order delivery (TCP)
 - congestion control
 - flow control
 - connection setup
- unreliable, unordered delivery: UDP
 - no-frills extension of "best-effort" IP
- services not available:
 - delay guarantees
 - bandwidth guarantees



Transport Layer 3-5

5

Chapter 3 outline

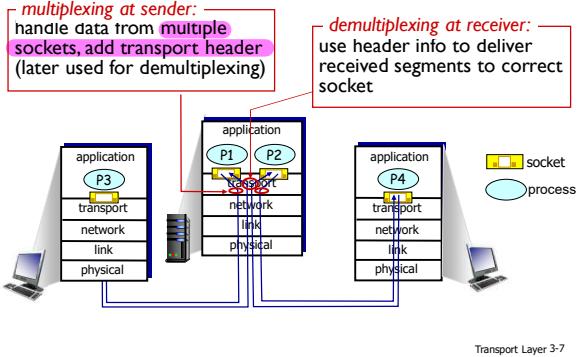
- 3.1 transport-layer services
- 3.2 multiplexing and demultiplexing
- 3.3 connectionless transport: UDP
- 3.4 principles of reliable data transfer
- 3.5 connection-oriented transport: TCP
 - segment structure
 - reliable data transfer
 - flow control
 - connection management
- 3.6 principles of congestion control
- 3.7 TCP congestion control

Transport Layer 3-6

6

1

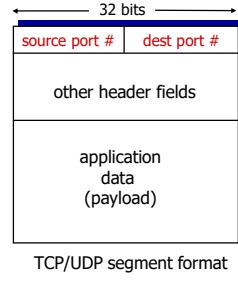
Multiplexing/demultiplexing



7

How demultiplexing works

- host receives IP datagrams
 - each datagram has source IP address, destination IP address
 - each datagram carries one transport-layer segment
 - each segment has source, destination port number
- host uses **IP addresses & port numbers** to direct segment to appropriate socket



Transport Layer 3-8

8

用 ip 地址不同的来区分 App
因为 Socket programming 会用 ip address

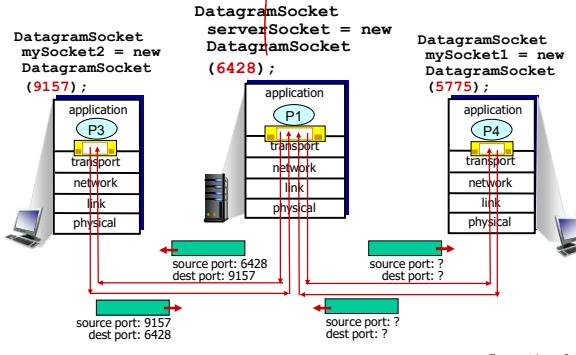
Connectionless demultiplexing

- recall: created socket has host-local port #:
`DatagramSocket mySocket1 = new DatagramSocket(12534);`
- recall: when creating datagram to send into UDP socket, must specify
 - destination IP address
 - destination port #
- when host receives UDP segment:
 - checks destination port # in segment
 - directs UDP segment to socket with that port #
- IP datagrams with **same dest. port #**, but different source IP addresses and/or source port numbers will be directed to **same socket** at dest

Transport Layer 3-9

9

Connectionless demux: example



10

Connection-oriented demux

- 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

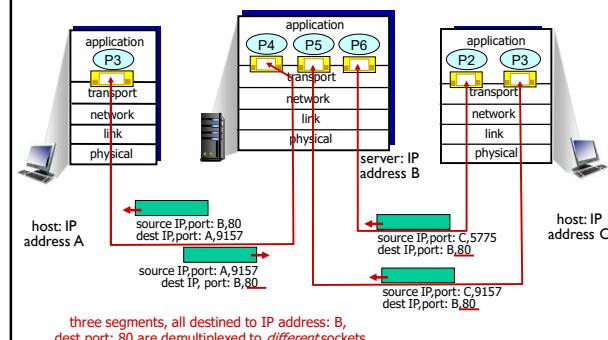
Transport Layer 3-11

11

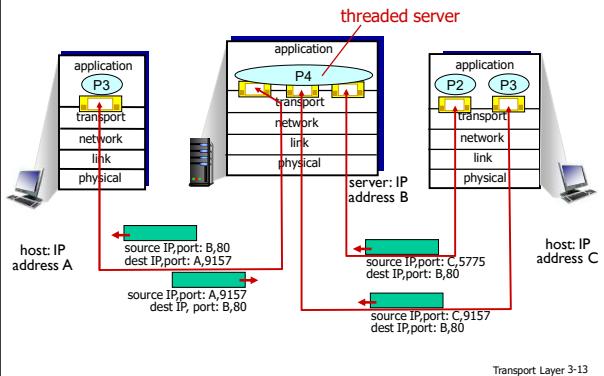
g mail has 3 tabs:
journey ip -> host ip -> des port ->
source port ->

12

Connection-oriented demux: example



Connection-oriented demux: example



13

★ TCP: bidirectional
UDP: unidirection

UDP: User Datagram Protocol [RFC 768]

- “no frills,” “bare bones” Internet transport protocol
- “best effort” service, UDP segments may be:
 - lost
 - delivered out-of-order to app
- **connectionless:**
 - no handshaking between UDP sender, receiver
 - each UDP segment handled independently of others
- UDP use:
 - streaming multimedia apps (loss tolerant, rate sensitive)
 - DNS
 - SNMP
- reliable transfer over UDP:
 - add reliability at application layer
 - application-specific error recovery!

Transport Layer 3-15

15

UDP checksum

Goal: detect “errors” (e.g., flipped bits) in transmitted segment

sender:

- treat segment contents, including header fields, as sequence of 16-bit integers
- checksum: addition (one's complement sum) of segment contents
- sender puts checksum value into UDP checksum field

receiver:

- compute checksum of received segment
- check if computed checksum equals checksum field value:
 - NO - error detected
 - YES - no error detected. But maybe errors nonetheless? More later

Transport Layer 3-17

17

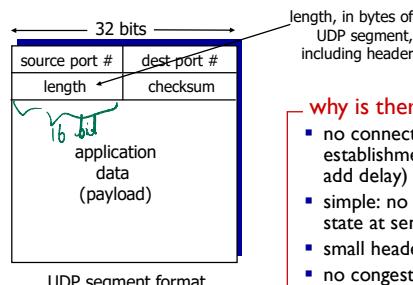
Chapter 3 outline

- 3.1 transport-layer services
- 3.2 multiplexing and demultiplexing
- 3.3 connectionless transport: UDP
- 3.4 principles of reliable data transfer
- 3.5 connection-oriented transport: TCP
 - segment structure
 - reliable data transfer
 - flow control
 - connection management
- 3.6 principles of congestion control
- 3.7 TCP congestion control

Transport Layer 3-14

14

UDP: segment header



why is there a UDP?

- no connection establishment (which can add delay)
- simple: no connection state at sender, receiver
- small header size
- no congestion control: UDP can blast away as fast as desired

Transport Layer 3-16

16

Internet checksum: example

example: add two 16-bit integers

1	1	1	1	0	0	1	1	0	0	1	1	0	0	1	1
1	1	0	1	0	1	0	1	0	1	0	1	0	1	0	1
wraparound															
1	0	1	1	1	0	1	1	1	0	1	1	1	0	1	1
sum	1	0	1	1	1	0	1	1	1	0	1	1	1	1	0
checksum	0	1	0	0	0	1	0	0	0	1	0	0	0	0	1

Note: when adding numbers, a carryout from the most significant bit needs to be added to the result

* Check out the online interactive exercises for more examples: http://gaia.cs.umass.edu/kurose_ross/interactive/

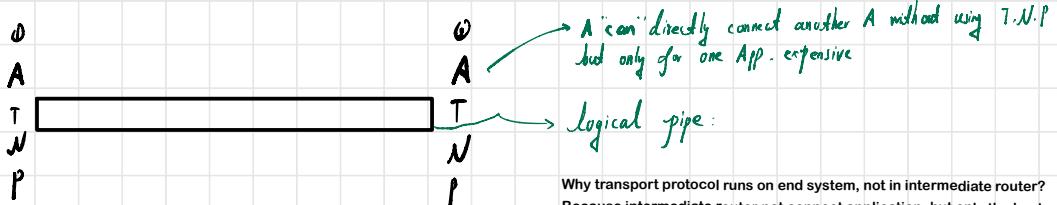
Transport Layer 3-18

18

Chapter 3

{ TCP: complicated, slow, reliable.
UDP: easy, fast, unreliable

"Jack steak and kingdom"



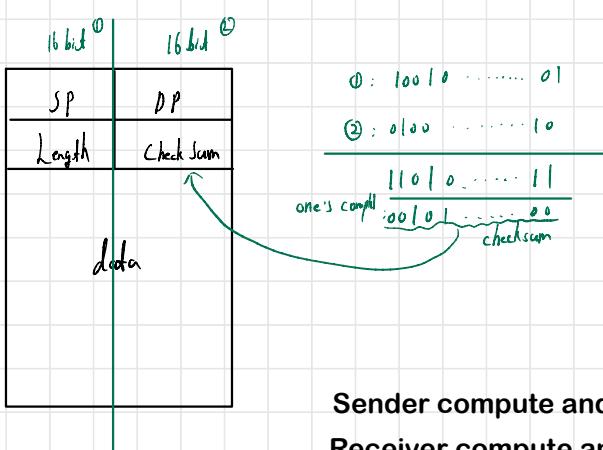
Transport: reassemble and chop up data

Stream delivery = sequence number delivery

TCP requires a hand shake.

Transport: App to app delivery.

A diagram showing a 1 GB file being processed by TCP and UDP. A large square box labeled "1 GB" has two arrows pointing away from it. One arrow points to the right with the text "TCP: chop it up for you (reasembly and segmentation) good for large size". The other arrow points to the right with the text "UDP: We manually chop up good for small size".



① : 10010.....01

② : 0100.....10

11010.....11

one's compd
0010.....00

checksum

Sender compute and populate Check sum
Receiver compute and compare the checksum

一直加多余的1直到只有十六个bits

Chapter 3 outline

3.1 transport-layer services

3.2 multiplexing and demultiplexing

3.3 connectionless transport: UDP

3.4 principles of reliable data transfer

3.5 connection-oriented transport: TCP

- segment structure
- reliable data transfer
- flow control
- connection management

3.6 principles of congestion control

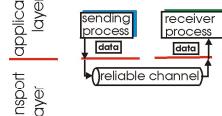
3.7 TCP congestion control

Transport Layer 3-19

19

Principles of reliable data transfer

- important in application, transport, link layers
- top-10 list of important networking topics!



(a) provided service

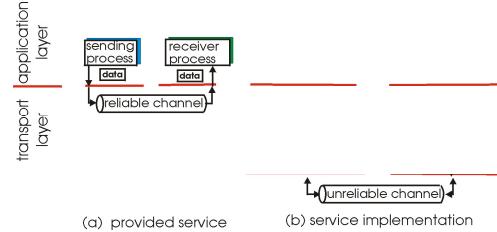
- characteristics of unreliable channel will determine complexity of reliable data transfer protocol (rdt)

Transport Layer 3-20

20

Principles of reliable data transfer

- important in application, transport, link layers
- top-10 list of important networking topics!



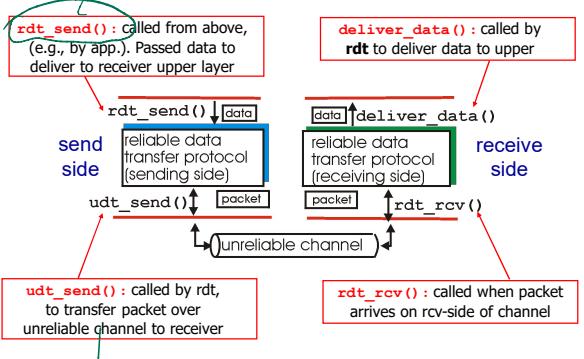
- characteristics of unreliable channel will determine complexity of reliable data transfer protocol (rdt)

Transport Layer 3-21

21

rdt_send , called by only A layer resides in T layer

Reliable data transfer: getting started



Transport Layer 3-23

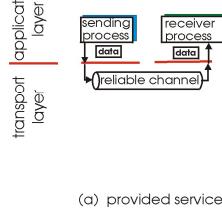
23

can be called by rdt in A

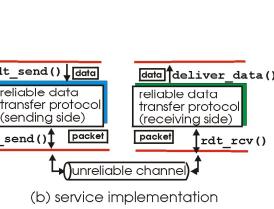
udt_send is in network layer - called by 1 layer

Principles of reliable data transfer

- important in application, transport, link layers
- top-10 list of important networking topics!



(a) provided service



(b) service implementation

- characteristics of unreliable channel will determine complexity of reliable data transfer protocol (rdt)

Transport Layer 3-22

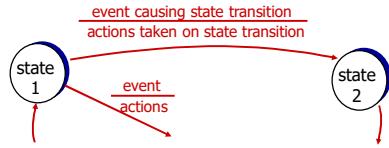
22

Reliable data transfer: getting started

we'll:

- incrementally develop sender, receiver sides of reliable data transfer protocol (rdt)
- consider only unidirectional data transfer
 - but control info will flow on both directions!
- use finite state machines (FSM) to specify sender, receiver

state: when in this "state" next state uniquely determined by next event



Transport Layer 3-24

24

02/22/2022

Why do we have reliable system on non-reliable network.

1: speed.

2: network guarantee something. Like TCP

3: simplicity. Do not want reliability on top of network layer.

Reliable network:

Reliability: A guarantee the packet will be delivered to B

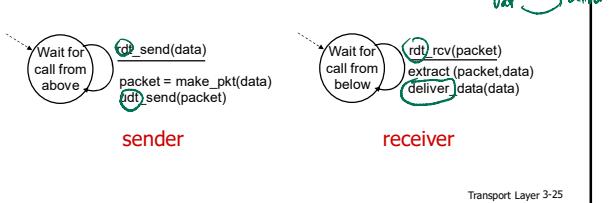
Integrity: If a data will be deliver to B it will be exact the A send.

check sum mechanism.

A → talk to only 7 layer, cont with N.P
T (Reliability/unreliability both control by 7 layer)
TCP UDP
N
P

rdt1.0: reliable transfer over a reliable channel

- underlying channel perfectly reliable
 - no bit errors
 - no loss of packets
- separate FSMs for sender, receiver:
 - sender sends data into underlying channel
 - receiver reads data from underlying channel



25

27 28

rdt2.0: channel with bit errors

- underlying channel may flip bits in packet
 - checksum to detect bit errors
- the question: how to recover from errors:

How do humans recover from “errors” during conversation?

Transport Layer 3-26

26

rdt2.0: channel with bit errors

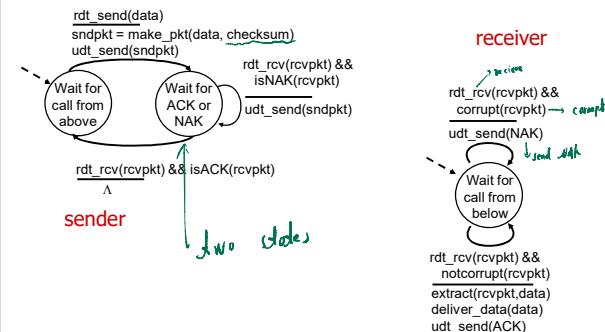
- underlying channel may flip bits in packet
 - checksum to detect bit errors
- the question: how to recover from errors:
 - acknowledgements (ACKs)**: receiver explicitly tells sender that pkt received OK
 - negative acknowledgements (NAKs)**: receiver explicitly tells sender that pkt had errors
 - sender retransmits pkt on receipt of NAK
- new mechanisms in rdt2.0 (beyond rdt1.0):
 - error detection
 - feedback: control msgs (ACK,NAK) from receiver to sender

Transport Layer 3-27

Transport Layer 3-28

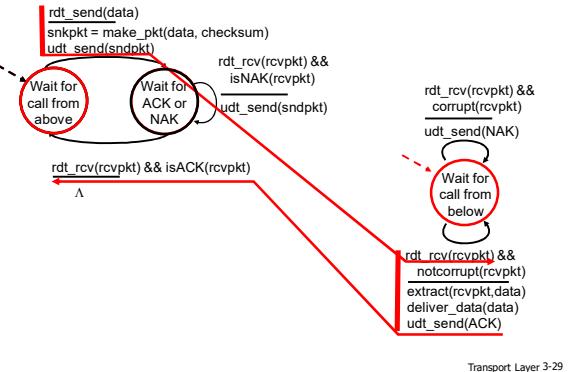
27

rdt2.0: FSM specification



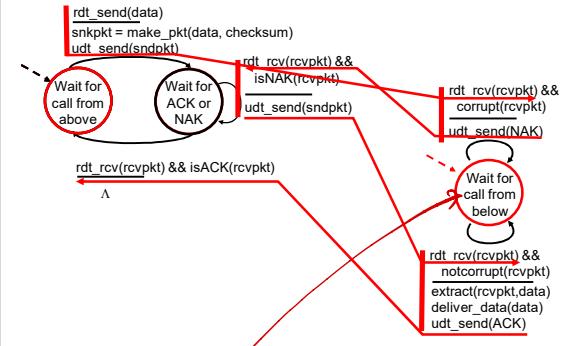
28

rdt2.0: operation with no errors



29

rdt2.0: error scenario



30

why only one state?
In either way (good pkt, bad pkt)
it will always listen.

rdt2.0 has a fatal flaw!

what happens if
ACK/NAK corrupted?

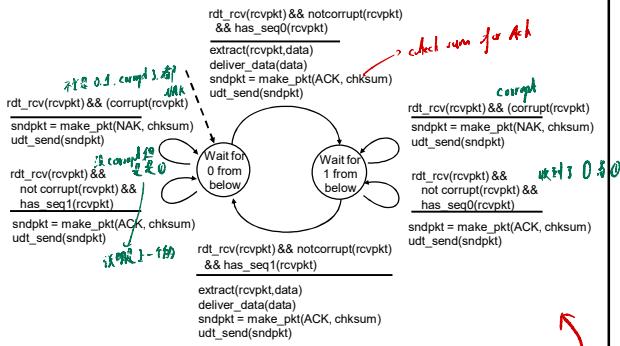
- sender doesn't know what happened at receiver!
 - can't just retransmit: possible duplicate

- stop and wait
sender sends one packet,
then waits for receiver
response

Transport Layer 3-31

31 ✓ sender and receiver will wait forever!!!!

rdt2.1: receiver, handles garbled ACK/NAKs



You slow!!!

33

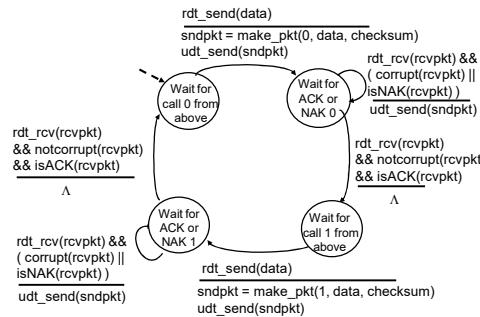
rdt2.2: a NAK-free protocol

- same functionality as rd_{2.1}, using ACKs only
 - instead of NAK, receiver sends ACK for last pkt received OK
 - receiver must explicitly include seq # of pkt being ACKed
 - duplicate ACK at sender results in same action as NAK: retransmit current pkt

Transport Layer 3-35

stop & wait protocol

rdt2.1: sender, handles garbled ACK/NAKs



Transport Layer 3-32

32 1. d. sequence 10. w. enoxyd.

rdt2.1: discussion

sender:

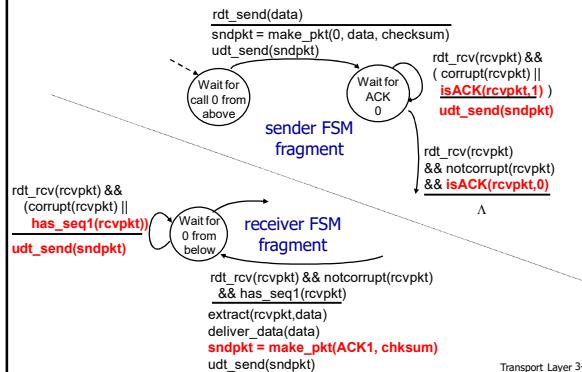
- seq # added to pkt
 - two seq. #'s (0,1) will suffice. Why?
 - must check if received ACK/NAK corrupted
 - twice as many states
 - state must "remember" whether "expected" pkt should have seq # of 0 or 1
 - must check if received packet is duplicate
 - state indicates whether 0 or 1 is expected pkt seq #
 - note: receiver can not know if its last ACK/NAK received OK at sender

Transport Layer 3-34

34

foster

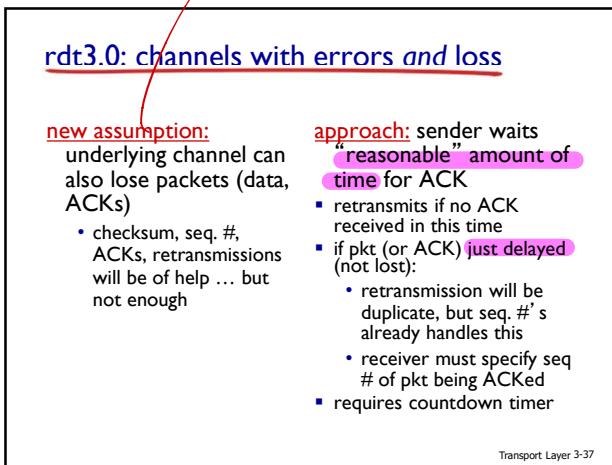
rdt2.2: sender, receiver fragments



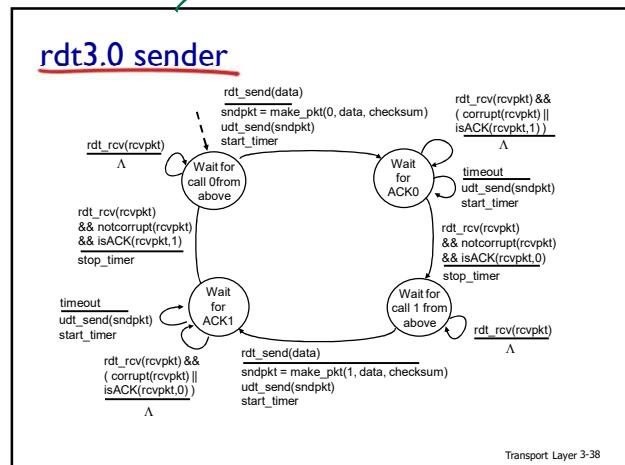
Transport Layer 3-36

36

已自动元了 drop
且 intermediate router 不是 sender

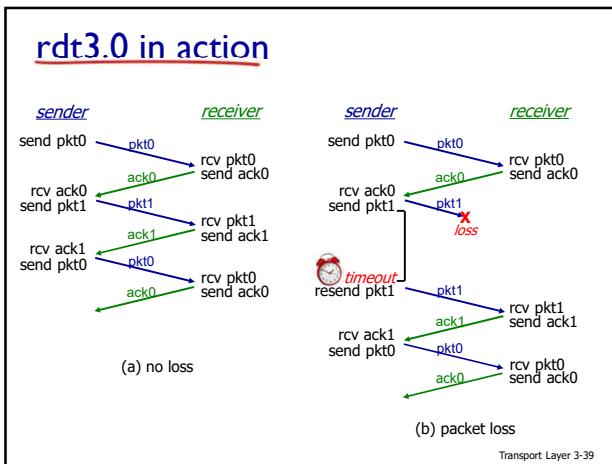


37

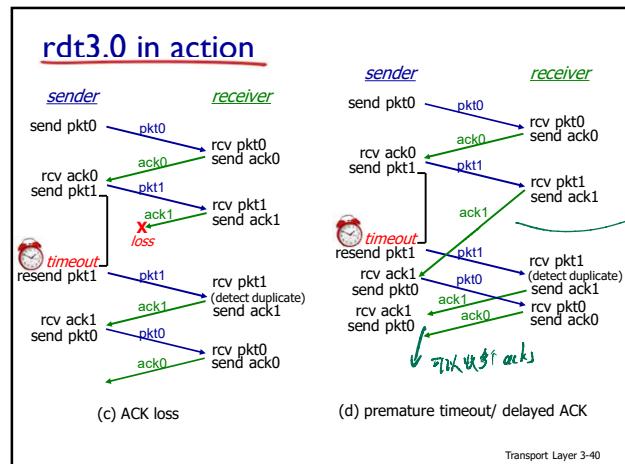


38

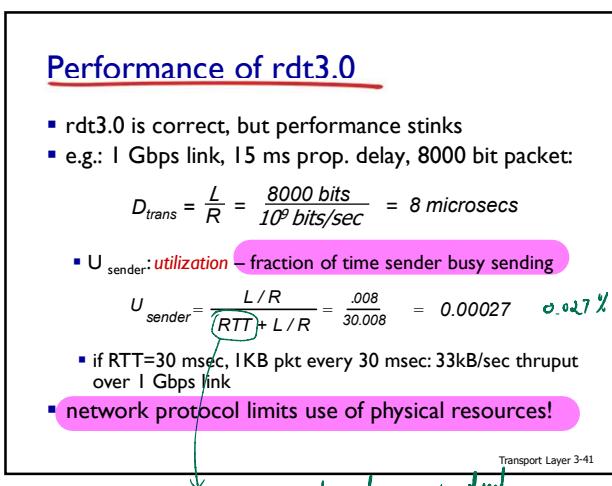
这里有几个 scenario 可以看



39

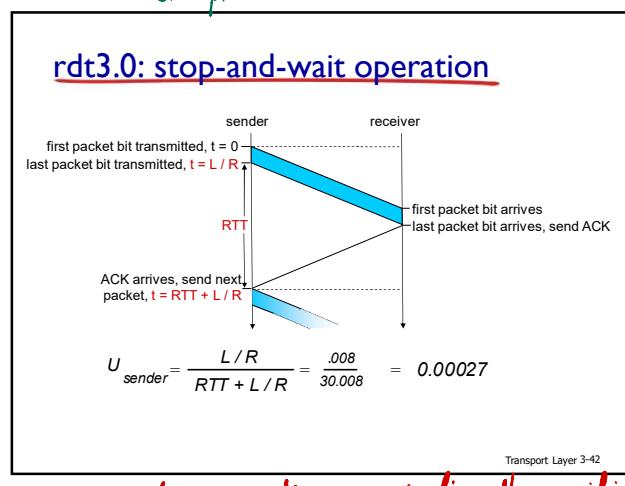


40



41

Transport layer protocol
is limiting the performance



42

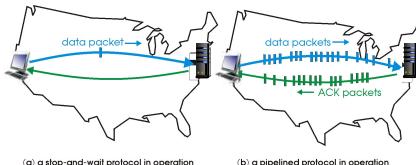
spend more time on sending than waiting
thus, higher utilization.

pipeline: multiple in-flight pkts

Pipelined protocols

pipelining: sender allows multiple, "in-flight", yet-to-be-acknowledged pkts

- range of sequence numbers must be increased
- buffering at sender and/or receiver



- two generic forms of pipelined protocols: go-Back-N, selective repeat

Transport Layer 3-43

43

pipelined protocols: { go-Back-N
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**

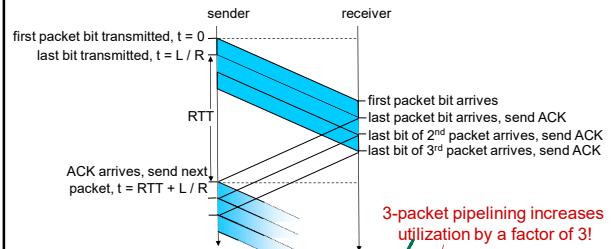
Selective Repeat:

- 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**

Transport Layer 3-45

45

Pipelining: increased utilization



$$U_{\text{sender}} = \frac{3L/R}{RTT + L/R} = \frac{.0024}{30.008} = 0.00081$$

Transport Layer 3-44

44

✓ sending time increased
has not back send yet

Go-Back-N: sender

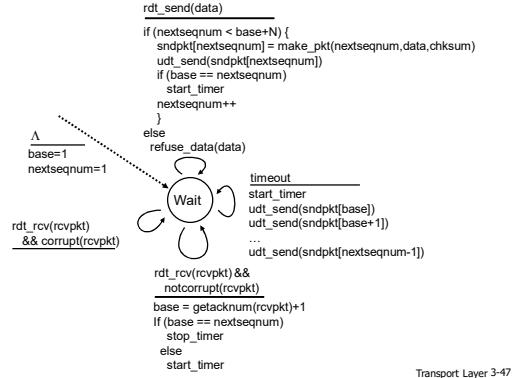
- k-bit seq # in pkt header
- "window" of up to N, consecutive unack'ed pkts allowed



Transport Layer 3-46

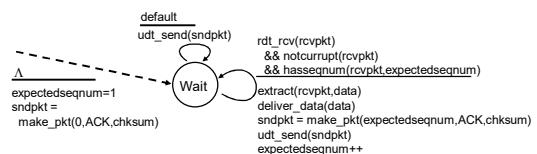
46 Sender required buffering

GBN: sender extended FSM



47

GBN: receiver extended FSM



ACK-only: always send ACK for correctly-received pkt with highest **in-order** seq #

- may generate duplicate ACKs
- need only remember **expectedseqnum**
- out-of-order pkt:
 - discard (don't buffer): **no receiver buffering!**
 - re-ACK pkt with highest in-order seq #

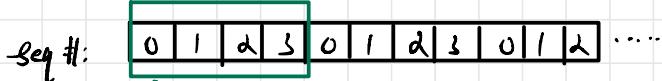
Transport Layer 3-48

48

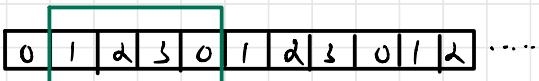
for receiver
usually is window-size

02/24/2022

Go-back-N: Window size: 4



oldest packet timer, if timeout, send all packed in the window



0 is acked, then slide window to next



d is acked, then slide the window
since ack is cumulative

Problem:

① what determines seq# : window size?

② issue of go-back-N

③ The seq# in GBN is

1	2	3	4	5	1	2	3	4	5
---	---	---	---	---	---	---	---	---	---

or

1	2	3	4	5	6	7	8
---	---	---	---	---	---	---	---

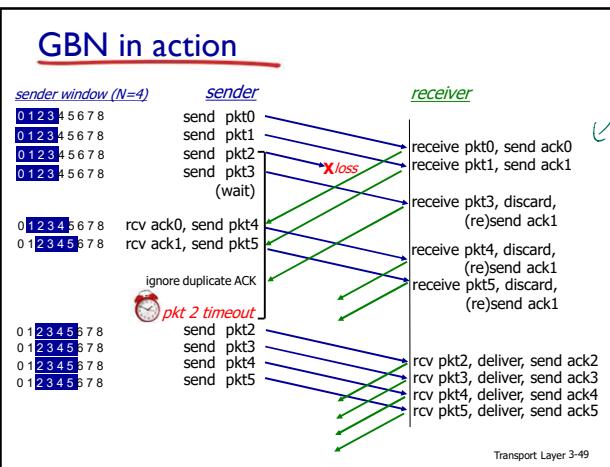
when window slice, will seq# change?

size of window should $\frac{1}{2}$ say a layer smaller.

The size of window is dynamic The more congested, the smaller window size

The size of window is at least half of the range of sequence numbers

Window size has threshold:



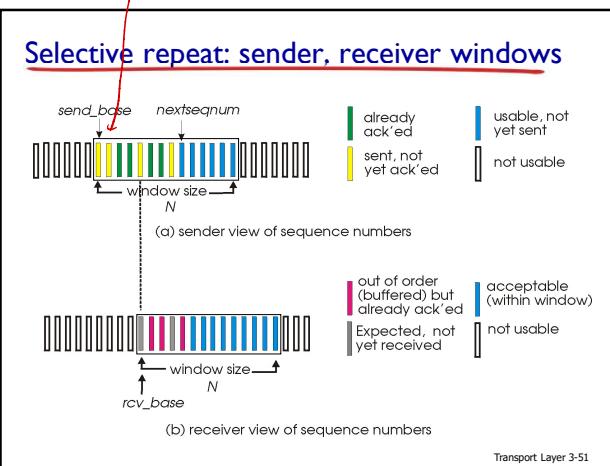
49

Selective repeat

- receiver **individually acknowledges** all correctly received pkts
 - buffers pkts, as needed, for eventual in-order delivery to upper layer
- sender only resends pkts for which ACK not received
 - sender timer for each unACKed pkt
- sender window
 - N consecutive seq #'s
 - limits seq #'s of sent, unACKed pkts

Transport Layer 3-50

50



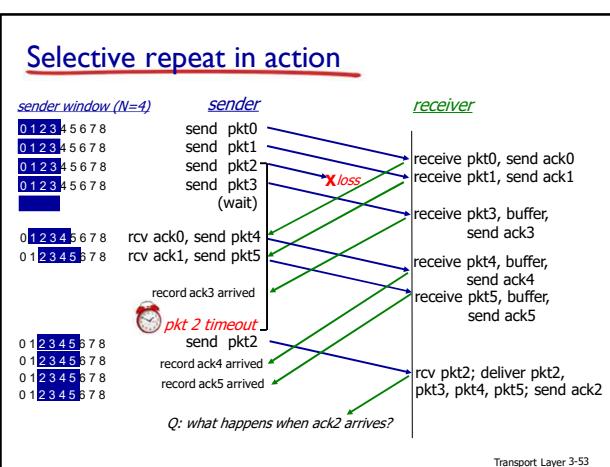
51

Selective repeat

- | | |
|---|---|
| sender
data from above: <ul style="list-style-type: none"> if next available seq # in window, send pkt timeout(n): <ul style="list-style-type: none"> resend pkt n, restart timer ACK(n) in [sendbase, sendbase+N]: <ul style="list-style-type: none"> mark pkt n as received if n smallest unACKed pkt, advance window base to next unACKed seq # | receiver
pkt n in [rcvbase, rcvbase+N-1] <ul style="list-style-type: none"> send ACK(n) out-of-order: buffer in-order: deliver (also deliver buffered, in-order pkts), advance window to next not-yet-received pkt pkt n in [rcvbase-N, rcvbase-1] <ul style="list-style-type: none"> ACK(n) otherwise: <ul style="list-style-type: none"> ignore |
|---|---|

Transport Layer 3-52

52



53

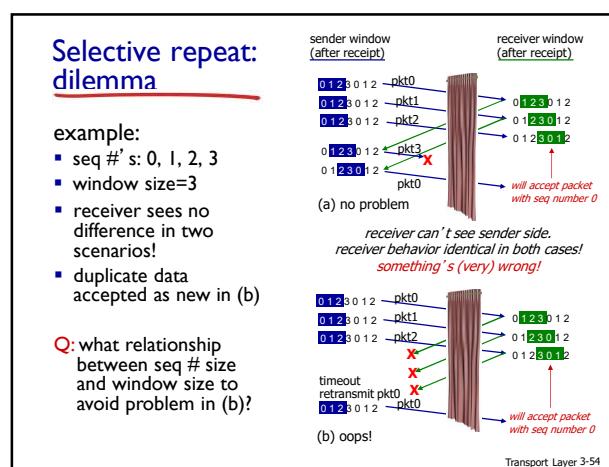
Selective repeat: dilemma

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

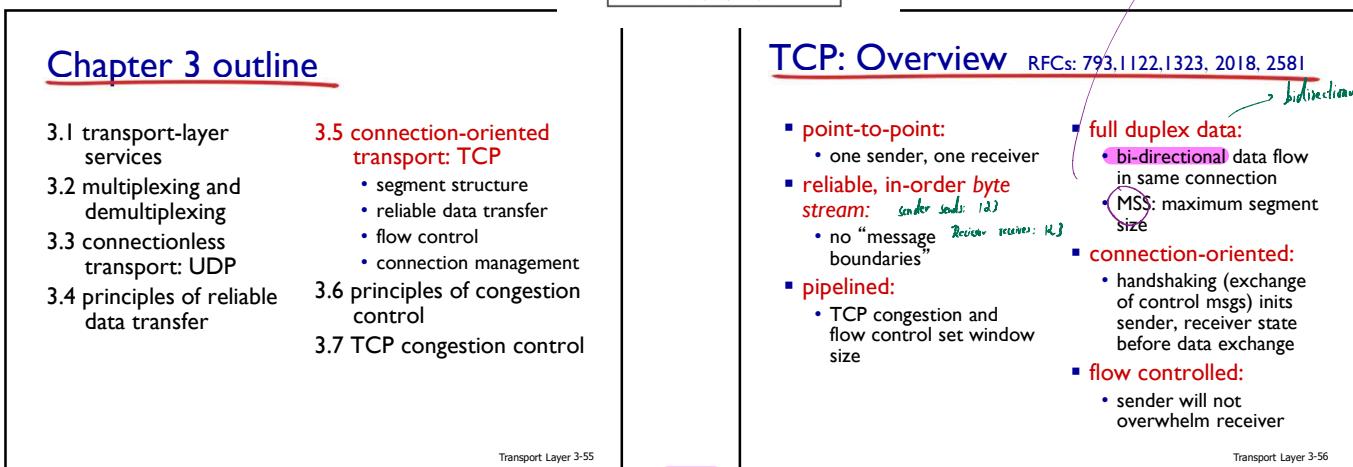
Q: what relationship between seq # size and window size to avoid problem in (b)?

54

both selective repeat and GBN have same problems.



receiver has 0
from sender it's old 0



55

(1) need source dest port since TCP is a T protocol

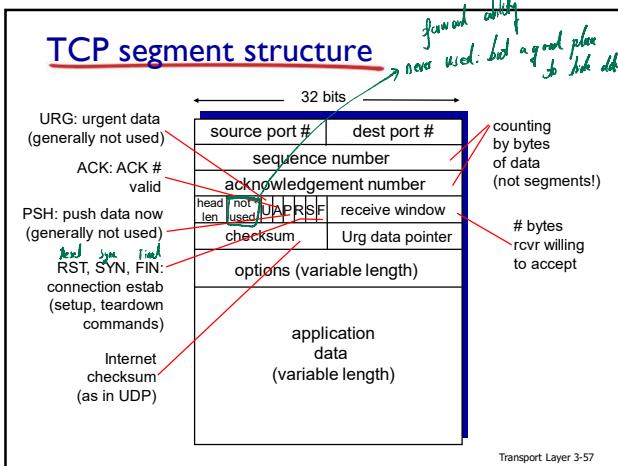
TCP: Overview RFCs: 793, 1122, 1323, 2018, 2581

- point-to-point:
 - one sender, one receiver
- reliable, in-order byte stream: *sender sends: 123, receiver receives: 123*
 - no “message boundaries”
- pipelined:
 - TCP congestion and flow control set window size
- full duplex data:
 - bi-directional data flow in same connection
 - MSS: maximum segment size
- connection-oriented:
 - handshaking (exchange of control msgs) init sender, receiver state before data exchange
- flow controlled:
 - sender will not overwhelm receiver

Transport Layer 3-56

56

→ ABN: selective segt of packets
乱序包

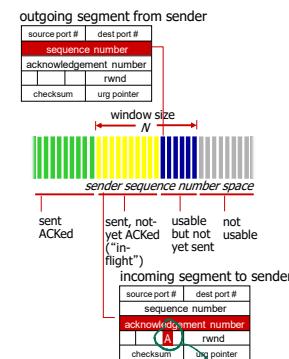


57

both sides could be sender and receiver

TCP seq. numbers, ACKs

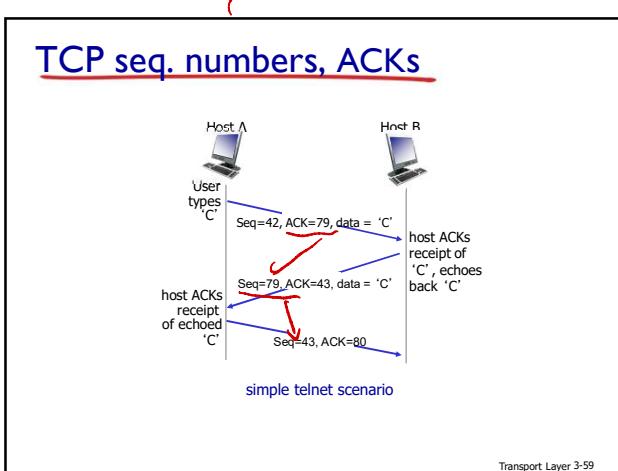
- sequence numbers:**
- byte stream “number” of first byte in segment’s data
- acknowledgements:**
- seq # of next byte expected from other side
 - cumulative ACK
- Q:** how receiver handles out-of-order segments
- A: TCP spec doesn’t say, up to implementor



58

Tcp: Ack: 10, 1024, 1024
G&N: Ack: 10, 1023

Flag: A, for acknowledgement



59

TCP round trip time, timeout

- Q:** how to set TCP timeout value?
- longer than RTT
 - but RTT varies
 - too short: premature timeout, unnecessary retransmissions
 - too long: slow reaction to segment loss

- Q:** how to estimate RTT?
- **SampleRTT:** measured time from segment transmission until ACK receipt
 - ignore retransmissions
 - SampleRTT will vary, want estimated RTT “smoother”
 - average several recent measurements, not just current SampleRTT

60

Transport Layer 3-60

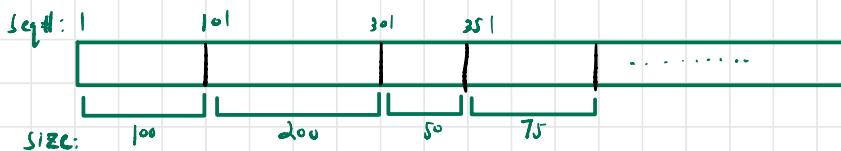
Sequence #

GBN: seq# for packets

TCP: seq# for byte-stream

TCP assumes data is continuous not discrete. That's why called stream.

The size of segment in TCP is different.



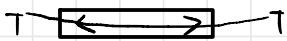
Why continuous:

①

②

A

1



Q: how receiver handles out-of-order segments

- A: TCP spec doesn't say,
- up to implementor

TCP: like Int. cars. wash machine
may not have luxury memory. (don't buffer)
either way
PC
(buffer)

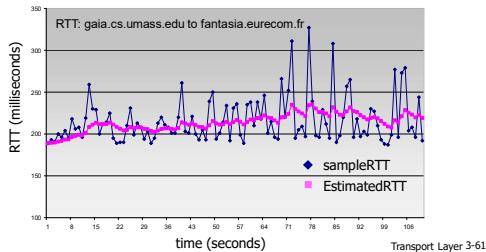
⚠ why not measure instantaneous RTT to determine time-out?

Smooth increase, not only influence by bursty data only, but also the previous EstimatedRTT, the history

TimeOut should based on smoother RTT.

TCP round trip time, timeout

- $\text{EstimatedRTT} = (1 - \alpha) * \text{EstimatedRTT} + \alpha * \text{SampleRTT}$
- exponential weighted moving average
 - influence of past sample decreases exponentially fast
 - typical value: $\alpha = 0.125$



Some of RTTs, not all of them, otherwise you're killing the system.

TCP round trip time, timeout

- timeout interval: EstimatedRTT plus "safety margin"
 - large variation in EstimatedRTT \rightarrow larger safety margin
- estimate SampleRTT deviation from EstimatedRTT:

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

(typically, $\beta = 0.25$)

$$\text{TimeoutInterval} = \text{EstimatedRTT} + 4 * \text{DevRTT}$$

estimated RTT "safety margin"

* Check out the online interactive exercises for more examples: http://gaia.cs.umass.edu/kurose_ross/interactive/

Transport Layer 3-62

61

62

Chapter 3 outline

- | | |
|--|--|
| 3.1 transport-layer services | 3.5 connection-oriented transport: TCP |
| 3.2 multiplexing and demultiplexing | <ul style="list-style-type: none"> • segment structure • reliable data transfer • flow control • connection management |
| 3.3 connectionless transport: UDP | 3.6 principles of congestion control |
| 3.4 principles of reliable data transfer | 3.7 TCP congestion control |

Transport Layer 3-63

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

Transport Layer 3-64

63

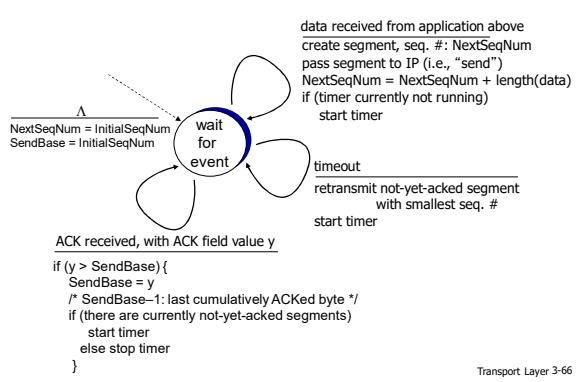
64

TCP sender events:

- data rcvd from app:*
- create segment with seq #
 - seq # is byte-stream number of first data byte in segment
 - start timer if not already running
 - think of timer as for oldest unacked segment
 - expiration interval: TimeOutInterval
- timeout:*
- retransmit segment that caused timeout
 - restart timer
- ack rcvd:*
- if ack acknowledges previously unacked segments
 - update what is known to be ACKed
 - start timer if there are still unacked segments

Transport Layer 3-65

TCP sender (simplified)



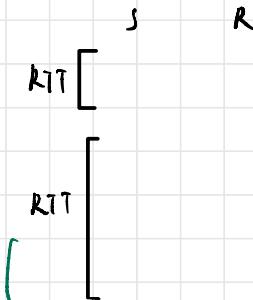
Transport Layer 3-66

65

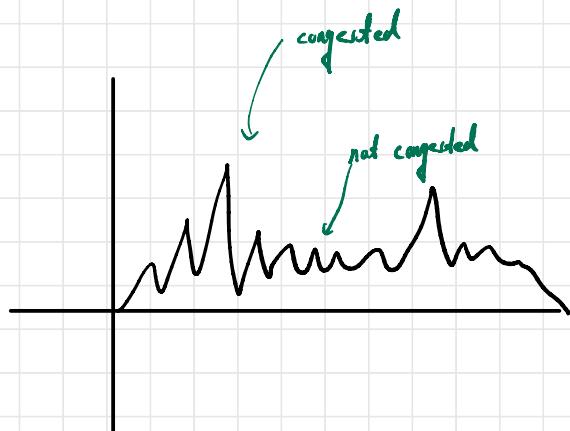
66

03/01/2022

internet is bursty data.



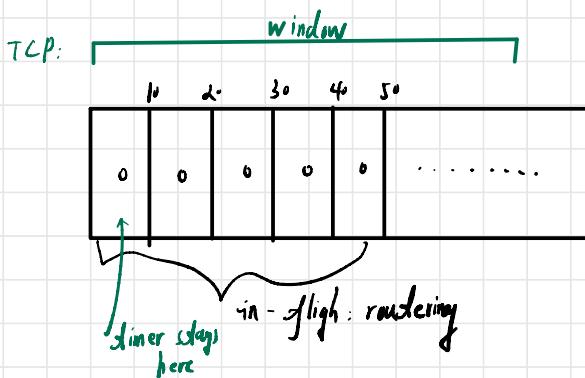
if RTT is shorter,
then internet is fast.
Longer, then slower.



If rtt is shorter, then we may have timer wait less time. If rtt is longer or bigger, then we should have the timer wait more time.

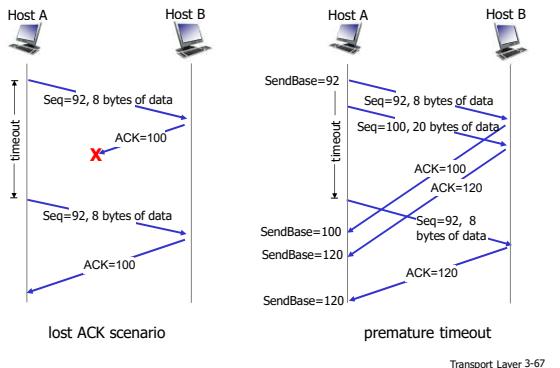
RTT should not be larger than time-out

It is good to spend time on sending rather than waiting to increase utilization.



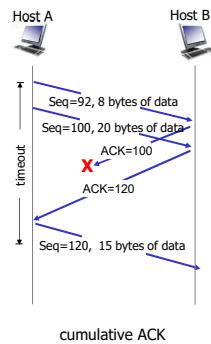
If time out, only send one segment, and whole window, then restart timer.

TCP: retransmission scenarios



67

TCP: retransmission scenarios



68

03/01

TCP ACK generation [RFC 1122, RFC 2581]

event at receiver	TCP receiver action
arrival of in-order segment with expected seq #. All data up to expected seq # already ACKed	delayed ACK. Wait up to 500ms for next segment. If no next segment, send ACK
arrival of in-order segment with expected seq #. One other segment has ACK pending	immediately send single cumulative ACK, ACKing both in-order segments
arrival of out-of-order segment higher-than-expect seq. #. Gap detected	immediately send <u>duplicate ACK</u> , indicating seq. # of next expected byte
arrival of segment that partially or completely fills gap	immediate send ACK, provided that segment starts at lower end of gap

Transport Layer 3-69

69

this can cause fast-retransmission

TCP fast retransmit

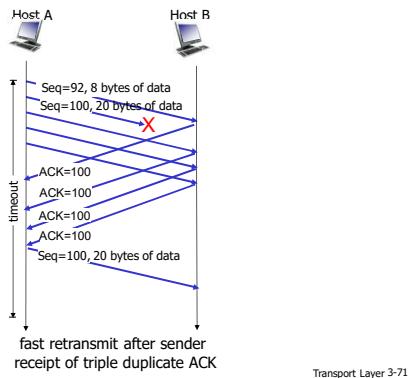
- time-out period often relatively long:
 - long delay before resending lost packet
- detect lost segments via duplicate ACKs.
 - sender often sends many segments back-to-back
 - if segment is lost, there will likely be many duplicate ACKs.

TCP fast retransmit
if sender receives 3 ACKs for same data ("triple duplicate ACKs"), resend unacked segment with smallest seq #
likely that unacked segment lost, so don't wait for timeout

Transport Layer 3-70

70

TCP fast retransmit



71

Chapter 3 outline

- 3.1 transport-layer services
- 3.2 multiplexing and demultiplexing
- 3.3 connectionless transport: UDP
- 3.4 principles of reliable data transfer
- 3.5 connection-oriented transport: TCP
 - segment structure
 - reliable data transfer
 - flow control**
 - connection management
- 3.6 principles of congestion control
- 3.7 TCP congestion control

Transport Layer 3-72

72

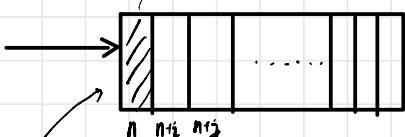
12

03/03/2022

color of using

optimistic to wait up to 500ms if the n received in order

①



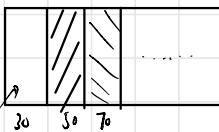
Last ack will
be ack n if n
is expected.

②



pending ack n if i arrives then immediately send ack i.

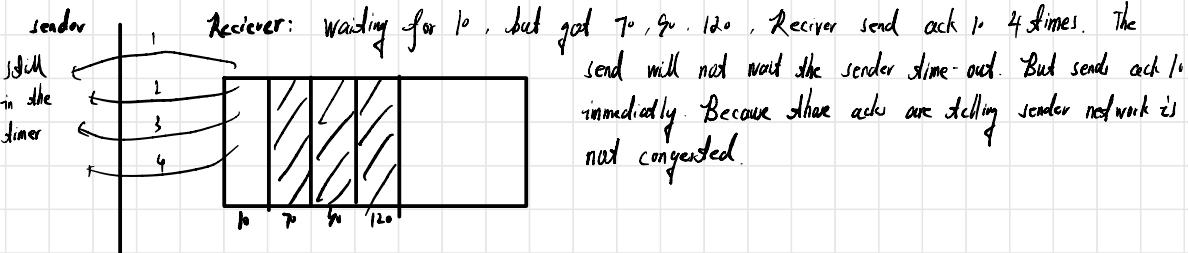
③



expected 30 but 50, 70 received in out of order.

Then immediately send duplicate of ack 30
since we assume 30 is lost.

Fast-retransmission:



03/03/2022

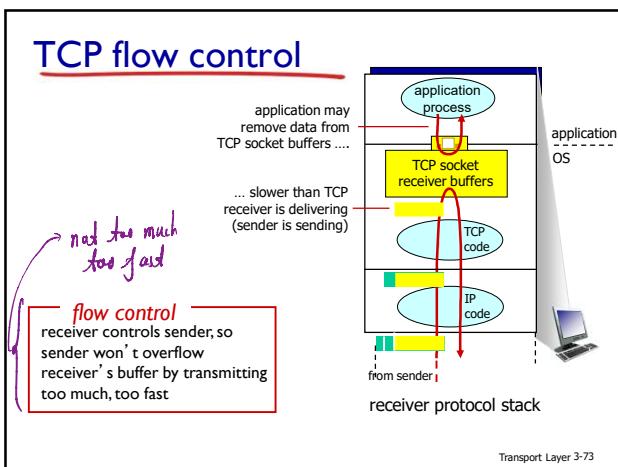
Sender

exam:

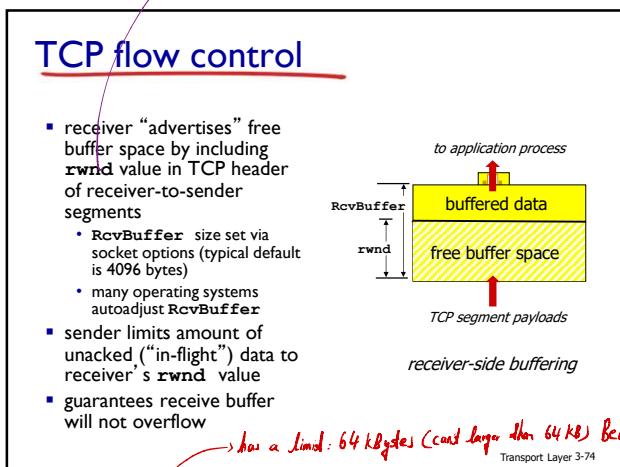
Receiver

If the timer is time-out and receives nothing
That means network is congest.

If receives duplicated acks the it means network is
fine but the sent segment is lost. And not need to
wait timer then fast retransmit.

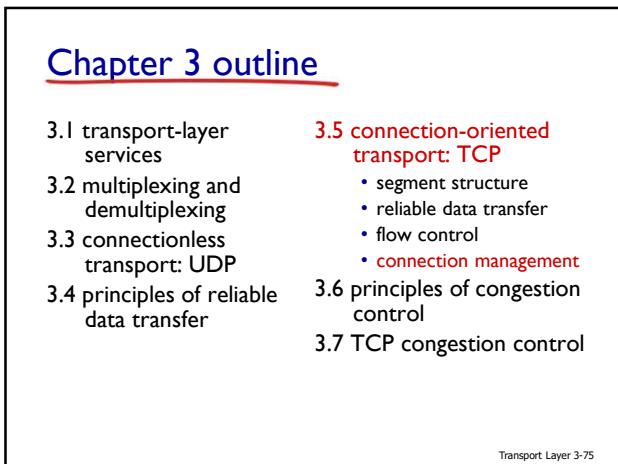


73

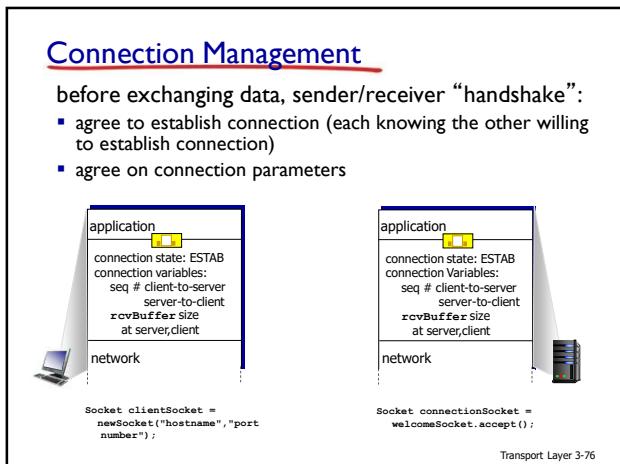


74

rwnd: is a value based on RcvBuffer. Network property
RcvBuffer: is system property
rwnd could be smaller than RcvBuffer

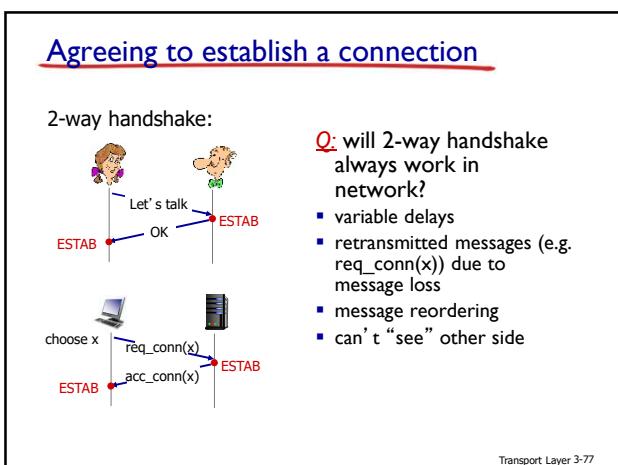


75

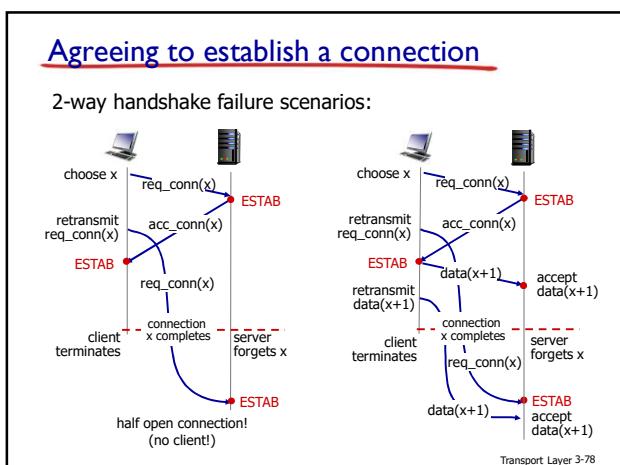


76

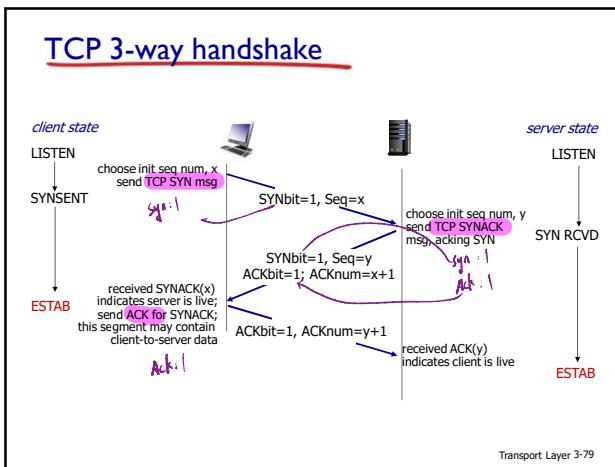
快过，无细讲



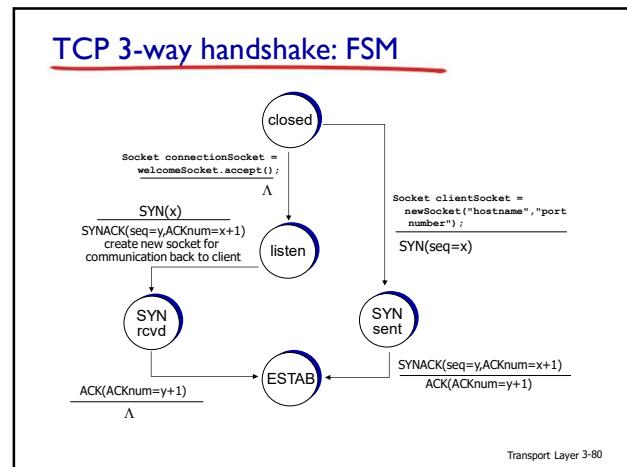
77



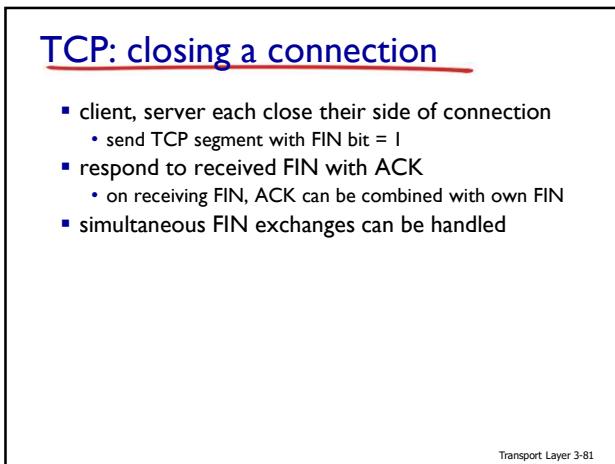
78



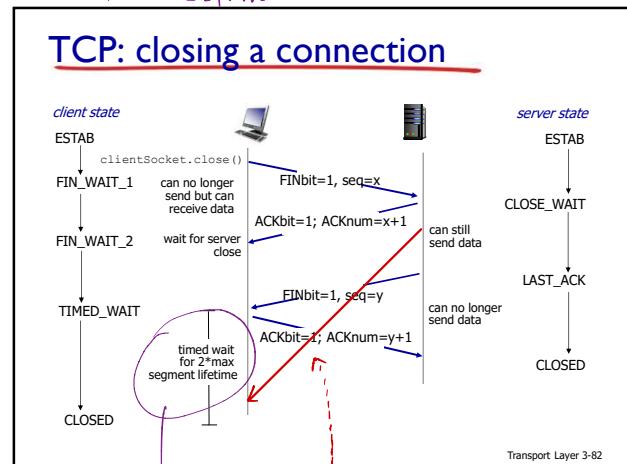
79



80



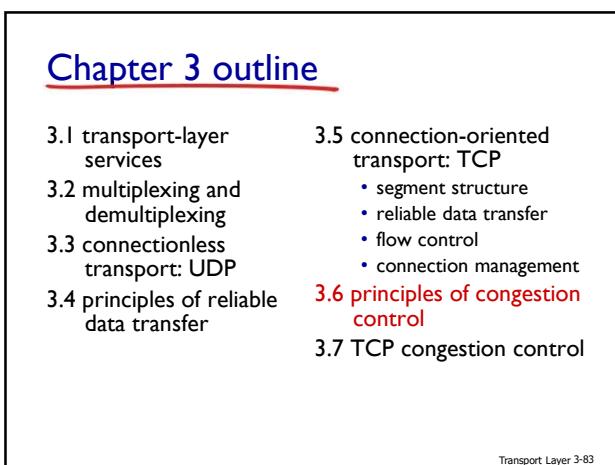
81



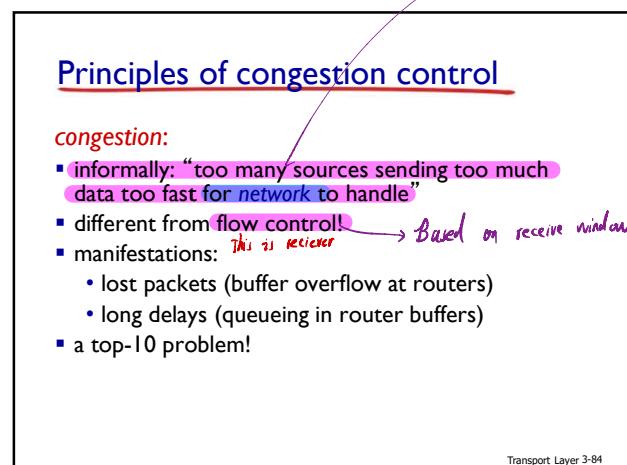
82

why wait? In case some data packet is delayed.

Based on estimated RTT



83

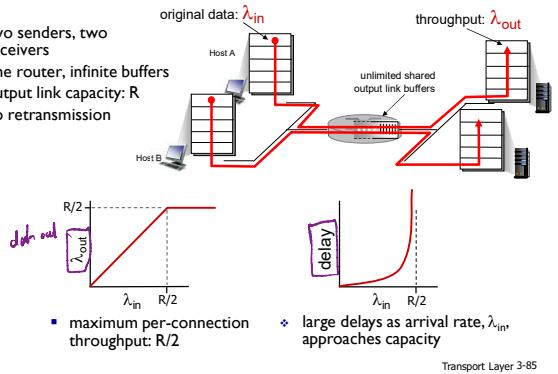


84

03/03/2022

Causes/costs of congestion: scenario 1

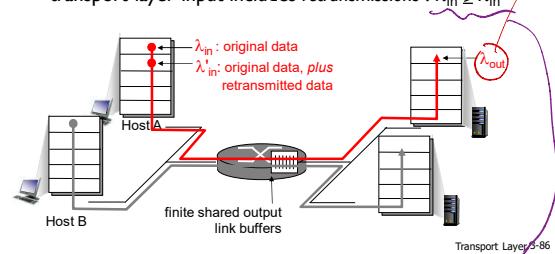
- two senders, two receivers
- one router, infinite buffers
- output link capacity: R
- no retransmission



85

Causes/costs of congestion: scenario 2

- one router, **finite** buffers
- sender retransmission of timed-out packet
- application-layer input = application-layer output: $\lambda_{in} = \lambda_{out}$
- transport-layer input includes retransmissions: $\lambda'_{in} \geq \lambda_{in}$



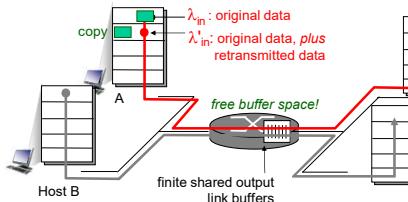
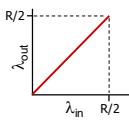
86

$$\lambda'_{in} = \lambda_{in} + \text{dormant packets}.$$

Causes/costs of congestion: scenario 2

idealization: perfect knowledge

- sender sends only when router buffers available



87

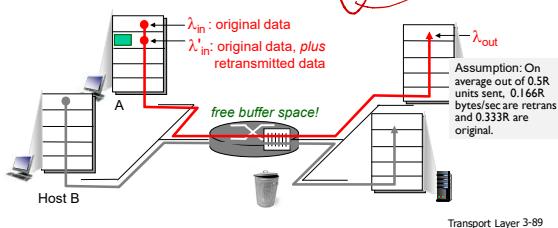
$$\lambda_{out} = \lambda'_{in} - \text{retransm. packet}$$

从转发点上有效率的极限
 先发3.已发送
 重传3.重传
 重传3.重传

Causes/costs of congestion: scenario 2

Idealization: known loss
 packets can be lost, dropped at router due to full buffers

- sender only resends if packet known to be lost

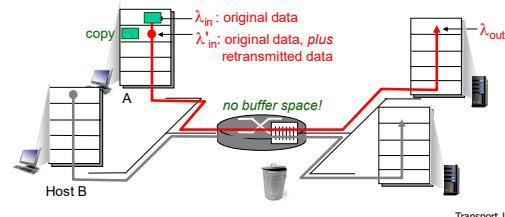


89

Causes/costs of congestion: scenario 2

Idealization: known loss
 packets can be lost, dropped at router due to full buffers
 $\lambda'_{in} > \lambda_{in}$

- sender only resends if packet known to be lost

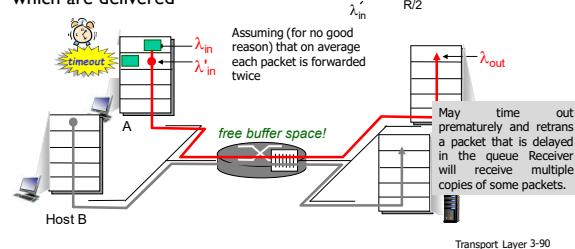


88

Causes/costs of congestion: scenario 2

Realistic: duplicates

- packets can be lost, dropped at router due to full buffers
- sender times out prematurely, sending two copies, both of which are delivered

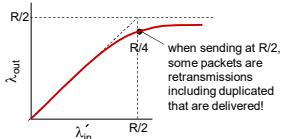


90

Causes/costs of congestion: scenario 2

Realistic: duplicates

- packets can be lost, dropped at router due to full buffers
- sender times out prematurely, sending **two** copies, both of which are delivered



"costs" of congestion:

- more work (retrans) for given "goodput"
- unnecessary retransmissions: link carries multiple copies of pkt
 - decreasing goodput

Transport Layer 3-91

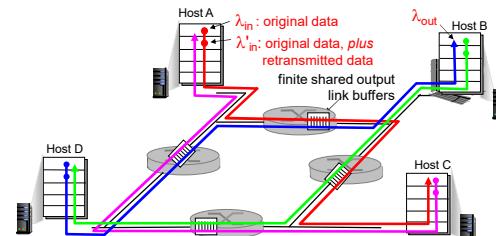
91

Causes/costs of congestion: scenario 3

- four senders
- multihop paths
- timeout/retransmit

Q: what happens as λ_{in} and λ_{in}' increase?

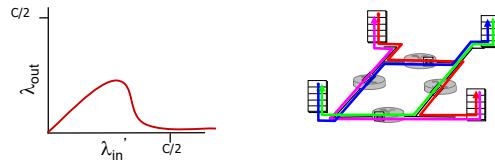
A: as red λ_{in} increases, all arriving blue pkts at upper queue are dropped, blue throughput $\rightarrow 0$



Transport Layer 3-92

92

Causes/costs of congestion: scenario 3



another "cost" of congestion:

- when packet dropped, any "upstream transmission capacity used for that packet was wasted!"

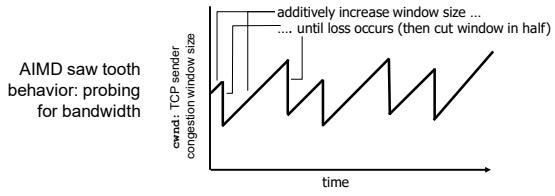
Transport Layer 3-93

93

rwnd and cwnd controls size of sliding window. And the minimal one of them decide will decides the size.

TCP congestion control: additive increase multiplicative decrease

- approach:** sender increases transmission rate (window size), probing for usable bandwidth, until loss occurs
 - additive increase:** increase **cwnd** (**congestion window**) by 1 MSS every RTT until loss detected
 - multiplicative decrease:** cut **cwnd** in half after loss



94

Chapter 3 outline

- 3.1 transport-layer services
- 3.2 multiplexing and demultiplexing
- 3.3 connectionless transport: UDP
- 3.4 principles of reliable data transfer

- 3.5 connection-oriented transport: TCP
 - segment structure
 - reliable data transfer
 - flow control
 - connection management
- 3.6 principles of congestion control
- 3.7 TCP congestion control

Transport Layer 3-94

TCP Congestion Control: details

- TCP sending rate:**
 - roughly: send **cwnd** bytes, wait RTT for ACKS, then send more bytes
- sender limits transmission:**

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

$$\text{LastByteSent} - \text{LastByteAcked} \leq \text{cwnd}$$
- cwnd is dynamic, function of perceived network congestion**

$$\text{LastByteSent} - \text{LastByteAcked} \leq \min(\text{cwnd}, \text{rwnd})$$

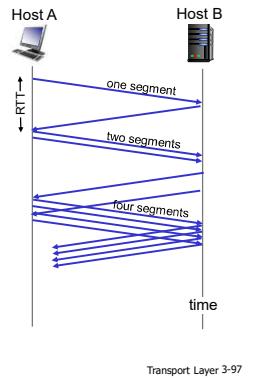
Transport Layer 3-96

95

96

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



97

七

TCP: detecting, reacting to loss

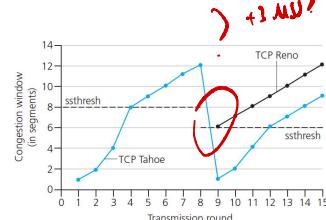
- loss indicated by timeout:
 - cwnd set to 1 MSS,
 - window then grows exponentially (as in slow start) to threshold, then grows linearly
 - loss indicated by 3 duplicate ACKs, TCP RENO
 - dup ACKs indicate network capable of delivering some segments ?
 - cwnd is cut in half window then grows linearly
 - TCP Tahoe always sets cwnd to 1 (timeout or 3 duplicate acks)

Transport Layer 3-98

到底加不加了呢?
±看書

TCP: switching from slow start to CA (Congestion Avoidance)

- Q: when should the exponential increase switch to linear?
 - A: when $cwnd$ gets to 1/2 of its value before timeout.



99

TCP throughput

- avg. TCP thruput as function of window size, RTT?
 - ignore slow start, assume always data to send
 - W: window size (measured in bytes) where loss occurs
 - avg. window size (# in-flight bytes) is $\frac{3}{4}W$
 - avg. thruput is $\frac{3}{4}W$ per RTT

$$\text{avg TCP thruput} = \frac{3}{4} \frac{W}{RTT} \text{ bytes/sec}$$



Transport Layer 3-10

101

TCP Futures: TCP over “long, fat pipes”

- example: 1500 byte segments, 100ms RTT, want 10 Gbps throughput
 - requires $W = 83,333$ in-flight segments
 - throughput in terms of segment loss probability, L [Mathis 1997]:

$$\text{TCP throughput} = \frac{1.22 \cdot \text{MSS}}{\text{RTT} \sqrt{L}}$$

- to achieve 10 Gbps throughput, need a loss rate of $L = 2 \cdot 10^{-10}$ – *a very small loss rate!*

- new versions of TCP for high-speed

Transport Layer 3-102

102

03/22/2022

color :

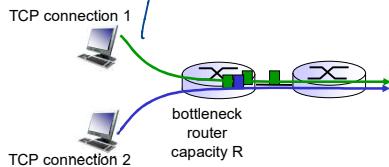
congestion control is defined by what?

Review questions

如果有 10↑ connections
那就 bandwidth / 10

TCP Fairness

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



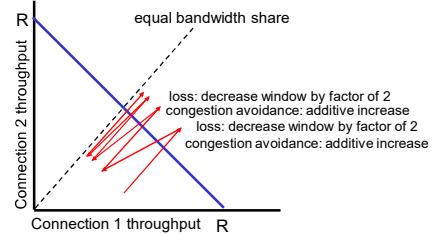
Transport Layer 3-103

103

Why is TCP fair?

two competing sessions:

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



Transport Layer 3-104

104

Fairness (more)

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

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 $R/2$

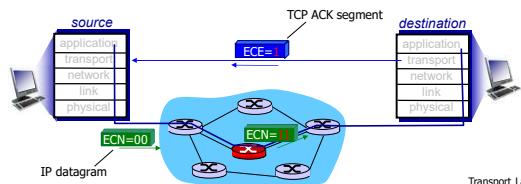
Transport Layer 3-105

105

Explicit Congestion Notification (ECN)

network-assisted congestion control:

- two bits in IP header (ToS field, type of service) marked by network router to indicate congestion
- congestion indication carried to receiving host
- receiver (seeing congestion indication in IP datagram) sets ECE bit on receiver-to-sender ACK segment to notify sender of congestion



Transport Layer 3-106

106

Chapter 3: summary

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

next:

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

Transport Layer 3-107

107

$$1500 - 20x = 1440$$

1500

x

500

$$100 - 20x = 440$$

100

$$15 \quad 1500 \quad 11$$

$$\begin{array}{|c|c|} \hline 440 & \\ \hline 0 & 0 \\ \hline \end{array}$$

$$\boxed{440}$$

$$\boxed{440}$$

0

1320

1320

$$(80)$$

Chapter 4

Chapter 4: outline

4.1 Overview of Network layer

- data plane
- control plane

4.2 What's inside a router

4.3 IP: Internet Protocol

- datagram format
- fragmentation
- IPv4 addressing
- network address translation
- IPv6

4.4 Generalized Forward and SDN

- match
- action
- OpenFlow examples of match-plus-action in action

Network Layer: Data Plane 4-1

1

Chapter 4: network layer

chapter goals:

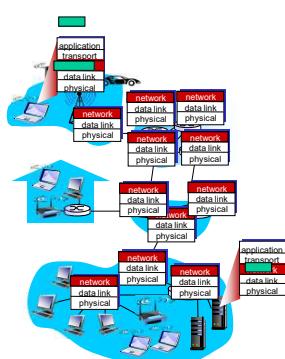
- understand principles behind network layer services, focusing on data plane:
 - network layer service models
 - forwarding versus routing
 - how a router works
 - generalized forwarding
- instantiation, implementation in the Internet

Network Layer: Data Plane 4-2

2

Network layer

- transport segment from sending to receiving host
- on sending side encapsulates segments into datagrams
- on receiving side, delivers segments to transport layer
- network layer protocols in **every** host, router
- router examines header fields in all IP datagrams passing through it



Network Layer: Data Plane 4-3

3

Two key network-layer functions

network-layer functions:

- **forwarding:** move packets from router's input to appropriate router output
- **routing:** determine route taken by packets from source to destination
 - routing algorithms

analogy: taking a trip

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

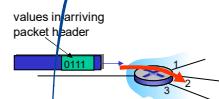
Network Layer: Data Plane 4-4

4

Network layer: data plane, control plane

Data plane

- local, per-router function
- determines how datagram arriving on router input port is forwarded to router output port
- **forwarding function**



Control plane

- **network-wide logic**
- determines how datagram is routed among routers along end-end path from source host to destination host
- two control-plane approaches:
 - **traditional routing algorithms:** implemented in routers
 - **software-defined networking (SDN):** implemented in (remote) servers

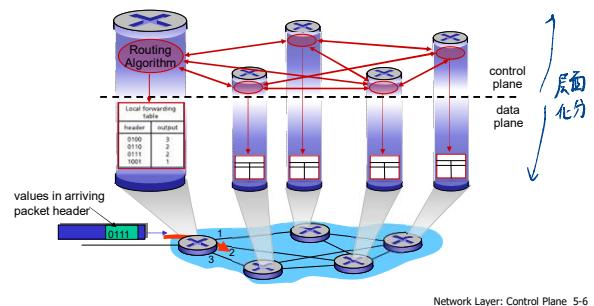
Network Layer: Data Plane 4-5

5

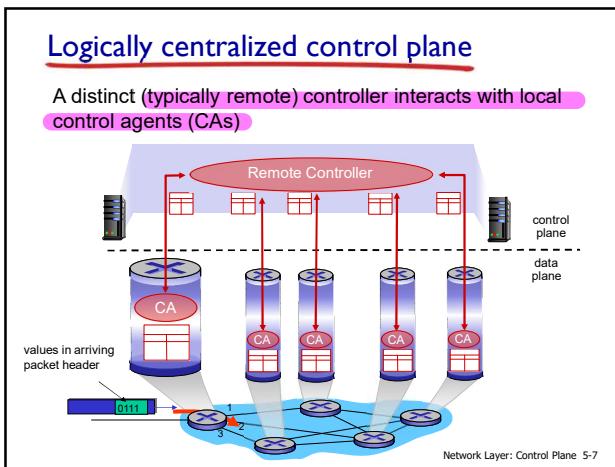
forwarding *routing*

Per-router control plane

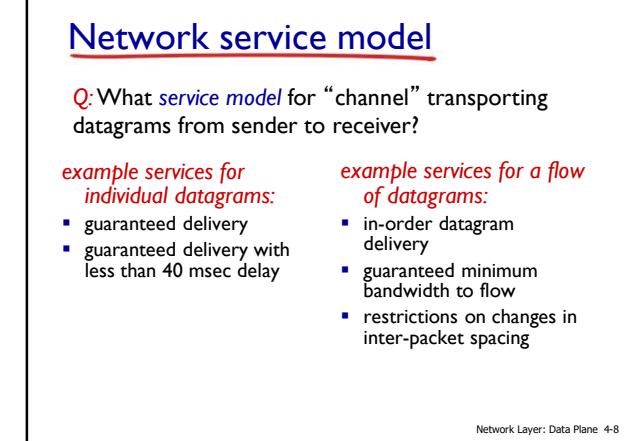
Individual routing algorithm components **in each and every router** interact in the control plane



6



7



8

Network layer service models:

Network Architecture	Service Model	Guarantees ?				Congestion feedback
		Bandwidth	Loss	Order	Timing	
Internet	best effort	none	no	no	no	no (inferred via loss)
ATM	CBR	constant rate	yes	yes	yes	no congestion
ATM	VBR	guaranteed rate	yes	yes	yes	no congestion
ATM	ABR	guaranteed minimum	no	yes	no	yes
ATM	UBR	none	no	yes	no	no

Network Layer: Data Plane 4-9

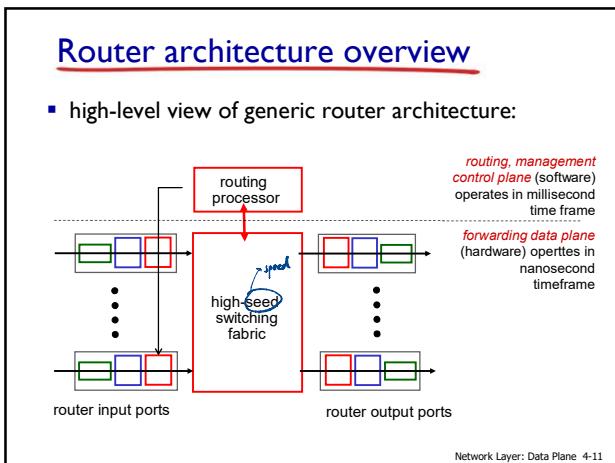
9

Chapter 4: outline

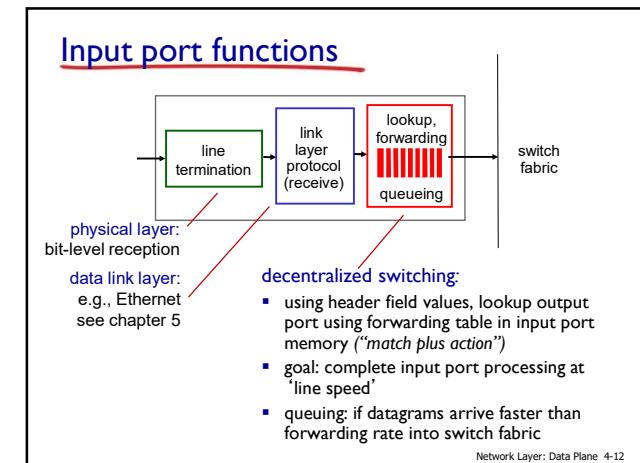
4.1 Overview of Network layer <ul style="list-style-type: none"> • data plane • control plane 	4.4 Generalized Forward and SDN <ul style="list-style-type: none"> • match • action • OpenFlow examples of match-plus-action in action
4.2 What's inside a router	4.3 IP: Internet Protocol <ul style="list-style-type: none"> • datagram format • fragmentation • IPv4 addressing • network address translation • IPv6

Network Layer: Data Plane 4-10

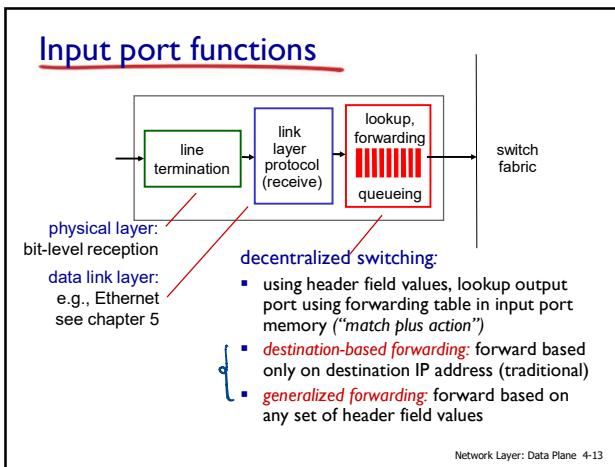
10



11



12



13

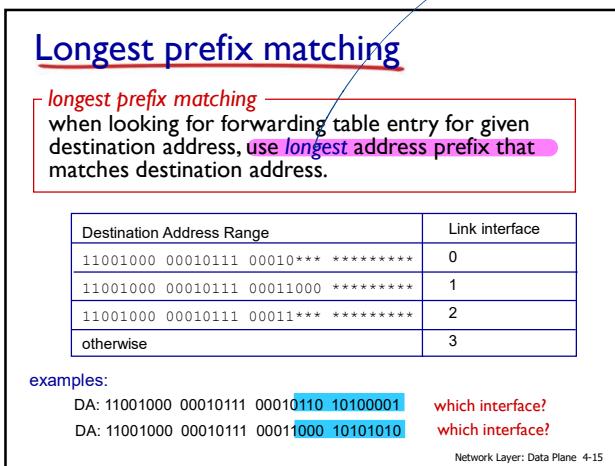
Destination-based forwarding

forwarding table	
Destination Address Range	Link Interface
11001000 00010111 00010000 00000000 through 11001000 00010111 00010111 11111111	0
11001000 00010111 00011000 00000000 through 11001000 00010111 00011000 11111111	1
11001000 00010111 00011001 00000000 through 11001000 00010111 00011111 11111111	2
otherwise	3

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

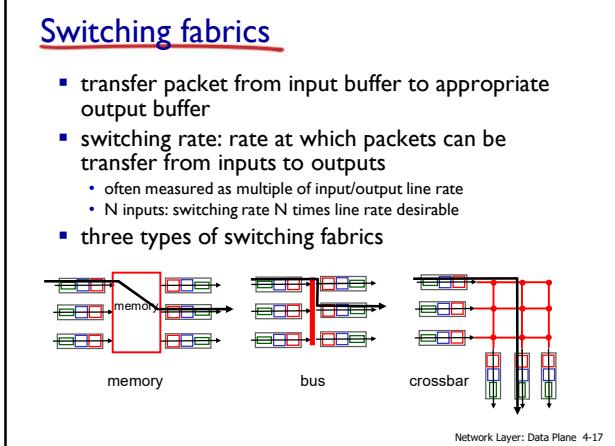
Network Layer: Data Plane 4-14

14



15 Router must has more than two interfaces

03/24/2022 starts:



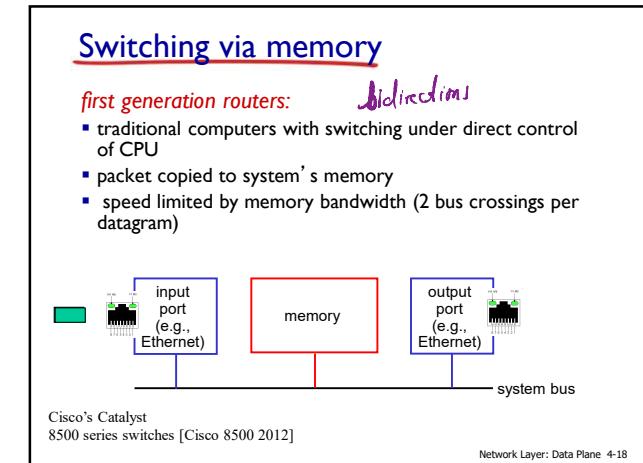
17

Longest prefix matching

- we'll see **why** longest prefix matching is used shortly, when we study addressing
- longest prefix matching: often performed using ternary content addressable memories (TCAMs)
 - content addressable:** present address to TCAM: retrieve address in one clock cycle, regardless of table size
 - Cisco Catalyst: can up ~1M routing table entries in TCAM

Network Layer: Data Plane 4-16

16



18

Switching via a bus

→ Wires

- datagram from input port memory to output port memory via a shared bus
- bus contention:** switching speed limited by bus bandwidth
- 32 Gbps bus, Cisco 5600: sufficient speed for access and enterprise routers

The Cisco 5600
[Cisco Switches 2012]

Network Layer: Data Plane 4-19

19

Switching via interconnection network

- overcome bus bandwidth limitations
- banyan networks, crossbar, other interconnection nets initially developed to connect processors in multiprocessor
- advanced design: fragmenting datagram into fixed length cells, switch cells through the fabric.
- Cisco 12000: switches 60 Gbps through the interconnection network

conductivity makes it slow

Network Layer: Data Plane 4-20

20

Input port queuing

queued up

- fabric slower than input ports combined -> queuing may occur at input queues
 - queueing delay and loss due to input buffer overflow!
- Head-of-the-Line (HOL) blocking:** queued datagram at front of queue prevents others in queue from moving forward

HOL:
eg: + 仔裡去待色的 output port
前面有红色的前向, PINK2X3.

Other example: 电影灯红灯亮

output port contention: only one red datagram can be transferred, lower red packet is blocked

one packet time later: green packet experiences HOL blocking

Network Layer: Data Plane 4-21

21

Output port queuing

at t , packets more from input to output

- buffering when arrival rate via switch exceeds output line speed
- queueing (delay) and loss due to output port buffer overflow!

Network Layer: Data Plane 4-23

23

Output ports

This slide is HUGELY important!

switch fabric → datagram buffer (queuing) → link layer protocol (send) → line termination

- buffering** required from fabric faster Datagram (packets) can be lost due to congestion, lack of buffers rate
- scheduling discipline** Priority scheduling – who gets best performance, network neutrality

Network Layer: Data Plane 4-22

22

How much buffering?

- RFC 3439 rule of thumb: average buffering equal to "typical" RTT (say 250 msec) times link capacity C
 - e.g., C = 10 Gbps link: 2.5 Gbit buffer
- recent recommendation: with N flows, buffering equal to

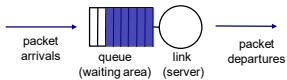
$$\frac{RTT \cdot C}{N}$$

Network Layer: Data Plane 4-24

24

Scheduling mechanisms

- scheduling:** choose next packet to send on link
- FIFO (first in first out) scheduling:** send in order of arrival to queue
 - real-world example?
 - discard policy:** if packet arrives to full queue: who to discard?
 - tail drop:** drop arriving packet
 - priority:** drop/remove on priority basis
 - random:** drop/remove randomly

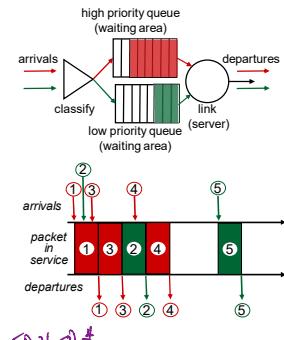


Network Layer: Data Plane 4-25

25

Scheduling policies: priority

- priority scheduling:** send highest priority queued packet
- multiple classes, with different priorities
 - class may depend on marking or other header info, e.g. IP source/dest, port numbers, etc.
 - real world example?



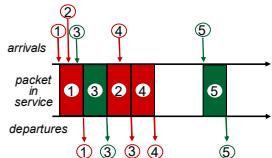
Network Layer: Data Plane 4-26

26

Scheduling policies: still more

Round Robin (RR) scheduling:

- multiple classes
- cyclically scan class queues, sending one complete packet from each class (if available)
- real world example?



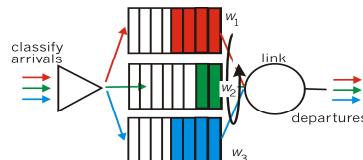
Network Layer: Data Plane 4-27

27

Scheduling policies: still more

Weighted Fair Queuing (WFQ):

- generalized Round Robin
- each class gets **weighted amount** of service in each cycle
- real-world example?
3%从红色出来, 20%从绿色出来



Network Layer: Data Plane 4-28

28

Chapter 4: outline

4.1 Overview of Network layer

- data plane
- control plane

4.2 What's inside a router

- #### 4.3 IP: Internet Protocol
- datagram format
 - fragmentation
 - IPv4 addressing
 - network address translation
 - IPv6

4.4 Generalized Forward and SDN

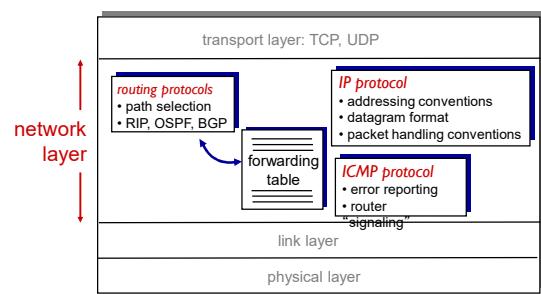
- match
- action
- OpenFlow examples of match-plus-action in action

Network Layer: Data Plane 4-29

29

The Internet network layer

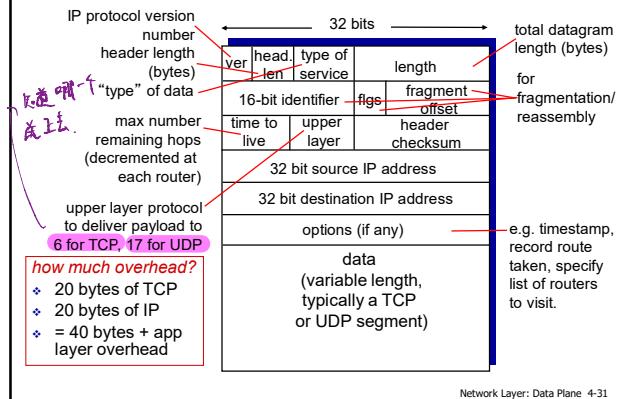
host, router network layer functions:



Network Layer: Data Plane 4-30

30

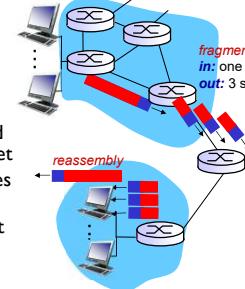
IP datagram format



31

IP fragmentation, reassembly

- network links have MTU (max.transfer size) - largest possible link-level frame
 - different link types, different MTUs
- large IP datagram divided ("fragmented") within net
 - one datagram becomes several datagrams
 - "reassembled" only at final destination
 - IP header bits used to identify, order related fragments



32

03/29/2020: color: start from ↓

Chapter 4: outline

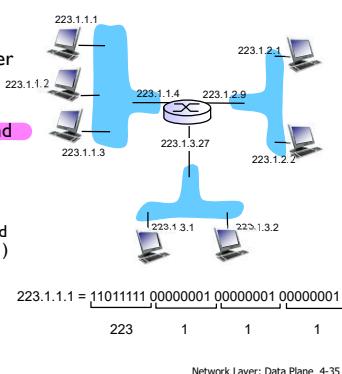
- 4.1 Overview of Network layer
 - data plane
 - control plane
- 4.2 What's inside a router
- 4.3 IP: Internet Protocol
 - datagram format
 - fragmentation
 - IPv4 addressing
 - network address translation
 - IPv6
- 4.4 Generalized Forward and SDN
 - match
 - action
 - OpenFlow examples of match-plus-action in action

Network Layer: Data Plane 4-34

34

IP addressing: introduction

- IP address:** 32-bit identifier for host, router interface
- interface:** connection between host/router and physical link
 - router's typically have multiple interfaces
 - host typically has one or two interfaces (e.g., wired Ethernet, wireless 802.11)
- IP addresses associated with each interface**



35

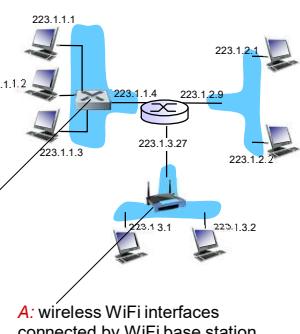
IP addressing: introduction

Q: how are interfaces actually connected?

A: we'll learn about that in chapter 5, 6.

A: wired Ethernet interfaces connected by Ethernet switches

For now: don't need to worry about how one interface is connected to another (with no intervening router)



36

03/24/2022

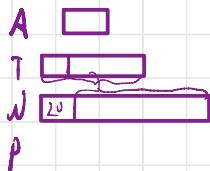
- Why we need queue at input port.
- Because we need queue to buffer the incoming packet if switch speed is slower than incoming packet rate. And if that happens without queue, then packet will loss
- Why we need queue at output port?
- If output port packet rate is faster than switch output going port speed, then packet may get lost without queue.

Higher priority packets: the packet is **router information** and generated by router, not users. Government packet for **security** from FBI, NSL.

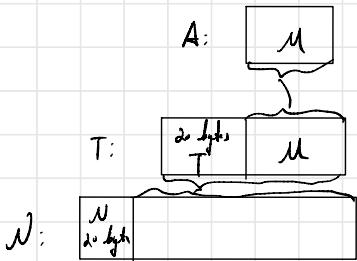
★ The priority queue will be on the exam. Carefully review.

Why there is no packet length in transport layer?

Because Network layer has the packet length, and use it minus 20 bytes get the payload, which is the length of TCP packet length.



$$\text{Length of } N \text{ packets} - 20 \text{ bytes} = \text{Segment Length}$$

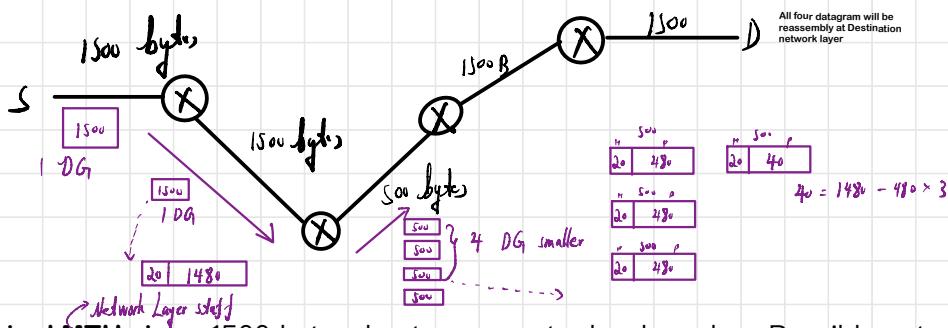


女：考題

Should I expected fragmentation at ending system?

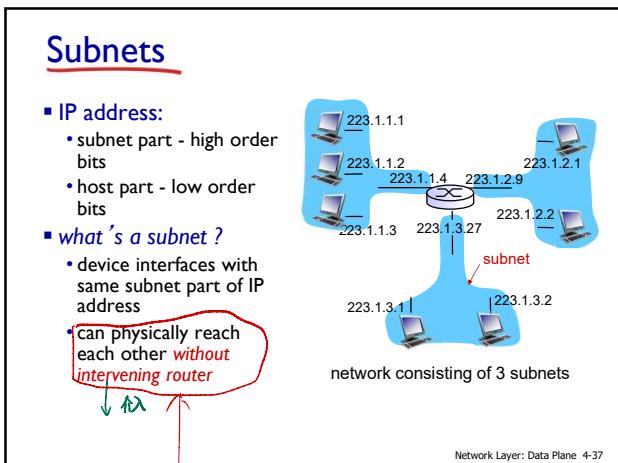
Chop up is never be the job of network layer, it is transport layer work. but it is capable to do it. Network layer job is host to host delivery, and we should give the job to transport layer. And much pressure will be on end system.

Fragmentation and reassembly: The below is network layer links

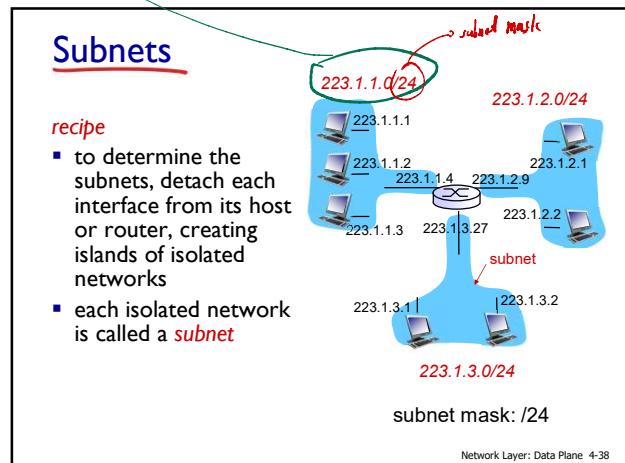


Typical MTU size: 1500 bytes, best guess not a hard number. Possible network layer packet has to be chopped up. (网时是整个路上最小的取值)

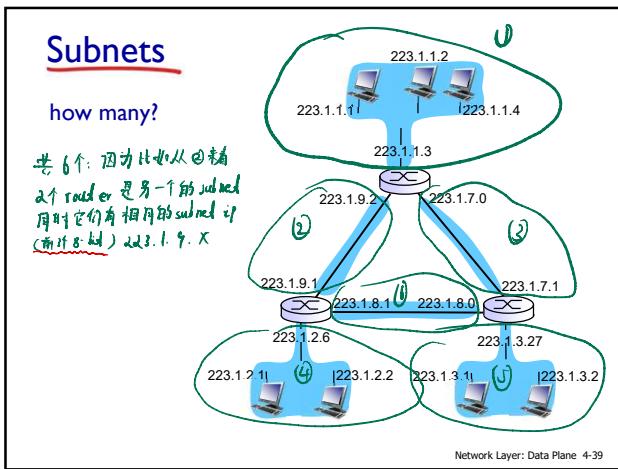
Offset calculation refer to slides.



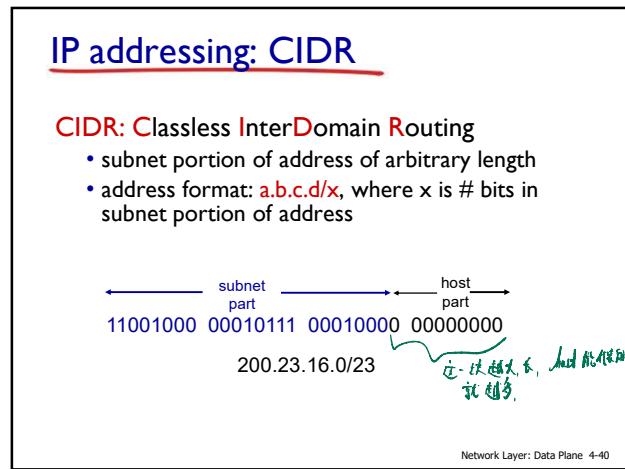
37

definition of subnet

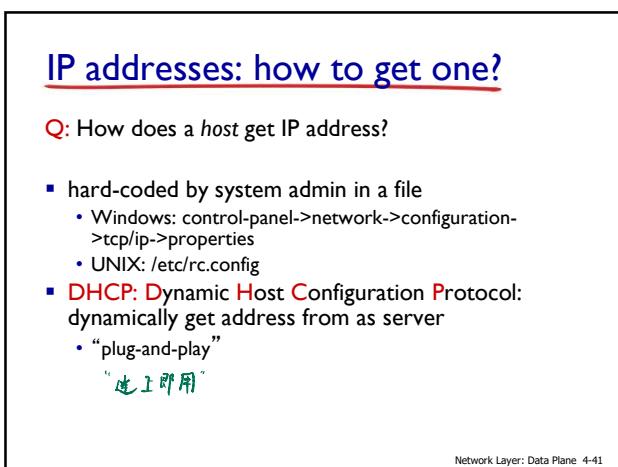
38



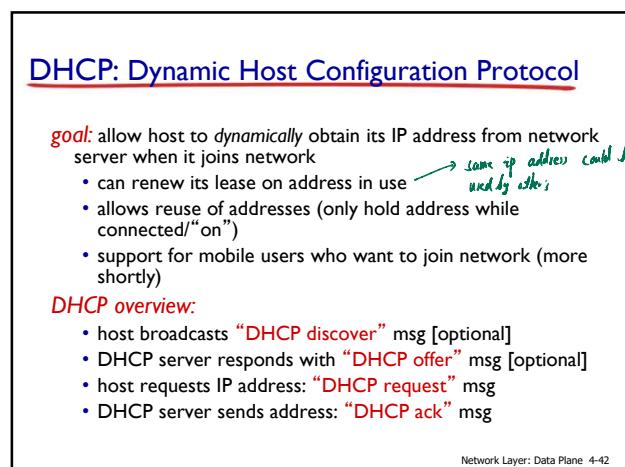
39



40

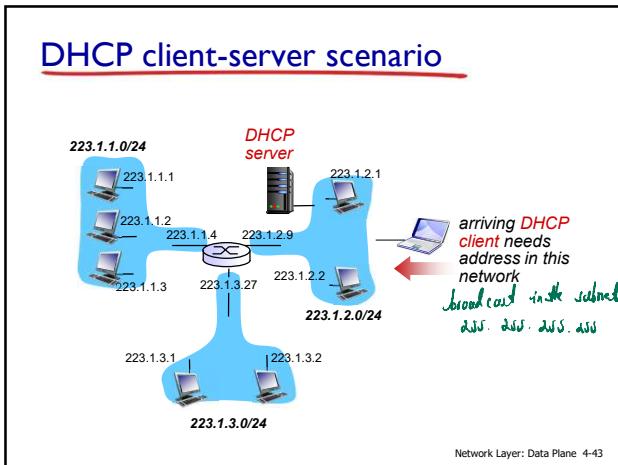


41



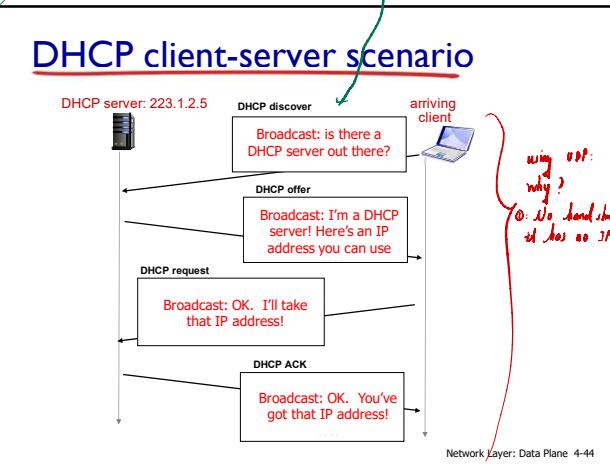
42

Scenario: Lease problem. DHCP gives you a lease to use the IP address. But it may reject you if the ip has been given to others

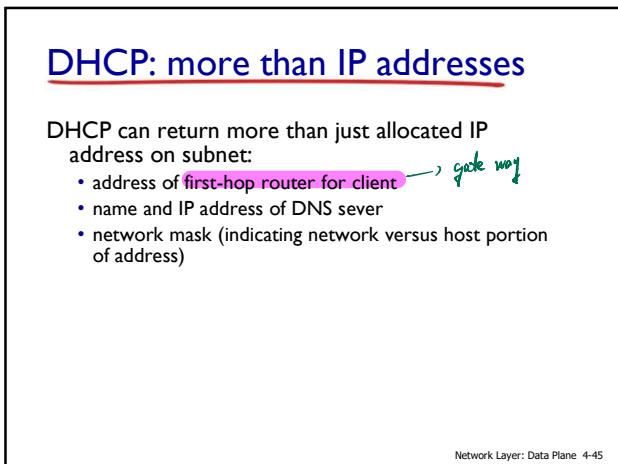


43

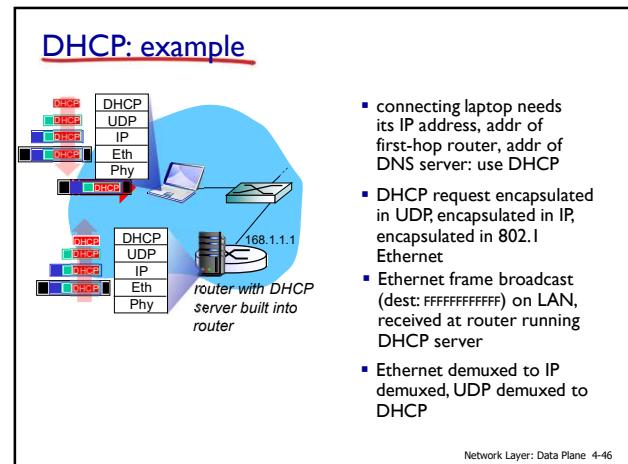
Cache 数量和gateway相同



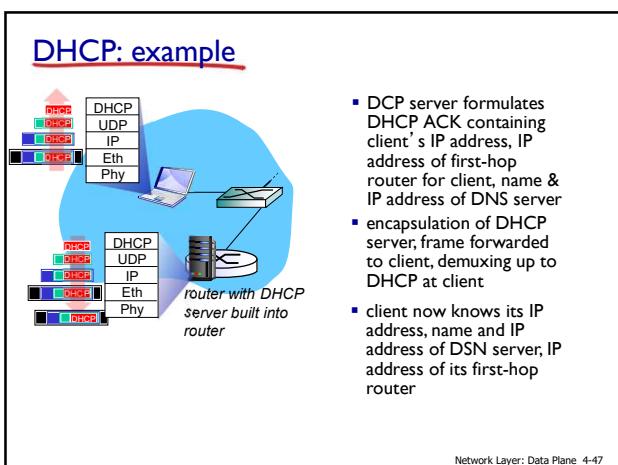
44



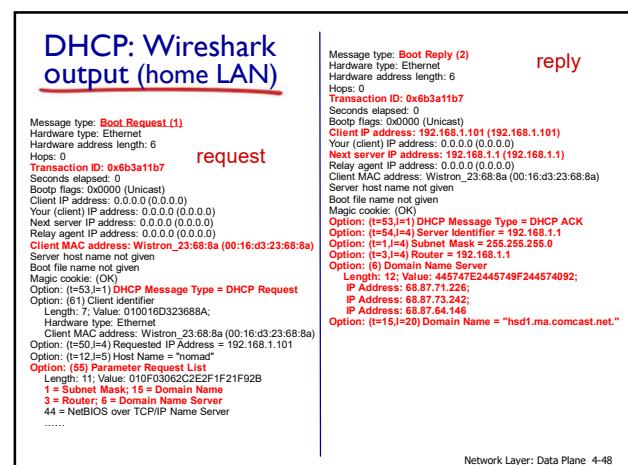
45



46



47



48

8

IP addresses: how to get one?

Q: how does network get subnet part of IP addr?

A: gets allocated portion of its provider ISP's address space

e.g. How many IP just bought?
2¹⁶

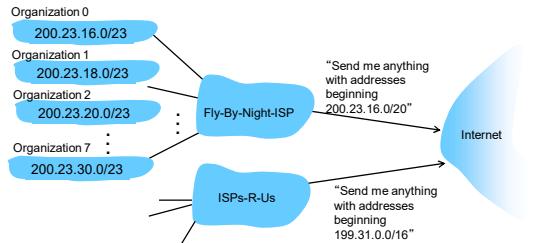
ISP's block	11001000 00010111 00010000 00000000	200.23.16.0/20
Organization 0	11001000 00010111 00010000 00000000	200.23.16.0/23
Organization 1	11001000 00010111 00010010 00000000	200.23.18.0/23
Organization 2	11001000 00010111 00010100 00000000	200.23.20.0/23
...
Organization 7	11001000 00010111 00011110 00000000	200.23.30.0/23

Network Layer: Data Plane 4-49

49

Hierarchical addressing: route aggregation

hierarchical addressing allows efficient advertisement of routing information:

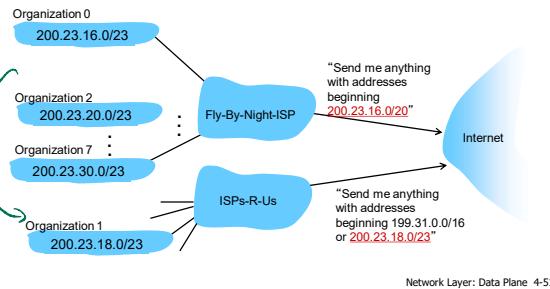


Network Layer: Data Plane 4-50

50

Hierarchical addressing: more specific routes

ISPs-R-U has a more specific route to Organization 1



Network Layer: Data Plane 4-51

51

如果 org 1 移到了下面的 ISP, 那么 Local prefix Match
会首先命中到上面的 org. 1 的 longer prefix, model B/ 下面的 org
因为 1 是 .../16, 但是 .../23

IP addressing: the last word...

Q: how does an ISP get block of addresses?

A: ICANN: Internet Corporation for Assigned

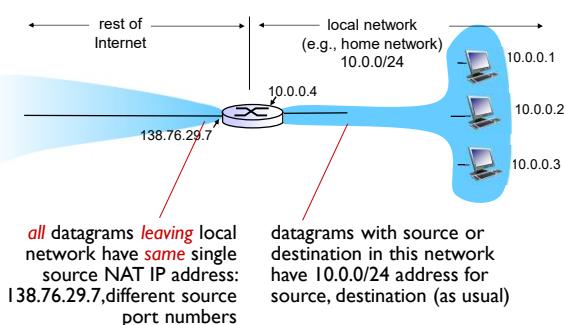
Names and Numbers <http://www.icann.org/>

- allocates addresses
- manages DNS
- assigns domain names, resolves disputes

Network Layer: Data Plane 4-52

52

NAT: network address translation



Network Layer: Data Plane 4-53

53

NAT: network address translation

motivation: local network uses just one IP address as far as outside world is concerned:

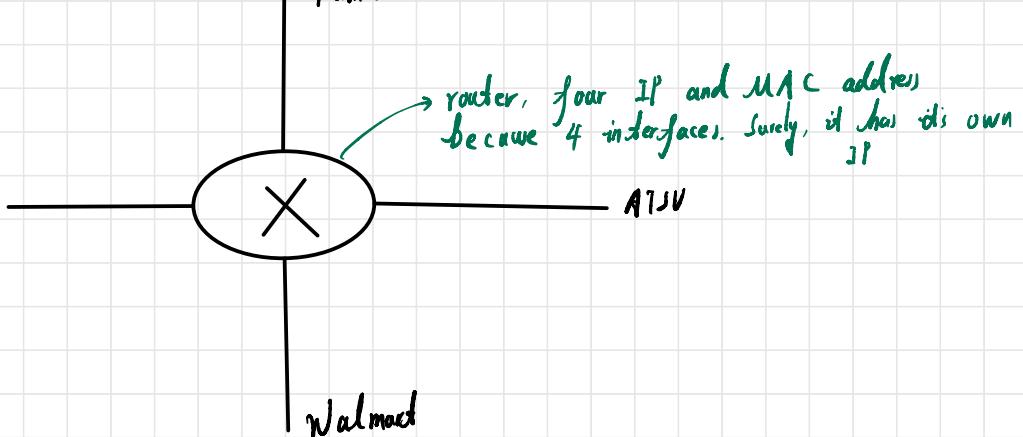
- range of addresses not needed from ISP: just one IP address for all devices
- can change addresses of devices in local network without notifying outside world
- can change ISP without changing addresses of devices in local network
- devices inside local net not explicitly addressable, visible by outside world (a security plus)

Network Layer: Data Plane 4-54

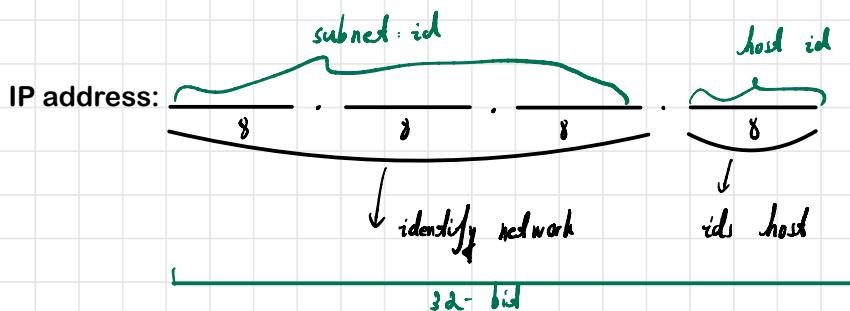
54

03/29/2022

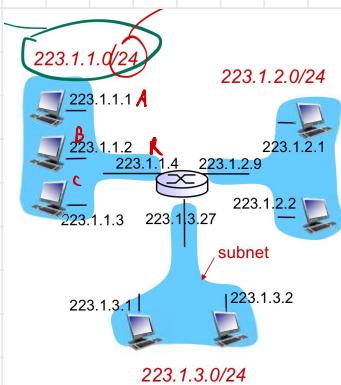
Truman



Subnet mask:



eg: 192.6.223.9 / 24, 前3块是 IP of subnet
后一块(8 bit)是 subnet mask



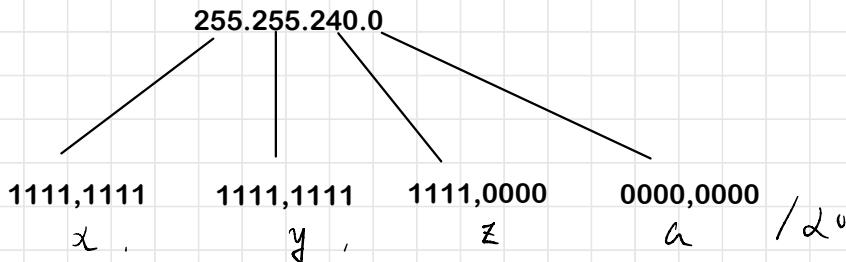
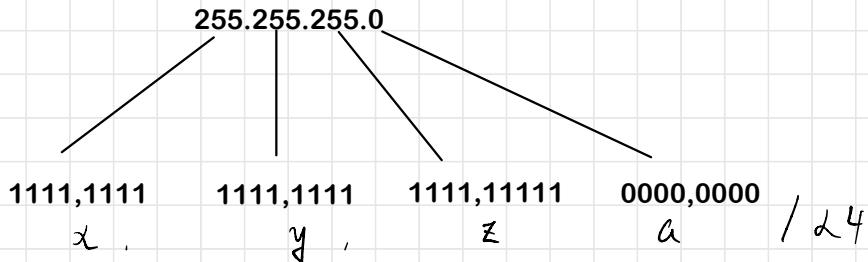
A can directly talk to B, C, D, because they have the same subnet IP



We need IP address to analysis if two pc are under same subnet, but only use MAC address to communicate.

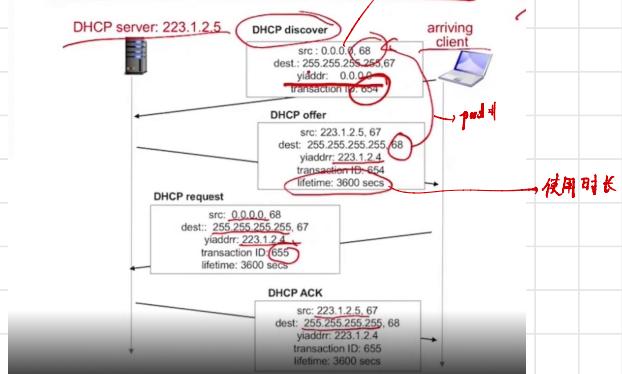
When it sends data by MAC address it, all node in the Wi-Fi receives it but only the pc with the MAC address accept it.

CIDR: real example for subnet mask



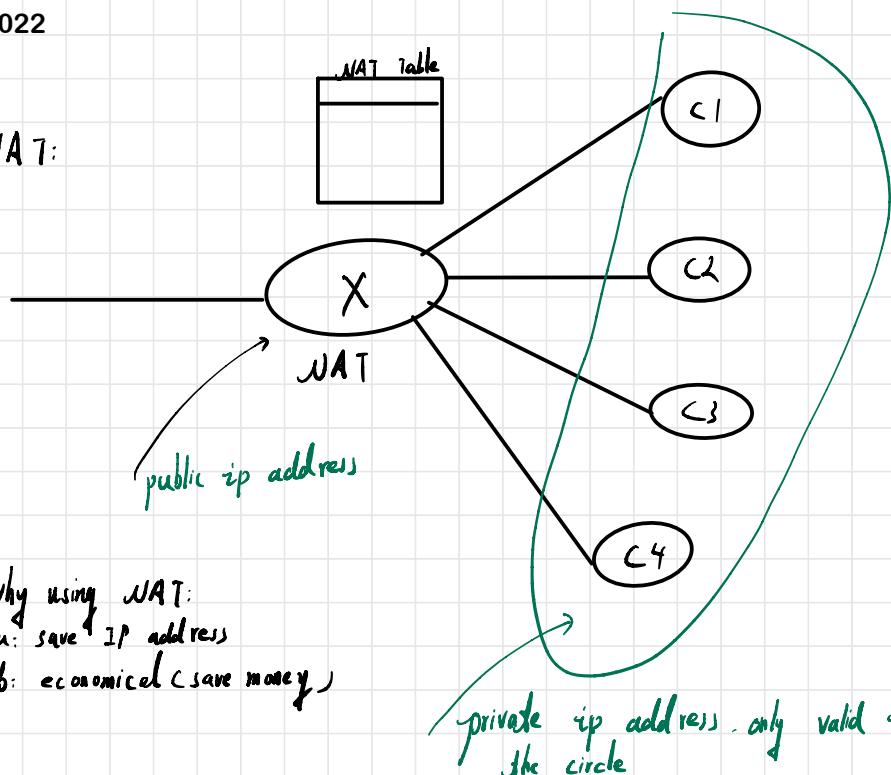
DHCP:

DHCP client-server scenario



03/29/2022

NAT:



Q: Why using NAT:

- a: save IP address
- b: economical (save money)

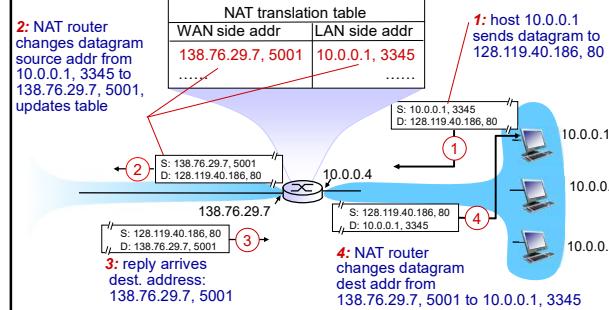
NAT: network address translation

implementation: NAT router must:

- **outgoing datagrams:** replace (source IP address, port #) of every outgoing datagram to (NAT IP address, new port #)
... remote clients/servers will respond using (NAT IP address, new port #) as destination addr
- remember (in **NAT translation table**) every (source IP address, port #) to (NAT IP address, new port #) translation pair
- **incoming datagrams:** replace (NAT IP address, new port #) in dest fields of every incoming datagram with corresponding (source IP address, port #) stored in NAT table

Network Layer: Data Plane 4-55

NAT: network address translation



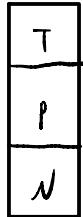
Network Layer: Data Plane 4-56

56

NAT: network address translation

- 16-bit port-number field:
 - 60,000 simultaneous connections with a single LAN-side address!
- NAT is controversial:
 - routers should only process up to layer 3 (**CP**)
 - address shortage should be solved by IPv6
 - violates end-to-end argument
 - NAT possibility must be taken into account by app designers, e.g., P2P applications
 - NAT traversal: what if client wants to connect to server behind NAT?

Network Layer: Data Plane 4-57



57

↓ part for forwarding

Chapter 4: outline

- | | |
|---|---|
| 4.1 Overview of Network layer <ul style="list-style-type: none"> • data plane • control plane 4.2 What's inside a router <ul style="list-style-type: none"> • datagram format • fragmentation • IPv4 addressing • network address translation • IPv6 | 4.4 Generalized Forward and SDN <ul style="list-style-type: none"> • match • action • OpenFlow examples of match-plus-action in action |
|---|---|

Network Layer: Data Plane 4-58

58

IPv6: motivation

- **initial motivation:** 32-bit address space soon to be completely allocated.
- additional motivation:
 - header format helps speed processing/forwarding
 - header changes to facilitate QoS

IPv6 datagram format:

- fixed-length 40 byte header
- no fragmentation allowed

Network Layer: Data Plane 4-59

59

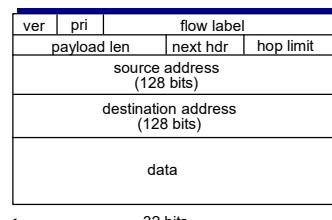
IPv6 datagram format

priority: identify priority among datagrams in flow

flow Label: identify datagrams in same "flow."

(concept of "flow" not well defined).

next header: identify upper layer protocol for data



Network Layer: Data Plane 4-60

60

Other changes from IPv4

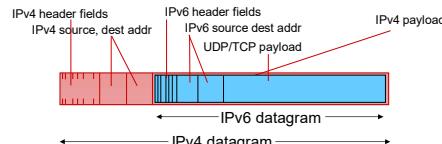
- **checksum**: removed entirely to reduce processing time at each hop
- **options**: allowed, but outside of header, indicated by “Next Header” field
- **ICMPv6**: new version of ICMP
 - additional message types, e.g. “Packet Too Big”
 - multicast group management functions

Network Layer: Data Plane 4-61

61

Transition from IPv4 to IPv6

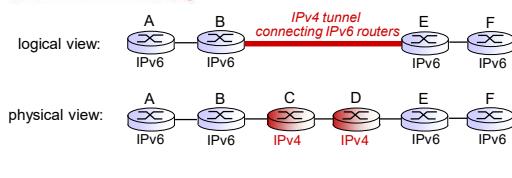
- not all routers can be upgraded simultaneously
 - no “flag days”
 - how will network operate with mixed IPv4 and IPv6 routers?
- **tunneling**: IPv6 datagram carried as payload in IPv4 datagram among IPv4 routers



Network Layer: Data Plane 4-62

62

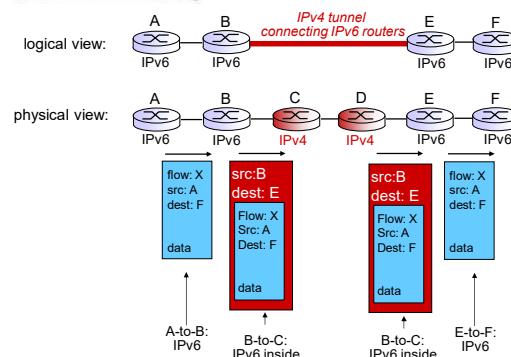
Tunneling



Network Layer: Data Plane 4-63

63

Tunneling



Network Layer: Data Plane 4-64

64

IPv6: adoption

- Google: 8% of clients access services via IPv6
- NIST: 1/3 of all US government domains are IPv6 capable
- **Long (long!) time for deployment, use**
 - 20 years and counting!
 - think of application-level changes in last 20 years: WWW, Facebook, streaming media, Skype, ...
 - **Why?**

Network Layer: Data Plane 4-65

65

IPv6: adoption



Network Layer: Data Plane 4-66

66

Chapter 4: outline

4.1 Overview of Network layer

- data plane
- control plane

4.2 What's inside a router

4.3 IP: Internet Protocol

- datagram format
- fragmentation
- IPv4 addressing
- network address translation
- IPv6

4.4 Generalized Forward and SDN

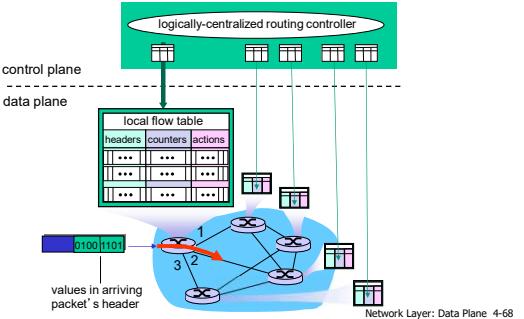
- match
- action
- OpenFlow examples of match-plus-action in action

Network Layer: Data Plane 4-67

67

Generalized Forwarding and SDN

Each router contains a *flow table* that is computed and distributed by a *logically centralized routing controller*



68

OpenFlow data plane abstraction

- **flow:** defined by header fields
- **generalized forwarding:** simple packet-handling rules
 - **Pattern:** match values in packet header fields
 - **Actions:** for matched packet: drop, forward, modify, matched packet or send matched packet to controller
 - **Priority:** disambiguate overlapping patterns
 - **Counters:** #bytes and #packets



Flow table in a router (computed and distributed by controller) define router's match+action rules

Network Layer: Data Plane 4-69

69

OpenFlow data plane abstraction

- **flow:** defined by header fields
- **generalized forwarding:** simple packet-handling rules
 - **Pattern:** match values in packet header fields
 - **Actions:** for matched packet: drop, forward, modify, matched packet or send matched packet to controller
 - **Priority:** disambiguate overlapping patterns
 - **Counters:** #bytes and #packets

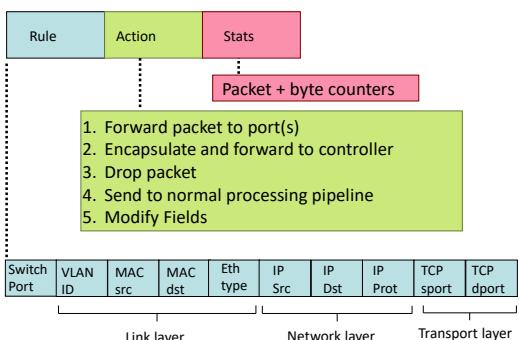


* : wildcard

1. src=1.2.*.* , dest=3.4.5.* → drop
2. src = *.*.*.* , dest=3.4.*.* → forward(2)
3. src=10.1.2.3 , dest=*.*.*.* → send to controller

70

OpenFlow: Flow Table Entries



71

Examples

Destination-based forwarding:

Switch Port	MAC src	MAC dst	Eth type	VLAN ID	IP Src	IP Dst	IP Prot	TCP sport	TCP dport	Action
*	*	*	*	*	*	51.6.0.8	*	*	*	port6

IP datagrams destined to IP address 51.6.0.8 should be forwarded to router output port 6

Firewall:

Switch Port	MAC src	MAC dst	Eth type	VLAN ID	IP Src	IP Dst	IP Prot	TCP sport	TCP dport	Forward
*	*	*	*	*	*	*	*	*	*	22 drop

do not forward (block) all datagrams destined to TCP port 22

Switch Port	MAC src	MAC dst	Eth type	VLAN ID	IP Src	IP Dst	IP Prot	TCP sport	TCP dport	Action
*	*	*	*	*	128.119.1.1	*	*	*	*	drop

do not forward (block) all datagrams sent by host 128.119.1.1

72

12

Examples

Destination-based layer 2 (switch) forwarding:

Switch Port	MAC src	MAC dst	Eth type	VLAN ID	IP Src	IP Dst	TCP Prot	TCP sport	TCP dport	Action
*	22:A7:23: 11:E1:02	*	*	*	*	*	*	*	*	port3

layer 2 frames from MAC address 22:A7:23:11:E1:02
should be forwarded to output port 6

Network Layer: Data Plane 4-73

73

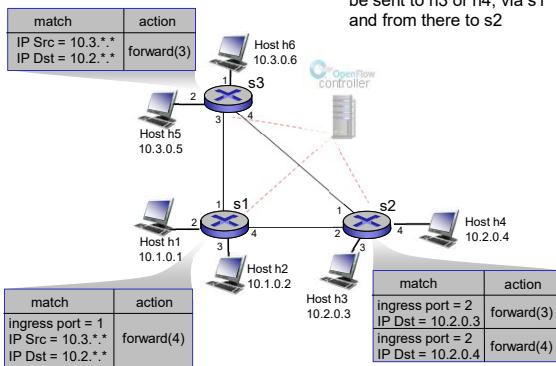
OpenFlow abstraction

- **match+action:** unifies different kinds of devices
- Router
 - **match:** longest destination IP prefix
 - **action:** forward out a link
- Switch
 - **match:** destination MAC address
 - **action:** forward or flood
- Firewall
 - **match:** IP addresses and TCP/UDP port numbers
 - **action:** permit or deny
- NAT
 - **match:** IP address and port
 - **action:** rewrite address and port

Network Layer: Data Plane 4-74

74

OpenFlow example



75

Chapter 4: done!

- 4.1 Overview of Network layer: data plane and control plane
- 4.2 What's inside a router
- 4.3 IP: Internet Protocol
 - datagram format
 - fragmentation
 - IPv4 addressing
 - NAT
 - IPv6

4.4 Generalized Forward and SDN

- match plus action
- OpenFlow example

Question: how do forwarding tables (destination-based forwarding) or flow tables (generalized forwarding) computed?

Answer: by the control plane (next chapter)

Network Layer: Data Plane 4-76

76

Chapter-5

Chapter 5: network layer control plane

chapter goals: understand principles behind network control plane

- traditional routing algorithms
- SDN controllers
- Internet Control Message Protocol
- network management

and their instantiation, implementation in the Internet:

- OSPF, BGP, OpenFlow, ODL and ONOS controllers, ICMP, SNMP

Network Layer: Control Plane 5-1

1

Chapter 5: outline

- | | |
|--|---|
| 5.1 introduction
5.2 routing protocols | 5.5 The SDN control plane
5.6 ICMP: The Internet Control Message Protocol
5.7 Network management and SNMP |
| ▪ link state
▪ distance vector
5.3 intra-AS routing in the Internet: OSPF
5.4 routing among the ISPs: BGP | |

Network Layer: Control Plane 5-2

2

Network-layer functions

Recall: two network-layer functions:

- **forwarding:** move packets from router's input to appropriate router output *data plane*
- **routing:** determine route taken by packets from source to destination *control plane*

Two approaches to structuring network control plane:

- {
- per-router control (traditional)
- logically centralized control (software defined networking)

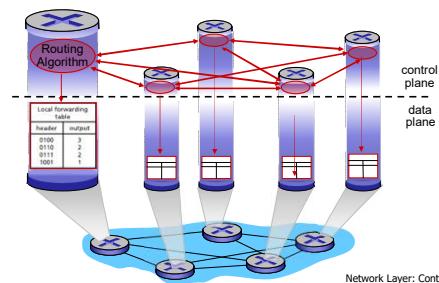
Network Layer: Control Plane 5-3

3

network control planes:
per-router
logically centralized control

Per-router control plane

Individual routing algorithm components *in each and every router* interact with each other in control plane to compute forwarding tables

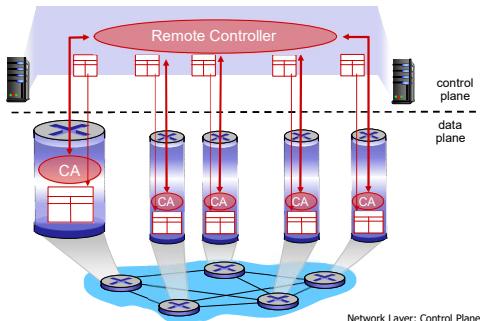


Network Layer: Control Plane 5-4

4

Logically centralized control plane or SDN

A distinct (typically remote) controller interacts with local control agents (CAs) in routers to compute forwarding tables



Network Layer: Control Plane 5-5

5

All the tables will be consistent.

Chapter 5: outline

- | | |
|--|---|
| 5.1 introduction
5.2 routing protocols | 5.5 The SDN control plane
5.6 ICMP: The Internet Control Message Protocol
5.7 Network management and SNMP |
| ▪ link state
▪ distance vector
5.3 intra-AS routing in the Internet: OSPF
5.4 routing among the ISPs: BGP | |

Network Layer: Control Plane 5-6

6

Routing protocols

Routing protocol goal: determine “good” paths (equivalently, routes), from sending hosts to receiving host, through network of routers

- path: sequence of routers packets will traverse in going from given initial source host to given final destination host
- “good”: least “cost”, “fastest”, “least congested”
- routing: a “top-10” networking challenge!

→ less cost

Network Layer: Control Plane 5-7

7

Graph abstraction of the network

graph: $G = (N, E)$
 $N = \{u, v, w, x, y, z\}$
 $E = \{(u, v), (u, x), (v, x), (v, w), (w, z), (x, w), (x, y), (y, z), (w, y), (w, z)\}$

aside: graph abstraction is useful in other network contexts, e.g., P2P, where N is set of peers and E is set of TCP connections

Network Layer: Control Plane 5-8

8

Graph abstraction: costs

$c(x, x') = \text{cost of link } (x, x')$
e.g., $c(w, z) = 5$

cost could always be 1, or inversely related to bandwidth, or inversely related to congestion

cost of path $(x_1, x_2, x_3, \dots, x_p) = c(x_1, x_2) + c(x_2, x_3) + \dots + c(x_{p-1}, x_p)$

key question: what is the least-cost path between u and z?
routing algorithm: algorithm that finds that least cost path

Network Layer: Control Plane 5-9

9

Routing algorithm classification

Q: global or decentralized information?

- **global:**
 - all routers have complete topology, link cost info
 - “link state” algorithms
- **decentralized:**
 - router knows physically-connected neighbors, link costs to neighbors
 - iterative process of computation, exchange of info with neighbors
 - “distance vector” algorithms

Q: static or dynamic?

- **static:**
 - routes change slowly over time
- **dynamic:**
 - routes change more quickly
 - periodic update
 - in response to link cost changes

Network Layer: Control Plane 5-10

10

Chapter 5: outline

5.1 introduction	5.5 The SDN control plane
5.2 routing protocols	5.6 ICMP: The Internet Control Message Protocol
▪ link state	5.7 Network management and SNMP
▪ distance vector	
5.3 intra-AS routing in the Internet: OSPF	
5.4 routing among the ISPs: BGP	

Network Layer: Control Plane 5-11

11

A link-state routing algorithm

Dijkstra's algorithm

- net topology, link costs known to all nodes
 - accomplished via “link state broadcast”
 - all nodes have same info
- computes least cost paths from one node (‘source’) to all other nodes
 - gives **forwarding table** for that node
- iterative: after k iterations, know least cost path to k dest.’s

notation:

- $c(x, y)$: link cost from node x to y ; $= \infty$ if not direct neighbors
- $D(v)$: current value of cost of path from source to dest. v
- $p(v)$: predecessor node along path from source to v
- N^* : set of nodes whose least cost path definitively known

Network Layer: Control Plane 5-12

12

2

Dijkstra's algorithm

```

1 Initialization:
2 N' = {u}
3 for all nodes v
4   if v adjacent to u
5     then D(v) = c(u,v)
6   else D(v) = ∞
7
8 Loop
9  find w not in N' such that D(w) is a minimum
10 add w to N'
11 update D(v) for all v adjacent to w and not in N':
12   D(v) = min( D(v), D(w) + c(w,v) )
13 /* new cost to v is either old cost to v or known
14 shortest path cost to w plus cost from w to v */
15 until all nodes in N'

```

Network Layer: Control Plane 5-13

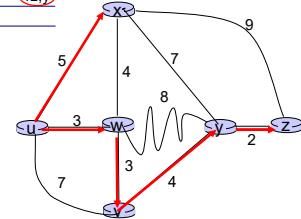
13

Dijkstra's algorithm: example

Step	N'	D(v)	D(w)	D(x)	D(y)	D(z)
		p(v)	p(w)	p(x)	p(y)	p(z)
0	u	7,u	3,u	5,u	∞	∞
1	uw	6,w	5,u	11,w	∞	∞
2	uwx	6,w	5,u	11,w	14,x	∞
3	uwvx	6,w	5,u	10,y	14,x	∞
4	uwvxy	6,w	5,u	12,y	14,x	∞
5	uwvxyz	6,w	5,u	12,y	14,x	9

notes:

- construct shortest path tree by tracing predecessor nodes
- ties can exist (can be broken arbitrarily)

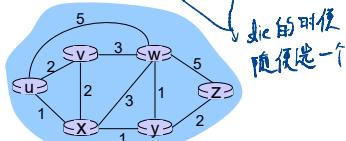


Network Layer: Control Plane 5-14

14

Dijkstra's algorithm: another example

Step	N'	D(v),p(v)	D(w),p(w)	D(x),p(x)	D(y),p(y)	D(z),p(z)
0	u	2,u	5,u	1,u	∞	∞
1	ux	2,u	4,x	2,y	∞	∞
2	uxy	2,u	3,y	3,y	4,y	∞
3	uxyy	2,u	3,y	3,y	4,y	∞
4	uxyvw	2,u	3,y	4,y	4,y	∞
5	uxyvwz	2,u	3,y	4,y	4,y	∞



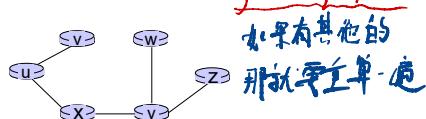
* Check out the online interactive exercises for more examples: http://gaia.cs.umass.edu/kurose_ross/interactive/

Network Layer: Control Plane 5-15

15

Dijkstra's algorithm: example (2)

resulting shortest-path tree from u:



resulting forwarding table in u:

destination	link
v	(u,v)
x	(u,x)
y	(u,x)
w	(u,x)
z	(u,x)

Network Layer: Control Plane 5-16

16

even y is not directed connecting to u

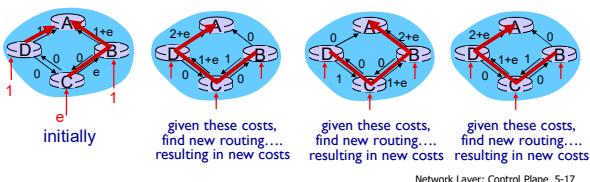
Dijkstra's algorithm, discussion

algorithm complexity: n nodes

- each iteration: need to check all nodes, w, not in N
- $n(n+1)/2$ comparisons: $O(n^2)$
- more efficient implementations possible: $O(n \log n)$

oscillations possible:

- e.g., support link cost equals amount of carried traffic:



Network Layer: Control Plane 5-17

17

based on current path
not previous one, and calculate
it before updating it.

Chapter 5: outline

- 5.1 introduction
- 5.2 routing protocols
 - link state
 - distance vector
- 5.3 intra-AS routing in the Internet: OSPF
- 5.4 routing among the ISPs: BGP
- 5.5 The SDN control plane
- 5.6 ICMP: The Internet Control Message Protocol
- 5.7 Network management and SNMP

18

Network Layer: Control Plane 5-18

Distance vector algorithm

Bellman-Ford equation (dynamic programming)

let
 $d_x(y) := \text{cost of least-cost path from } x \text{ to } y$
then

$$d_x(y) = \min_v \{ c(x,v) + d_v(y) \}$$

↓ cost from neighbor v to destination y
 ↓ cost to neighbor v
 min taken over all neighbors v of x

Network Layer: Control Plane 5-19

19

if $z \rightarrow u$, y will still be w , x , not u .
 Since x, w are y 's direct neighbor.
 If new least cost path is found, then update the table.
 If not, do nothing.

Bellman-Ford example

clearly, $d_v(z) = 5$, $d_x(z) = 3$, $d_w(z) = 3$

B-F equation says:

$$\begin{aligned} d_u(z) &= \min \{ c(u,v) + d_v(z), \\ &\quad c(u,x) + d_x(z), \\ &\quad c(u,w) + d_w(z) \} \\ &= \min \{ 2+5, \\ &\quad 1+3, \\ &\quad 5+3 \} = 4 \end{aligned}$$

node achieving minimum is next hop in shortest path, used in forwarding table

Network Layer: Control Plane 5-20

20

Distance vector algorithm

- $D_x(y) = \text{estimate of least cost from } x \text{ to } y$
 - x maintains distance vector $D_x = [D_x(y): y \in N]$
- node x :
 - knows cost to each neighbor v : $c(x,v)$
 - maintains its neighbors' distance vectors. For each neighbor v , x maintains $D_v = [D_v(y): y \in N]$

Network Layer: Control Plane 5-21

21

Distance vector algorithm

key idea:

- from time-to-time, each node sends its own distance vector estimate to neighbors
- when x receives new DV estimate from neighbor, it updates its own DV using B-F equation:

$$D_x(y) \leftarrow \min_v \{ c(x,v) + D_v(y) \} \text{ for each node } y \in N$$

- under minor, natural conditions, the estimate $D_x(y)$ converge to the actual least cost $d_x(y)$

Network Layer: Control Plane 5-22

22

Distance vector algorithm

iterative, asynchronous:
 each local iteration caused by:

- local link cost change
- DV update message from neighbor

distributed:
 each node notifies neighbors *only* when its DV changes

- neighbors then notify their neighbors if necessary

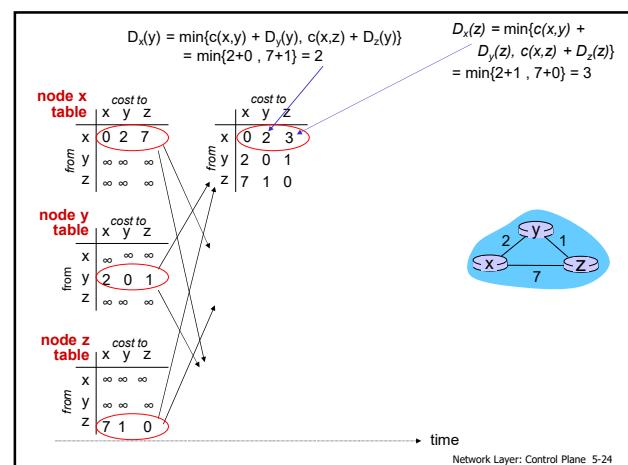
each node:

```

  wait for (change in local link cost or msg from neighbor)
  recomputes estimates
  if DV to any dest has changed, notify neighbors
  
```

Network Layer: Control Plane 5-23

23



24

- Routing is a P2P network.

Logically centralized control: pros: tables are consistent, less routing data transferring in network. Cons: central control is down.

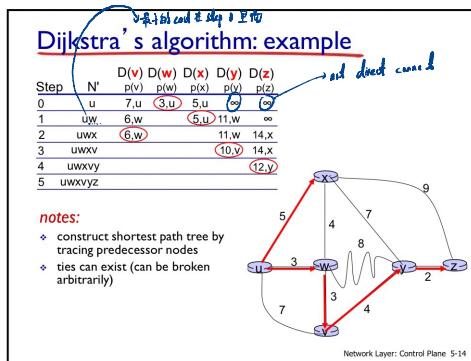
- Routing algorithms
 - Global: linked state, you know the entire topology.
 - Good for: small or medium size of network.
 - Broadcast for each router causing too much traffic. Router knows all others.
 - Good for static networks, since no too much traffic*
 - Decentralized: distance vectors, you only know neighbors
 - Good for: large size network
 - Good for dynamic network.*
- Dijkstra's algorithm — Global, linked state

notation:

- $c(x,y)$: link cost from node x to y ; $= \infty$ if not direct neighbors
- $D(v)$: current value of cost of path from source to dest. v
- $p(v)$: predecessor node along path from source to v
- N^* : set of nodes whose least cost path definitively known

ex:

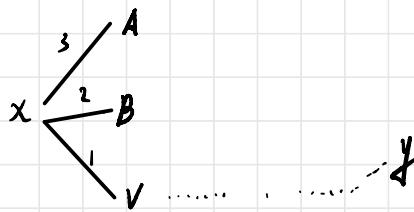
其他他每一次都是在找least cost的link 把所有的nodes连接起来



04/05/2022

Bellman-Ford

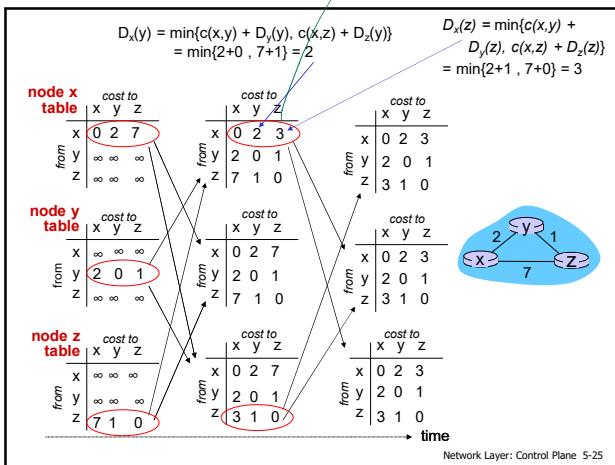
concept:



选中x至v, 至于怎么让他到y, 那是v node的事情, 这也是和Linked state的区别, linked state知道所有的node信息, dan sh

04/07/2022 starts here

逐 x 收到新的 link cost
更新其 link cost and cost path
如果更新 cost 则更新 (x & x node ts)

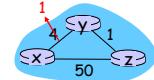


25

Distance vector: link cost changes

link cost changes:

- node detects local link cost change
- updates routing info, recalculates distance vector
- if DV changes, notify neighbors



"good news travels fast"

t_0 : y detects link-cost change, updates its DV, informs its neighbors.

t_1 : z receives update from y, updates its table, computes new least cost to x, sends its neighbors its DV.

t_2 : y receives z's update, updates its distance table. y's least costs do not change, so y does not send a message to z.

* Check out the online interactive exercises for more examples: http://gaia.cs.umass.edu/kurose_ross/interactive/

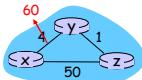
Network Layer: Control Plane 5-26

26

Distance vector: link cost changes

link cost changes:

- node detects local link cost change
- bad news travels slow - "count to infinity" problem!
- 44 iterations before algorithm stabilizes: see text



What would have happened if the link cost $c(y,x)$ had changed from 4 to 10,000 and the cost $c(z,x)$ had been 9,999? Because of such scenarios, the problem we have seen is sometimes referred to as the count-to-infinity problem.

poisoned reverse:

- If Z routes through Y to get to X:
 - Z tells Y its (Z's) distance to X is infinite (so Y won't route to X via Z)
- will this completely solve count to infinity problem?

count - infinity problem

Network Layer: Control Plane 5-27

27

Comparison of LS and DV algorithms

message complexity

- LS:** with n nodes, E links, $O(nE)$ msgs sent
- DV:** exchange between neighbors only
 - convergence time varies

robustness: what happens if router malfunctions?

LS:

- node can advertise incorrect link cost
- each node computes only its own table

DV:

- DV node can advertise incorrect path cost
- each node's table used by others
 - error propagate thru network

In 1997, a malfunctioning router in a small ISP provided national backbone routers with erroneous routing information. This caused other routers to flood the malfunctioning router with traffic and caused large portions of the Internet to become disconnected for up to several hours.

Network Layer: Control Plane 5-28

28

04/12/2022 starts here

Chapter 5: outline

- 5.1 introduction
- 5.2 routing protocols
 - link state
 - distance vector
- 5.3 intra-AS routing in the Internet: OSPF
- 5.4 routing among the ISPs: BGP
- 5.5 The SDN control plane
- 5.6 ICMP: The Internet Control Message Protocol
- 5.7 Network management and SNMP

Network Layer: Control Plane 5-29

29

Making routing scalable

our routing study thus far - idealized

- all routers identical
- network "flat"
- ... not true in practice

eg: AT&T, spark light

scale: with billions of destinations:

- can't store all destinations in routing tables!
- routing table exchange would swamp links!

administrative autonomy

- internet = network of networks
- each network admin may want to control routing in its own network

30

Link state store all ip address in the table, and the update may cause large scale information exchange. butterfly effect.

Network Layer: Control Plane 5-30

linked stats is not good for larger network (quiz question)

- Higher frequency update, higher frequency of results sent, more number of broadcast.

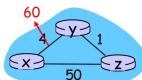
When does router in bellmen ford sends information?

Regular synchronization like each 12 hr sent information to make sure all router has synced information. And a new costs path is found and need to be updated.

Distance vector: link cost changes

link cost changes:

- node detects local link cost change
- **bad news travels slow** - "count to infinity" problem!
- 44 iterations before algorithm stabilizes: see text



What would have happened if the link cost $c(y,x)$ had changed from 4 to 10,000 and the cost $c(z,x)$ had been 9,999? Because of such scenarios, the problem we have seen is sometimes referred to as the count-to-infinity problem.

poisoned reverse:

- If Z routes through Y to get to X:
 - Z tells Y its (Z's) distance to X is infinite (so Y won't route to X via Z)
- will this completely solve count to infinity problem?

重点：由于在未改变costs 4之前，y知道z到x的least costs是4+1，然后在改变了xy的costs 到60之后，y就认为我需要从z送包是最远的。因为他是5，总比60好，随后他的新路线就是1+原来z的least costs=5=6。问题就在于y不知道的是原来z的least cost路线是经过我自己的，傻逼了。好了这下z自己也要更新table了，一看y给的数据，y到x的least costs是6，那就把包给y然后这样的话我z到x的least costs path就是6 +1 = 7. 然后又把这个数据告诉了y. 随后他们两个就不停的增加1然后update table。直到发现z到x的least costs有一条路是50，然后就停止了相互增加。

node x table			node y table			node z table			
cost to			cost to			cost to			
from	x	y	z	x	y	z	x	y	z
x	0	4	1	x	0	51	50		
y	4	0	1	y	51	0			
z	5	1	0	z	5	1	0		

poisoned reverse t_0

node x table			node y table			node z table			
cost to			cost to			cost to			
from	x	y	z	x	y	z	x	y	z
x	0	4	1	x	0	4	5		
y	4	0	1	y	51	0			
z	5	1	0	z	5	1	0		

$d_z(x) = \min(c(y,x) + d_z(y), c(y,z) + d_z(y)) = \min(60+0, 1+50) = 60$

Poly Huang, NTU EE

node x table			node y table			node z table			
cost to			cost to			cost to			
from	x	y	z	x	y	z	x	y	z
x	0	4	5	x	0	4	5		
y	4	0	1	y	51	0			
z	5	1	0	z	5	1	0		

poisoned reverse t_1 and t_2

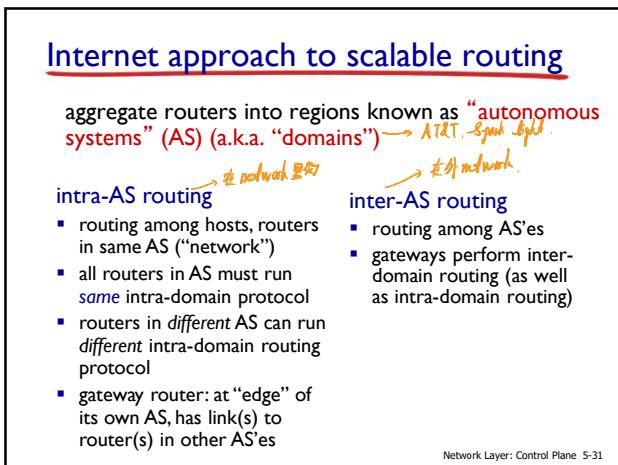
$d_z(x) = \min(c(y,x) + d_z(y), c(y,z) + d_z(y)) = \min(60+0, 1+50) = 51$

$d_z(x) = \min(c(z,x) + d_z(z), c(z,y) + d_z(y)) = \min(50+0, 1+60) = 50$

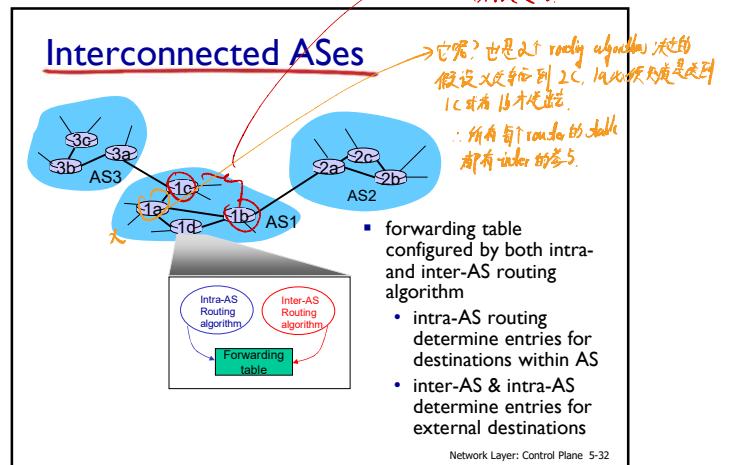
Poly Huang, NTU EE

<https://www.youtube.com/watch?v=UcTgNZwGmRU>

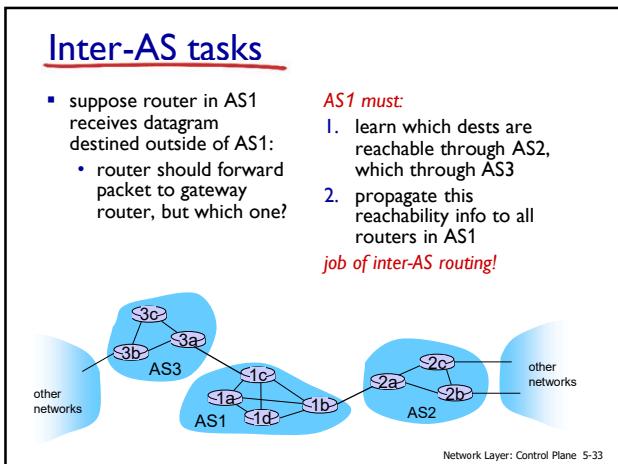
HW: poisoned reverse 不能解决问题，因3个路由器，更新不行。
read book



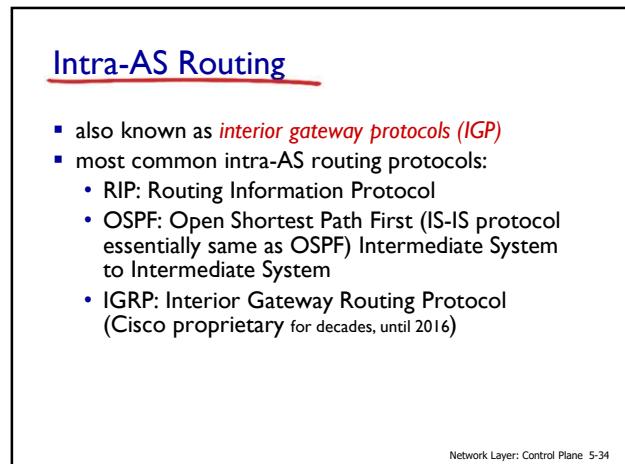
31



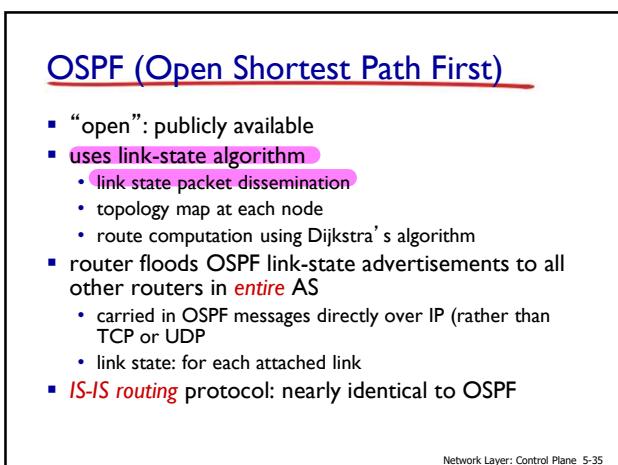
32



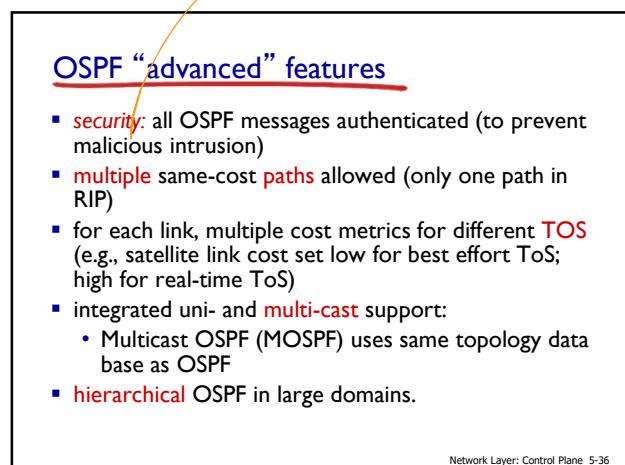
33



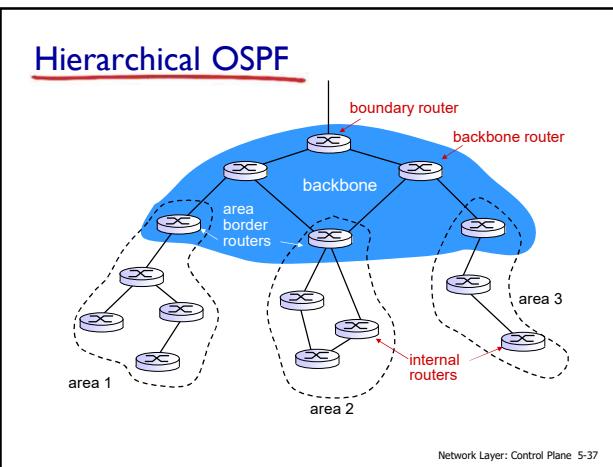
34



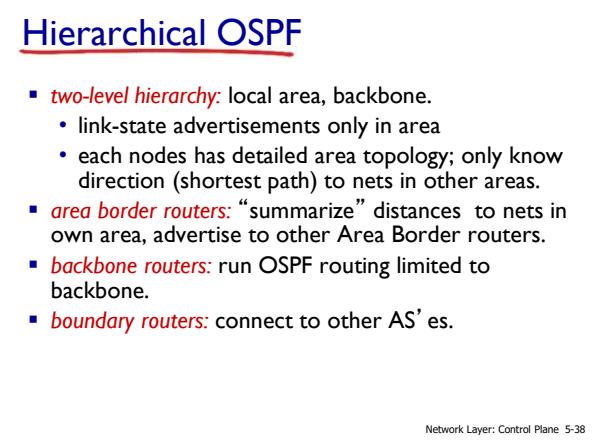
35



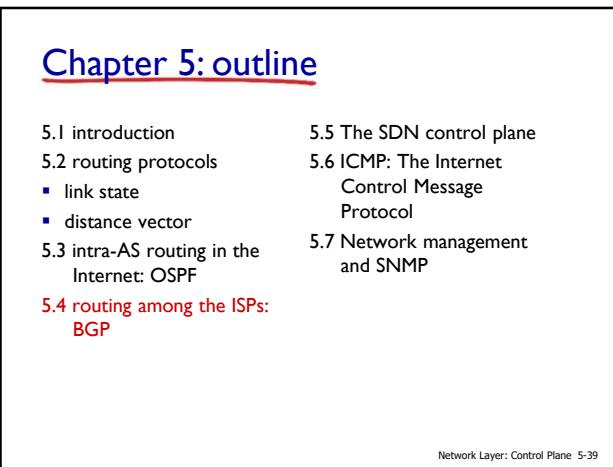
36



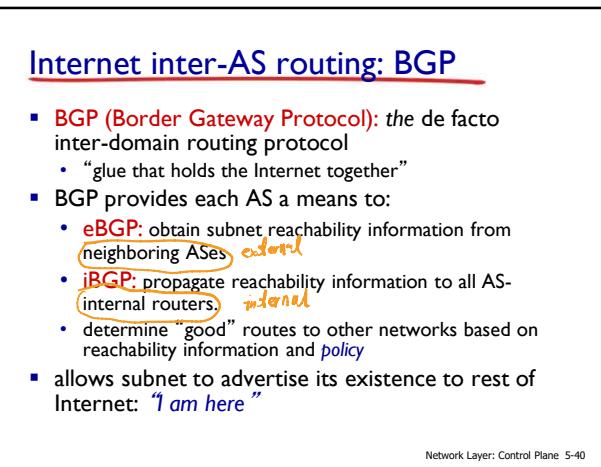
37



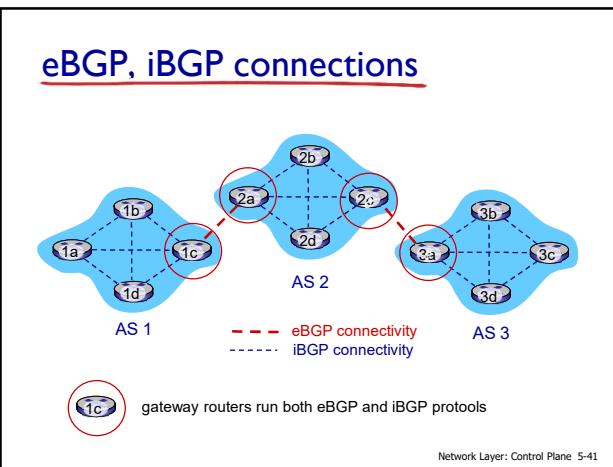
38



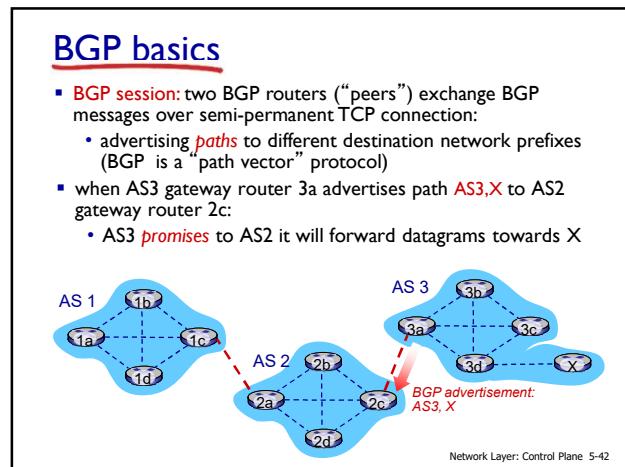
39



40



41



42

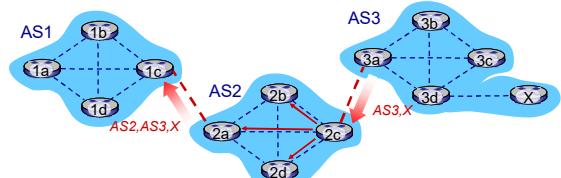
Path attributes and BGP routes

- advertised prefix includes BGP attributes
 - prefix + attributes = “route”
- two important attributes:
 - **AS-PATH**: list of ASes through which prefix advertisement has passed
 - **NEXT-HOP**: indicates specific internal-AS router to next-hop AS
- **Policy-based routing**:
 - gateway receiving route advertisement uses **import policy** to accept/decline path (e.g., never route through AS Y).
 - AS policy also determines whether to **advertise** path to other other neighboring ASes

Network Layer: Control Plane 5-43

43

BGP path advertisement

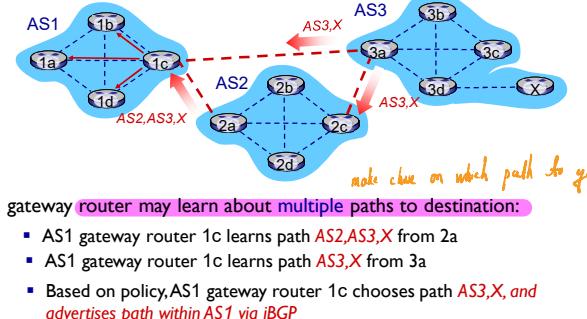


- AS2 router 2c receives path advertisement **AS3,X** (via eBGP) from AS3 router 3a
- Based on AS2 policy, AS2 router 2c accepts path **AS3,X**, propagates (via iBGP) to all AS2 routers
- Based on AS2 policy, AS2 router 2a advertises (via eBGP) path **AS2, AS3,X** to AS1 router 1c

Network Layer: Control Plane 5-44

44

BGP path advertisement



Network Layer: Control Plane 5-45

45

BGP messages

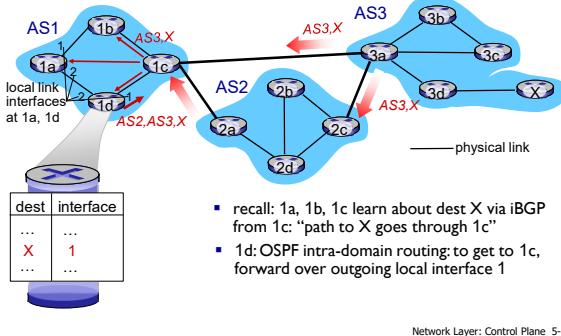
- BGP messages exchanged between peers over TCP connection
- BGP messages:
 - **OPEN**: opens TCP connection to remote BGP peer and authenticates sending BGP peer
 - **UPDATE**: advertises new path (or withdraws old)
 - **KEEPALIVE**: keeps connection alive in absence of UPDATES; also ACKs OPEN request
 - **NOTIFICATION**: reports errors in previous msg; also used to close connection

Network Layer: Control Plane 5-46

46

BGP, OSPF, forwarding table entries

Q: how does router set forwarding table entry to distant prefix?

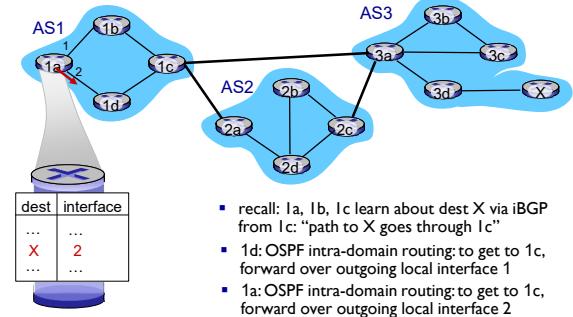


Network Layer: Control Plane 5-47

47

BGP, OSPF, forwarding table entries

Q: how does router set forwarding table entry to distant prefix?

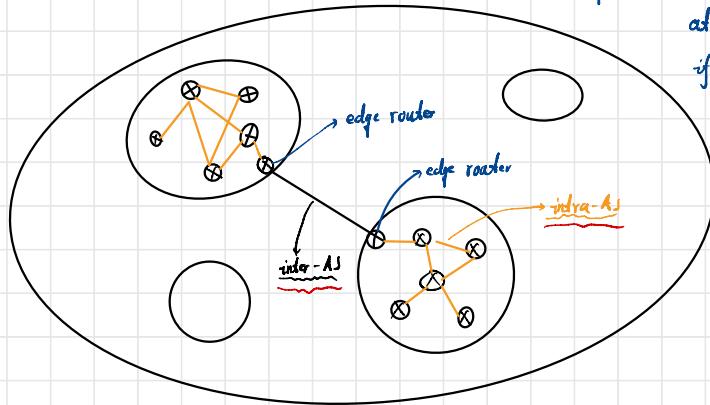


Network Layer: Control Plane 5-48

48

04/12/2022

^{routing}
edge router: at least 2 algorithms: 1 intra, 1 inter
at most 3 algorithms: 1 intra, 2 inter
if more auto-link connecting to the router.



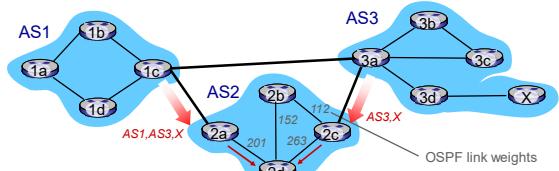
BGP route selection

- router may learn about more than one route to destination AS, selects route based on:
 1. local preference value attribute: policy decision
 2. shortest AS-PATH
 3. closest NEXT-HOP router: hot potato routing
 4. additional criteria

Network Layer: Control Plane 5-49

49

Hot Potato Routing

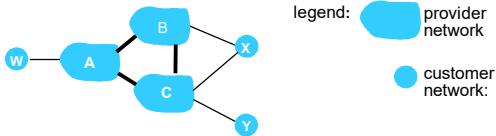


- 2d learns (via iBGP) it can route to X via 2a or 2c
- **hot potato routing**: choose local gateway that has least intra-domain cost (e.g., 2d chooses 2a, even though more AS hops to X): don't worry about inter-domain cost!

Network Layer: Control Plane 5-50

50

BGP: achieving policy via advertisements



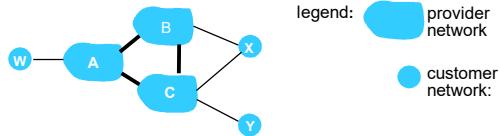
Suppose an ISP only wants to route traffic to/from its customer networks (does not want to carry transit traffic between other ISPs)

- A advertises path Aw to B and to C
- **B chooses not to advertise BAw to C:** *选择不告诉其他人*
- B gets no “revenue” for routing CBAw, since none of C, A, w are B’s customers
- C does not learn about CBAw path
- C will route CAw (not using B) to get to w

Network Layer: Control Plane 5-51

51

BGP: achieving policy via advertisements



Suppose an ISP only wants to route traffic to/from its customer networks (does not want to carry transit traffic between other ISPs)

- A,B,C are **provider networks**
- X,W,Y are customer (of provider networks)
- X is **dual-homed**: attached to two networks
- **policy to enforce**: X does not want to route from B to C via X
 - ... so X will not advertise to B a route to C

Network Layer: Control Plane 5-52

52

Why different Intra-, Inter-AS routing ?

policy:

- inter-AS: admin wants control over how its traffic routed, who routes through its net.
- intra-AS: single admin, so no policy decisions needed

scale:

- hierarchical routing saves table size, reduced update traffic

performance:

- **intra-AS**: can focus on performance
- **inter-AS**: policy may dominate over performance

Network Layer: Control Plane 5-53

53

Chapter 5: outline

- | | |
|---|---|
| 5.1 introduction
5.2 routing protocols
<ul style="list-style-type: none"> ▪ link state ▪ distance vector 5.3 intra-AS routing in the Internet: OSPF
5.4 routing among the ISPs:
BGP | 5.5 The SDN control plane
5.6 ICMP: The Internet Control Message Protocol
5.7 Network management and SNMP |
|---|---|

Network Layer: Control Plane 5-54

54

Software defined networking (SDN)

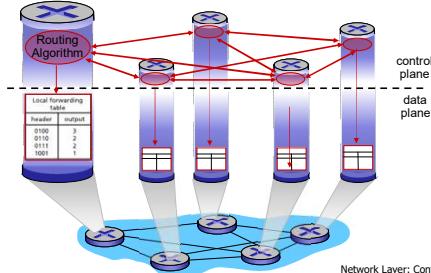
- Internet network layer: historically has been implemented via distributed, per-router approach
 - monolithic* router contains switching hardware, runs proprietary implementation of Internet standard protocols (IP, RIP, IS-IS, OSPF, BGP) in proprietary router OS (e.g., Cisco IOS)
 - different “middleboxes” for different network layer functions: firewalls, load balancers, NAT boxes, ..
- ~2005: renewed interest in rethinking network control plane

Network Layer: Control Plane 5-55

55

Recall: per-router control plane

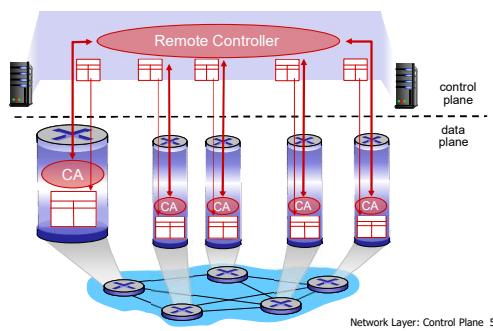
Individual routing algorithm components *in each and every router* interact with each other in control plane to compute forwarding tables



56

Recall: logically centralized control plane

A distinct (typically remote) controller interacts with local control agents (CAs) in routers to compute forwarding tables



Network Layer: Control Plane 5-57

57

Software defined networking (SDN)

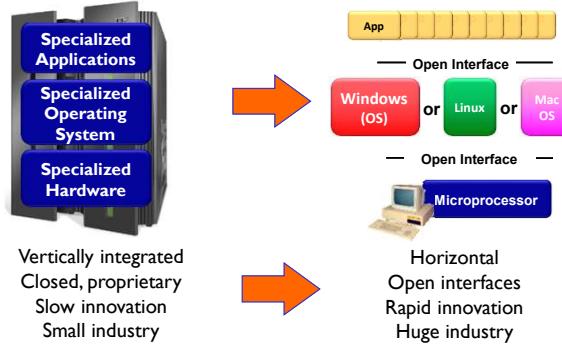
Why a *logically centralized* control plane?

- easier network management: avoid router misconfigurations, greater flexibility of traffic flows
- table-based forwarding (recall OpenFlow API) allows “programming” routers
 - centralized “programming” easier: compute tables centrally and distribute
 - distributed “programming”: more difficult: compute tables as result of distributed algorithm (protocol) implemented in each and every router
- open (non-proprietary) implementation of control plane

Network Layer: Control Plane 5-58

58

Analogy: mainframe to PC evolution:

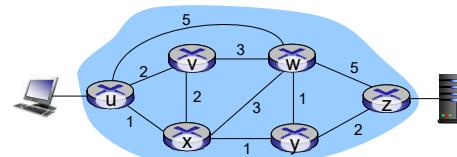


Network Layer: Control Plane 5-59

* Slide courtesy: N. McKeown

59

Traffic engineering: difficult traditional routing



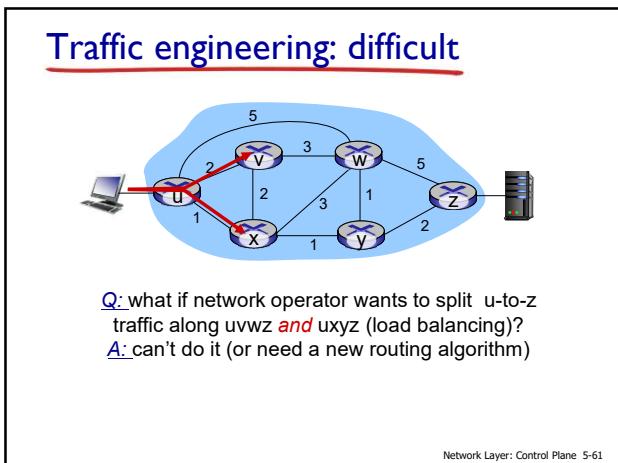
Q: what if network operator wants u-to-z traffic to flow along *uvwz*, x-to-z traffic to flow *xwyz*?

A: need to define link weights so traffic routing algorithm computes routes accordingly (or need a new routing algorithm)!

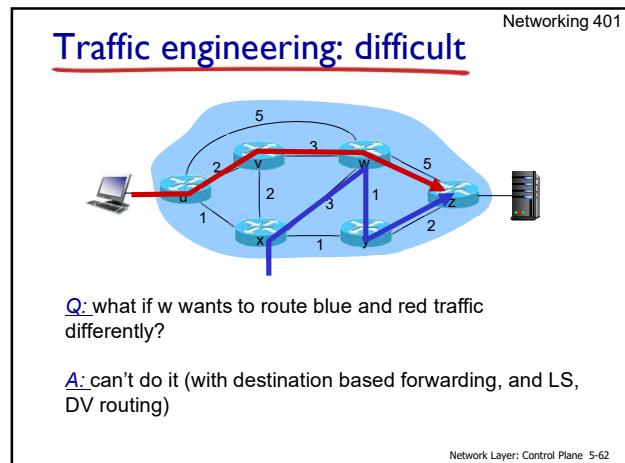
Link weights are only control “knobs”: wrong!

Network Layer: Control Plane 5-60

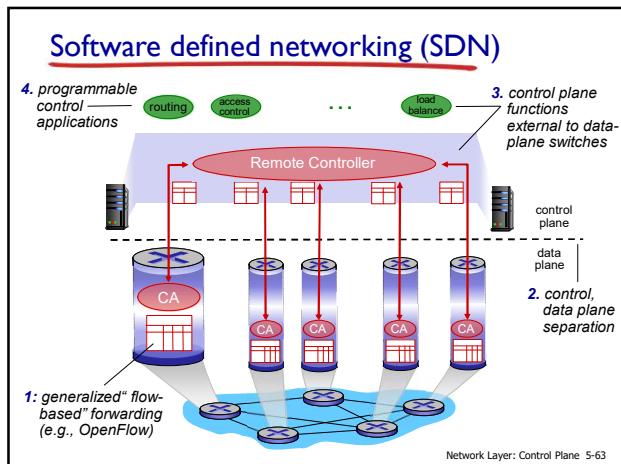
60



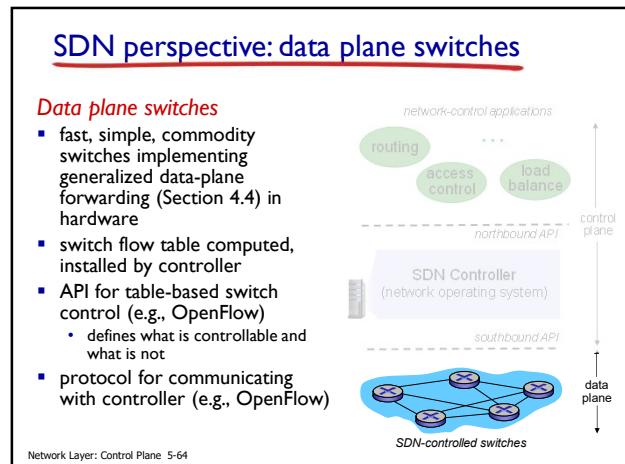
61



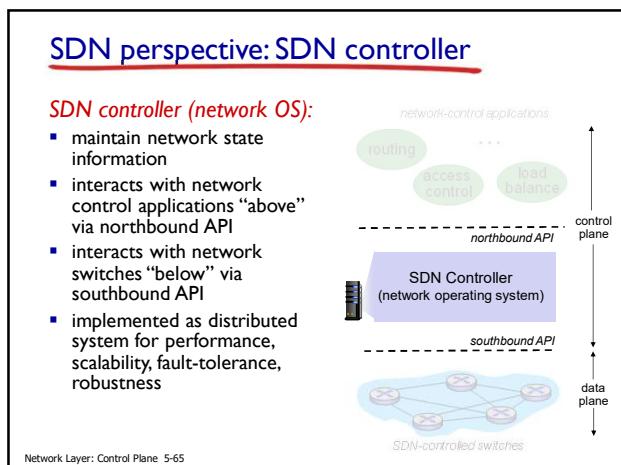
62



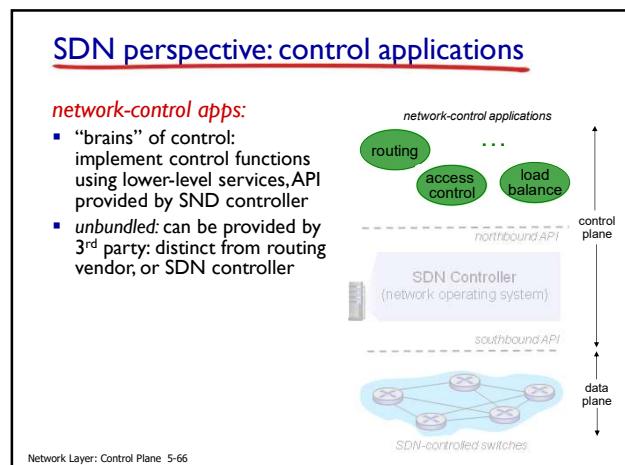
63



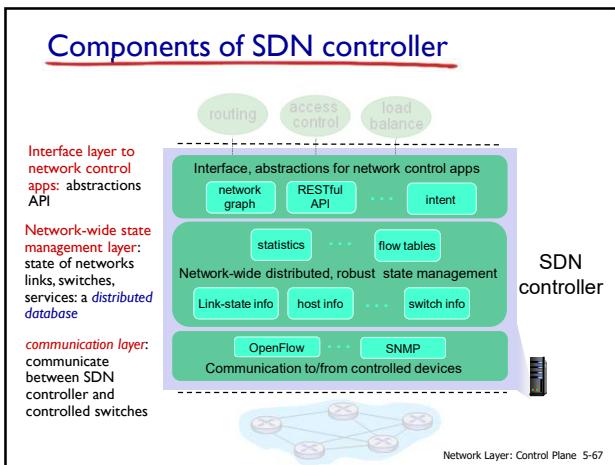
64



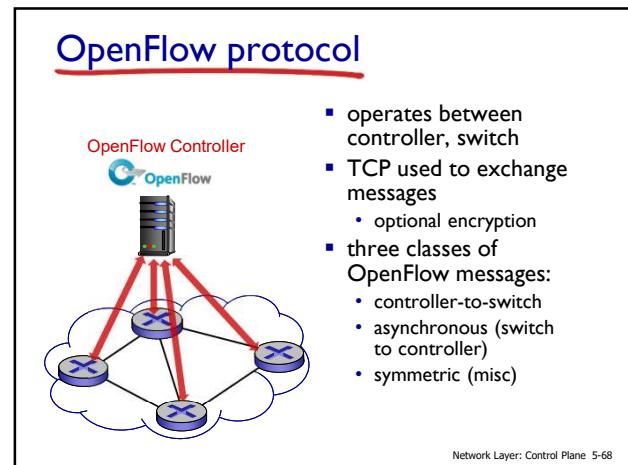
65



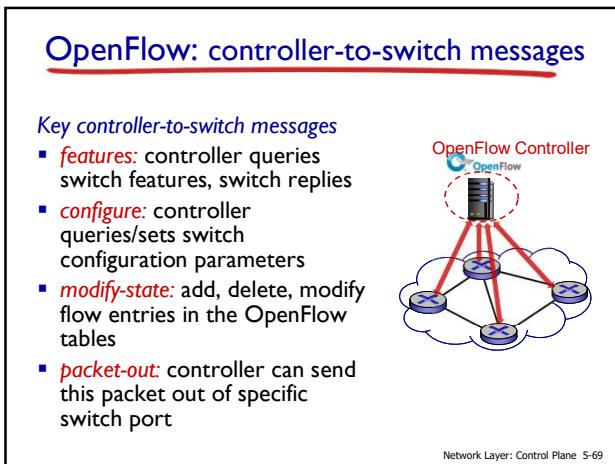
66



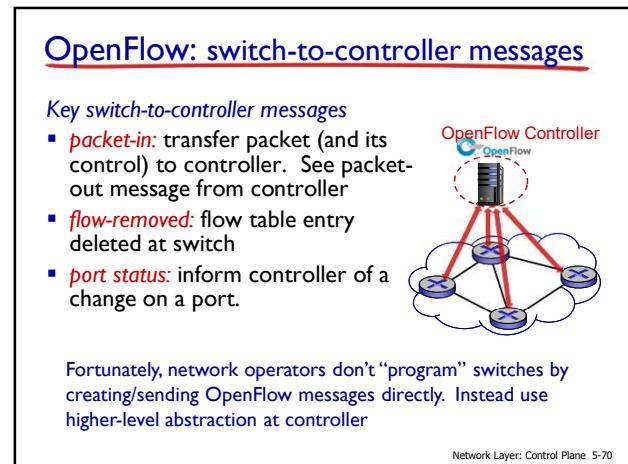
67



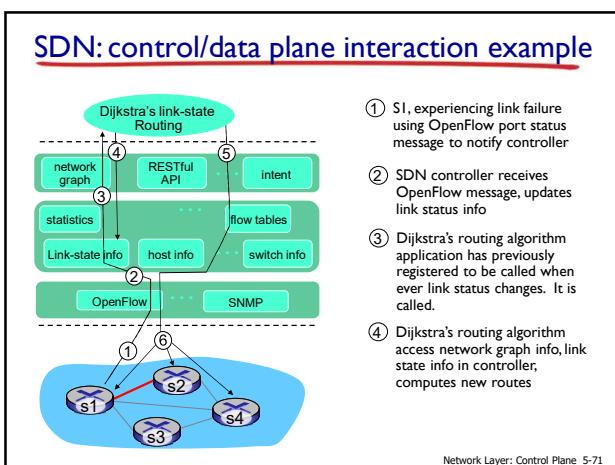
68



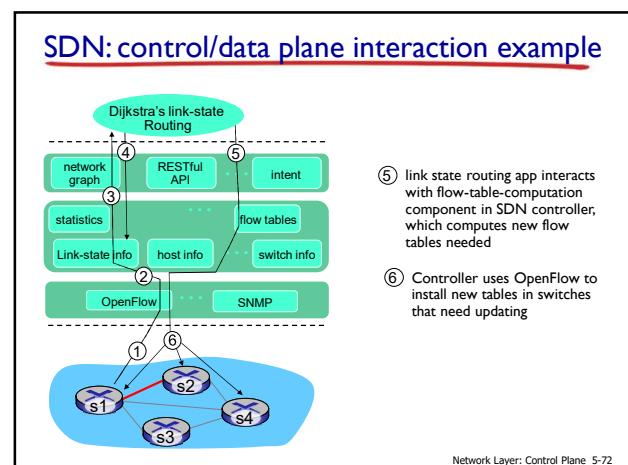
69



70

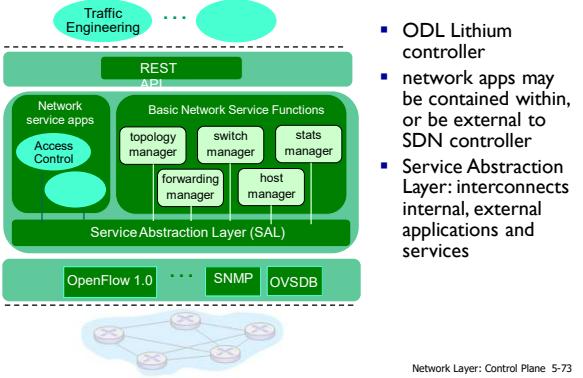


71



72

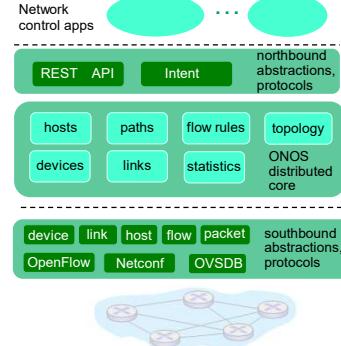
OpenDaylight (ODL) controller



- ODL Lithium controller
- network apps may be contained within, or be external to SDN controller
- Service Abstraction Layer: interconnects internal, external applications and services

73

ONOS controller



- control apps separate from controller
- intent framework: high-level specification of service: what rather than how
- considerable emphasis on distributed core: service reliability, replication, performance scaling

74

前面的大概略过3.

SDN: selected challenges

- hardening the control plane: dependable, reliable, performance-scalable, secure distributed system
 - robustness to failures: leverage strong theory of reliable distributed system for control plane
 - dependability, security: "baked in" from day one?
- networks, protocols meeting mission-specific requirements
 - e.g., real-time, ultra-reliable, ultra-secure
- Internet-scaling

Network Layer: Control Plane 5-75

75

Chapter 5: outline

- | | |
|--|---|
| 5.1 introduction | 5.5 The SDN control plane |
| 5.2 routing protocols | 5.6 ICMP: The Internet Control Message Protocol |
| ▪ link state | 5.7 Network management and SNMP |
| ▪ distance vector | |
| 5.3 intra-AS routing in the Internet: OSPF | |
| 5.4 routing among the ISPs: BGP | |

Network Layer: Control Plane 5-76

76

Traceroute exploits ICMP?
Application layer

ICMP: internet control message protocol

- used by hosts & routers to communicate network-level information
 - error reporting: unreachable host, network, port, protocol
 - echo request/reply (used by ping)
- network-layer "above" IP:
 - ICMP msgs carried in IP datagrams
- ICMP message: type, code plus first 8 bytes of IP datagram causing error

Type	Code	description
0	0	echo reply (ping)
3	0	dest. network unreachable
3	1	dest host unreachable
3	2	dest protocol unreachable
3	3	dest port unreachable
3	6	dest network unknown
3	7	dest host unknown
4	0	source quench (congestion control - not used)
8	0	echo request (ping)
9	0	route advertisement
10	0	router discovery
11	0	TTL expired
12	0	bad IP header

Network Layer: Control Plane 5-77

77

implemented on top network layer
but not transport layer (middle between N and T)
use IP but not use TCP/UDP *

Traceroute and ICMP

- source sends series of UDP segments to destination
 - first set has TTL = 1
 - second set has TTL = 2, etc.
 - unlikely port number
- when datagram in nth set arrives to nth router:
 - router discards datagram and sends source ICMP message (type 11, code 0)
 - ICMP message include name of router & IP address



- when ICMP message arrives, source records RTTs
- stopping criteria:
 - UDP segment eventually arrives at destination host
 - destination returns ICMP "port unreachable" message (type 3, code 3)
 - source stops

Network Layer: Control Plane 5-78

78

we send good # on purpose. not to be read by destination host.

Chapter 5: outline

- 5.1 introduction
- 5.2 routing protocols
 - link state
 - distance vector
- 5.3 intra-AS routing in the Internet: OSPF
- 5.4 routing among the ISPs:
 - BGP
- 5.5 The SDN control plane
- 5.6 ICMP: The Internet Control Message Protocol
- 5.7 Network management and SNMP**

Network Layer: Control Plane 5-79

79

What is network management?

- **autonomous systems (aka “network”):** 1000s of interacting hardware/software components
- other complex systems requiring monitoring, control:
 - jet airplane
 - nuclear power plant
 - others?



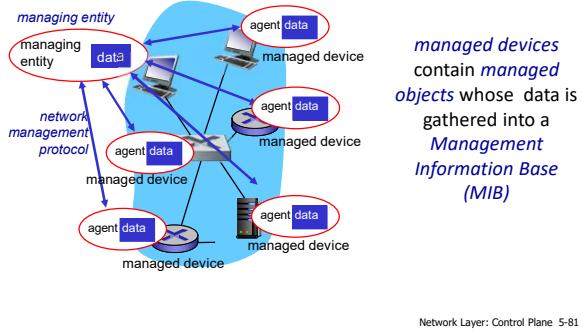
“Network management includes the deployment, integration and coordination of the hardware, software, and human elements to monitor, test, poll, configure, analyze, evaluate, and control the network and element resources to meet the real-time, operational performance, and Quality of Service requirements at a reasonable cost.”

Network Layer: Control Plane 5-80

80

Infrastructure for network management

definitions:

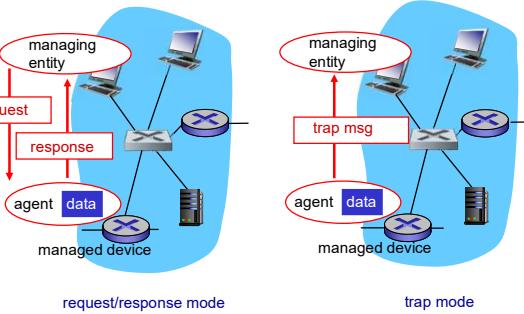


Network Layer: Control Plane 5-81

81

SNMP protocol

Two ways to convey MIB info, commands:



Network Layer: Control Plane 5-82

82

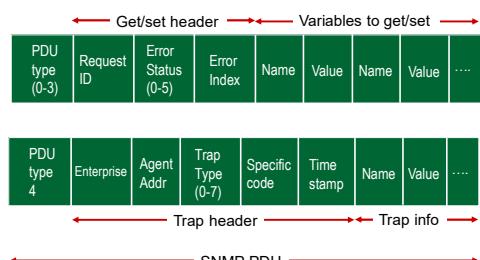
SNMP protocol: message types

Message type	Function
GetRequest	manager-to-agent: “get me data”
GetNextRequest	(data instance, next data in list, block of data)
GetBulkRequest	
InformRequest	manager-to-manager: here’s MIB value
SetRequest	manager-to-agent: set MIB value
Response	Agent-to-manager: value, response to Request
Trap	Agent-to-manager: inform manager of exceptional event

Network Layer: Control Plane 5-83

83

SNMP protocol: message formats



More on network management: see earlier editions of text!

Network Layer: Control Plane 5-84

84

Chapter 5: summary

we've learned a lot!

- approaches to network control plane
 - per-router control (traditional)
 - logically centralized control (software defined networking)
- traditional routing algorithms
 - implementation in Internet: OSPF, BGP
- SDN controllers
 - implementation in practice: ODL, ONOS
- Internet Control Message Protocol
- network management

next stop: link layer!

Network Layer: Control Plane 5-85

Chapter 6

Chapter 6: Link layer and LANs

our goals:

- understand principles behind link layer services:
 - error detection, correction
 - sharing a broadcast channel: multiple access
 - link layer addressing
 - local area networks: Ethernet, VLANs
- instantiation, implementation of various link layer technologies

Link Layer and LANs 6-1

1

Link layer, LANs: outline

- | | |
|--|---|
| 6.1 introduction, services | 6.5 link virtualization: MPLS |
| 6.2 error detection, correction | 6.6 data center networking |
| 6.3 multiple access protocols | 6.7 a day in the life of a web request |
| 6.4 LANs <ul style="list-style-type: none"> • addressing, ARP • Ethernet • switches • VLANs | |

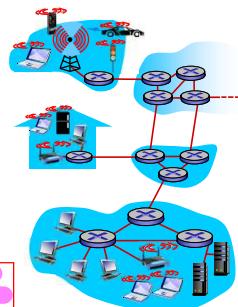
Link Layer and LANs 6-2

2

Link layer: introduction

terminology:

- hosts and routers: **nodes**
- communication channels that connect adjacent nodes along communication path: **links**
 - wired links
 - wireless links
 - LANs
- layer-2 packet: **frame**, encapsulates datagram



Link Layer and LANs 6-3

3

Link layer: context

- datagram transferred by different link protocols over different links:
 - e.g., Ethernet on first link, frame relay on intermediate links, 802.11 on last link
- each link protocol provides different services
 - e.g., may or may not provide rdt over link

- transportation analogy:*
- trip from Princeton to Lausanne
 - limo: Princeton to JFK
 - plane: JFK to Geneva
 - train: Geneva to Lausanne
 - tourist = **datagram**
 - transport segment = **communication link**
 - transportation mode = **link layer protocol**
 - travel agent = **routing algorithm**

Link Layer and LANs 6-4

4

Link layer services

- **framing, link access:**
 - encapsulate datagram into frame, adding header, trailer
 - channel access if shared medium
 - “MAC” addresses used in frame headers to identify source, destination
 - different from IP address!
- **reliable delivery between adjacent nodes**
 - we learned how to do this already (chapter 3)!
 - seldom used on low bit-error link (fiber, some twisted pair)
 - wireless links: high error rates (uses Forward Error Correction, FEC)
 - Q: why both link-level and end-end reliability?

Link Layer and LANs 6-5

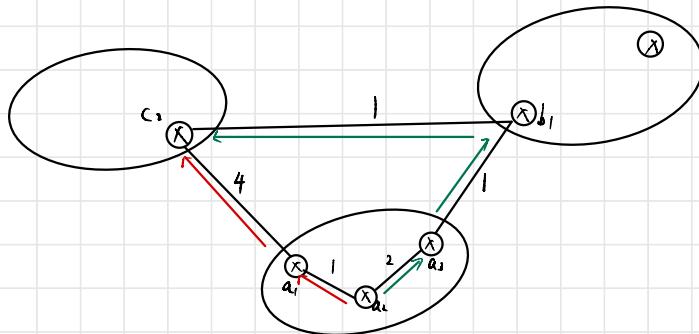
5

Link layer services (more)

- **flow control:**
 - pacing between adjacent sending and receiving nodes
- **error detection:**
 - errors caused by signal attenuation, noise.
 - receiver detects presence of errors:
 - signals sender for retransmission or drops frame
- **error correction:**
 - receiver identifies and corrects bit error(s) without resorting to retransmission
- **half-duplex and full-duplex**
 - with half duplex, nodes at both ends of link can transmit, but not at same time
 - 单边传 双边同时传

Link Layer and LANs 6-6

6



hot potato routing

从 A_2 出来到 C_1 , 用的是 红色路径, 而不是 绿色. 那怕总的 cost 等于红线. 因为在 A_2 中我想要

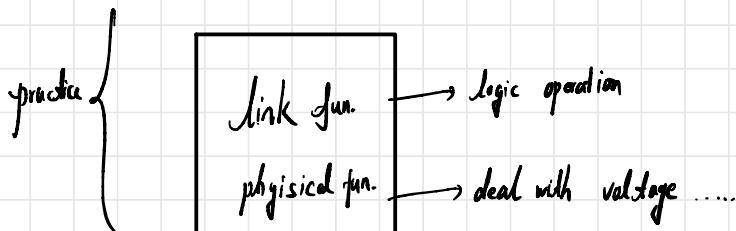
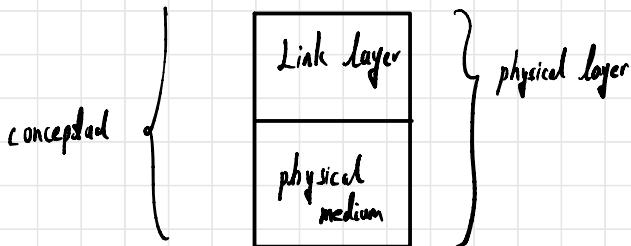
- hot potato routing: choose local gateway that has least intra-domain cost (e.g., 2d chooses 2a, even though more AS hops to X): don't worry about inter-domain cost!

{
inter-As: dominated by policy
intra-As: dominated by performance

ch: 6:

edge $\left[\begin{array}{l} T: \text{TCP/UDP protocol} \\ N: \text{IP protocol} \\ p \end{array} \right]$ core

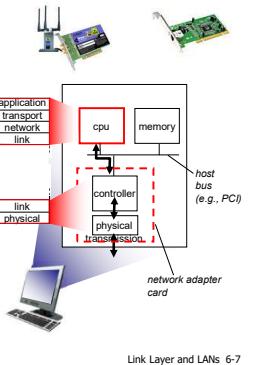
Introduce reliability on physical layer is debatable: for example, fiber has low error rate, the reliability may slow it down, and Wi-Fi has relatively high errors rate and the reliability may slow the speed down further.



one chip, two functions;

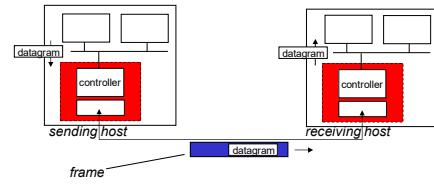
Where is the link layer implemented?

- in each and every host
- link layer implemented in “adaptor” (aka **network interface card** NIC) or on a chip
 - Ethernet card, 802.11 card; Ethernet chipset
 - implements link, physical layer
- attaches into host’s system buses
- combination of hardware, software, firmware



7

Adaptors communicating



- sending side:**
 - encapsulates datagram in frame
 - adds error checking bits, rdt, flow control, etc.
- receiving side:**
 - looks for errors, rdt, flow control, etc.
 - extracts datagram, passes to upper layer at receiving side

Link Layer and LANs 6-8

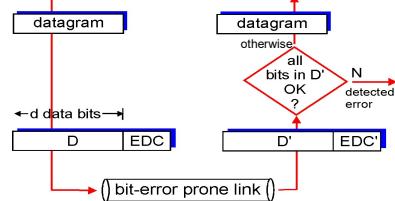
8

04/19/2022 starts here

Error detection

EDC= Error Detection and Correction bits (redundancy)
D = Data protected by error checking, may include header fields

- Error detection not 100% reliable!
 - protocol may miss some errors, but rarely
 - larger EDC field yields better detection and correction



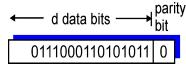
Link Layer and LANs 6-10

10

Parity checking

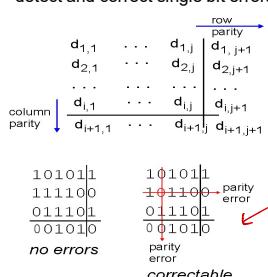
single bit parity:

- detect single bit errors



two-dimensional bit parity:

- detect and correct single bit errors



Link Layer and LANs 6-11

11

Internet checksum (review)

goal: detect “errors” (e.g., flipped bits) in transmitted packet
(note: used at transport layer only)

sender:

- treat segment contents as sequence of 16-bit integers
- checksum: addition ($1's$ complement sum) of segment contents
- sender puts checksum value into UDP checksum field

receiver:

- compute checksum of received segment
- check if computed checksum equals checksum field value:
 - NO - error detected
 - YES - no error detected. *But maybe errors nonetheless!*

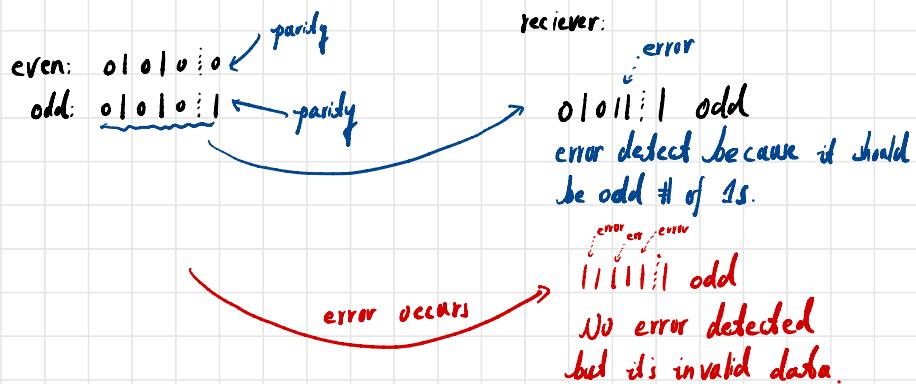
Link Layer and LANs 6-12

12

2

Parity check

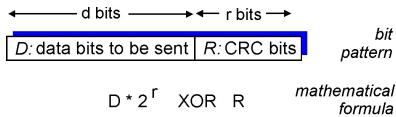
Even parity and odd parity



* 本讲题 材人
参考 w3school 的例

Cyclic redundancy check

- more powerful error-detection coding
- view data bits, D , as a binary number
- choose $r+1$ bit pattern (generator), G
- goal: choose r CRC bits, R , such that
 - $\langle D, R \rangle$ exactly divisible by G (modulo 2)
 - receiver knows G , divides $\langle D, R \rangle$ by G . If non-zero remainder: error detected!
 - can detect all burst errors less than $r+1$ bits
- widely used in practice (Ethernet, 802.11 WiFi, ATM)



Link Layer and LANs 6-13

13

CRC example

want:

$$D \cdot 2^r \text{ XOR } R = nG$$

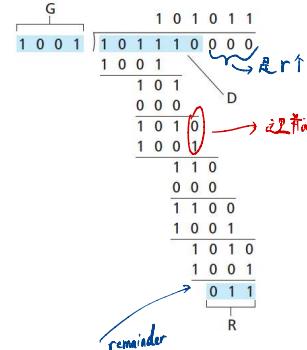
equivalently:

$$D \cdot 2^r = nG \text{ XOR } R$$

equivalently:

if we divide $D \cdot 2^r$ by G , want remainder R to satisfy:

$$R = \text{remainder} \left[\frac{D \cdot 2^r}{G} \right]$$



* Check out the online interactive exercises for more examples: http://gaia.cs.umass.edu/kurose_ross/interactive/

Link Layer and LANs 6-14

14

Link layer, LANs: outline

- | | |
|--------------------------------------|--|
| 6.1 introduction, services | 6.5 link virtualization: MPLS |
| 6.2 error detection, correction | 6.6 data center networking |
| 6.3 multiple access protocols | 6.7 a day in the life of a web request |
| 6.4 LANs | |
| • addressing, ARP | |
| • Ethernet | |
| • switches | |
| • VLANs | |

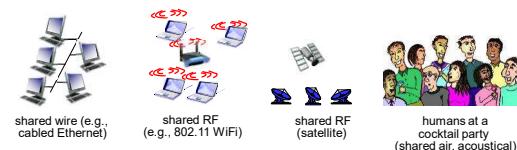
Link Layer and LANs 6-15

15

Multiple access links, protocols

two types of "links":

- point-to-point**
 - PPP for dial-up access
 - point-to-point link between Ethernet switch, host
- broadcast (shared wire or medium)**
 - old-fashioned Ethernet
 - upstream HFC
 - 802.11 wireless LAN



Link Layer and LANs 6-16

16

Multiple access protocols

- single shared broadcast channel
- two or more simultaneous transmissions by nodes: interference
 - collision** if node receives two or more signals at the same time

multiple access protocol

- distributed algorithm that determines how nodes share channel, i.e., determine when node can transmit
- communication about channel sharing must use channel itself!
 - no out-of-band channel for coordination

Link Layer and LANs 6-17

17

An ideal multiple access protocol

given: broadcast channel of rate R bps

desiderata:

- when one node wants to transmit, it can send at rate R .
- when M nodes want to transmit, each can send at average rate R/M
- fully decentralized:
 - no special node to coordinate transmissions
 - no synchronization of clocks, slots
- simple

Link Layer and LANs 6-18

18

MAC protocols: taxonomy

three broad classes:

- **channel partitioning**
 - divide channel into smaller "pieces" (time slots, frequency, code)
 - allocate piece to node for exclusive use
- **random access**
 - channel not divided, allow collisions
 - "recover" from collisions
- **"taking turns"**
 - nodes take turns, but nodes with more to send can take longer turns

34 MAC protocols

Link Layer and LANs 6-19

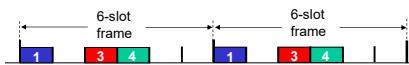
19

Fair, but not efficient.

Channel partitioning MAC protocols: TDMA

TDMA: time division multiple access

- access to channel in "rounds"
- each station gets fixed length slot (length = packet transmission time) in each round
- unused slots go idle
- example: 6-station LAN, 1,3,4 have packets to send, slots 2,5,6 idle



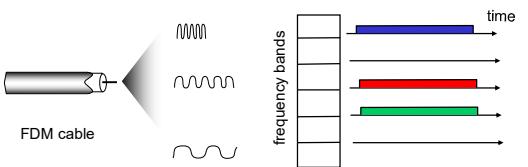
Link Layer and LANs 6-20

20

Channel partitioning MAC protocols: FDMA

FDMA: frequency division multiple access

- channel spectrum divided into frequency bands
- each station assigned fixed frequency band
- unused transmission time in frequency bands go idle
- example: 6-station LAN, 1,3,4 have packet to send, frequency bands 2,5,6 idle



Link Layer and LANs 6-21

21

Fair, not efficient.
idle bands may occur.

Random access protocols

- when node has packet to send
 - transmit at full channel data rate R.
 - no *a priori* coordination among nodes
- two or more transmitting nodes → "collision"
- **random access MAC protocol** specifies:
 - how to detect collisions
 - how to recover from collisions (e.g., via delayed retransmissions)
- examples of random access MAC protocols:
 - slotted ALOHA
 - ALOHA
 - CSMA, CSMA/CD, CSMA/CA

Link Layer and LANs 6-22

22

Slotted ALOHA

assumptions:

- all frames same size
- time divided into equal size slots (time to transmit 1 frame) *at least long enough for one frame*
- nodes start to transmit only slot beginning
- **nodes are synchronized**
- if 2 or more nodes transmit in slot, all nodes detect collision

must be time synchronized

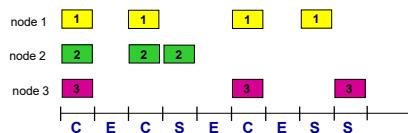
operation:

- when node obtains fresh frame, transmits in next slot
 - if no collision: node can send new frame in next slot
 - if collision: node retransmits frame in each subsequent slot with prob. p until success

Link Layer and LANs 6-23

23

Slotted ALOHA



Pros: *只有一個 active node* Cons: *只有一個 active node*

- single active node can continuously transmit at full rate of channel
- highly decentralized: only slots in nodes need to be in sync
- simple
- collisions, wasting slots
- idle slots
- nodes may be able to detect collision in less than time to transmit packet
- clock synchronization

Link Layer and LANs 6-24

24

more collision the node has seen, lower probability it does

Slotted ALOHA: efficiency

efficiency: long-run fraction of successful slots (many nodes, all with many frames to send)

- suppose: N nodes with many frames to send, each transmits in slot with probability p
- prob that given node has success in a slot = $p(1-p)^{N-1}$
- prob that *any* node has a success = $Np(1-p)^{N-1}$

- max efficiency: find p^* that maximizes $Np(1-p)^{N-1}$
- for many nodes, take limit of $Np^*(1-p^*)^{N-1}$ as N goes to infinity, gives:
 $\text{max efficiency} = 1/e = .37$

at best: channel used for useful transmissions 37% of time!

Link Layer and LANs 6-25

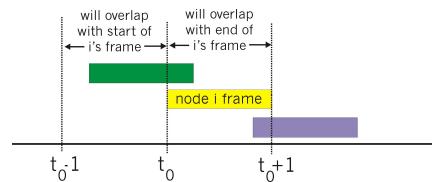
只有 37%

25

Pure (unslotted) ALOHA

- unslotted Aloha: simpler, no synchronization
- when frame first arrives
 - transmit immediately
- collision probability increases:
 - frame sent at t_0 collides with other frames sent in $[t_0-1, t_0+1]$

18% success.



Link Layer and LANs 6-26

26

Pure ALOHA efficiency

$$\begin{aligned} P(\text{success by given node}) &= P(\text{node transmits}) \cdot \\ &\quad P(\text{no other node transmits in } [t_0-1, t_0]) \cdot \\ &\quad P(\text{no other node transmits in } [t_0, t_0+1]) \\ &= p \cdot (1-p)^{N-1} \cdot (1-p)^{N-1} \\ &= p \cdot (1-p)^{2(N-1)} \\ \dots \text{choosing optimum } p \text{ and then letting } n \rightarrow \infty \\ &= 1/(2e) = .18 \end{aligned}$$

even worse than slotted Aloha!

Link Layer and LANs 6-27

27

CSMA (carrier sense multiple access)

CSMA: listen before transmit:

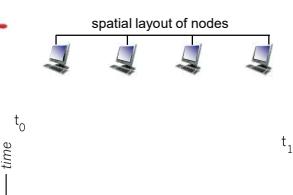
- if channel sensed idle: transmit entire frame
- if channel sensed busy, defer transmission
- human analogy: don't interrupt others!

Link Layer and LANs 6-28

28

CSMA collisions

- collisions can still occur: propagation delay means two nodes may not hear each other's transmission
- collision: entire packet transmission time wasted
 - distance & propagation delay play role in determining collision probability



Link Layer and LANs 6-29

29

CSMA/CD (collision detection)

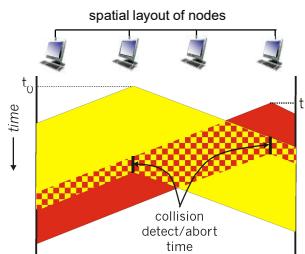
CSMA/CD: carrier sensing, deferral as in CSMA

- collisions detected within short time
- colliding transmissions aborted, reducing channel wastage
- collision detection:
 - easy in wired LANs: measure signal strengths, compare transmitted, received signals
 - difficult in wireless LANs: received signal strength overwhelmed by local transmission strength
- human analogy: the polite conversationalist

Link Layer and LANs 6-30

30

CSMA/CD (collision detection)



Link Layer and LANs 6-31

31

CSMA/CD efficiency

- T_{prop} = max prop delay between 2 nodes in LAN
- t_{trans} = time to transmit max-size frame

$$\text{efficiency} = \frac{1}{1 + 5t_{prop}/t_{trans}}$$

Longer data you have
 smaller network, less t_{prop} .
 Because if you are sending no one else will send

- efficiency goes to 1
 - as t_{prop} goes to 0
 - as t_{trans} goes to infinity
- better performance than ALOHA: and simple, cheap, decentralized!

Link Layer and LANs 6-33

33

Ethernet CSMA/CD algorithm

1. NIC receives datagram from network layer, creates frame
2. If NIC senses channel idle, starts frame transmission. If NIC senses channel busy, waits until channel idle, then transmits.
3. If NIC transmits entire frame without detecting another transmission, NIC is done with frame !
4. If NIC detects another transmission while transmitting, aborts and sends jam signal *(collisions)*
5. After aborting, NIC enters *binary (exponential) backoff*
 - after m th collision, NIC chooses K at random from $\{0, 1, 2, \dots, 2^m - 1\}$. NIC waits $K \cdot 512$ bit times, returns to Step 2
 - longer backoff interval with more collisions

Link Layer and LANs 6-32

32

more collisions, longer time to wait, back off

"Taking turns" MAC protocols

channel partitioning MAC protocols:

- share channel efficiently and fairly at high load
- inefficient at low load: delay in channel access, I/N bandwidth allocated even if only 1 active node!

random access MAC protocols

- efficient at low load: single node can fully utilize channel
- high load: collision overhead

"taking turns" protocols

look for best of both worlds!

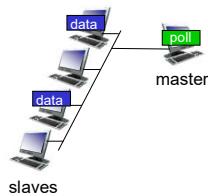
Link Layer and LANs 6-34

34

"Taking turns" MAC protocols

polling:

- master node "invites" slave nodes to transmit in turn
- typically used with "dumb" slave devices
- concerns:
 - polling overhead
 - latency
 - single point of failure (master)



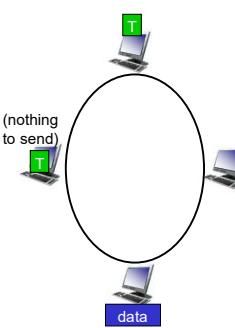
Link Layer and LANs 6-35

35

"Taking turns" MAC protocols

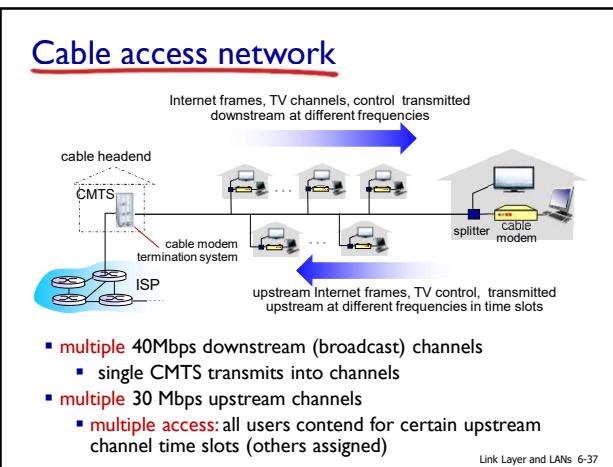
token passing:

- control *token* passed from one node to next sequentially.
- token message
- concerns:
 - token overhead
 - latency
 - single point of failure (token)

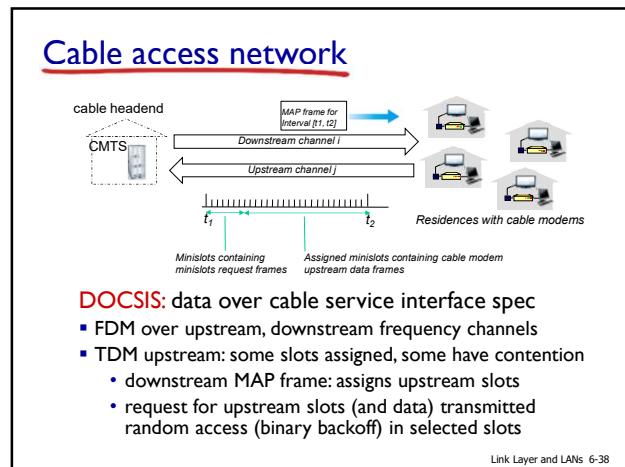


Link Layer and LANs 6-36

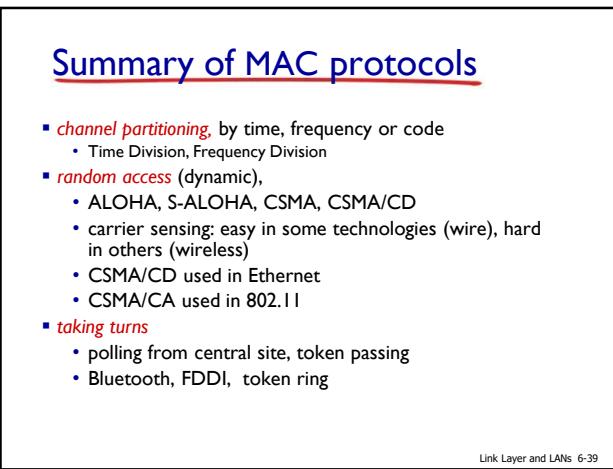
36



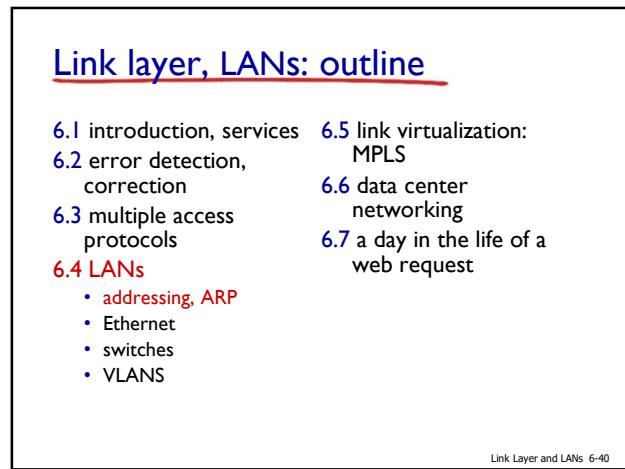
37



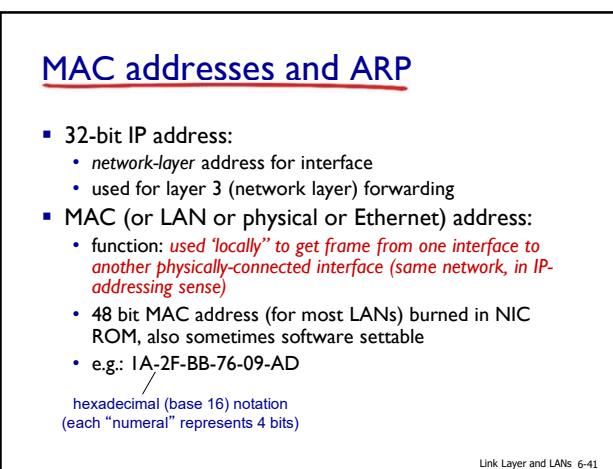
38



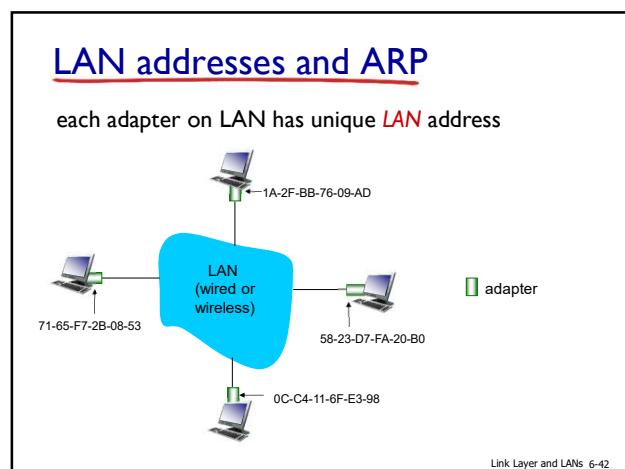
39



40



41



42

LAN addresses (more)

- MAC address allocation administered by IEEE
- manufacturer buys portion of MAC address space (to assure uniqueness)
- analogy:
 - MAC address: like Social Security Number
 - IP address: like postal address
- MAC flat address → portability
 - can move LAN card from one LAN to another
- IP hierarchical address *not* portable
 - address depends on IP subnet to which node is attached

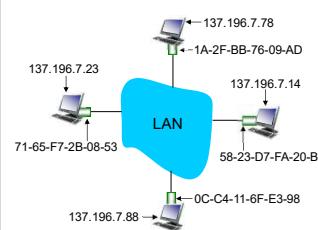
Link Layer and LANs 6-43

43

ARP: address resolution protocol

Question: how to determine interface's MAC address, knowing its IP address?

ARP table: each IP node (host, router) on LAN has table



- IP/MAC address mappings for some LAN nodes:
< IP address; MAC address; TTL>
- TTL (Time To Live): time after which address mapping will be forgotten (typically 20 min)

Link Layer and LANs 6-44

44

ARP protocol: same LAN

- A wants to send datagram to B
 - B's MAC address not in A's ARP table.
- A broadcasts ARP query packet, containing B's IP address
 - destination MAC address = FF-FF-FF-FF-FF-FF
 - all nodes on LAN receive ARP query
- B receives ARP packet, replies to A with its (B's) MAC address
 - frame sent to A's MAC address (unicast)
- A caches (saves) IP-to-MAC address pair in its ARP table until information becomes old (times out)
 - soft state: information that times out (goes away) unless refreshed
- ARP is "plug-and-play":
 - nodes create their ARP tables *without intervention from net administrator*

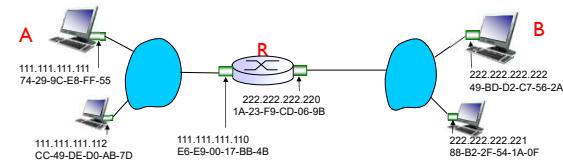
Link Layer and LANs 6-45

45

Addressing: routing to another LAN

walkthrough: send datagram from A to B via R

- focus on addressing – at IP (datagram) and MAC layer (frame)
- assume A knows B's IP address
- assume A knows IP address of first hop router, R (how?)
- assume A knows R's MAC address (how?)

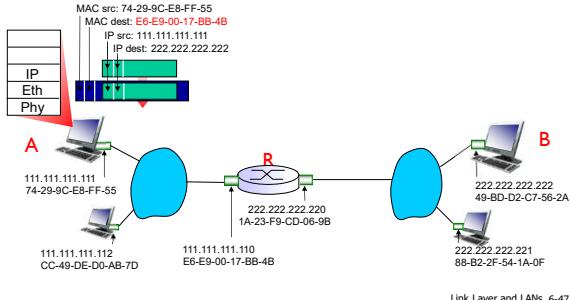


Link Layer and LANs 6-46

46

Addressing: routing to another LAN

- A creates IP datagram with IP source A, destination B
- A creates link-layer frame with R's MAC address as destination address, frame contains A-to-B IP datagram

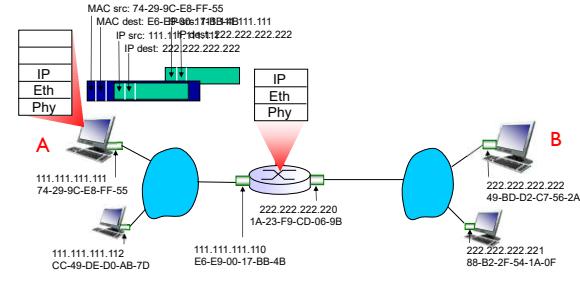


Link Layer and LANs 6-47

47

Addressing: routing to another LAN

- frame sent from A to R
- frame received at R, datagram removed, passed up to IP

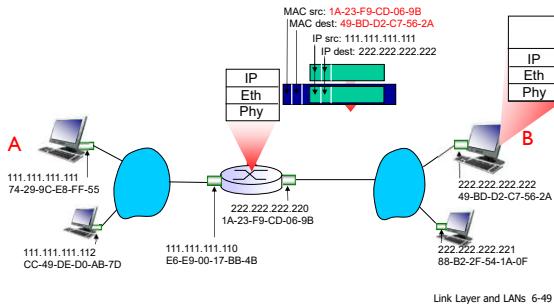


Link Layer and LANs 6-48

48

Addressing: routing to another LAN

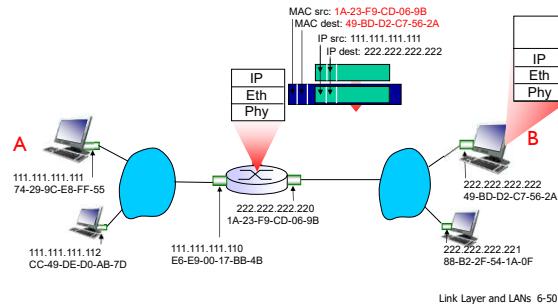
- R forwards datagram with IP source A, destination B
- R creates link-layer frame with B's MAC address as destination address, frame contains A-to-B IP datagram



49

Addressing: routing to another LAN

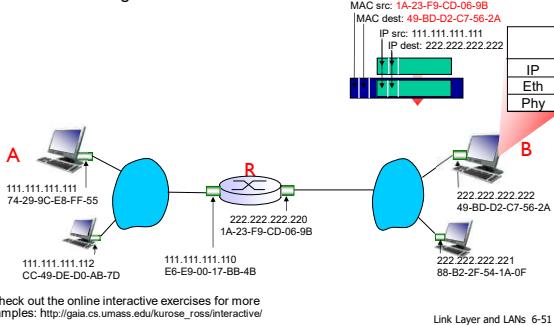
- R forwards datagram with IP source A, destination B
- R creates link-layer frame with B's MAC address as destination address, frame contains A-to-B IP datagram



50

Addressing: routing to another LAN

- R forwards datagram with IP source A, destination B
- R creates link-layer frame with B's MAC address as dest, frame contains A-to-B IP datagram



51

Link layer, LANs: outline

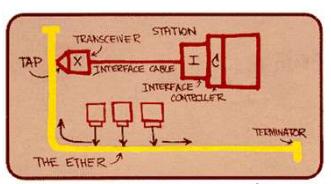
- 6.1 introduction, services
- 6.2 error detection, correction
- 6.3 multiple access protocols
- 6.4 LANs
 - addressing, ARP
 - Ethernet
 - switches
 - VLANs
- 6.5 link virtualization: MPLS
- 6.6 data center networking
- 6.7 a day in the life of a web request

Link Layer and LANs 6-52

Ethernet

“dominant” wired LAN technology:

- single chip, multiple speeds (e.g., Broadcom BCM5761)
- first widely used LAN technology
- simpler, cheap
- kept up with speed race: 10 Mbps – 10 Gbps

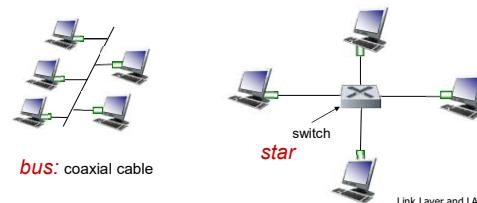


Link Layer and LANs 6-53

53

Ethernet: physical topology

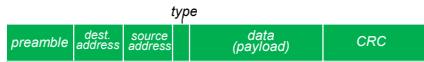
- bus:** popular through mid 90s
 - all nodes in same collision domain (can collide with each other)
- star:** prevails today
 - active **switch** in center
 - each “spoke” runs a (separate) Ethernet protocol (nodes do not collide with each other)



54

Ethernet frame structure

sending adapter encapsulates IP datagram (or other network layer protocol packet) in **Ethernet frame**



preamble:

- 7 bytes with pattern 10101010 followed by one byte with pattern 10101011
- used to synchronize receiver, sender clock rates

Link Layer and LANs 6-55

55

Ethernet frame structure (more)

- **addresses:** 6 byte source, destination MAC addresses

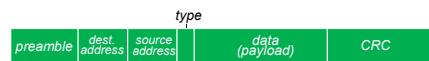
- if adapter receives frame with matching destination address, or with broadcast address (e.g. ARP packet), it passes data in frame to network layer protocol

- otherwise, adapter discards frame

- **type:** indicates higher layer protocol (mostly IP but others possible, e.g., Novell IPX, AppleTalk)

- **CRC:** cyclic redundancy check at receiver

- error detected: frame is dropped



Link Layer and LANs 6-56

56

Ethernet: unreliable, connectionless

- **connectionless:** no handshaking between sending and receiving NICs
- **unreliable:** receiving NIC doesn't send acks or nacks to sending NIC
 - data in dropped frames recovered only if initial sender uses higher layer rdt (e.g., TCP), otherwise dropped data lost
- Ethernet's MAC protocol: unslotted **CSMA/CD with binary backoff**

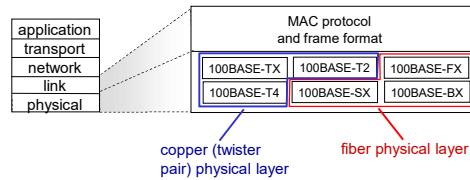
Link Layer and LANs 6-57

57

802.3 Ethernet standards: link & physical layers

- **many** different Ethernet standards

- common MAC protocol and frame format
- different speeds: 2 Mbps, 10 Mbps, 100 Mbps, 1Gbps, 10 Gbps, 40 Gbps
- different physical layer media: fiber, cable



Link Layer and LANs 6-58

58

Link layer, LANs: outline

- 6.1 introduction, services
- 6.2 error detection, correction
- 6.3 multiple access protocols
- 6.4 LANs
 - addressing, ARP
 - Ethernet
 - switches
 - VLANs
- 6.5 link virtualization: MPLS
- 6.6 data center networking
- 6.7 a day in the life of a web request

Link Layer and LANs 6-59

59

Ethernet switch

- **link-layer device: takes an active role**

- store, forward Ethernet frames
- examine incoming frame's MAC address, **selectively** forward frame to one-or-more outgoing links when frame is to be forwarded on segment, uses CSMA/CD to access segment

- **transparent**

- hosts are unaware of presence of switches

- **plug-and-play, self-learning**

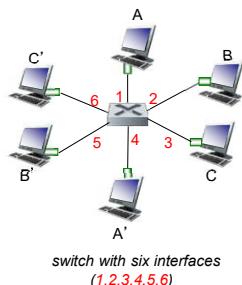
- switches do not need to be configured

Link Layer and LANs 6-60

60

Switch: multiple simultaneous transmissions

- hosts have dedicated, direct connection to switch
- switches buffer packets
- Ethernet protocol used on each incoming link, but no collisions; full duplex
 - each link is its own collision domain
- switching:** A-to-A' and B-to-B' can transmit simultaneously, without collisions



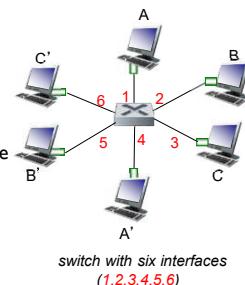
Link Layer and LANs 6-61

61

Switch forwarding table

- Q:** how does switch know A' reachable via interface 4, B' reachable via interface 5?

- A:** each switch has a **switch table**, each entry:
 - (MAC address of host, interface to reach host, time stamp)
- looks like a routing table!



Link Layer and LANs 6-62

62

Switch: self-learning

- switch **learns** which hosts can be reached through which interfaces
 - when frame received, switch "learns" location of sender: incoming LAN segment
 - records sender/location pair in switch table

MAC addr	interface	TTL
A	1	60

Switch table
(initially empty)

Link Layer and LANs 6-63

63

Switch: frame filtering/forwarding

when frame received at switch:

- record incoming link, MAC address of sending host
- index switch table using MAC destination address
- if entry found for destination
 - if destination on segment from which frame arrived
 - then drop frame
 - else forward frame on interface indicated by entry
- else flood /* forward on all interfaces except arriving interface */

Link Layer and LANs 6-64

64

Self-learning, forwarding: example

- frame destination, A', location unknown: **flood**
- destination A location known: **selectively send on just one link**

MAC addr	interface	TTL
A	1	60
A'	4	60

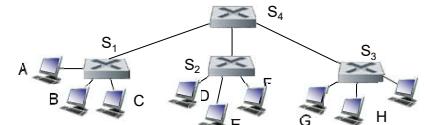
switch table
(initially empty)

Link Layer and LANs 6-65

65

Interconnecting switches

self-learning switches can be connected together:



- Q:** sending from A to G - how does S1 know to forward frame destined to G via S4 and S3?

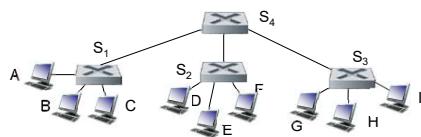
- A:** self learning! (works exactly the same as in single-switch case!)

Link Layer and LANs 6-66

66

Self-learning multi-switch example

Suppose C sends frame to I, I responds to C

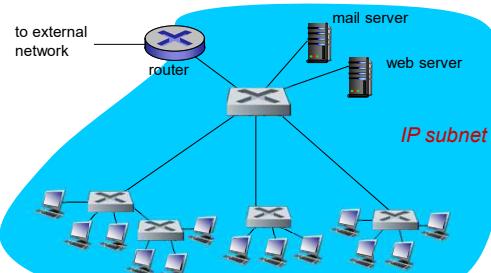


- Q: show switch tables and packet forwarding in S₁, S₂, S₃, S₄

Link Layer and LANs 6-67

67

Institutional network



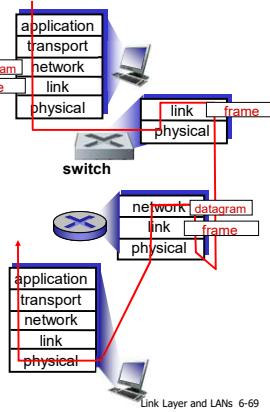
Link Layer and LANs 6-68

68

Switches vs. routers

both are store-and-forward:

- routers:** network-layer devices (examine network-layer headers)
- switches:** link-layer devices (examine link-layer headers)



Link Layer and LANs 6-69

69

VLANs: motivation

consider:

- CS user moves office to EE, but wants connect to CS switch?
- single broadcast domain:
 - all layer-2 broadcast traffic (ARP, DHCP, unknown location of destination MAC address) must cross entire LAN
 - security/privacy, efficiency issues

Link Layer and LANs 6-70

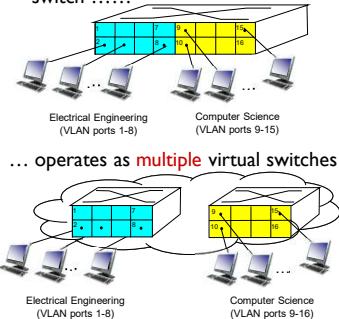
70

VLANs

Virtual Local Area Network

switch(es) supporting VLAN capabilities can be configured to define multiple **virtual** LANs over single physical LAN infrastructure.

port-based VLAN: switch ports grouped (by switch management software) so that **single** physical switch



Link Layer and LANs 6-71

71

Port-based VLAN

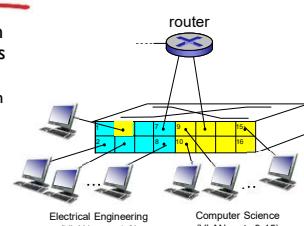
traffic isolation: frames to/from ports 1-8 can only reach ports 1-8

- can also define VLAN based on MAC addresses of endpoints, rather than switch port

dynamic membership: ports can be dynamically assigned among VLANs

forwarding between VLANs: done via routing (just as with separate switches)

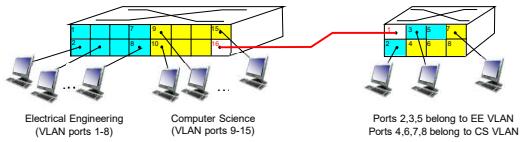
- in practice vendors sell combined switches plus routers



Link Layer and LANs 6-72

72

VLANs spanning multiple switches

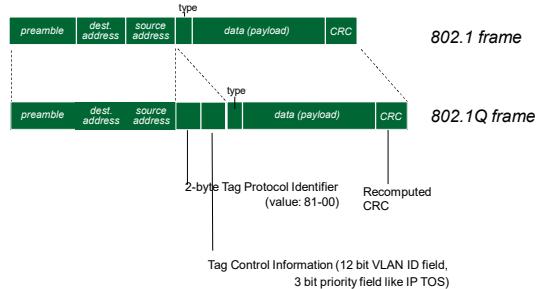


- **trunk port:** carries frames between VLANs defined over multiple physical switches
 - frames forwarded within VLAN between switches can't be vanilla 802.1 frames (must carry VLAN ID info)
 - 802.1Q protocol adds/removed additional header fields for frames forwarded between trunk ports

Link Layer and LANs 6-73

73

802.1Q VLAN frame format



Link Layer and LANs 6-74

74

Link layer, LANs: outline

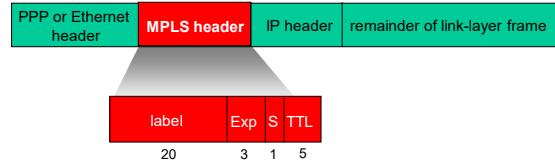
- | | |
|---|--|
| 6.1 introduction, services
6.2 error detection, correction
6.3 multiple access protocols
6.4 LANs <ul style="list-style-type: none"> • addressing, ARP • Ethernet • switches • VLANs | 6.5 link virtualization: MPLS
6.6 data center networking
6.7 a day in the life of a web request |
|---|--|

Link Layer and LANs 6-75

75

Multiprotocol label switching (MPLS)

- initial goal: high-speed IP forwarding using fixed length label (instead of IP address)
 - fast lookup using fixed length identifier (rather than shortest prefix matching)
 - borrowing ideas from Virtual Circuit (VC) approach
 - but IP datagram still keeps IP address!



Link Layer and LANs 6-76

76

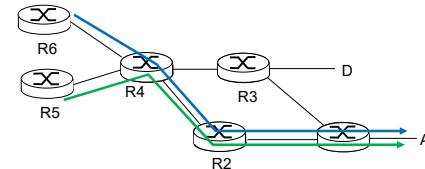
MPLS capable routers

- a.k.a. label-switched router
- forward packets to outgoing interface based only on label value (*don't inspect IP address*)
 - MPLS forwarding table distinct from IP forwarding tables
- **flexibility:** MPLS forwarding decisions can differ from those of IP
 - use destination *and* source addresses to route flows to same destination differently (traffic engineering)
 - re-route flows quickly if link fails: pre-computed backup paths (useful for VoIP)

Link Layer and LANs 6-77

77

MPLS versus IP paths

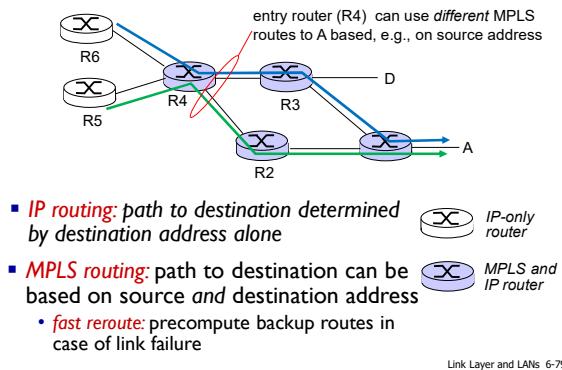


- **IP routing:** path to destination determined by destination address alone

Link Layer and LANs 6-78

78

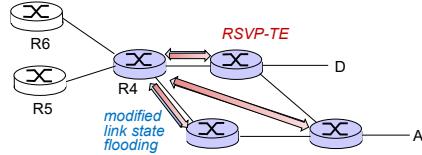
MPLS versus IP paths



79

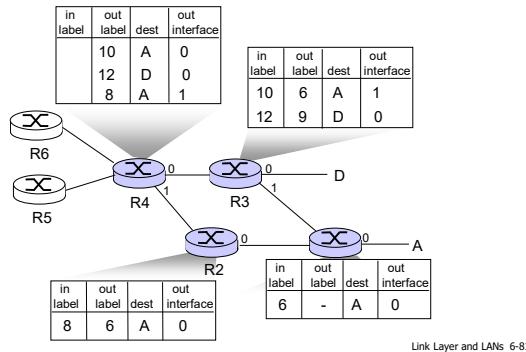
MPLS signaling

- modify OSPF, IS-IS link-state flooding protocols to carry info used by MPLS routing,
 - e.g., link bandwidth, amount of “reserved” link bandwidth
- **entry MPLS router** uses **RSVP-TE signaling protocol** to set up **MPLS forwarding at downstream routers**



80

MPLS forwarding tables



81

Link layer, LANs: outline

- 6.1 introduction, services
- 6.2 error detection, correction
- 6.3 multiple access protocols
- 6.4 LANs
 - addressing, ARP
 - Ethernet
 - switches
 - VLANS
- 6.5 link virtualization: MPLS
- 6.6 data center networking
- 6.7 a day in the life of a web request

Link Layer and LANs 6-82

82

Data center networks

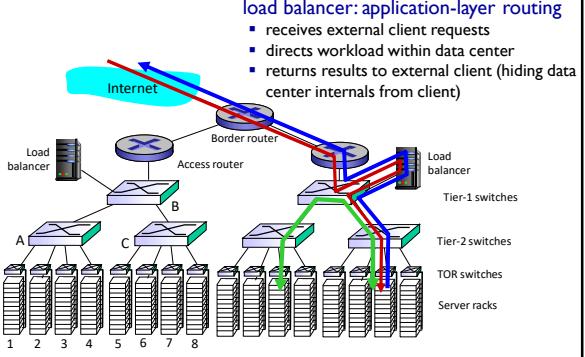
- 10's to 100's of thousands of hosts, often closely coupled, in close proximity:
 - e-business (e.g. Amazon)
 - content-servers (e.g., YouTube, Akamai, Apple, Microsoft)
 - search engines, data mining (e.g., Google)
- challenges:
 - multiple applications, each serving massive numbers of clients
 - managing/balancing load, avoiding processing, networking, data bottlenecks



Link Layer and LANs 6-83

83

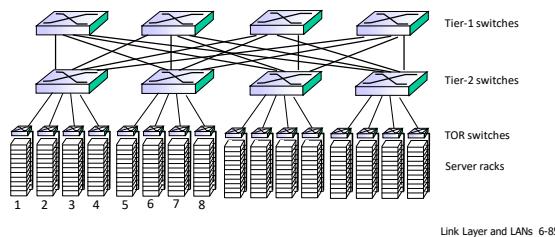
Data center networks



84

Data center networks

- rich interconnection among switches, racks:
 - increased throughput between racks (multiple routing paths possible)
 - increased reliability via redundancy



Link Layer and LANs 6-85

85

Link layer, LANs: outline

- | | |
|---|--|
| 6.1 introduction, services
6.2 error detection, correction
6.3 multiple access protocols
6.4 LANs <ul style="list-style-type: none"> • addressing, ARP • Ethernet • switches • VLANs | 6.5 link virtualization: MPLS
6.6 data center networking
6.7 a day in the life of a web request |
|---|--|

Link Layer and LANs 6-86

86

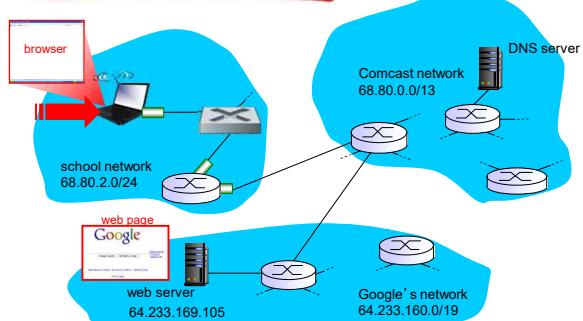
Synthesis: a day in the life of a web request

- journey down protocol stack complete!
 - application, transport, network, link
- putting-it-all-together: synthesis!
 - **goal:** identify, review, understand protocols (at all layers) involved in seemingly simple scenario: requesting www page
 - **scenario:** student attaches laptop to campus network, requests/receives www.google.com

Link Layer and LANs 6-87

87

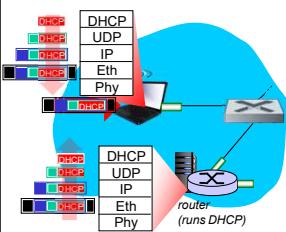
A day in the life: scenario



Link Layer and LANs 6-88

88

A day in the life... connecting to the Internet

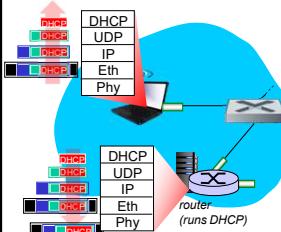


- connecting laptop needs to get its own IP address, addr of first-hop router, addr of DNS server: use **DHCP**
- DHCP request **encapsulated** in **UDP**, encapsulated in **IP**, encapsulated in **802.3** Ethernet
- Ethernet frame **broadcast** (dest: FFFFFFFFFFFF) on LAN, received at router running **DHCP** server
- Ethernet **demuxed** to IP demuxed, UDP demuxed to DHCP

Link Layer and LANs 6-89

89

A day in the life... connecting to the Internet



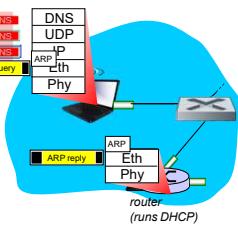
- DHCP server formulates **DHCP ACK** containing client's IP address, IP address of first-hop router for client, name & IP address of DNS server
- encapsulation at DHCP server, frame forwarded (**switch learning**) through LAN, demultiplexing at client
- DHCP client receives DHCP ACK reply

Client now has IP address, knows name & addr of DNS server, IP address of its first-hop router

Link Layer and LANs 6-90

90

A day in the life... ARP (before DNS, before HTTP)

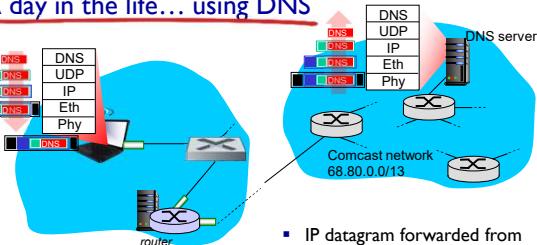


- before sending **HTTP** request, need IP address of www.google.com: **DNS**
- DNS query created, encapsulated in UDP, encapsulated in IP, encapsulated in Eth. To send frame to router, need MAC address of router interface: **ARP**
- **ARP query** broadcast, received by router, which replies with **ARP reply** giving MAC address of router interface
- client now knows MAC address of first hop router; so can now send frame containing DNS query

Link Layer and LANs 6-91

91

A day in the life... using DNS

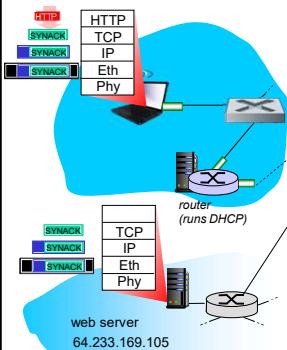


- IP datagram forwarded from campus network into Comcast network, routed (tables created by **RIP, OSPF, IS-IS** and/or **BGP** routing protocols) to DNS server
- demuxed to DNS server
- DNS server replies to client with IP address of www.google.com

Link Layer and LANs 6-92

92

A day in the life... TCP connection carrying HTTP

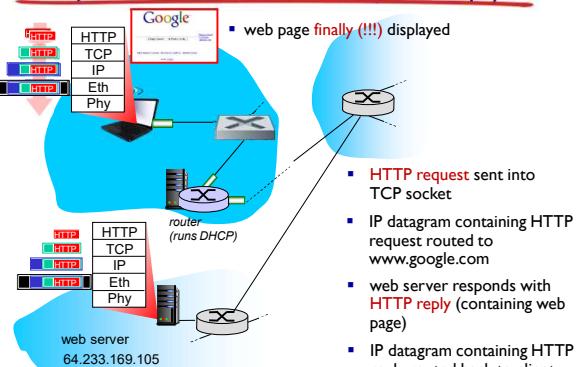


- to send HTTP request, client first opens **TCP socket** to web server
- TCP **SYN segment** (step 1 in 3-way handshake) inter-domain routed to web server
- web server responds with **TCP SYNACK** (step 2 in 3-way handshake)
- **TCP connection established!**

Link Layer and LANs 6-93

93

A day in the life... HTTP request/reply



- **HTTP request sent into TCP socket**
- IP datagram containing HTTP request routed to www.google.com
- web server responds with **HTTP reply** (containing web page)
- IP datagram containing HTTP reply routed back to client

Link Layer and LANs 6-94

94

Chapter 6: Summary

- principles behind data link layer services:
 - error detection, correction
 - sharing a broadcast channel: multiple access
 - link layer addressing
- instantiation and implementation of various link layer technologies
 - Ethernet
 - switched LANS, VLANs
 - virtualized networks as a link layer: MPLS
- synthesis: a day in the life of a web request

Link Layer and LANs 6-95

95

Chapter 6: let's take a breath

- journey down protocol stack **complete** (except PHY)
- solid understanding of networking principles, practice
- could stop here but **lots** of interesting topics!
 - wireless
 - multimedia
 - security

Link Layer and LANs 6-96

96

4/19/2022

6.1 Link layer serves: ppt (review)

6.2 error checksum

Parity Checking: ppt (~6:51)



(~1d:17) one error is fine, multiple error is no good.

CRC example

→: (21:11) exam Q:

(23:15) Mac protocols: (PPT)

channel partitioning Mac protocol:

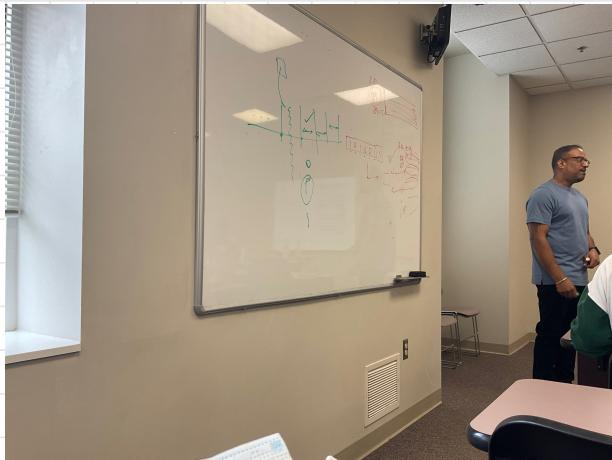
fair, not efficient (26:4)



Random Access Protocol (30:30~)

slotted ALOHA

time synchronized.



P ↑ , ↓ collision, (~38:00)

slotted ALOHA : efficiency 1/P(1 - e^{-P}) (23:00~)

pure ALLOHA

: lower even

CSMA collisions

双向 collision

CIMA / CD