

# Computer Networking

## Slide Set 1

Andrew W. Moore

[Andrew.Moore@cl.cam.ac.uk](mailto:Andrew.Moore@cl.cam.ac.uk)

2024-2025 update

## Topic 1 Foundation

- Administrivia
- Networks
- Channels
- Multiplexing
- Performance: loss, delay, throughput

2

## Course Administration

### Commonly Available Texts

- ❑ Computer Networks: A Systems Approach  
Peterson and Davie  
<https://book.systemsapproach.org>  
<https://github.com/SystemsApproach/book>
- ❑ Computer Networking : Principles, Protocols and Practice  
Olivier Bonaventure (and friends)  
Less GitHub but more practical exercises  
<https://www.computer-networking.info/>

Other textbooks are available, details on the materials tab of webpage.

3

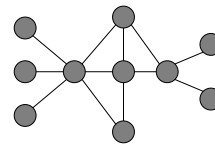
## Thanks

- Slides are a fusion of material from  
to Stephen Strowes, Tilman Wolf & Mike Zink, Ashish Padalkar , Evangelia Kalyvianaki, Brad Smith, Ian Leslie, Richard Black, Jim Kurose, Keith Ross, Larry Peterson, Bruce Davie, Jen Rexford, Ion Stoica, Vern Paxson, Scott Shenker, Frank Kelly, Stefan Savage, Jon Crowcroft , Mark Handley, Sylvia Ratnasamy, Adam Greenhalgh, and Anastasia Courtney.
- Supervision material is drawn from  
Stephen Kell, Andy Rice, and the TA teams of 144 and 168
- Finally thanks to the fantastic past Part 1b students and Andrew Rice for all the tremendous feedback.

4

## What is a network?

- A system of “links” that interconnect “nodes” in order to move “information” between nodes

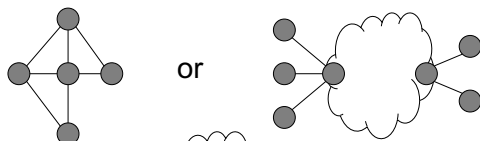


- Yes, this is all rather abstract

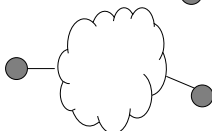
5

## What is a network?

- We also talk about



or even



- Yes, abstract, vague, and under-defined....

6

## There are *many* different types of networks

- Internet
- Telephone network
- Transportation networks
- Cellular networks
- Supervisory control and data acquisition networks
- Optical networks
- Sensor networks
- Satellite networks

We will focus almost exclusively on the Internet

7



## The Internet has transformed everything

- The way we do business
  - E-commerce, advertising, cloud-computing
- The way we have relationships
  - Facebook friends, E-mail, IM, virtual worlds
- The way we learn
  - Wikipedia, search engines
- The way we govern and view law
  - E-voting, censorship, copyright, cyber-attacks

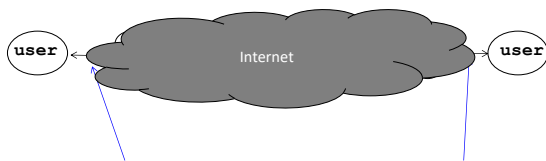
8

## A few defining characteristics of the Internet

9

## A federated system

- The Internet ties together different networks
  - >22,000 ISP networks (the definition is fuzzy)



Tied together by IP -- the “Internet Protocol” : a single common interface between users and the network and between networks

10

## A federated system

- The Internet ties together different networks
  - >22,000 ISP networks
- A single, common interface is great for interoperability...
- ...but tricky for business
- Why does this matter?
  - ease of interoperability is the Internet's most important goal
  - practical realities of incentives, economics and real-world trust, drive topology, route selection and service evolution

11

## Enormous diversity and dynamic range

- Communication latency: nanoseconds to seconds ( $10^9$ )
- Bandwidth: 100bits/second to 1.600 Terabits/second ( $10^{12}$ )
- Packet loss: 0 – 90%
- Technology: optical, wireless, satellite, copper
- **Endpoint devices**: from sensors and cell phones to datacenters and supercomputers
- **Applications**: social networking, file transfer, skype, live TV, gaming, remote medicine, backup, IM
- **Users**: the governing, governed, operators, **malicious**, naïve, savvy, embarrassed, paranoid, addicted, cheap ...

14

## Constant Evolution

1970s:

- 56kilobits/second “backbone” links
- <100 computers, a handful of sites in the US (and one UK)
- Telnet and file transfer were the “killer” applications

Today

- 400+Gigabits/second backbone links
- 40B+ devices, all over the globe
  - 27B+ IoT devices alone

15

## Asynchronous Operation

- Fundamental constraint: **speed of light**
- Consider:
  - How many cycles does your 3GHz CPU in Cambridge execute before it can possibly get a response from a message it sends to a server in Palo Alto?
    - Cambridge to Palo Alto: 8,609 km
    - Traveling at 300,000 km/s: 28.70 milliseconds
    - Then back to Cambridge:  $2 \times 28.70 = 57.39$  milliseconds
    - $3,000,000,000 \text{ cycles/sec} \times 0.05739 = 172,179,999$  cycles!
- Thus, communication feedback is always *dated*

How much can change with 172 Million instructions

16

## Prone to Failure

- To send a message, **all** components along a path must function correctly
  - software, wireless access point, firewall, links, network interface cards, switches,...
  - Including **human operators**
- Consider: 50 components in a system, each working correctly 99% of time → 39.5% chance communication will fail
- Plus, recall
  - scale → lots of components
  - asynchrony → takes a long time to hear (bad) news
  - federation (**internet**) → hard to identify fault or assign blame

17

## Recap: The Internet is...

- A complex federation
- Of enormous scale
- Dynamic range
- Diversity
- Constantly evolving
- Asynchronous in operation
- Failure prone
- Constrained by what's practical to engineer
- Too complex for (simple) theoretical models
- "Working code" doesn't mean much
- Performance benchmarks are too narrow

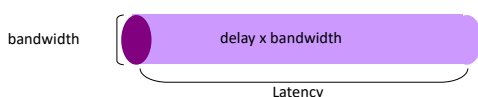
18

## An Engineered System

- Constrained by what technology is practical
  - Link bandwidths
  - Switch port counts
  - Bit error rates
  - **Cost**
  - ...

19

## Properties of Links (Channels)



- Bandwidth (capacity): "width" of the links
  - number of bits sent (or received) per unit time (bits/sec or bps)
- Latency (delay): "length" of the link
  - propagation time for data to travel along the link (seconds)
- Bandwidth-Delay Product (BDP): "volume" of the link
  - amount of data that can be "in flight" at any time
  - propagation delay  $\times$  bits/time = total bits in link

21

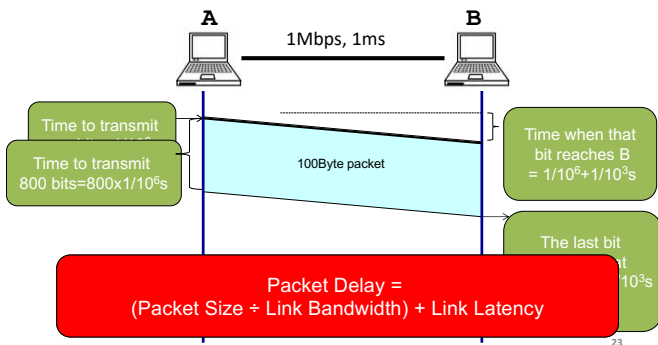
## Examples of Bandwidth-Delay

- Same city over a slow link:
  - BW~100Mbps
  - Latency~10msec
  - BDP  $\sim 10^6$ bits  $\sim 125$ KBytes
  - $17\text{km} / c = 56\mu\text{s} \ll 10\text{ms}$
- Intra Datacenter:
  - BW~100Gbps
  - Latency~30usec
  - BDP  $\sim 10^6$ bits  $\sim 375$ KBytes
  - $750\text{m} / c = 56\mu\text{s} \approx 30\mu\text{s}$
- To California over a fast link:
  - BW~10Gbps
  - Latency~140msec
  - BDP  $\sim 1.4 \times 10^9$ bits  $\sim 175$ MBytes
  - $9708\text{km} / c = 32\text{ms} \ll 140\text{ms}$
- Intra (inside) Host:
  - BW~800Gbps
  - Latency~16nsec
  - BDP  $\sim 12 \times 10^3$ bits  $\sim 5$ KBytes
  - $25\text{cm} / c = 83\text{ps} \ll 16\text{ns}$

22

## Packet Delay

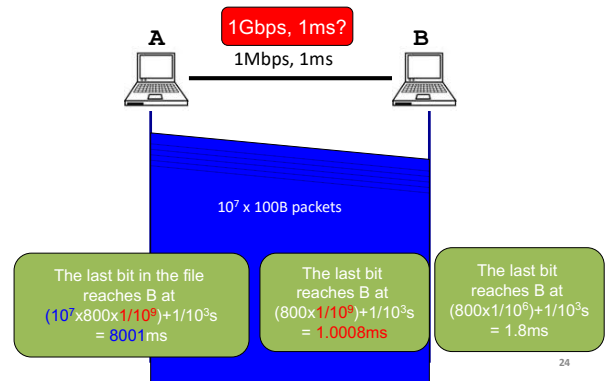
Sending a 100B packet from A to B?



23

## Packet Delay

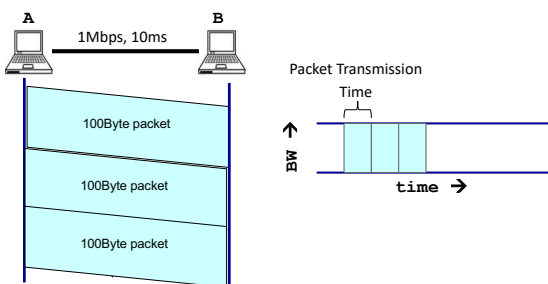
Sending a 100B packet from A to B?



24

## Packet Delay: The “pipe” view

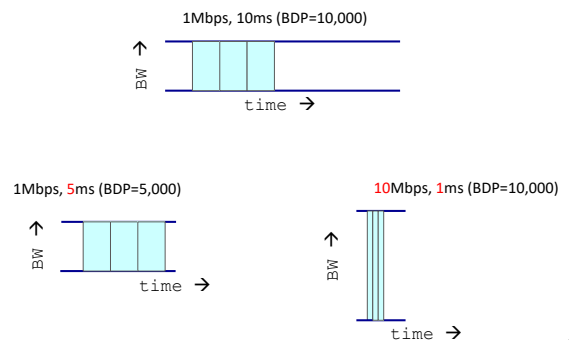
Sending 100B packets from A to B?



25

## Packet Delay: The “pipe” view

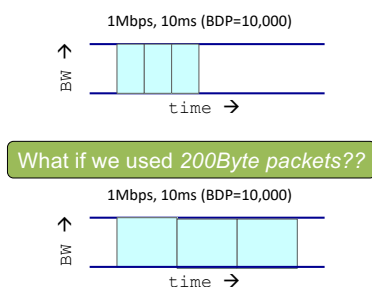
Sending 100B packets from A to B?



26

## Packet Delay: The “pipe” view

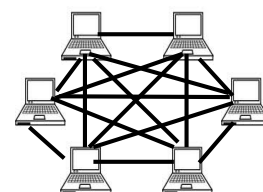
Sending 100B packets from A to B?



27

## What if we have more nodes?

One link for every node?

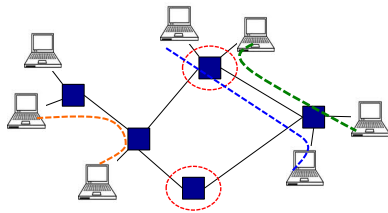


Need a scalable way to interconnect nodes

29

## Solution: A *switched* network

Nodes share network link resources

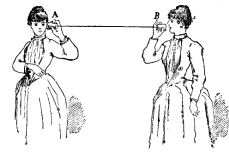


How is this sharing implemented?

30

## Two examples of switched networks

- Circuit switching (used in the *POTS*: Plain Old Telephone system) emphasis on old



- Packet switching (used in the Internet)

31

## Circuit switching



Telephone



Exchange



Exchange

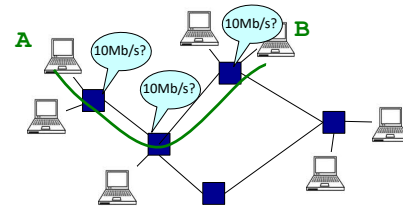


Telephone



## Circuit switching

Idea: source **reserves** network capacity along a path



- (1) Node A sends a reservation request
- (2) Interior switches establish a connection -- i.e., "circuit"
- (3) A starts sending data
- (4) A sends a "teardown circuit" message

33

## Multiplexing



Sharing makes things efficient (cost less)

- One airplane/train for 100's of people
- One telephone for many calls
- One lecture theatre for many classes
- One computer for many tasks
- One network for many computers
- One datacenter many applications

34

## Multiplexing

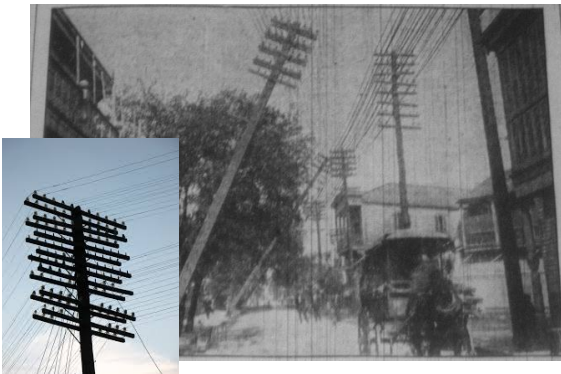


Sharing makes things efficient (cost less)

- One airplane/train for 100's of people
- One telephone for many calls
- One ~~Lecturer?~~ lecture theatre for many classes
- One computer for many tasks
- One network for many computers
- One datacenter many applications

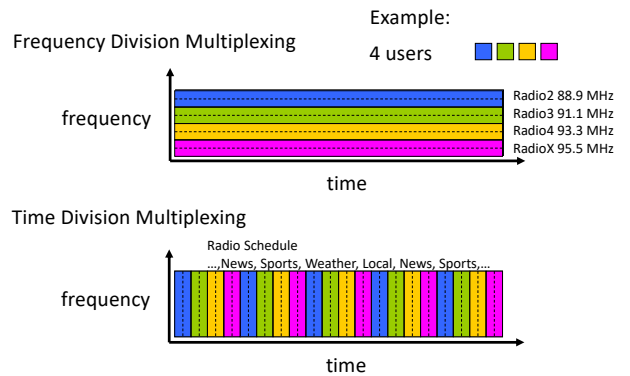
35

## Old Time Multiplexing



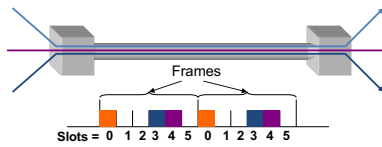
36

## Sharing Circuit Switching: FDM and TDM



37

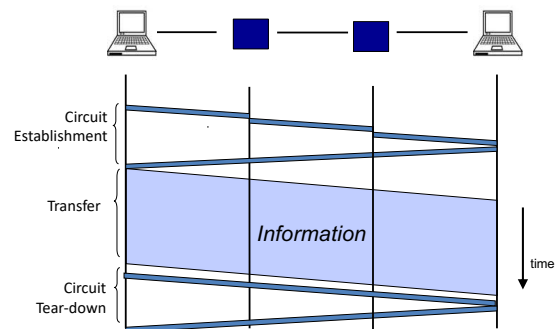
## Time-Division Multiplexing/Demultiplexing



- Time divided into frames; frames into slots
- Relative slot position inside a frame determines to which conversation data belongs
  - e.g., slot 0 belongs to **orange** conversation
- Slots are reserved (released) during circuit setup (teardown)
- If a conversation does not use its circuit **capacity is lost!**

38

## Timing in Circuit Switching



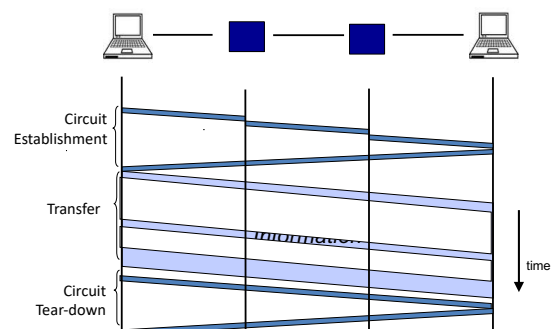
39

## Circuit switching: pros and cons

- Pros
  - guaranteed performance
  - fast transfer (once circuit is established)
- Cons

40

## Timing in Circuit Switching



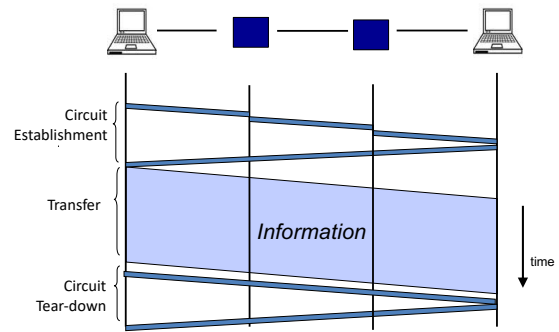
41

## Circuit switching: pros and cons

- Pros
  - guaranteed performance
  - fast transfer (once circuit is established)
- Cons
  - **wastes bandwidth if traffic is “bursty”**

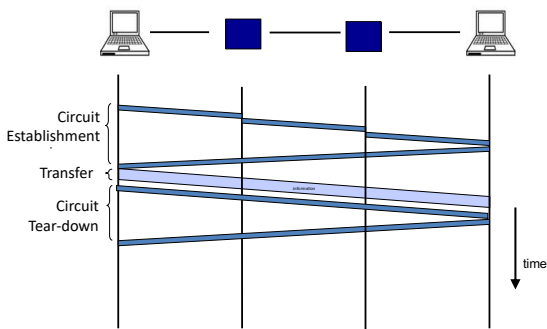
42

## Timing in Circuit Switching



43

## Timing in Circuit Switching



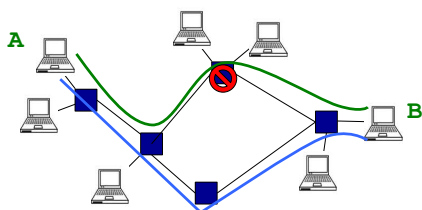
44

## Circuit switching: pros and cons

- Pros
  - guaranteed performance
  - fast transfers (once circuit is established)
- Cons
  - wastes bandwidth if traffic is “bursty”
  - **connection setup time is overhead**

45

## Circuit switching



Circuit switching doesn't “route around failure”

46

## Circuit switching: pros and cons

- Pros
  - guaranteed performance
  - fast transfers (once circuit is established)
- Cons
  - wastes bandwidth if traffic is “bursty”
  - connection setup time is overhead
  - **recovery from failure is slow**

47

## Numerical example

- How long does it take to send a file of 640,000 bits from host A to host B over a circuit-switched network?
  - All links are 1.536 Mbps
  - Each link uses TDM with 24 slots/sec
  - 500 msec to establish end-to-end circuit

Let's work it out!

48

## Two examples of switched networks

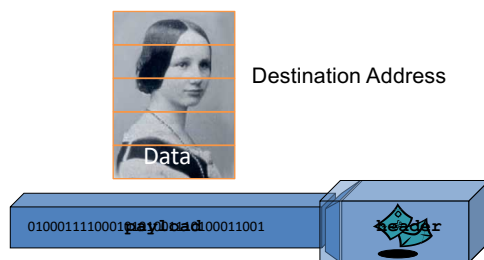
- Circuit switching (used in the *POTS*: Plain Old Telephone system)
- Packet switching (used in the Internet)



49

## Packet Switching

- Data is sent as chunks of formatted bits (**Packets**)
- Packets consist of a “**header**” and “**payload**”\*



After Nick McKeown © 2006

51

## Packet Switching

- Data is sent as chunks of formatted bits (**Packets**)
- Packets consist of a “**header**” and “**payload**”\*
  - payload is the data being carried
  - header holds instructions to the network for how to handle packet (think of the header as an API)
    - In this example, the header has a destination address
    - More complex headers may include
      - How this traffic should be handled? (first class, second class, etc)
      - Do I acknowledge this? Who signed for it?
      - Were the contents ok?

52

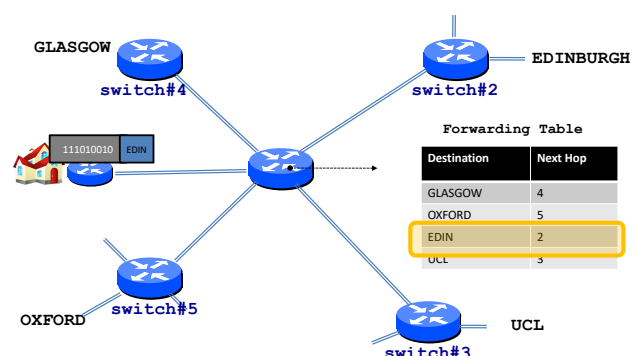
## Packet Switching

- Data is sent as chunks of formatted bits (Packets)
- Packets consist of a “header” and “payload”
- Switches “**forward**” packets based on their headers

*A switch looks at the header and immediately decides which physical port*  
***In a switch: address maps to port***

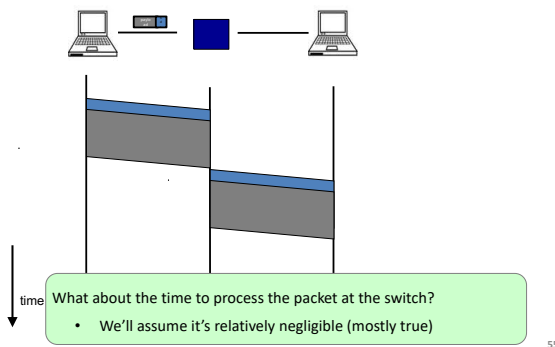
53

## Switches forward packets

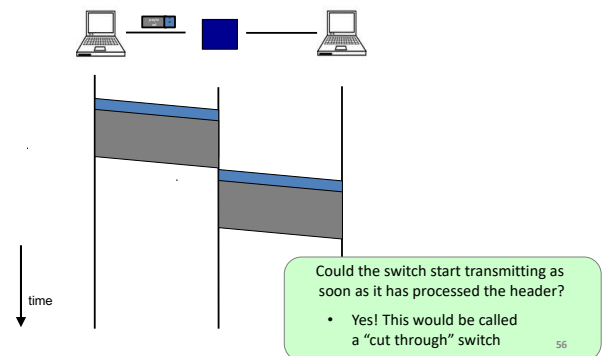


54

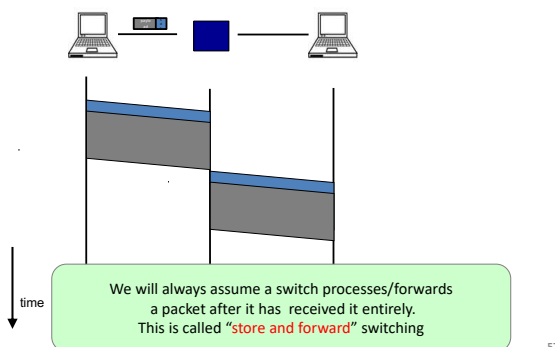
## Timing in Packet Switching



## Timing in Packet Switching



## Timing in Packet Switching



## Packet Switching

- Data is sent as chunks of formatted bits (Packets)
- Packets consist of a "header" and "payload"
- Switches "**forward**" packets based on their headers

58

## Packet Switching

- Data is sent as chunks of formatted bits (Packets)
- Packets consist of a "header" and "payload"
- Switches "**forward**" packets based on their headers
- Each packet travels independently
  - no notion of packets belonging to a "circuit"

59

## Packet Switching

- Data is sent as chunks of formatted bits (Packets)
- Packets consist of a "header" and "payload"
- Switches "**forward**" packets based on their headers
- Each packet travels independently
- No link resources are reserved in advance. Instead packet switching leverages **statistical multiplexing** (stat muxing)

60



## Multiplexing

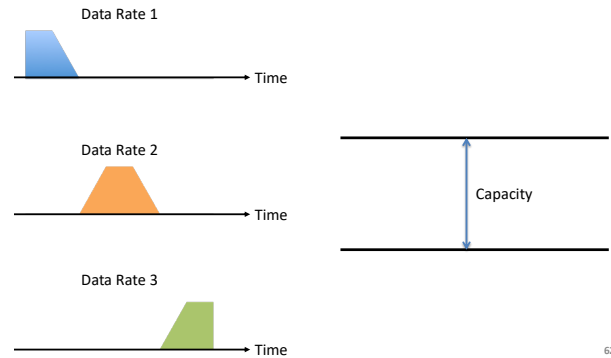


Sharing makes things efficient (cost less)

- One airplane/train for 100's of people
- One telephone for many calls
- One lecture theatre for many classes
- One computer for many tasks
- One network for many computers
- One datacenter many applications

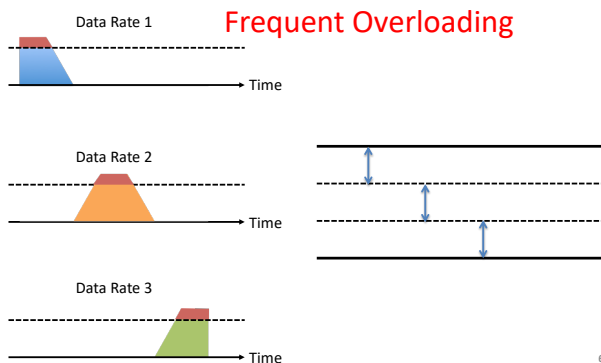
61

## Three Flows with Bursty Traffic



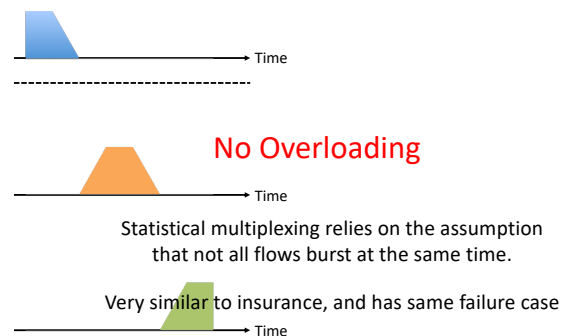
62

## When Each Flow Gets 1/3<sup>rd</sup> of Capacity



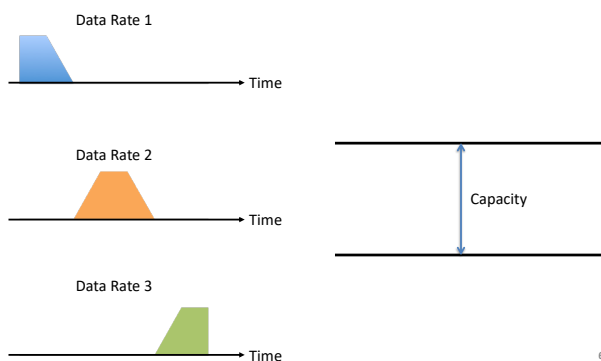
63

## When Flows Share Total Capacity



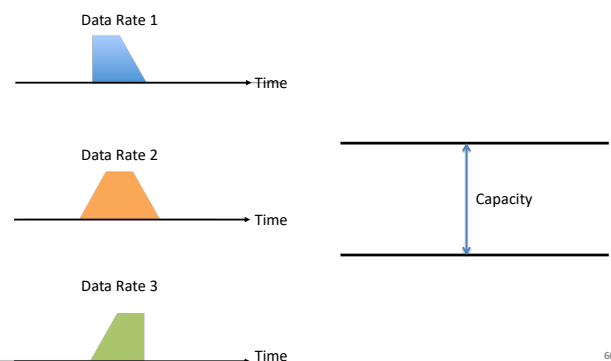
64

## Three Flows with Bursty Traffic



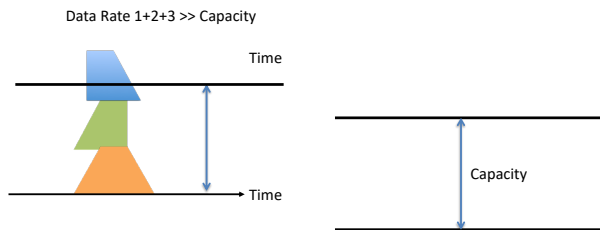
65

## Three Flows with Bursty Traffic



66

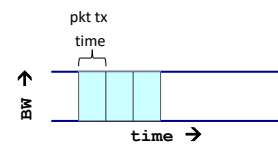
## Three Flows with Bursty Traffic



What do we do under overload?

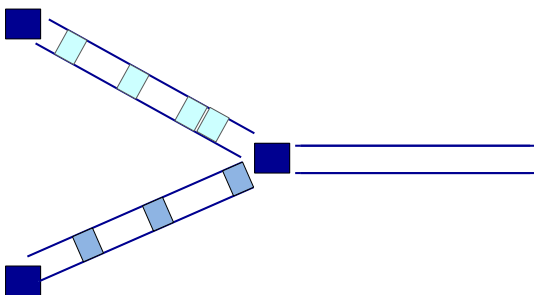
67

## Statistical multiplexing: pipe view



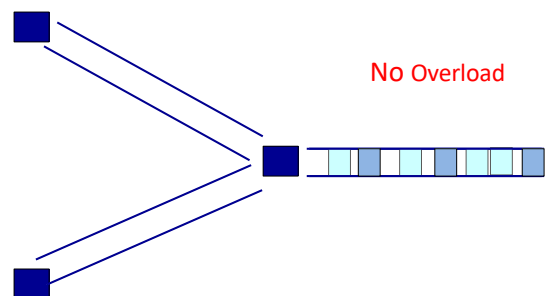
68

## Statistical multiplexing: pipe view



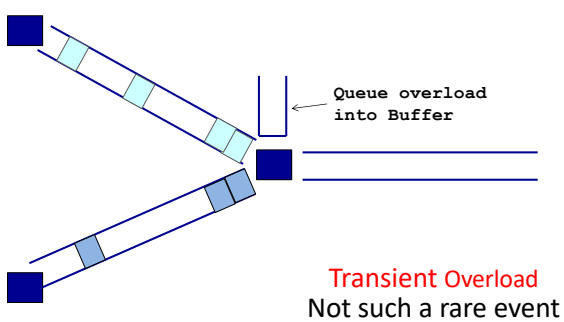
69

## Statistical multiplexing: pipe view



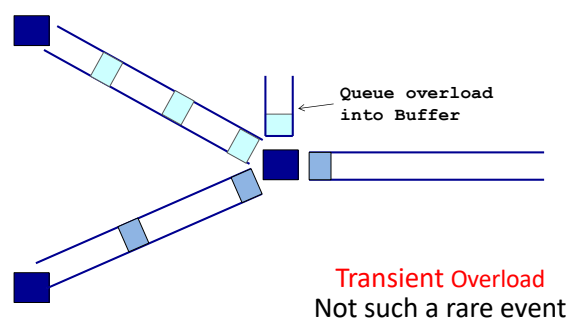
70

## Statistical multiplexing: pipe view



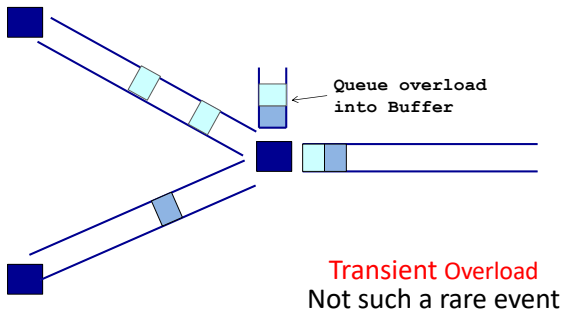
71

## Statistical multiplexing: pipe view



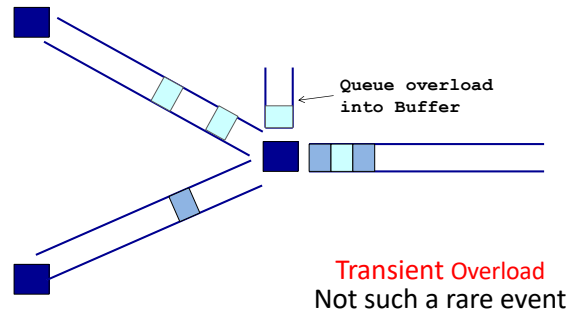
72

## Statistical multiplexing: pipe view



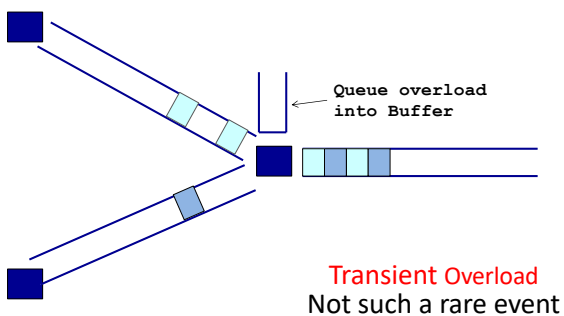
73

## Statistical multiplexing: pipe view



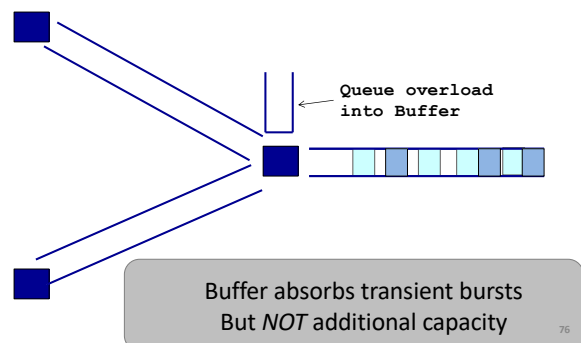
74

## Statistical multiplexing: pipe view



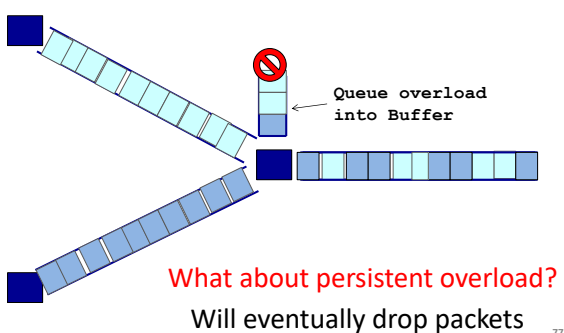
75

## Statistical multiplexing: pipe view



76

## Statistical multiplexing: pipe view



77

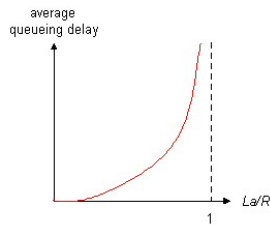
## Queues introduce queuing delays

- Recall,
 
$$\text{packet delay} = \text{transmission delay} + \text{propagation delay} (*)$$
  - With queues (statistical multiplexing)
 
$$\text{packet delay} = \text{transmission delay} + \text{propagation delay} + \text{queuing delay} (*)$$
  - Queuing delay caused by "packet interference"
  - Made worse at high load
    - less "idle time" to absorb bursts
    - think about traffic jams at rush hour or rail network failure
- (\*) plus per-hop processing delay that we define as negligible)

78

## Queuing delay extremes

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



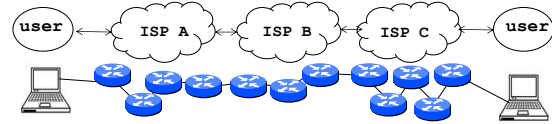
traffic intensity =  $La/R$

- $La/R \approx 0$ : average queuing delay small
- $La/R \rightarrow 1$ : delays become large
- $La/R > 1$ : more "work" arriving than can be serviced, average delay infinite – or data is lost (dropped).

79

## Recall the Internet *federation*

- The Internet ties together different networks
  - >20,000 ISP networks



We can see (hints) of the nodes and links using traceroute...

80

## "Real" Internet delays and routes

traceroute: department ssh server to melbourneisp.com (Melbourne)  
(tracepath on windows is similar)

```

awm22@svr-ssh-0:~$ traceroute melbourneisp.com
traceroute to melbourneisp.com (116.206.130.24), 30 hops max, 60 byte packets
 1 vian398.gatwick.net.c1.cam.ac.uk (128.232.64.2) 10.299 ms 10.593 ms 11.864 ms
 2 cl-mgb.d-mw.net.cam.ac.uk (193.60.89.5) 0.743 ms 0.789 ms 0.506 ms
 3 d-mw-c-ce.net.cam.ac.uk (131.111.6.53) 0.927 ms 1.123 ms 0.986 ms
 4 c-ce-b-jc.net.cam.ac.uk (131.111.6.82) 0.795 ms 0.771 ms 0.781 ms
 5 ips-out-b-jc.net.cam.ac.uk (131.111.7.217) 1.096 ms 0.906 ms 1.032 ms
 6 ae0.lowds-ban1.ja.net (146.97.41.37) 3.404 ms 3.019 ms 3.076 ms
 7 ae26.lowds-sbr1.ja.net (146.97.35.245) 3.740 ms 3.430 ms 3.374 ms
 8 ae31.londtw-sbr2.ja.net (146.97.33.38) 7.034 ms 6.631 ms 6.962 ms
 9 ae28.londtt-sbr1.ja.net (146.97.33.61) 8.820 ms 8.476 ms 10.954 ms
10 ae0.londtt-ban2.ja.net (146.97.35.194) 7.320 ms 6.467 ms 6.387 ms
11 ldn-b11-link.ip.twelve99.net (62.115.175.186) 6.476 ms 6.234 ms 6.585 ms
12 ldn-b01-link.ip.twelve99.net (62.115.138.268) 2.479 ms * 2.718 ms
13 nyl-b02-link.ip.twelve99.net (62.115.139.246) 76.127 ms nyl-b05-link.ip.twelve99.net (62.115.139.244) 75.315 ms *
14 * chi-bb2-link.ip.twelve99.net (62.115.132.135) 140.582 ms 140.836 ms
15 * den-bb2-link.ip.twelve99.net (62.115.115.76) 115.630 ms 114.872 ms
16 den-bb2-link.ip.twelve99.net (62.115.137.114) 113.977 ms * 113.656 ms
17 palo-bb2-link.ip.twelve99.net (62.115.139.112) 143.238 ms 144.527 ms *
18 * tpg-ic-387776.ip.twelve99-cust.net (62.115.188.5) 295.626 ms 296.153 ms
19 tpg-ic-387776.ip.twelve99-cust.net (62.115.188.5) 295.143 ms syd-apt-ros-crt3-be-100.tpgi.com.au (203.29.134.43) 291
20 syd-apt-ros-crt3-be-100.tpgi.com.au (203.29.134.43) 295.545 ms syd-sot-ken-crs2-fe-0-0-9.tpgi.com.au (203.26.22.12)
21 syd-sot-ken-crs2-fe-0-0-9.tpgi.com.au (203.26.22.12) 300.416 ms AU-VI-1015-IP0-221-Bundle-Ether1.tpgi.com.au (27.
22 AU-VI-4901-IP0-01-Eth-Trunk21.tpgi.com.au (203.220.216.30) 301.396 ms AU-VI-1015-IP0-221-Bundle-Ether2.tpgi.com.au (
23 AU-VI-4901-IP0-01-Eth-Trunk21.tpgi.com.au (203.220.216.30) 300.724 ms 14-202-130-170.static.tpgi.com.au (14.202.130.1
24 14-202-130-170.static.tpgi.com.au (14.202.130.170) 303.742 ms * 303.963 ms
25 *
26 *
27 *
28 *
29 *
30 *
awm22@svr-ssh-0:~$
    
```

Three delay measurements from  
rio.cl.cam.ac.uk to London JaNet gateway

Crossing the Atlantic Ocean

Crossing the Pacific Ocean

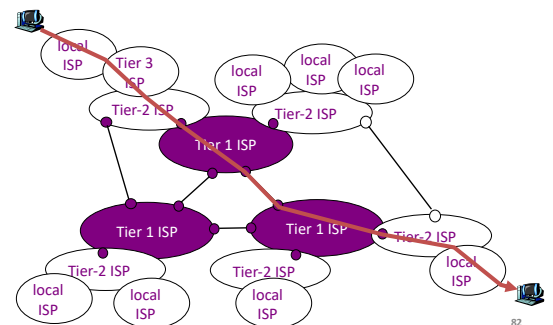
\* means no response (probe or reply lost, router not replying)

300ms RTT, 150ms one way Internet, 59.4ms by photon, 42ms by neutron

81

## Internet structure: network of networks

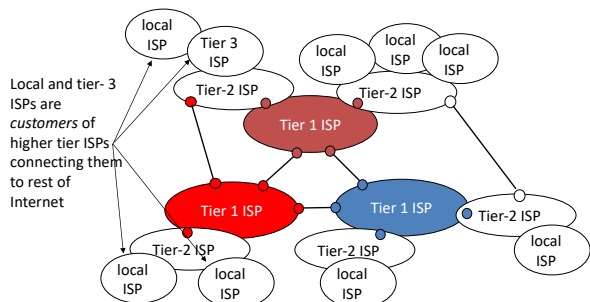
- a packet passes through many networks!



82

## Internet structure: network of networks

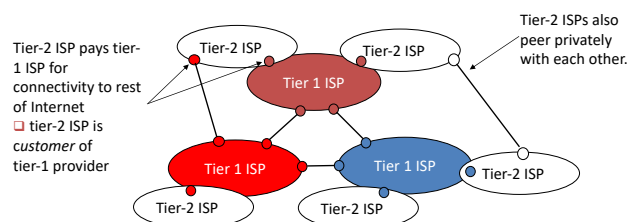
- "Tier-3" ISPs and local ISPs
  - last hop ("access") network (closest to end systems)



83

## Internet structure: network of networks

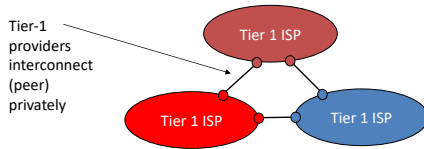
- "Tier-2" ISPs: smaller (often regional) ISPs
  - Connect to one or more tier-1 ISPs, possibly other tier-2 ISPs



84

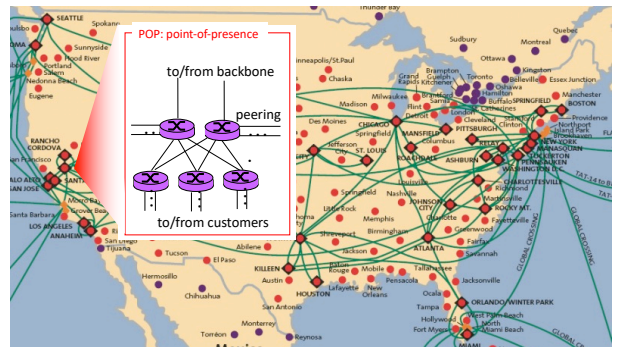
## Internet structure: network of networks

- roughly hierarchical
- **at center:** “tier-1” ISPs (e.g., Verizon, Sprint, AT&T, Cable and Wireless), national/international coverage
  - treat each other as equals



85

## Tier-1 ISP: e.g., Sprint



86

## Packet Switching

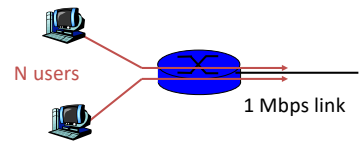
- Data is sent as chunks of formatted bits (Packets)
- Packets consist of a “header” and “payload”
- Switches “forward” packets based on their headers
- Each packet travels independently
- No link resources are reserved in advance. Instead, packet switching depends on **statistical multiplexing**
  - allows efficient use of resources
  - but introduces queues and queuing delays

87

## Packet switching versus circuit switching

*Packet switching may (does!) allow more users to use network*

- 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



Q: how did we get value 0.0004?

88

## Packet switching versus circuit switching

Q: how did we get value 0.0004?

- 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

Let U be number of users active  
N the total users  
P is 0.1 in our example to get 0.0004

$$\begin{aligned} P(U = k) &= \binom{n}{k} p^k (1-p)^{n-k} \\ \therefore P(U \leq K) &= \sum_{k=0}^K \binom{n}{k} p^k (1-p)^{n-k} \quad \left[ P(U > K) = 1 - \sum_{k=0}^K \binom{n}{k} p^k (1-p)^{n-k} \right] \\ \text{for } n=35, K=10 \\ P(U \leq 10) &= \sum_{k=0}^{10} \binom{35}{k} p^k (1-p)^{35-k} \\ \text{where } p=0.1: \\ P(U \leq 10) &= 0.99958 \\ \therefore P(U > 10) &= 0.00042 \end{aligned}$$

90

## Circuit switching: pros and cons

- Pros
  - guaranteed performance
  - fast transfers (once circuit is established)
- Cons
  - wastes bandwidth if traffic is “bursty”
  - connection setup adds delay
  - recovery from failure is slow

92

## Packet switching: pros and cons

- Pros
  - efficient use of bandwidth (stat. muxing)
  - no overhead due to connection setup
  - resilient -- can ‘route around trouble’
- Cons
  - no guaranteed performance
  - header overhead per packet
  - queues and queuing delays

93

## Summary

- A sense of how the basic ‘plumbing’ works
  - links and switches
  - packet delays = transmission + propagation + queuing + (negligible) per-switch processing
  - statistical multiplexing and queues
  - circuit vs. packet switching

94

## Topic 2 – Architecture and Philosophy

- Abstraction
- Layering
- Layers and Communications
- Entities and Peers
- What is a protocol?
- Protocol Standardization
- The architects process
  - How to break system into modules
  - Where modules are implemented
  - Where is state stored
- Internet Philosophy and Tensions

1

## Abstraction Concept

A mechanism for breaking down a problem

*what not how*

- eg Specification *versus* implementation
- eg Modules in programs

Allows replacement of implementations without affecting system behavior

*Vertical versus Horizontal*

“*Vertical*” what happens in a box “How does it attach to the network?”

“*Horizontal*” the communications paths running through the system

**Hint:** paths are built (“layered”) on top of other paths

3

## Computer System Modularity

Partition system into modules & abstractions:

- Well-defined interfaces give flexibility
  - **Hides** implementation - can be freely changed
  - Extend functionality of system by adding new modules
- E.g., libraries encapsulating set of functionality
- E.g., programming language + compiler abstracts away how the particular CPU works ...

4

## Computer System Modularity (cnt’ d)

- Well-defined interfaces hide information
  - Isolate **assumptions**
  - Present high-level **abstractions**
- **But can impair performance!**
- Ease of implementation vs worse performance

5

## Network System Modularity

Like software modularity, but:

- Implementation is distributed across many machines (routers and hosts)
- Must decide:
  - How to break system into modules
    - **Layering**
  - Where modules are implemented
    - **End-to-End Principle**
  - Where state is stored
    - **Fate-sharing**

6

## Layering Concept

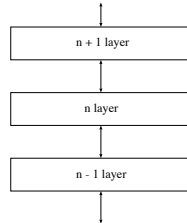
- A restricted form of abstraction: system functions are divided into layers, one built upon another
- Often called a *stack*; but **not** a data structure!

speaking 1	thoughts
speaking 2	words
speaking 3	phonemes
D/A, A/D	7 KHz analog voice
companding	8 K 12 bit samples per sec
multiplexing	8 KByte per sec stream
framing	Framed Byte Stream
modulation	Bitstream
	Analog signal

7

## Layers and Communications

- Interaction only between adjacent layers
- *layer n* uses services provided by *layer n-1*
- *layer n* provides service to *layer n+1*
- Bottom layer is physical media
- Top layer is application



8

## Entities and Peers

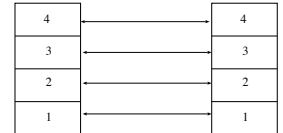
*Entity* – a *thing* (an independent existence)

Entities *interact* with the layers above and below

Entities *communicate* with *peer* entities

- same level but different place (eg different person, different box, different host)

Communications between peers is supported by entities at the lower layers



9

## Entities and Peers

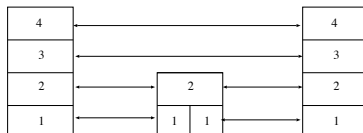
Entities usually do something useful

- Encryption – Error correction – Reliable Delivery
- Nothing at all is also reasonable

Not all communications is end-to-end

Examples for things in the middle

- IP Router – Mobile Phone Cell Tower
- Person translating French to English



10

## Layering and Embedding

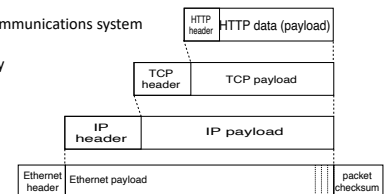
In Computer Networks we often see higher-layer information embedded within lower-layer information

- Such embedding can be considered a form of layering
- Higher layer information is generated by stripping off headers and trailers of the current layer
- eg an IP entity only looks at the IP headers

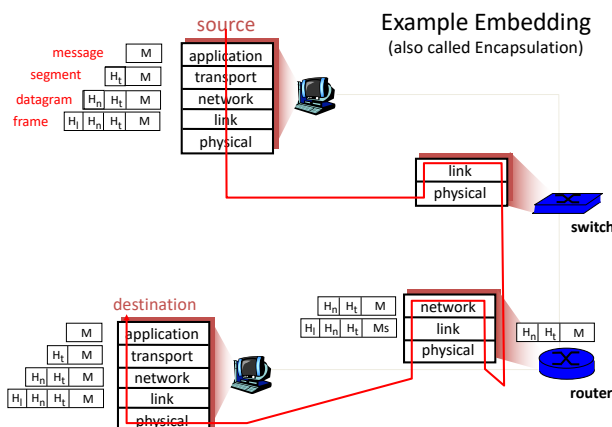
**BUT embedding is not the only form of layering**

Layering is to help understand a communications system

**NOT**  
determine implementation strategy

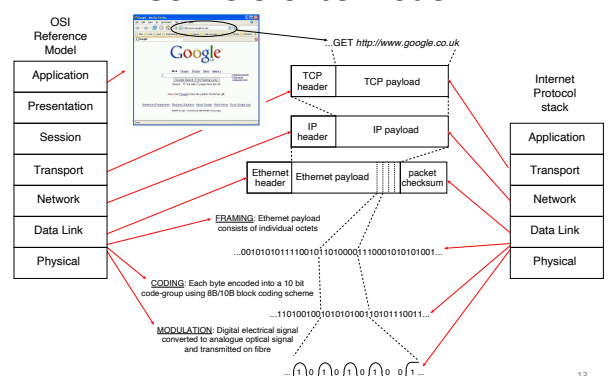


11



12

## Internet protocol stack versus OSI Reference Model

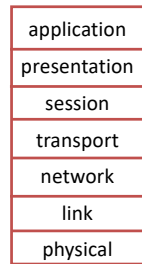


13



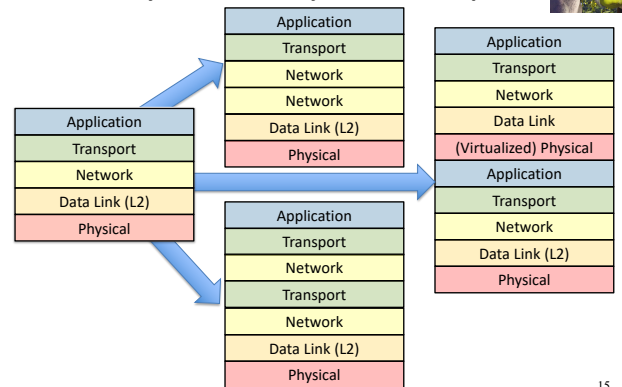
## ISO/OSI reference model

- **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



14

## Layers on Layers examples



15

## What is a protocol?

### human protocols:

- “what’s the time?”
- “I have a question”
- introductions

... specific msgs sent  
... specific actions taken  
when msgs received, or  
other events

### network protocols:

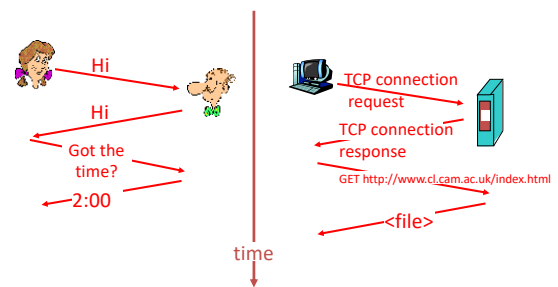
- machines rather than humans
- all communication activity in Internet governed by protocols

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

16

## What is a protocol?

a human protocol and a computer network protocol:

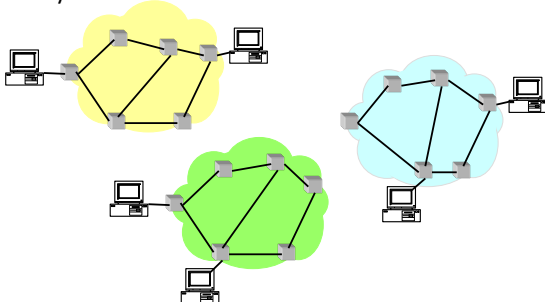


Q: Other human protocols?

17

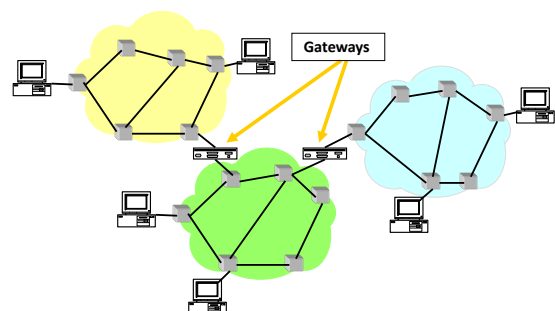
## So many Standards Problem

- Many different packet-switching networks
- Each with its own Protocol
- Only nodes on the same network could communicate



19

## INTERnet Solution



20

## Internet Design Goals (Clark '88)

- **Connect existing networks**
- Robust in face of failures
- Support multiple types of delivery services
- Accommodate a variety of networks
- Allow distributed management
- Easy host attachment
- Cost effective
- Allow resource accountability

21

## Real Goals

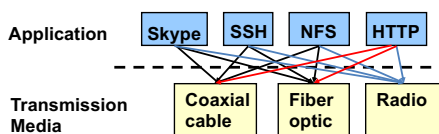
Internet Motto

*We reject kings , presidents, and voting. We believe in rough consensus and running code.* – David Clark

- **Build something that works!**
- Connect existing networks
- Robust in face of failures
- Support multiple types of delivery services
- Accommodate a variety of networks
- Allow distributed management
- Easy host attachment
- Cost effective
- ~~Allow resource accountability~~

22

## A Multitude of Apps Problem

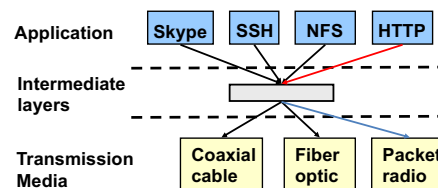


- Re-implement every application for every technology?
- No! But how does the Internet design avoid this?

23

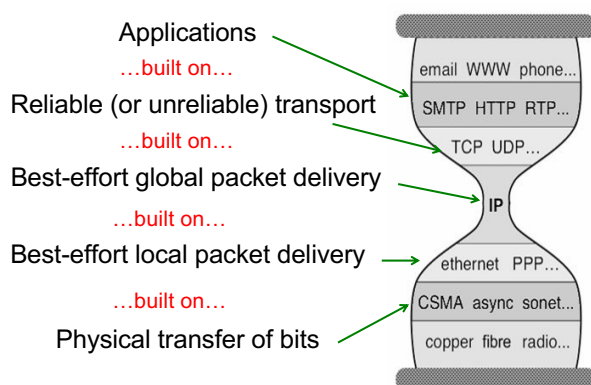
## Solution: Intermediate Layers

- Introduce intermediate layers that provide [set of abstractions](#) for various network functionality and technologies
  - A new app/media implemented only once
  - Variation on “add another level of indirection”



24

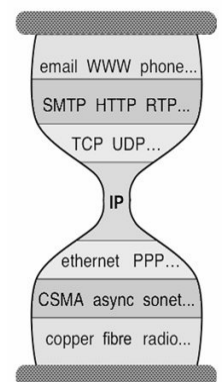
## In the context of the Internet



25

## Three Observations

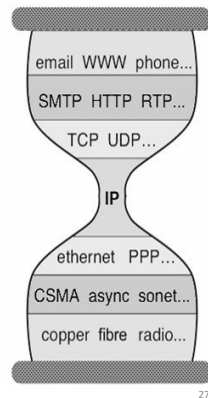
- Each layer:
  - Depends on layer below
  - Supports layer above
  - Independent of others
- Multiple versions in layer
  - Interfaces differ somewhat
  - Components pick which lower-level protocol to use
- But only one IP layer
  - Unifying protocol



26

## Layering Crucial to Internet's Success

- Reuse
- Hides underlying detail
- Innovation at each level can proceed in parallel
- Pursued by very different communities



27

What are some of the drawbacks of protocols and layering?

28

## Drawbacks of Layering

- Layer N may duplicate lower layer functionality
  - e.g., error recovery to retransmit lost data
- Information hiding may hurt performance
  - e.g., packet loss due to corruption vs. congestion
- Headers start to get really big
  - e.g., typical TCP+IP+Ethernet is 54 bytes
- Layer violations when the gains too great to resist
  - e.g., TCP-over-wireless
- Layer violations when network doesn't trust ends
  - e.g., firewalls

29

## Placing Network Functionality

- Hugely influential paper: "End-to-End Arguments in System Design" by Saltzer, Reed, and Clark ('84)
  - articulated as the "End-to-End Principle" (E2E)
- Endless debate over what it means
- Everyone cites it as supporting their position  
(regardless of the position!)

30

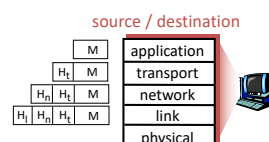
## Basic Observation

- Some application requirements can only be correctly implemented **end-to-end**
  - reliability, security, etc.
- Implementing these in the network is hard
  - every step along the way must be fail proof
- Hosts
  - **Can** satisfy the requirement without network's help
  - **Will/must** do so, since they can't rely on the network

31

## What Gets Implemented on Host?

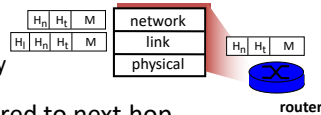
- Bits arrive on wire, must make it up to application
- Therefore, all layers must exist at the host



42

## What Gets Implemented on a Router?

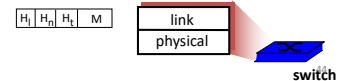
- Bits arrive on wire
  - Physical layer necessary
- Packets must be delivered to next-hop
  - Datalink layer necessary
- Routers participate in global delivery
  - Network layer necessary
- Routers don't support reliable delivery
  - Transport layer (and above) **not** supported



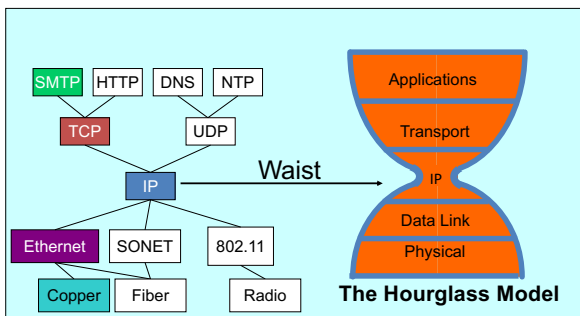
43

## What Gets Implemented on Switches?

- Switches do what routers do, except they don't participate in global delivery, just local delivery
- They only need to support Physical and Datalink
  - Don't need to support Network layer
- Won't focus on the router/switch distinction
  - Almost all boxes support network layer these days
  - Routers have switches but switches do not have routers



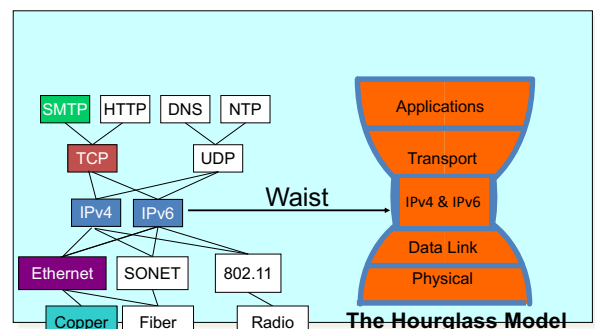
## The Internet Hourglass



There is just **one** network-layer protocol, **IP**.  
The "narrow waist" facilitates **interoperability**.

45

## The middle-age Internet Hourglass



There is just **TWO** network-layer protocol, **IPv4 + v6**  
The "narrow waist" facilitates **interoperability(???)**.

## Protocol Standardization

- All hosts must follow same protocol
  - Very small modifications can make a big difference
  - Or prevent it from working altogether
- This is why we have standards
  - Can have multiple implementations of protocol
- Internet Engineering Task Force (IETF)
  - Based on working groups that focus on specific issues
  - Produces "Request For Comments" (RFCs)
  - IETF Web site is <http://www.ietf.org>
  - RFCs archived at <http://www.rfc-editor.org>

47

## Alternative to Standardization?

- Have one implementation used by everyone
- Open-source projects
  - Which has had more impact, Linux or POSIX?
- Or just sole-sourced implementation
  - zoom, Signal, FaceTime, etc.
  - which in turn create monopolies...
  - for example, limit/remove interoperability

48

# Summary

- Engineering tools
  - Abstraction, Layering, Layers and Communications, Entities and Peers, Protocol as motivation, Examples of the *architects* process
  - Example: Internet Philosophy and Tensions
  - Decisions left un-made: Where to put stuff... What stuff needs putting...
- Standards
  - Political
  - Blind (and thus sometimes silly)
  - Good for user choice

## Topic 3a: The Physical Layer

### Our goals:

- Understand physical channel fundamentals
  - Physical channels can carry data in proportion to the signal and inversely in proportion to noise
  - Modulation represents Digital data in analog channels
  - Baseband vs. Broadband
  - Synchronous vs. Asynchronous

1

## Physical Channels / The Physical Layer

these example physical channels are also known as *Physical Media*

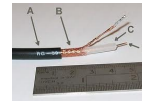
### Twisted Pair (TP)

- two insulated copper wires
  - Category 3: traditional phone wires, 10 Mbps Ethernet
  - Category 8: 25Gbps Ethernet
- Shielded (STP)
- Unshielded (UTP)



### Coaxial cable:

- two concentric copper conductors
- bidirectional
- baseband:
  - single channel on cable
  - legacy Ethernet
- broadband:
  - multiple channels on cable
  - HFC (Hybrid Fiber Coax)



### Fiber optic cable:

- high-speed operation
- point-to-point transmission
- (10' s-100' s Gbps)
- low error rate
- immune to electromagnetic noise



2

## More Physical media: Radio

- Bidirectional and multiple access
- propagation environment effects:
  - reflection
  - obstruction by objects
  - interference

### Radio link types:

- ❑ terrestrial microwave
  - ❖ e.g. 90 Mbps channels
- ❑ LAN (e.g., Wifi)
  - ❖ 11Mbps, 54 Mbps, 600 Mbps
- ❑ wide-area (e.g., cellular)
  - ❖ 5G cellular: ~ 40 Mbps - 10Gbps
- ❑ satellite
  - ❖ 27-50MHz typical bandwidth
  - ❖ geosynchronous versus low altitude
  - ❖ For geosync - 270 msec end-end delay to orbit



3

## Physical Channel Characteristics - Fundamental Limits -

**symbol type:** generally, an analog waveform — voltage, current, photo intensity etc.

**capacity:** bandwidth

**delay:** speed of light in medium and distance travelled

**fidelity:** signal to noise ratio

- measure of the range of frequencies of sinusoidal signal that channel supports
- E.g., a channel that supports sinusoids from 1 MHz to 1.1 MHz has a bandwidth of 100 KHz
- “supports” in this context means “comes out the other end of the channel”
- some frequencies supported better than others
- analysing what happens to an arbitrary waveform is done by examining what happens to its component sinusoids → Fourier analysis
- bandwidth is a resource

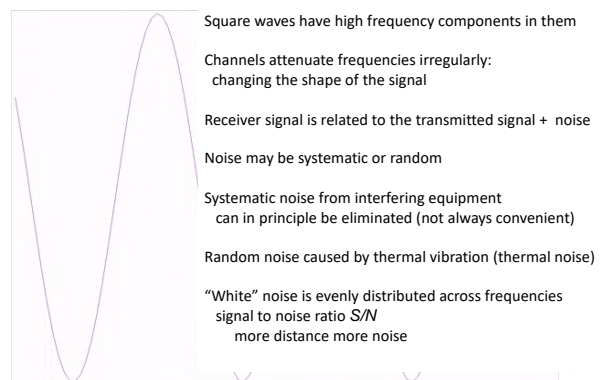
4

## Analog meet Digital



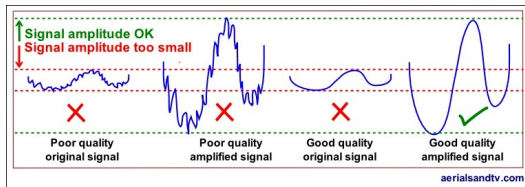
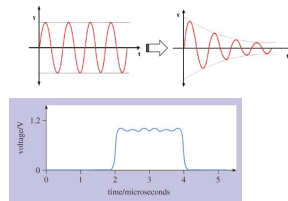
5

## Analog meet Digital



## Noise: Enemy of Communications

Attenuation, External Noise,  
Systematic, non-systematic,  
digitization, interference, reflection, ....



7

## Bandwidth vs Signal to Noise

what's better: high bandwidth or low signal to noise?

- channels subject to white noise have information capacity  $C$  measured in bits per second, of a channel

$$C = B \log_2(1 + S/N)$$

$B$  is the bandwidth of the channel  $S/N$  is the ratio of received signal power to received noise power.

- channels with no noise have information capacity determined only by bandwidth
- channels with any signal have nonzero information capacity
- channels with signal to noise ratio of unity have an information capacity in bits per second equal to its bandwidth in hertz
- (This is actually NOT the definition of information capacity; it is derived from the definition)

8

## (Digital) Channels

- Physical layer provides a channel
- Fixed rate for now
- Symbols are discrete values sent on the channel at fixed rate
- Symbols need not be binary
- Fidelity of the channel usually measured as a bit error rate — the probability that a bit sent as a 1 was interpreted as a 0 by the receiver or vice versa.
- Baud rate is the rate at which symbols can be transmitted
- Data rate (or bit rate) is the equivalent number of binary digits which can be sent
- E.g., if symbols represent with rate  $R$  then the data rate is  $2 \times R$ .

9

## Modulation

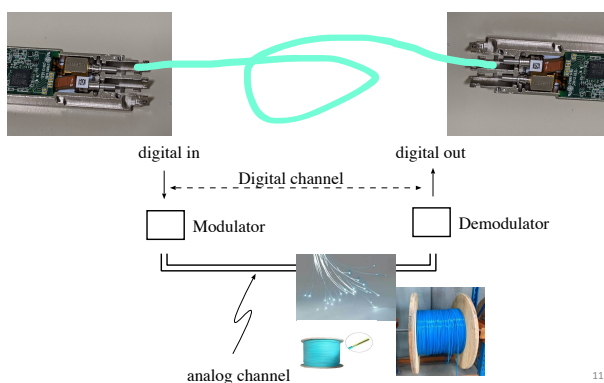
Two definitions:

- Transform an information signal into a signal more appropriate for transmission on a physical medium
- The systematic alteration of a carrier waveform by an information signal

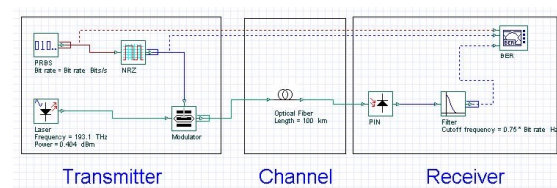
In general, we mean the first here (which encompasses the second).

10

## Communications



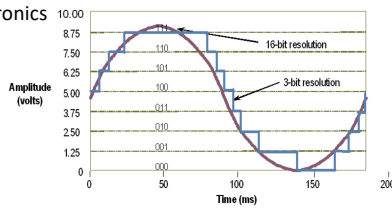
11



12

## Analog/Digital Digital/Analog

Recall from Digital Electronics



Conversion errors can occur in both directions

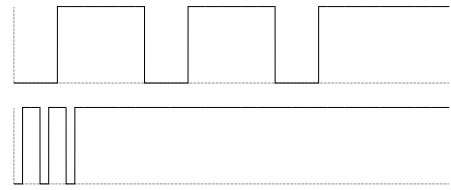
e.g.

Noise leads to incorrect digitization

Insufficient digitization resolution leads to information loss

13

## More Challenges



Where are the bits?

**WHEN** are the bits?

Bit boundaries can be asynchronous or synchronous

14

## Asynchronous versus Synchronous

- Transmission is sporadic, divided into frames
- Receiver and transmitter have oscillators which are close in frequency producing tx clocks and rx clock
- Receiver synchronises the phase of the rx clock with the tx clock by looking at one or more bit transitions
- RX clock drifts with respect to the tx clock but stays within a fraction of a bit of tx clock throughout the duration of a frame
- Transmission time is limited by accuracy of oscillators
- Transmission is continuous
- Receiver continually adjusts its frequency to track clock from incoming signal
- Requires bit transitions to inform clock
- Phase locked loop: rx clock predicts when incoming clock will change and corrects slightly when wrong.

15

## Asynchronous versus Synchronous

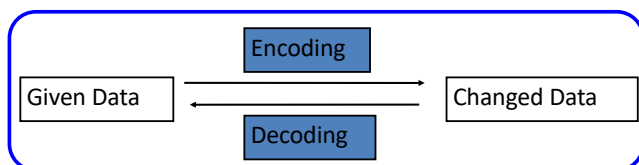
- Transmission is sporadic, divided into frames
- Receiver and transmitter have oscillators which are close in frequency producing tx clocks and rx clock
- Receiver synchronises the phase of the rx clock with the tx clock by looking at one or more bit transitions
- RX clock drifts with respect to the tx clock but stays within a fraction of a bit of tx clock throughout the duration of a frame
- Transmission time is limited by accuracy of oscillators
- Transmission is continuous
- Receiver continually adjusts its frequency to track clock from incoming signal
- Requires bit transitions to inform clock
- Phase locked loop: rx clock predicts when incoming clock will change and corrects slightly when wrong.

Bit transitions are critical

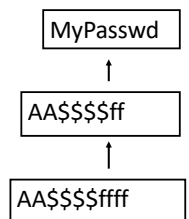
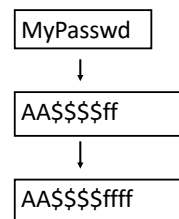
16

## Coding – a channel function

Change the representation of data.



17

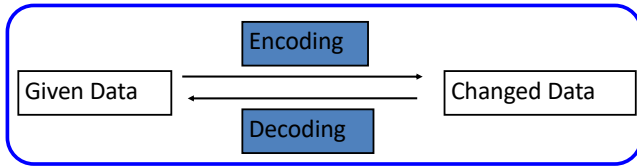


18



# Coding

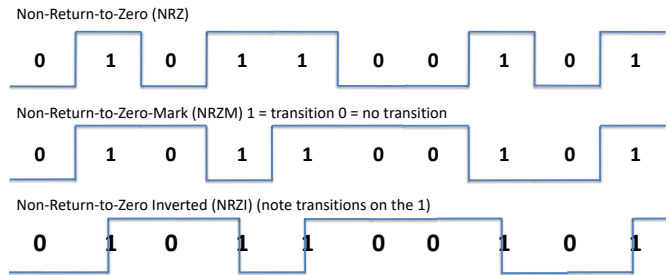
Change the representation of data.



1. **Encryption:** MyPasswd <-> AA\$\$\$f\$ff
2. **Error Detection:** AA\$\$\$f\$ff <-> AA\$\$\$f\$fff
3. **Compression:** AA\$\$\$f\$fff <-> A2\$4f4
4. **Analog:** A2\$4f4 <->

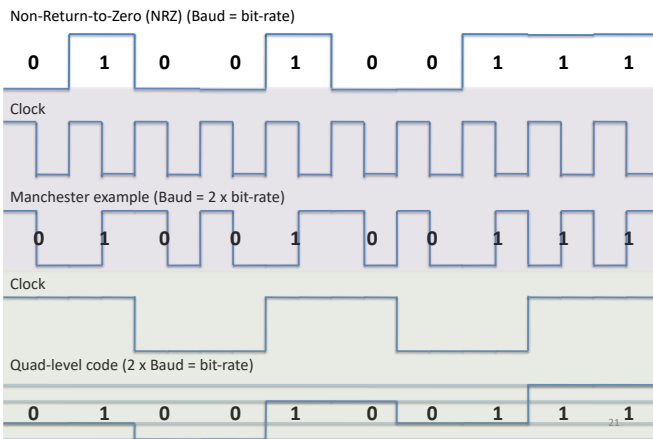
19

## Line Coding Examples where Baud=bit-rate

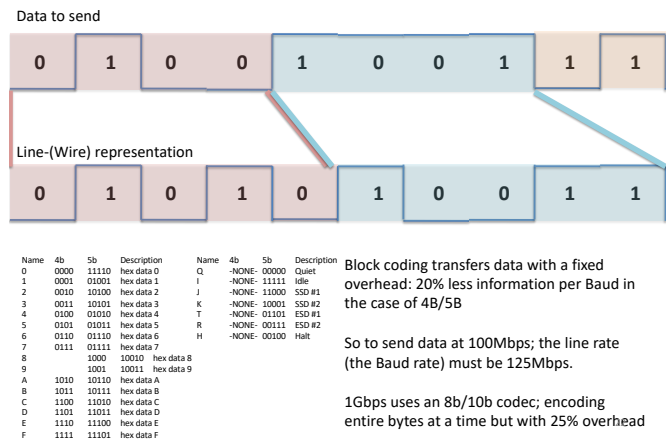


20

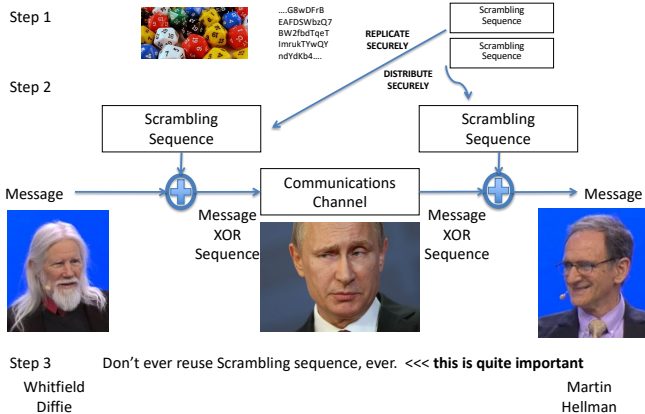
## Line Coding Examples



## Line Coding – Block Code example

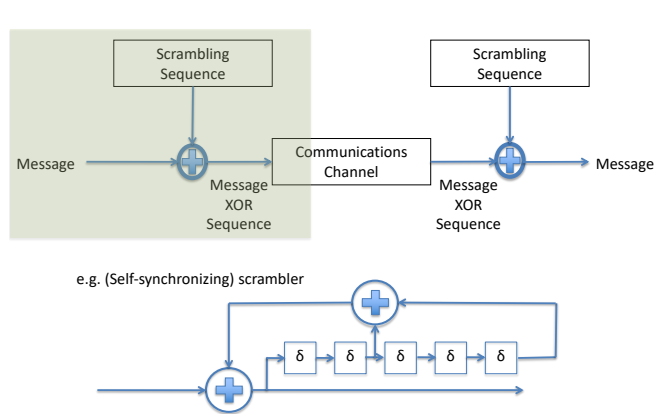


## Line Coding Scrambling – with secrecy



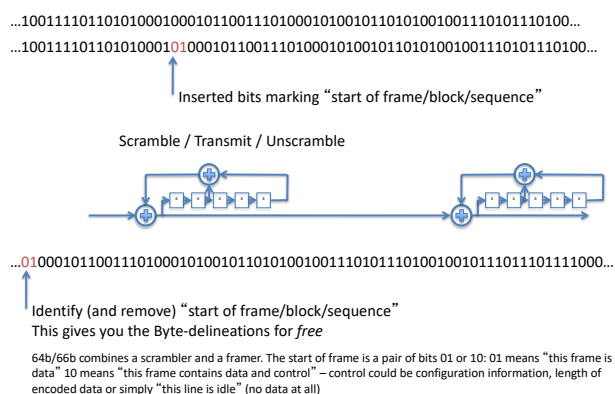
23

## Line Coding Scrambling– no secrecy

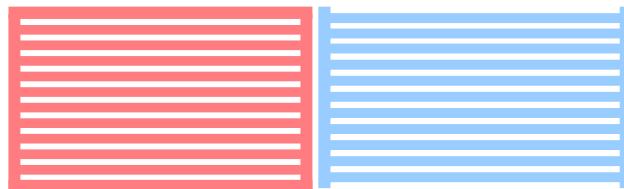
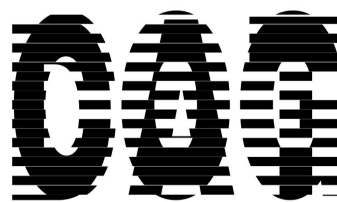


24

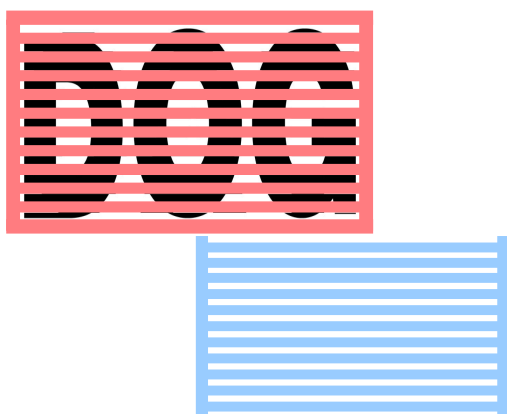
## Line Coding Examples (Hybrid)



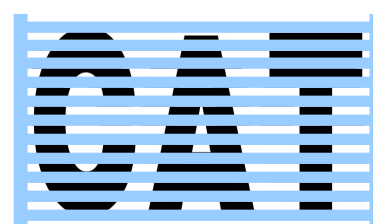
25



26



27



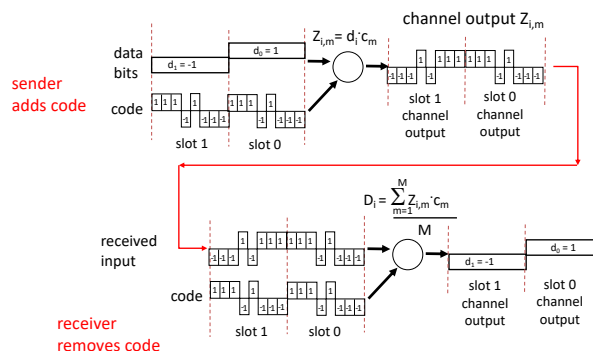
28

## Code Division Multiple Access (CDMA) (not to be confused with CSMA!)

- used in several wireless broadcast channels (cellular, satellite, etc) standards
- unique "code" assigned to each user; i.e., code set partitioning
- all users share same frequency, but each user has own "chipping" sequence (i.e., code) to encode data
- encoded signal** = (original data) XOR (chipping sequence)
- decoding**: inner-product of encoded signal and chipping sequence
- allows multiple users to "coexist" and transmit simultaneously with minimal interference (if codes are "orthogonal")

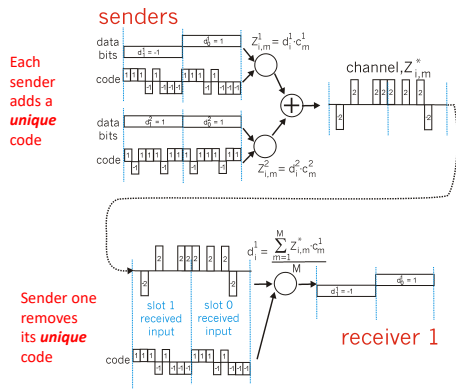
29

## CDMA Encode/Decode



30

## CDMA: two-sender interference



31

## Coding Examples summary

- Common Wired coding
  - Block codecs: table-lookups
    - fixed overhead, inline control signals
  - Scramblers: shift registers
    - overhead free

Like earlier coding schemes and error correction/detection; you can combine these

- e.g, 10Gb/s Ethernet may use a hybrid

### CDMA (Code Division Multiple Access)

- coping intelligently with competing sources
- Mobile phones

33

## Error Detection and Correction

Transmission media are not perfect and cause signal impairments:

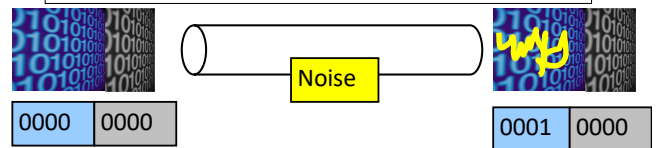
- Attenuation
  - Loss of energy to overcome medium's resistance
- Distortion
  - The signal changes its form or shape, caused in composite signals
- Noise
  - Thermal noise, induced noise, crosstalk, impulse noise

Interference can change the shape or timing of a signal:

$0 \rightarrow 1$  or  $1 \rightarrow 0$

## Error Detection and Correction

How to use coding to deal with errors in data communication?



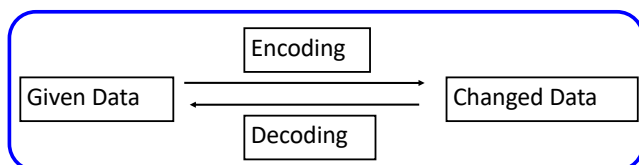
Basic Idea :

- Add additional information (redundancy) to a message.
  - Detect an error and discard
- Or, fix an error in the received message.

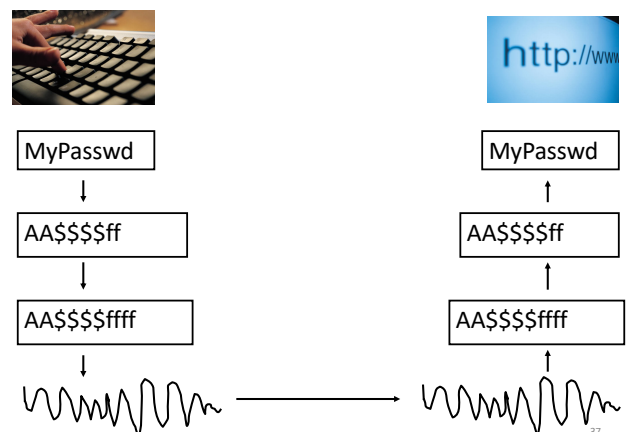
35

## Coding – a channel function

Change the representation of data.



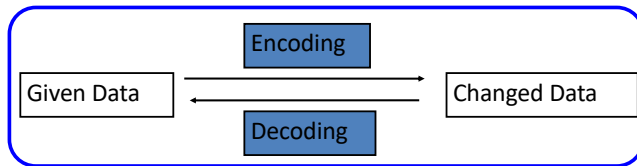
36



37

## Coding Examples

Change the representation of data.

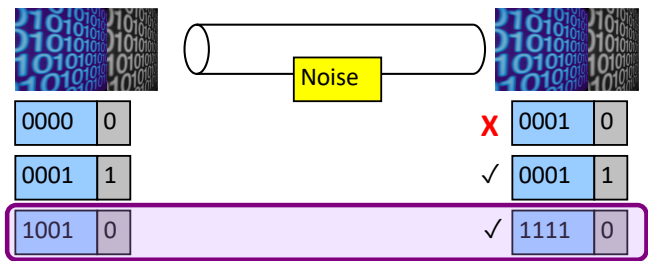


1. **Encryption:** MyPasswd  $\leftrightarrow$  AA\$\$\$f
2. **Error Detection:** AA\$\$\$f  $\leftrightarrow$  AA\$\$\$fff
3. **Compression:** AA\$\$\$fff  $\leftrightarrow$  A2\$4f4
4. **Analog:** A2\$4f4  $\leftrightarrow$

38

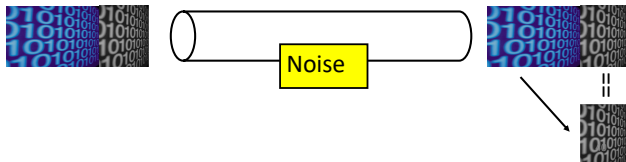
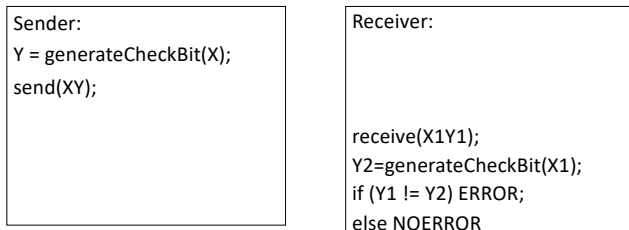
## Error Detection Code: Parity

Add one bit, such that the number of all 1's is even.



Problem: This simple parity cannot detect two-bit errors.

## Error Detection Code

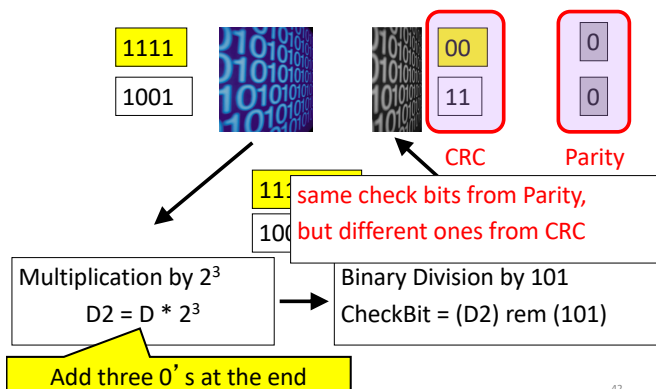


## Error Detection Code: CRC

- CRC means "Cyclic Redundancy Check".
- "A sequence of redundant bits, called CRC, is appended to the end of data so that the resulting data becomes exactly divisible by a second, predetermined binary number."
- $CRC = remainder (data \div predetermined\ divisor)$
- More powerful than parity.
  - It can detect various kinds of errors, including 2-bit errors.
- More complex: multiplication, binary division.
- Parameterized by n-bit divisor P.
  - Example: 3-bit divisor 101.
  - Choosing good P is crucial.

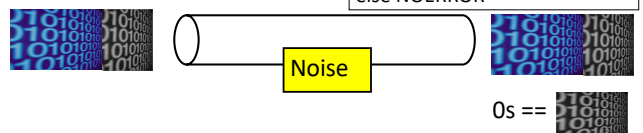
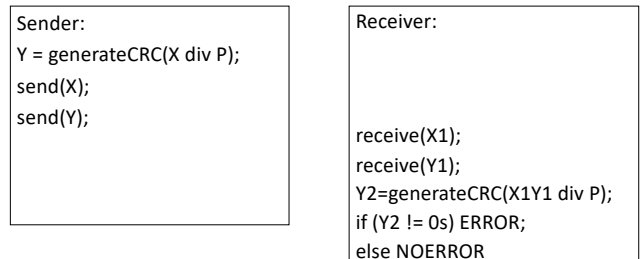
41

## CRC with 3-bit Divisor 101



42

## Error Detection Code



0s ==

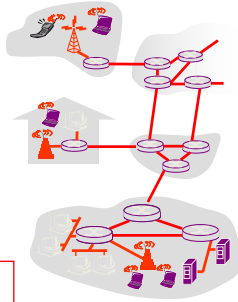


## Link Layer: Introduction

### Some reminder-terminology:

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

**data-link layer** has responsibility of transferring datagram from one node to adjacent node over a link



50

## Link Layer (Channel) Services - 1/2

- **framing, physical addressing:**
  - encapsulate datagram into frame, adding header, trailer
  - channel access if shared medium
  - “MAC” addresses used in frame headers to identify source, destination
    - This is **not** an IP address!
- **reliable delivery between adjacent nodes**
  - we revisit this again in the Transport Topic
  - seldom used on low bit-error link (fiber, some twisted pair)
  - wireless links: high error rates

51

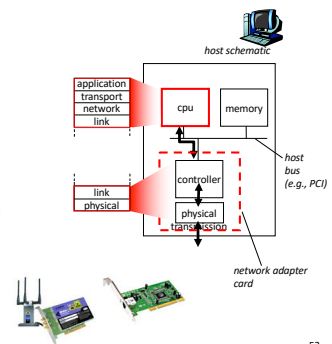
## Link Layer (Channel) Services – 2/2

- **flow control:**
  - pacing between adjacent sending and receiving nodes
- **error control:**
  - **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
- **access control: half-duplex and full-duplex**
  - with half duplex, nodes at both ends of link can transmit, but not at same time

52

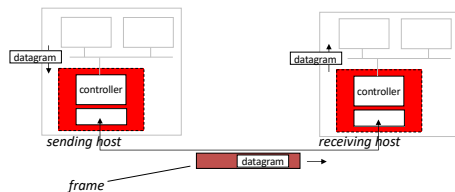
## Where is the link layer implemented?

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



53

## Adaptors Communicating



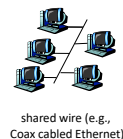
- **sending side:**
  - encapsulates datagram in frame
  - encodes data for the physical layer
  - adds error checking bits, provide reliability, flow control, etc.
- **receiving side**
  - decodes data from the physical layer
  - looks for errors, provide reliability, flow control, etc
  - extracts datagram, passes to upper layer at receiving side

54

## Multiple Access Links and Protocols

### Two types of “links”:

- **point-to-point**
  - point-to-point link between Ethernet switch and host
- **broadcast** (shared wire or medium)
  - old-fashioned wired Ethernet (*here be dinosaurs* – extinct)
  - upstream HFC (Hybrid Fiber-Coax – the Coax may be broadcast)
  - Home plug / Powerline networking
  - 802.11 wireless LAN



55

## 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

56

## Ideal Multiple Access Protocol

### Broadcast channel of rate $R$ bps

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

57

## MAC Protocols: a taxonomy

Three broad classes:

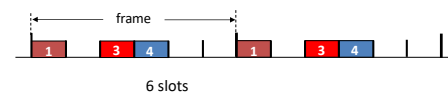
- **Channel Partitioning**
  - divide channel into smaller “pieces” (time slots, frequency, code)
  - allocate piece to node for exclusive use
- **Random Access**
  - channel not divided, allow collisions
  - “recover” from collisions
- **“Taking turns”**
  - nodes take turns, but nodes with more to send can take longer turns

58

## Channel Partitioning MAC protocols: TDMA (we discussed this earlier)

### TDMA: time division multiple access

- access to channel in “rounds”
- each station gets fixed length slot (length = pkt trans time) in each round
- unused slots go idle
- example: station LAN, 1,3,4 have pkt, slots 2,5,6 idle

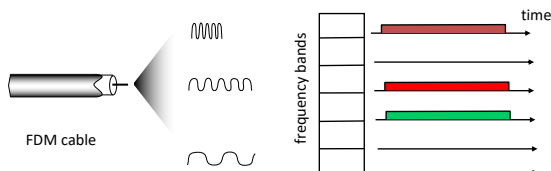


59

## Channel Partitioning MAC protocols: FDMA (we discussed this earlier)

### 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: station LAN, 1,3,4 have pkt, frequency bands 2,5,6 idle



60

## “Taking Turns” MAC protocols

### channel partitioning MAC protocols:

- share channel *efficiently* and *fairly* at high load
- inefficient at low load: delay in channel access,  $1/N$  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:

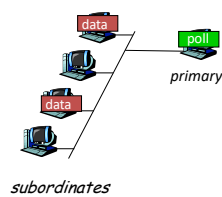
look for best of both worlds!

61

## “Taking Turns” MAC protocols

### Polling:

- Primary node “invites” subordinates nodes to transmit in turn
- typically used with simpler subordinate devices
- concerns:
  - polling overhead
  - latency
  - single point of failure (primary)

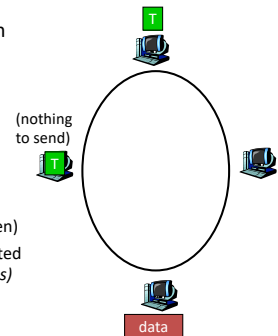


62

## “Taking Turns” MAC protocols

### Token passing:

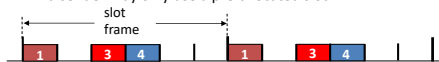
- control **token** passed from one node to next sequentially.
- token message
- concerns:
  - token overhead
  - latency
  - single point of failure (token)
- concerns fixed in part by a slotted ring (many simultaneous *tokens*)



63

## ATM

In TDM a sender may only use a pre-allocated slot



In ATM a sender transmits labeled cells whenever necessary



ATM = Asynchronous Transfer Mode – an ugly expression  
think of it as ATDM – Asynchronous Time Division Multiplexing

That's a variant of **PACKET SWITCHING** to the rest of us – just like Ethernet  
but using fixed length slots/packets/cells

Use the media when you need it, but  
ATM had virtual circuits and these needed setup....

64

## “Taking Turns” MAC protocols

### channel partitioning MAC protocols:

- share channel *efficiently* and *fairly* at high load
- inefficient at low load: delay in channel access,  $1/N$  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

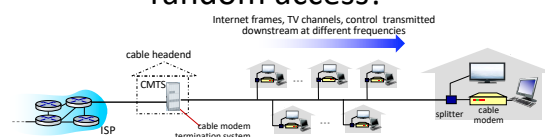
### “taking turns” protocols:

look for best of both worlds!

Recall.....

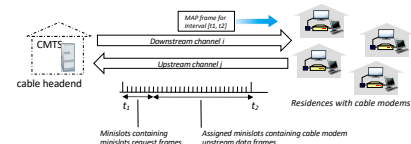
65

## Cable access network: FDM, TDM and random access!



- multiple** downstream (broadcast) FDM channels: up to 1.6 Gbps/channel
  - single CMTS transmits into channels
- multiple** upstream channels (up to 1 Gbps/channel)
  - multiple access**: all users contend (random access) for certain upstream channel time slots; others assigned TDM

## Cable access network:



### DOCSIS: data over cable service interface specification

- FDM over upstream, downstream frequency channels
- TDM upstream: some slots assigned, some have contention
  - downstream MAP frame: assigns upstream slots
  - request for upstream slots (and data) transmitted random access (binary backoff) in selected slots



## Random Access MAC Protocols

- When node has packet to send
  - Transmit at full channel data rate
  - No *a priori* coordination among nodes
- Two or more transmitting nodes  $\Rightarrow$  collision
  - Data lost
- Random access MAC protocol specifies:
  - How to detect collisions
  - How to recover from collisions
- Examples
  - ALOHA and Slotted ALOHA
  - CSMA, CSMA/CD, CSMA/CA (wireless)

68

## Key Ideas of Random Access

- **Carrier sense**
  - Listen before speaking, and don't interrupt
  - Checking if someone else is already sending data
  - ... and waiting till the other node is done
- **Collision detection**
  - If someone else starts talking at the same time, stop
  - Realizing when two nodes are transmitting at once
  - ...by detecting that the data on the wire is garbled
- **Randomness**
  - Don't start talking again right away
  - Waiting for a random time before trying again

69

## 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!
- Does this eliminate all collisions?
  - No, because of nonzero propagation delay

70

## CSMA Collisions

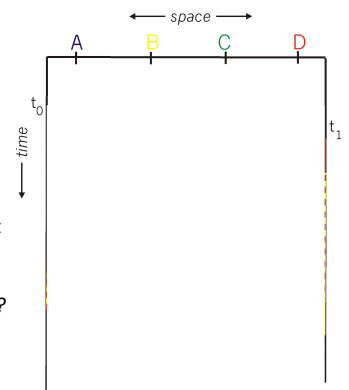
Propagation delay: two nodes may not hear each other's before sending.

Would slots hurt or help?

CSMA reduces but does not eliminate collisions

Biggest remaining problem?

Collisions still take full slot!  
How do you fix that?



## CSMA/CD (Collision Detection)

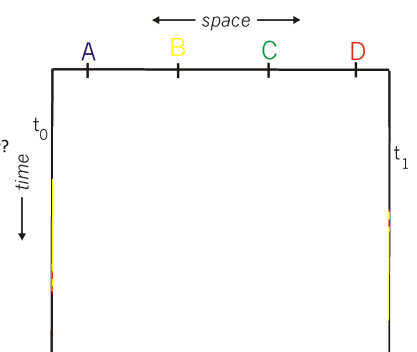
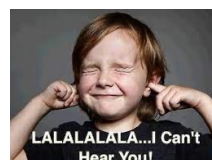
- CSMA/CD: carrier sensing, deferral as in CSMA
  - Collisions detected **within short time**
  - Colliding transmissions aborted, reducing wastage
- Collision detection easy in wired LANs:
  - Compare transmitted, received signals
- Collision detection difficult in wireless LANs:
  - Reception shut off while transmitting (well, perhaps not)
  - Not perfect broadcast (limited range) so collisions local
  - Leads to use of *collision avoidance* instead (later)

72

## CSMA/CD Collision Detection

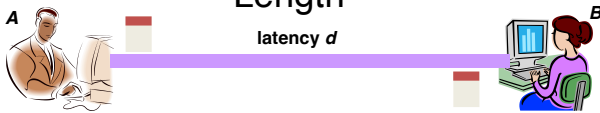
B and D can tell that collision occurred.

Note: for this to work, need restrictions on minimum frame size and maximum distance. Why?



73

## Limits on CSMA/CD Network Length



- Latency depends on physical length of link
  - Time to propagate a packet from one end to the other
- Suppose A sends a packet at time  $t$ 
  - And B sees an idle line at a time just before  $t+d$
  - ... so B happily starts transmitting a packet
- B detects a collision, and sends **jamming signal**
  - But A can't see collision until  $t+2d$

74

## Performance of CSMA/CD

- Time wasted in collisions
  - Proportional to distance  $d$
- Time spend transmitting a packet
  - Packet length  $p$  divided by bandwidth  $b$
- Rough estimate for efficiency ( $K$  some constant)

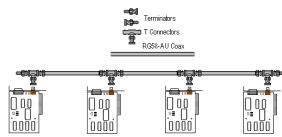
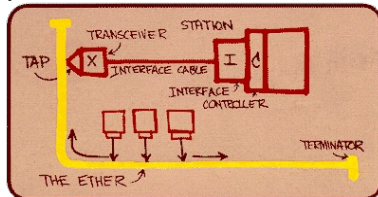
$$E \sim \frac{\frac{p}{b}}{\frac{p}{b} + Kd}$$

- Note:
  - For large packets, small distances,  $E \sim 1$
  - As bandwidth increases,  $E$  decreases
  - That is why high-speed LANs are all switched aka packets are sent via a switch - (any  $d$  is bad)

75

## Ethernet...

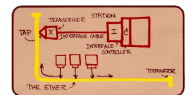
yet another product of XEROX/PARC



Preamble								Destination MAC						Source MAC						EtherType/Size		Payload				CRC							
1	2	3	4	5	6	7	8	1	2	3	4	5	6	1	2	3	4	5	6	1	2									1	2	3	4

77

## Ethernet: CSMA/CD Protocol



- **Carrier sense:** wait for link to be idle
- **Collision detection:** listen while transmitting
  - No collision: transmission is complete
  - Collision: abort transmission & send **jam** signal
- **Random access:** **binary exponential back-off**
  - After collision, wait a random time before trying again
  - After  $m^{\text{th}}$  collision, choose  $K$  randomly from  $\{0, \dots, 2^m-1\}$
  - ... and wait for  $K*512$  bit times before trying again
    - Using min packet size as "slot"
    - If transmission occurring when ready to send, wait until end of transmission (CSMA)

77

## STARVATION WARNING

- **Carrier sense:** wait for link to be idle
- **Collision detection:** listen while transmitting
  - No collision: transmission is complete
  - Collision: abort transmission & send jam signal
- **Random access:** **binary exponential back-off**
  - After collision, wait a random time before trying again
  - After  $m^{\text{th}}$  collision, choose  $K$  randomly from  $\{0, \dots, 2^m-1\}$
  - ... and wait for  $K*512$  bit times before trying again
    - Using min packet size as "slot"
    - If transmission occurring when ready to send, wait until end of transmission (CSMA)

78

## Benefits of Ethernet

- Easy to administer and maintain
- Inexpensive
- Increasingly higher speed
- Evolvable!

79

# Evolution of Ethernet

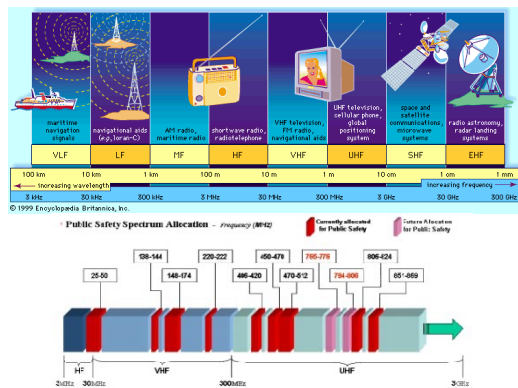
- Changed **everything** except the frame **format**
  - From single coaxial cable to hub-based star
  - From shared media to **switches**
  - From electrical signaling to optical
- **Lesson #1**
  - The right **interface** can accommodate many **changes**
  - Implementation is hidden behind interface
- **Lesson #2**
  - Really hard to displace the dominant technology
  - Slight performance improvements are not enough

80



81

## The Wireless Spectrum



82

## Metrics for evaluation / comparison of wireless technologies

- Bitrate or Bandwidth
- Range - PAN, LAN, MAN, WAN
- Two-way / One-way
- Multi-Access / Point-to-Point
- Digital / Analog
- Applications and industries
- Frequency – Affects most physical properties:
  - Distance (free-space loss)
  - Penetration, Reflection, Absorption
  - Energy proportionality
  - Policy: Licensed / Deregulated
  - Line of Sight (Fresnel zone)
  - Size of antenna
- Determined by wavelength –  $\lambda = \frac{v}{f}$

83

## Wireless Communication Standards

- Cellular (800/900/1700/1800/1900Mhz):
  - 2G: GSM / CDMA / GPRS / EDGE
  - 3G: CDMA2000/UMTS/HSDPA/EVDO
  - 4G: LTE, WiMax
- IEEE 802.11 (aka WiFi): (some examples)
  - b: 2.4Ghz band, 11Mbps (~4.5 Mbps operating rate)
  - g: 2.4Ghz, 54-108Mbps (~19 Mbps operating rate)
  - a: 5.0Ghz band, 54-108Mbps (~25 Mbps operating rate)
  - n: 2.4/5Ghz, 150-600Mbps (4x4 mimo)
  - ac: 2.4/5Ghz, 433-1300Mbps (improved coding 256-QAM)
  - ad: 60Ghz, 7Gbps
  - af: 54/790Mhz, 26-35Mbps (TV whitespace)
- IEEE 802.15 – lower power wireless:
  - 802.15.1: 2.4Ghz, 2.1 Mbps (Bluetooth)
  - 802.15.4: 2.4Ghz, 250 Kbps (Sensor Networks)

84

## What Makes Wireless Different?

- Broadcast and multi-access medium...
  - err, so....
- BUT, Signals sent by sender don't always end up at receiver intact
  - Complicated physics involved, which we won't discuss
  - But what can go wrong?

85

## Lets focus on 802.11

aka - WiFi ...

What makes it special?

Deregulation > Innovation > Adoption > Lower cost = Ubiquitous technology

JUST LIKE ETHERNET – not lovely but sufficient

86

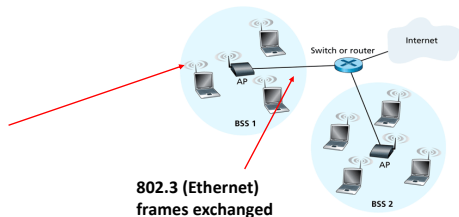
## IEEE 802.11 Wireless LAN

IEEE 802.11 standard	Year	Max data rate	Range	Frequency
802.11b	1999	11 Mbps	30 m	2.4 Ghz
802.11g	2003	54 Mbps	30m	2.4 Ghz
802.11n (WiFi 4)	2009	600	70m	2.4, 5 Ghz
802.11ac (WiFi 5)	2013	3.47Gbps	70m	5 Ghz
802.11ax (WiFi 6)	2020 (exp.)	14 Gbps	70m	2.4, 5 Ghz
802.11af	2014	35 – 560 Mbps	1 Km	unused TV bands (54-790 MHz)
802.11ah	2017	347Mbps	1 Km	900 Mhz

- all use CSMA/CA for multiple access, and have base-station and ad-hoc network versions

## 802.11 Architecture

802.11 frames exchanges



802.3 (Ethernet) frames exchanged

Figure 6.7 • IEEE 802.11 LAN architecture

- Designed for limited area
- AP's (Access Points) set to specific channel
- Broadcast beacon messages with SSID (Service Set Identifier) and MAC Address periodically
- Hosts scan all the channels to discover the AP's
  - Host associates with AP

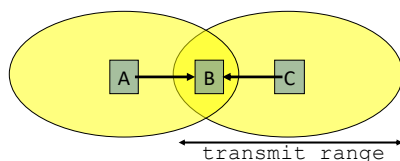
88

## Wireless Multiple Access Technique?

- Carrier Sense?
  - Sender can listen before sending
  - What does that tell the sender?
- Collision Detection?
  - Where do collisions occur?
  - How can you detect them?

89

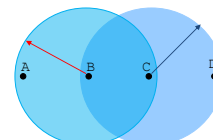
## Hidden Terminals



- A and C can both send to B but **can't hear each other**
  - A is a **hidden terminal** for C and vice versa
- Carrier Sense will be **ineffective**

90

## Exposed Terminals



- Exposed node:** B sends a packet to A; C hears this and decides not to send a packet to D (despite the fact that this will not cause interference)!
- Carrier sense would prevent a successful transmission.

91

## Key Points

- No concept of a global collision
  - Different receivers hear different signals
  - Different senders reach different receivers
- Collisions are at receiver, not sender
  - Only care if receiver can hear the sender clearly
  - It does not matter if sender can hear someone else
  - As long as that signal does not interfere with receiver
- Goal of protocol:
  - Detect if receiver can hear sender
  - Tell senders who might interfere with receiver to shut up

92

## Basic Collision Avoidance

- Since can't detect collisions, we try to **avoid** them
- Carrier sense:
  - When medium busy, choose random interval
  - Wait that many **idle** timeslots to pass before sending
- When a collision is inferred, retransmit with binary exponential backoff (like Ethernet)
  - Use **ACK** from receiver to infer "no collision"
  - Use exponential backoff to adapt contention window

93

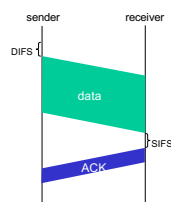
## IEEE 802.11 MAC Protocol: CSMA/CA

### 802.11 sender

- 1 if sense channel idle for **DIFS** then  
transmit entire frame (no CD)
- 2 if sense channel busy then  
start random backoff time  
timer counts down while channel idle  
transmit when timer expires  
if no ACK, increase random backoff interval, repeat 2

### 802.11 receiver

- if frame received OK  
return ACK after **SIFS** (ACK needed due to hidden terminal problem)

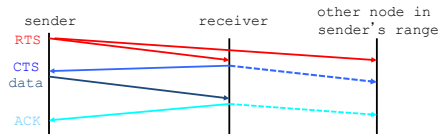


## Avoiding collisions

idea: sender "reserves" channel use for data frames using small reservation packets

- sender first transmits *small* request-to-send (RTS) packet to BS using CSMA
  - RTSs may still collide with each other (but they're short)
- BS broadcasts clear-to-send CTS in response to RTS
  - CTS heard by all nodes
  - sender transmits data frame
  - other stations defer transmissions

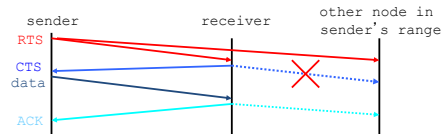
## CSMA/CA – and in this case RTS/CTS



- Before every data transmission
  - Sender sends a Request to Send (RTS) frame containing the length of the transmission
  - Receiver respond with a Clear to Send (CTS) frame
  - Sender sends data
  - Receiver sends an ACK; now another sender can send data
- When sender doesn't get a CTS back, it assumes collision

96

## CSMA/CA, con't

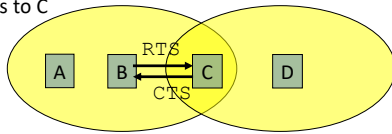


- If other nodes hear RTS, but not CTS: **send**
  - Presumably, destination for first sender is out of node's range ...
  - ... Can cause problems when a CTS is **lost**
- When you hear a CTS, you keep quiet until scheduled transmission is over (hear ACK)

97

## RTS / CTS Protocols (CSMA/CA)

B sends to C



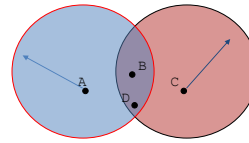
Overcome hidden terminal problems with contention-free protocol

1. B sends to C **Request To Send** (RTS)
2. A hears RTS and defers (to allow C to answer)
3. C replies to B with **Clear To Send** (CTS)
4. D hears CTS and defers to allow the data
5. B sends to C

98

## Preventing Collisions Altogether

- Frequency Spectrum partitioned into several channels
  - Nodes within interference range can use separate channels

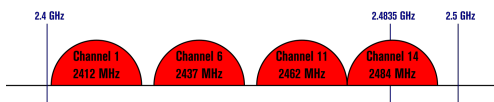


- Now A and C can send without any interference!
- Most cards have only 1 transceiver
  - **Not Full Duplex: Cannot send and receive at the same time**
- Aggregate Network throughput doubles

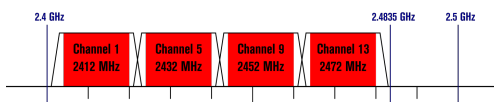
99

## Non-Overlapping Channels for 2.4 GHz WLAN

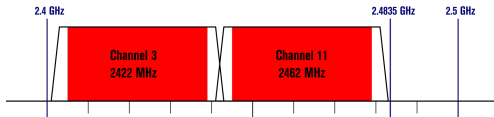
802.11b (DSSS) channel width 22 MHz



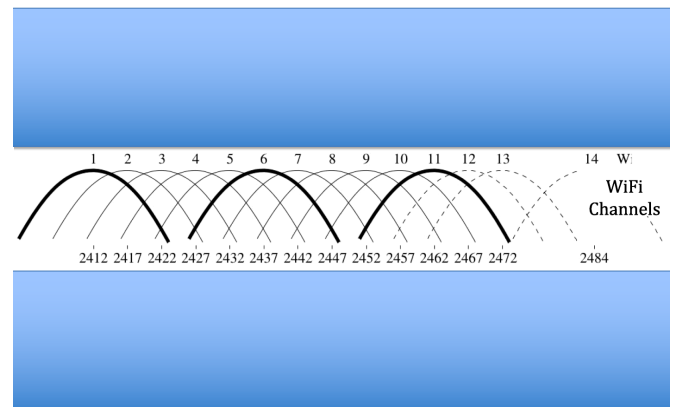
802.11g/n (OFDM) 20 MHz ch. width - 16.25 MHz used by sub-carriers



802.11n (OFDM) 40 MHz ch. width - 33.75 MHz used by sub-carriers

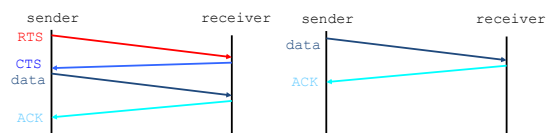


100



101

## CSMA/CA and RTS/CTS



RTS/CTS

- helps with hidden terminal
- good for high-traffic Access Points
- often turned on/off dynamically

Without RTS/CTS

- lower latency -> faster!
- reduces wasted b/w if the  $Pr(\text{collision})$  is low
- good for when net is small and not weird
  - eg no hidden/exposed terminals

103



Wifi has been evolving!

Using dual band (2.4GHz + 5GHz), multiple channels, MIMO, Meshing Wifi

Outside this introduction but the state of the art is very fast and very flexible

102

## CSMA/CD vs CSMA/CA (without RTS/CTS)

### CD Collision Detect

wired – listen and talk

1. Listen for others
2. Busy? goto 1.
3. Send message (and listen)
4. Collision?
  - a. JAM
  - b. increase your BEB
  - c. sleep
  - d. goto 1.

### CA Collision Avoidance

wireless – talk OR listen

1. Listen for others
2. Busy? goto 1.
3. Send message
4. Wait for ACK (MAC ACK)
5. Got No ACK from MAC?
  - a. increase your BEB
  - b. sleep
  - c. goto 1.

104

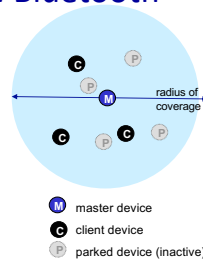
## 802.11: advanced capabilities

### power management

- node-to-AP: “I am going to sleep until next beacon frame”
  - AP knows not to transmit frames to this node
  - node wakes up before next beacon frame
- beacon frame: contains list of mobiles with AP-to-mobile frames waiting to be sent
  - node will stay awake if AP-to-mobile frames to be sent; otherwise sleep again until next beacon frame

## Personal area networks: Bluetooth

- TDM, 625 µsec sec. slot
- FDM: sender uses 79 frequency channels in known, pseudo-random order slot-to-slot (spread spectrum)
  - other devices/equipment not in piconet only interfere in some slots
- **parked mode**: clients can “go to sleep” (park) and later wakeup (to preserve battery)
- **bootstrapping**: nodes self-assemble (plug and play) into piconet



- **channel partitioning**, by time, frequency or code
  - Time Division (TDMA), Frequency Division (FDMA), Code Division (CDMA)
- **random access** (dynamic),
  - ALOHA, S-ALOHA, CSMA, CSMA/CD
  - carrier sensing: easy in some technologies (wire), hard in others (wireless)
  - CSMA/CD used in (old-style, coax) Ethernet, and PowerLine
  - CSMA/CA used in 802.11
- **taking turns**
  - polling from central site, token passing
  - Bluetooth, FDDI, IBM Token Ring

107

## MAC Addresses

- MAC (or LAN or physical or Ethernet) address:
  - function: *get frame from one interface to another physically-connected interface (same network)*
  - 48 bit MAC address (for most LANs)
    - *burned* in NIC ROM, nowadays usually software settable and set at boot time

```
dmz22@rio:~$ ifconfig eth0
eth0      Link encap:Ethernet  HWaddr 00:30:48:fe:c0:64
          inet addr:128.232.33.7  Bcast:128.232.47.255  Mask:255.255.240.0
          inet6 addr: fe80::230:48ff:fe:c064/64 Scope:Link
          UP BROADCAST RUNNING MULTICAST  MTU:1500  Metric:1
          RX packets:215084512 errors:252 dropped:25 overruns:0 frame:123
          TX packets:146711866 errors:0 dropped:0 overruns:0 carrier:0
          collisions:0 txqueuelen:1000
          RX bytes:170815941033 (170.8 GB)  TX bytes:86755864270 (86.7 GB)
          Memory: F0000000-F0020000
```

108

## LAN Address (more)

- MAC address allocation administered by IEEE
- manufacturer buys portion of MAC address space (to assure uniqueness)
- analogy:
  - (a) MAC address: like a National Insurance Number
  - (b) IP address: like a 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

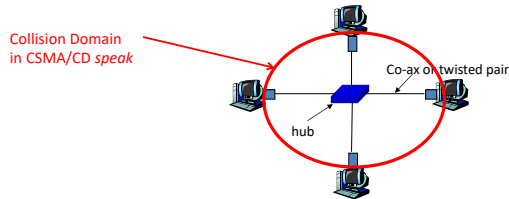
109

## Hubs



... physical-layer ("dumb") repeaters:

- bits coming in one link go out **all** other links at same rate
- all nodes connected to hub can collide with one another
- no frame buffering
- no CSMA/CD at hub: host NICs detect collisions



110

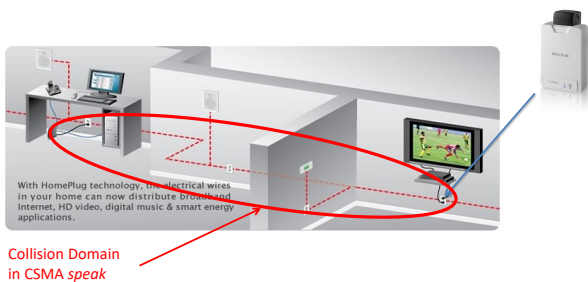
## CSMA in our home

### Home Plug Powerline Networking....



111

## Home Plug and similar Powerline Networking....



Collision Domain in CSMA speak

To secure network traffic on a specific HomePlug network, each set of adapters use an encryption key common to a specific HomePlug network

112

## Switch (example: Ethernet Switch)

- **link-layer device: smarter than hubs, take 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

If you want to connect different physical media (optical – copper – coax – wireless - ....)

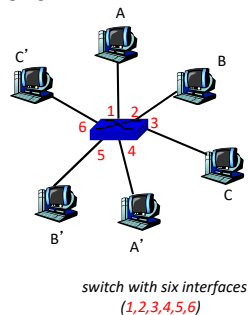
you **NEED** a switch.

Why? (Because each link, each media access protocol is specialised)

113

## Switch: allows *multiple* simultaneous transmissions

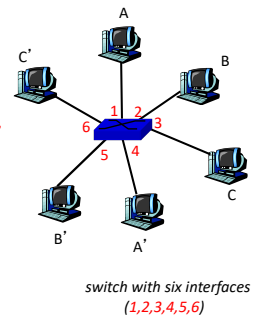
- hosts have dedicated, direct connection to switch
- switches buffer packets
- Ethernet protocol used on *each* incoming link, but no collisions; full duplex
  - each link is its own collision domain
- **switching:** A-to-A' and B-to-B' simultaneously, without collisions
  - not possible with dumb hub



114

## Switch Table

- **Q:** how does switch know that 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!
- **Q:** how are entries created, maintained in switch table?
  - something like a routing protocol?

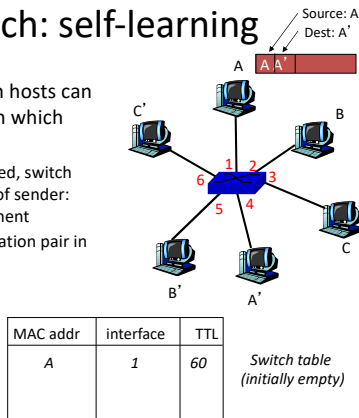


115



## 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



116

## Switch: frame filtering/forwarding

### When frame received:

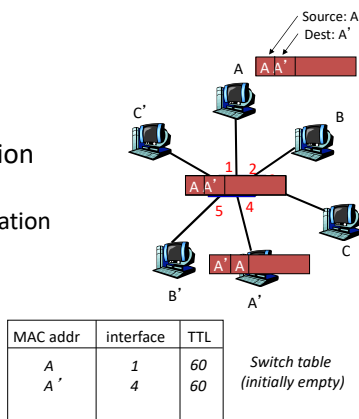
- record link associated with sending host
- index switch table using MAC dest address
- if entry found for destination
  - then {
    - if dest on segment from which frame arrived
      - then drop the frame
      - else forward the frame on interface indicated
- else flood

forward on all but the interface on which the frame arrived

117

## Self-learning, forwarding: example

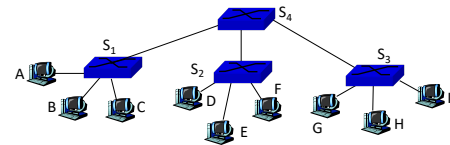
- frame destination unknown: **flood**
- destination A location known: **selective send**



118

## Interconnecting switches

- switches can be connected together

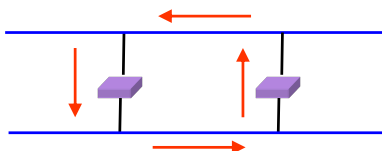


- Q: sending from A to G - how does S<sub>1</sub> know to forward frame destined to F via S<sub>4</sub> and S<sub>3</sub>?
- A: self learning! (works exactly the same as in single-switch case – flood/forward/drop)

119

## Flooding Can Lead to Loops

- Flooding can lead to **forwarding loops**
  - E.g., if the network contains a cycle of switches
  - “Broadcast storm”

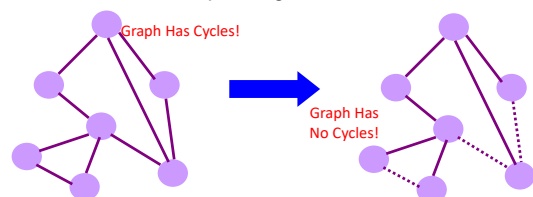


120



## Solution: Spanning Trees

- Ensure the forwarding **topology** has no loops
  - Avoid using some of the links when flooding
  - ... to prevent loop from forming
- Spanning tree**
  - Sub-graph that covers all vertices but *contains no cycles*
  - Links not in the spanning tree do not forward frames



121

## What Do We Know?

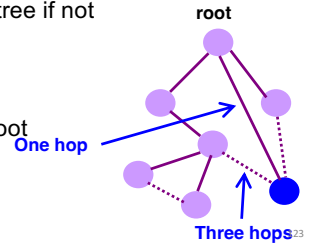
- “Spanning tree algorithm is an algorithm to create a tree out of a graph that includes all nodes with a minimum number of edges connecting to vertices.”
- Shortest paths to (or from) a node form a tree
- So, algorithm has two aspects :
  - Pick a root
  - Compute shortest paths to it
- Only keep the links on shortest-path

122

## Constructing a Spanning Tree

- Switches need to **elect** a **root**
  - The switch w/ smallest identifier (MAC addr)
- Each switch determines if each interface is on the **shortest path** from the root
  - Excludes it from the tree if not

- Messages (Y, d, X)
  - From node X
  - Proposing Y as the root
  - And the distance is d



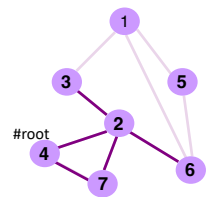
## Steps in Spanning Tree Algorithm

- Initially, each switch proposes itself as the root
  - Switch sends a message out every interface
  - ... proposing itself as the root with distance 0
  - Example: switch X announces (X, 0, X)
- Switches update their view of the root
  - Upon receiving message (Y, d, Z) from Z, check Y's id
  - If new id smaller, start viewing that switch as root
- Switches compute their distance from the root
  - Add 1 to the distance received from a neighbor
  - Identify interfaces not on shortest path to the root
  - ... and exclude them from the spanning tree
- If root or shortest distance to it **changed**, “flood” updated message (Y, d+1, X)

124

## Example From Switch #4's Viewpoint

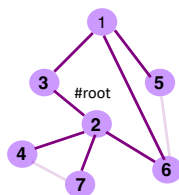
- Switch #4 thinks it is the root
  - Sends (4, 0, 4) message to 2 and 7
- Then, switch #4 hears from #2
  - Receives (2, 0, 2) message from 2
  - ... and thinks that #2 is the root
  - And realizes it is just one hop away
- Then, switch #4 hears from #7
  - Receives (2, 1, 7) from 7
  - And realizes this is a longer path
  - So, prefers its own one-hop path
  - And removes 4-7 link from the tree



125

## Example From Switch #4's Viewpoint

- Switch #2 hears about switch #1
  - Switch 2 hears (1, 1, 3) from 3
  - Switch 2 starts treating 1 as root
  - And sends (1, 2, 2) to neighbors
- Switch #4 hears from switch #2
  - Switch 4 starts treating 1 as root
  - And sends (1, 3, 4) to neighbors
- Switch #4 hears from switch #7
  - Switch 4 receives (1, 3, 7) from 7
  - And realizes this is a longer path
  - So, prefers its own three-hop path
  - And removes 4-7 link from the tree



126

## Robust Spanning Tree Algorithm

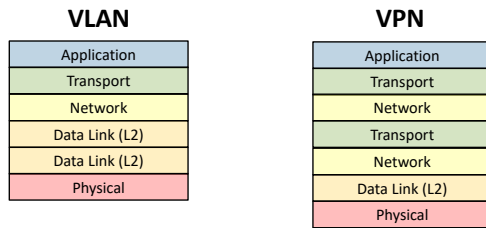
- Algorithm must react to **failures**
  - Failure of the root node
    - Need to elect a new root, with the next lowest identifier
  - Failure of other switches and links
    - Need to recompute the spanning tree
- Root switch continues sending messages
  - Periodically reannouncing itself as the root (1, 0, 1)
  - Other switches continue forwarding messages
- Detecting failures through timeout (**soft state**)
  - If no word from root, times out and claims to be the root
  - Delay in reestablishing spanning tree is **major problem**
  - Work on rapid spanning tree algorithms...

Given a switch-tree of a given size, link length, speed of computation, ...

How long does a failure take to rectify?

127

## Weirder “Data Link Layer” Networks



### Datacenter

“so you think your LAN has a lot of computers....”

128

## Datacenter 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)

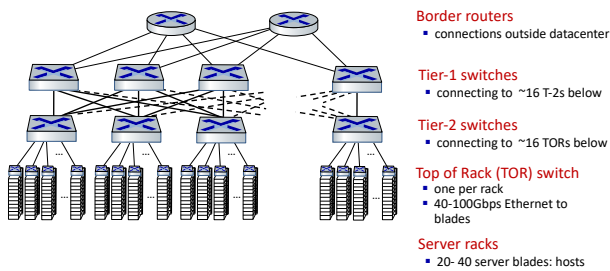
challenges:

- multiple applications, each serving massive numbers of clients
- reliability
- managing/balancing load, avoiding processing, networking, data bottlenecks



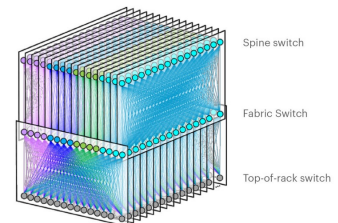
Inside a 40-ft Microsoft container, Chicago data center

## Datacenter networks: network elements



## Datacenter networks: network elements

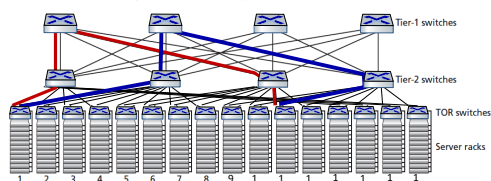
Facebook F16 data center network topology:



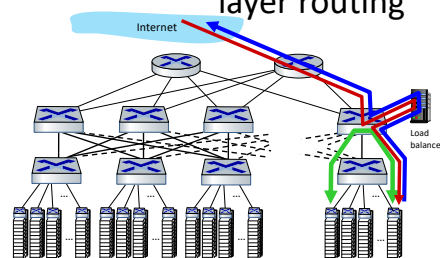
<https://engineering.fb.com/data-center-engineering/f16-mini-pack/> (posted 3/2019)

## Datacenter networks: multipath

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



## Datacenter networks: application-layer routing



**load balancer: application-layer routing**

- receives external client requests
- directs workload within data center
- returns results to external client (hiding data center internals from client)

## 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
  - WiFi
- algorithms
  - Binary Exponential Backoff
  - Spanning Tree

## Topic 4: Network Layer

### Our goals:

- understand principles behind network layer services:
  - network layer service models
  - forwarding versus routing (versus switching)
  - how a router works
  - routing (path selection)
  - IPv6 (one day IPv6 will be first but not this year)

For the most part, the Internet is our example – again.

1

Recall: Network layer is responsible for **GLOBAL** delivery

Name: a *something*

Address: Where is a *something*

Routing: How do I get to the *something*

Forwarding: What path do I take next to get to the *something*

2

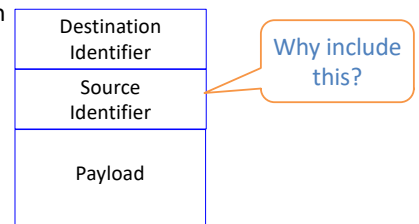
## Addressing (at a conceptual level)

- Assume all hosts have unique IDs
- No particular structure to those IDs
- Later in topic I will talk about real IP addressing
- Do I route on location or identifier?
- If a host moves, should its address change?
  - If not, how can you build scalable Internet?
  - If so, then what good is an address for identification?

3

## Packets (at a conceptual level)

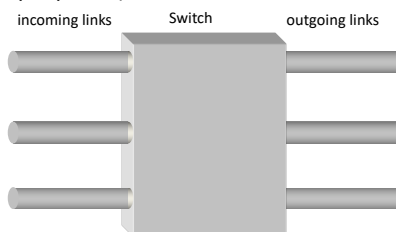
- Assume packet headers contain:
  - Source ID, Destination ID, and perhaps other information



4

## Switches/Routers

- Multiple ports (attached to other switches or hosts)



- Ports are typically duplex (incoming and outgoing)

5

## A Variety of *(Internet Protocol-based)* Networks

- ISPs: carriers
  - Backbone
  - Edge
  - Border (to other ISPs)
- Enterprises: companies, universities
  - Core
  - Edge
  - Border (to outside)
- Datacenters: massive collections of machines
  - Top-of-Rack
  - Aggregation and Core
  - Border (to outside)

6

## A Variety of (Internet Protocol-based) Routers

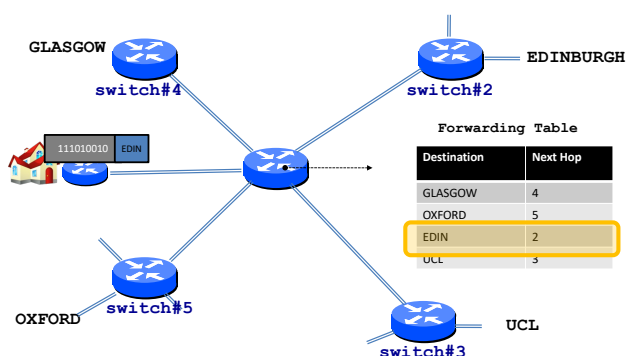
- ISPs: carriers
  - Backbone
  - Edge
  - Border (to other ISPs)
- Enterprises: companies, universities
  - Core
  - Edge
  - Border (to outside)
- Datacenters: massive collections of machines
  - Top-of-Rack
  - Aggregation and Core
  - Border (to outside)



7

Recall....

## Switches forward packets



8

## Forwarding Decisions

- When packet arrives..
  - Must decide which outgoing port to use
  - In single transmission time
  - Forwarding decisions must be *simple*
- Routing state dictates where to forward packets
  - Assume decisions are *deterministic*
- *Global routing state* is the collection of routing state in each of the routers
  - Will focus on where this routing state comes from
  - But first, a few preliminaries....

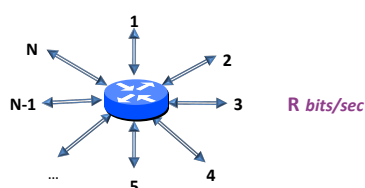
9

## Forwarding vs Routing

- Forwarding: “data plane”
  - Directing a data packet to an outgoing link
  - Individual router using routing state
- Routing: “control plane”
  - Computing paths the packets will follow
  - Routers talking amongst themselves
  - Jointly creating the routing state
- Two very different timescales....

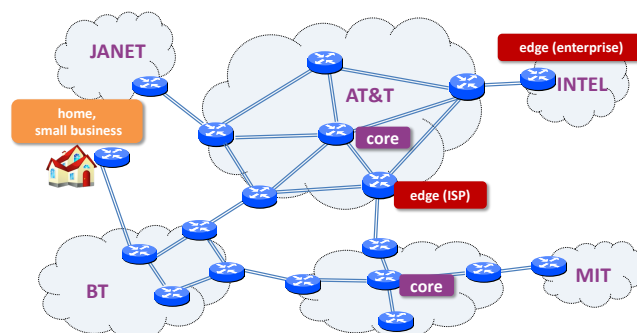
10

## Router definitions

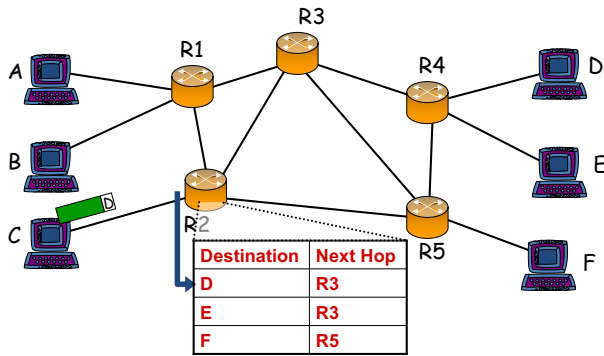


- $N$  = number of external router “ports”
- $R$  = speed (“line rate”) of a port
- Router capacity =  $N \times R$

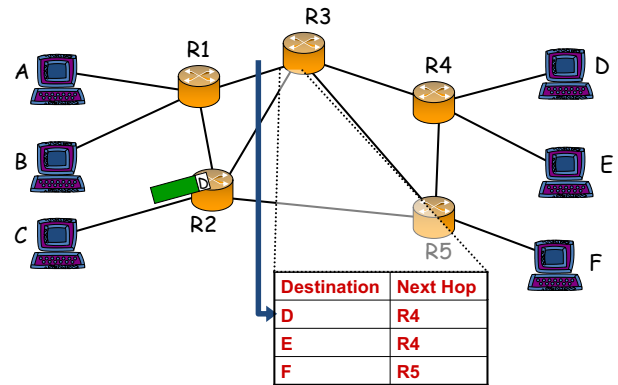
## Networks and routers



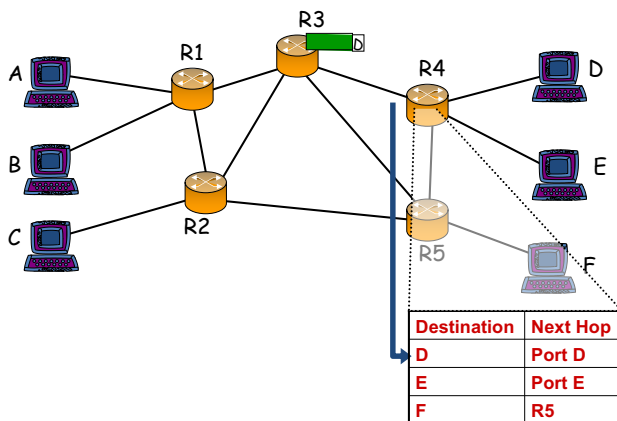
## Basic Operation of Router



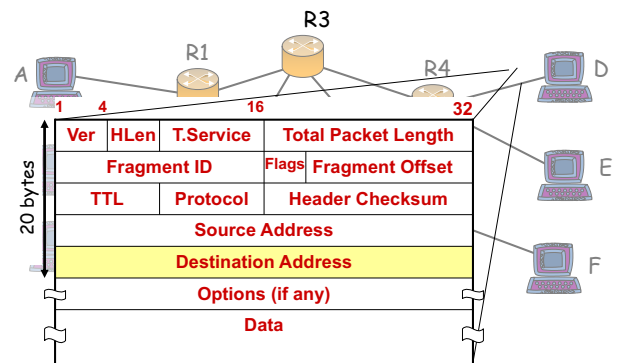
## Basic Operation of Router



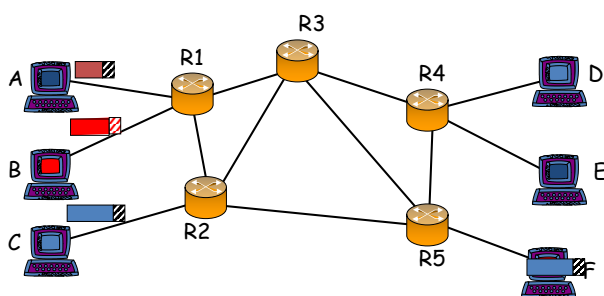
## Basic Operation of Router



## What does a router do?

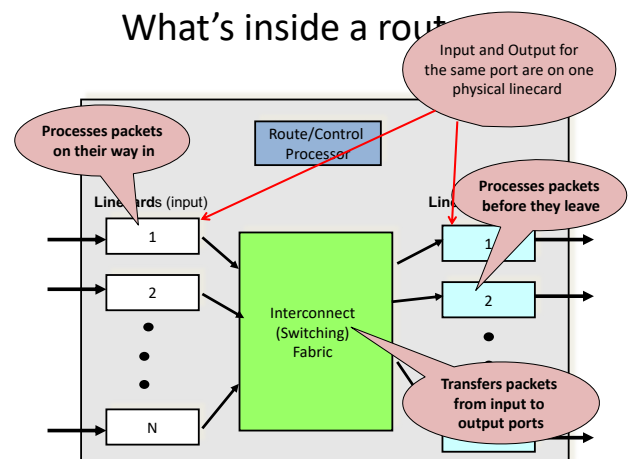


## What does a router do?

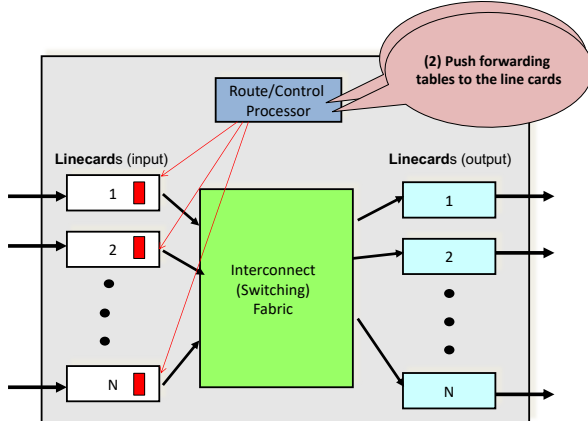


1. Every router performs a per-packet lookup for every packet
2. Each router performs a lookup in its local lookup table
3. Each router performs lookups (ENTIRELY) independently of every other router

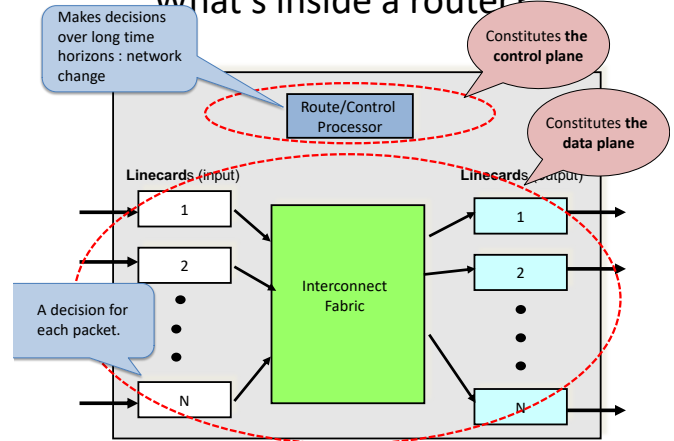
## What's inside a router



## What's inside a router?

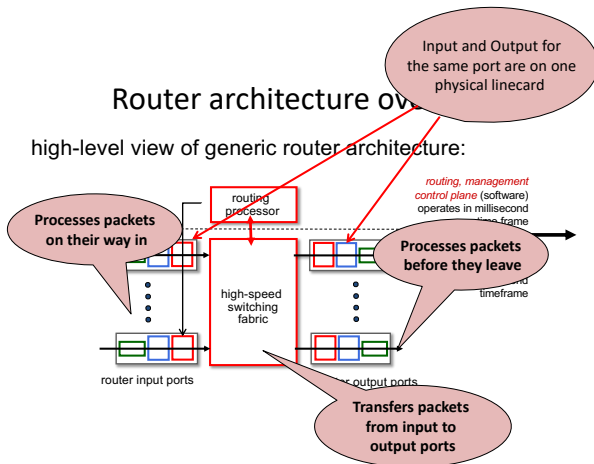


## What's inside a router?

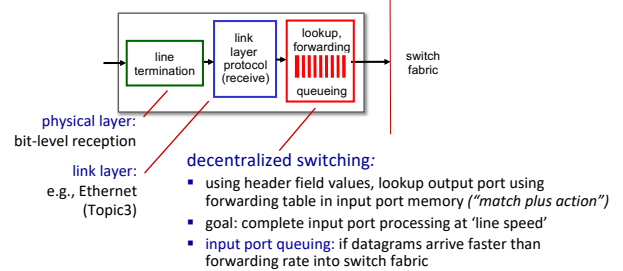


## Router architecture overview

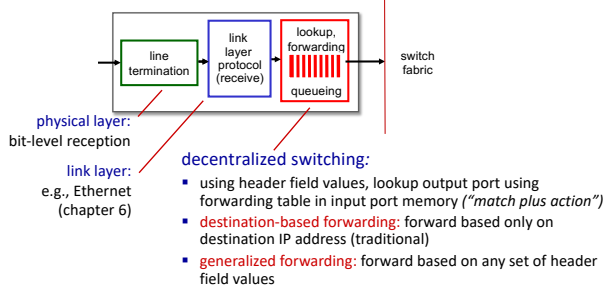
high-level view of generic router architecture:



## Input port functions

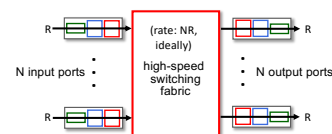


## Input port functions



## Switching fabrics

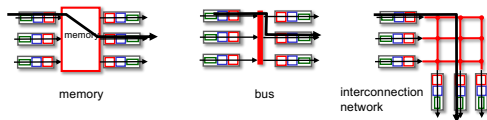
- transfer packet from input link to appropriate output link
- switching rate:** rate at which packets can be transfer from inputs to outputs
  - often measured as multiple of input/output line rate
  - N inputs: switching rate N times line rate desirable





## Switching fabrics

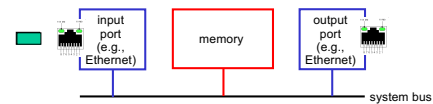
- transfer packet from input link to appropriate output link
- switching rate**: rate at which packets can be transfer from inputs to outputs
  - often measured as multiple of input/output line rate
  - N inputs: switching rate N times line rate desirable
- three major types of switching fabrics:



## Switching via memory

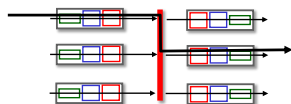
### first generation routers:

- traditional computers with switching under direct control of CPU
- packet copied to system's memory
- speed limited by memory bandwidth (2 bus crossings per datagram)



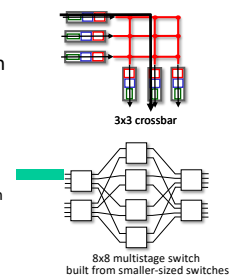
## Switching via a bus

- 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 routers



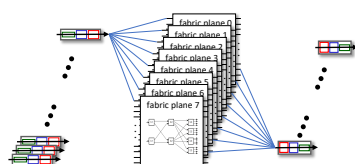
## Switching via interconnection network

- Crossbar, Clos networks, other interconnection nets initially developed to connect processors in multiprocessor
- multistage switch**:  $n \times n$  switch from multiple stages of smaller switches
- exploiting parallelism**:
  - fragment datagram into fixed length cells on entry
  - switch cells through the fabric, reassemble datagram at exit



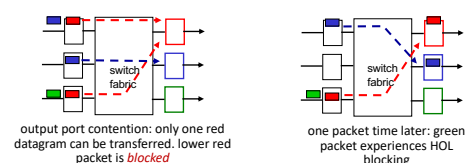
## Switching via interconnection network

- scaling, using multiple switching "planes" in parallel:
  - speedup, scaleup via parallelism
- Cisco CRS router:
  - basic unit: 8 switching planes
  - each plane: 3-stage interconnection network
  - up to 100's Tbps switching capacity

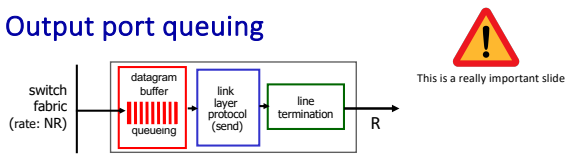


## Input port queuing

- If switch fabric slower than input ports combined -> queueing 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

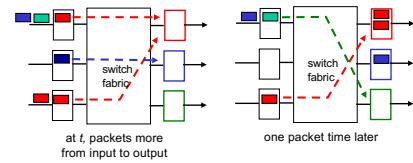


## Output port queuing



- **Buffering** required when datagrams arrive from fabric faster than link transmission rate. **Drop policy**: which datagrams to drop if no free buffers?
  - **Scheduling discipline** chooses among queued datagrams for transmission
- Datagrams can be lost due to congestion, lack of buffers
- Priority scheduling – who gets best performance, network neutrality

## Output port queuing

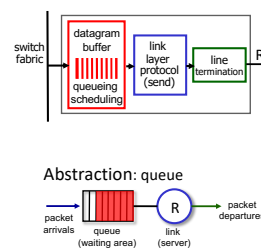


- buffering when arrival rate via switch exceeds output line speed
- **queueing (delay) and loss due to output port buffer overflow!**

## How much buffering? (related material in Topic 5)

- 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
- more recent recommendation: with  $N$  flows, buffering equal to  $\frac{RTT \cdot C}{\sqrt{N}}$
- but *too much* buffering can increase delays (particularly in home routers)
  - long RTTs: poor performance for realtime apps, sluggish TCP response
  - recall delay-based congestion control: “keep bottleneck link just full enough (busy) but no fuller”

## Buffer Management



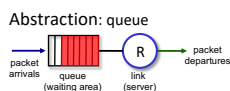
### buffer management:

- **drop**: which packet to add, drop when buffers are full
  - **tail drop**: drop arriving packet
  - **priority**: drop/remove on priority basis
- **marking**: which packets to mark to signal congestion (ECN, RED)

## Packet Scheduling: FCFS

**packet scheduling**: deciding which packet to send next on link

- first come, first served
- priority
- round robin
- weighted fair queueing



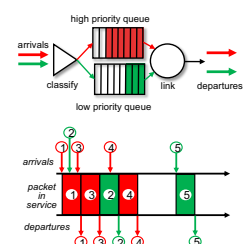
**FCFS**: packets transmitted in order of arrival to output port

- also known as: First-in-first-out (FIFO)
- real world examples?

## Scheduling policies: priority

### Priority scheduling:

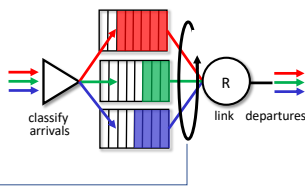
- arriving traffic classified, queued by class
  - any header fields can be used for classification
- send packet from highest priority queue that has buffered packets
  - FCFS within priority class



## Scheduling policies: round robin

### Round Robin (RR) scheduling:

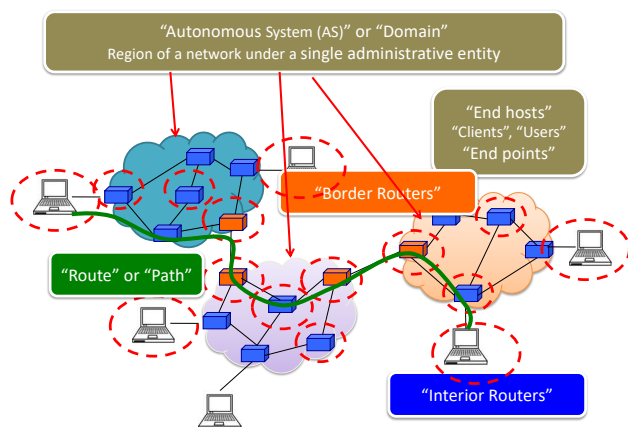
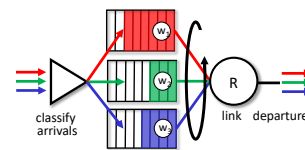
- arriving traffic classified, queued by class
  - any header fields can be used for classification
- server cyclically, repeatedly scans class queues, sending one complete packet from each class (if available) in turn



## Scheduling policies: weighted fair queueing

### Weighted Fair Queueing (WFQ):

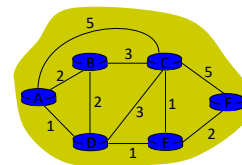
- generalized Round Robin
- each class,  $i$ , has weight,  $w_i$ , and gets weighted amount of service in each cycle:
 
$$\frac{w_i}{\sum_j w_j}$$
- minimum bandwidth guarantee (per-traffic-class)



39

## Routing Protocols

- Routing protocols implement the core function of a network
  - Establish paths between nodes
  - Part of the network's "control plane"
- Network modeled as a graph
  - Routers are graph vertices
  - Links are edges
  - Edges have an associated "cost"
    - e.g., distance, loss
- Goal: compute a "good" path from source to destination
  - "good" usually means the shortest (least cost) path



41

## Internet Routing

- Internet Routing works at two levels
- Each AS runs an **intra-domain** routing protocol that establishes routes within its domain
  - AS -- region of network under a single administrative entity
  - Link State, e.g., Open Shortest Path First (OSPF)
  - Distance Vector, e.g., Routing Information Protocol (RIP)
- ASes participate in an **inter-domain** routing protocol that establishes routes between domains
  - Path Vector, e.g., Border Gateway Protocol (BGP)

42

## Addressing (to date)

- a reminder -

Recall each host has a unique ID (address)

- No particular structure to those IDs (e.g. *Ethernet*)
- IP addressing – in contrast – has implicit structure (Why???)

43

## Outline

- Popular Routing Algorithms:
  - Link State Routing
  - Distance Vector Algorithm
- Routing: goals and metrics

## Link-State Routing

Examples:

Open Shortest Path First (**OSPF**) or  
Intermediate System to Intermediate System  
(written as **IS-IS** and pronounced eye-ess-eye-ess)

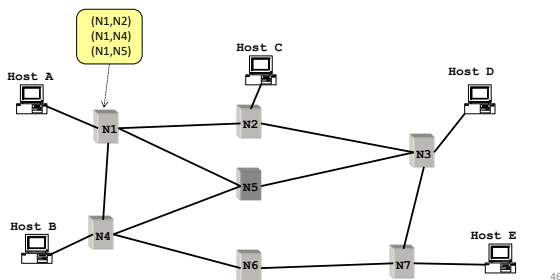
The two common Intradomain routing or  
interior gateway protocols (IGP)

44

45

## Link State Routing

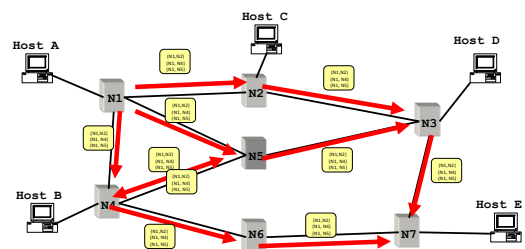
- Each node maintains its **local** "link state" (LS)
  - i.e., a list of its directly attached links and their costs



46

## Link State Routing

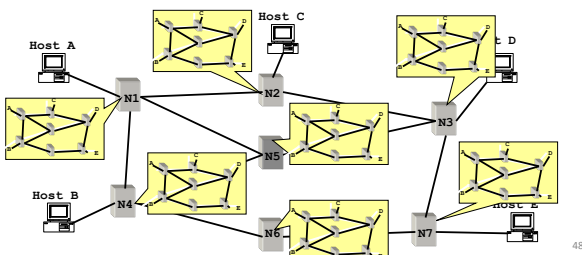
- Each node maintains its local "link state" (LS)
- Each node floods its local link state
  - on receiving a **new** LS message, a router forwards the message to all its neighbors other than the one it received the message from



47

## Link State Routing

- Each node maintains its local "link state" (LS)
- Each node floods its local link state
- Hence, each node learns the entire network topology
  - Can use Dijkstra's to compute the shortest paths between nodes



48

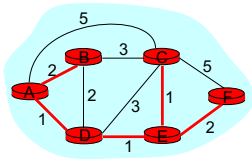
## Dijkstra's Shortest Path Algorithm

- **INPUT:**
  - Network topology (graph), with link costs
- **OUTPUT:**
  - Least cost paths from one node to all other nodes
- Iterative: after  $k$  iterations, a node knows the least cost path to its  $k$  closest neighbors
- This is covered in Algorithms

49

## The Forwarding Table

- Running Dijkstra at node A gives the shortest path from A to all destinations
- We then construct the *forwarding table*



Destination	Link
B	(A,B)
C	(A,D)
D	(A,D)
E	(A,D)
F	(A,D)

50

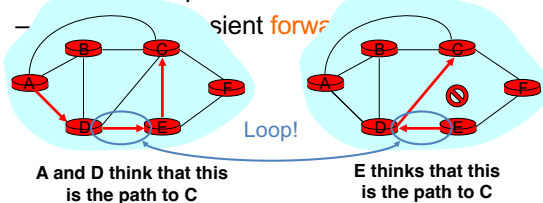
## Issue #1: Scalability

- How many messages needed to flood link state messages?
  - $O(N \times E)$ , where  $N$  is #nodes;  $E$  is #edges in graph
- Processing complexity for Dijkstra's algorithm?
  - $O(N^2)$ , because we check all nodes w not in  $S$  at each iteration and we have  $O(N)$  iterations
  - more efficient implementations:  $O(N \log(N))$
- How many entries in the LS topology database?  $O(E)$
- How many entries in the forwarding table?  $O(N)$

51

## Issue#2: Transient Disruptions

- Inconsistent link-state database
  - Some routers know about failure before others
  - The shortest paths are no longer consistent



52

## Distance Vector Routing

## Learn-By-Doing

Let's try to collectively develop distance-vector routing from first principles

## Experiment

- Your job: find the (route to) the youngest person in the room
- Ground Rules
  - You may not** leave your seat, nor shout loudly across the class
  - You may** talk with your immediate neighbors (N-S-E-W only) (hint: "exchange updates" with them)
- At the end of 5 minutes, I will pick a victim and ask:
  - who is the youngest person in the room? (date&name)
  - which one of your neighbors first told you this info.?

EQUIPMENT REQUIRED: PIECE OF PAPER and a PEN (or your emotional equivalent)

54

55

Go!

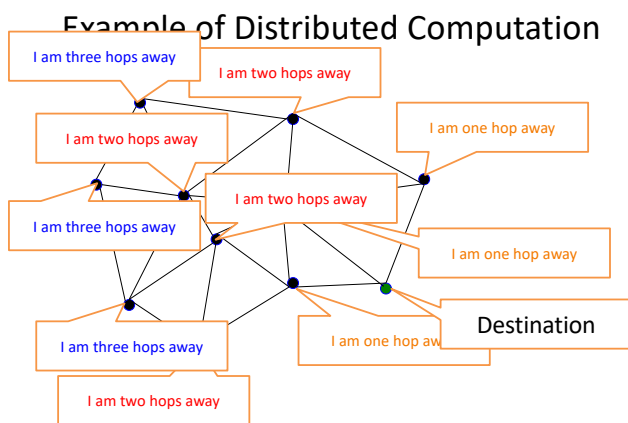
## Distance-Vector Routing

Example:

Routing Information Protocol (RIP)

56

57



58

## Distance Vector Routing

*Each router sends its knowledge about the “whole” network to its neighbors. Information sharing at regular intervals.*

- Each router knows the links to its neighbors
  - Does *not* flood this information to the whole network
- Each router has provisional “shortest path” to **every** other router
  - E.g.: Router A: “I can get to router B with cost 11”
- Routers exchange this **distance vector** information with their neighboring routers
  - Vector because one entry per destination
- Routers look over the set of options offered by their neighbors and select the best one
- Iterative process converges to set of shortest paths

59

## A few other inconvenient truths

- What if we use a non-additive metric?
  - E.g., maximal capacity
- What if routers don’t use the same metric?
  - I want low delay, you want low loss rate?
- What happens if nodes lie?

60

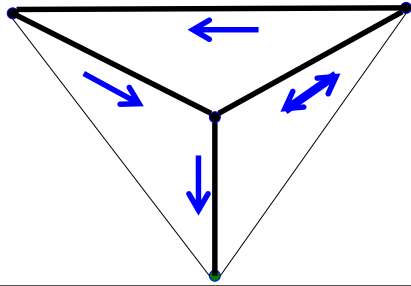
## Can You Use Any Metric?

- I said that we can pick any metric. Really?
- What about maximizing capacity?

61

## What Happens Here?

*Problem: "cost" does not change around loop*



*Additive measures avoid this problem!*

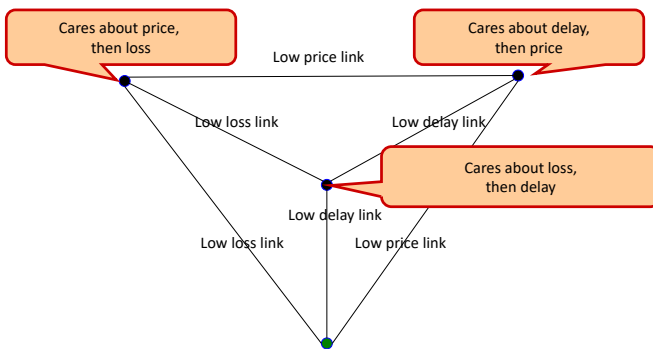
62

## No agreement on metrics?

- If the nodes choose their paths according to different criteria, then bad things might happen
- Example
  - Node A is minimizing latency
  - Node B is minimizing loss rate
  - Node C is minimizing price
- Any of those goals are fine, if globally adopted
  - Only a problem when nodes use different criteria
- Consider a routing algorithm where paths are described by delay, cost, loss

63

## What Happens Here?



64

## Must agree on loop-avoiding metric

- When all nodes minimize same metric
- And that metric increases around loops
- Then process is guaranteed to converge

65

## What happens when routers lie?

- What if a router claims a 1-hop path to everywhere?
- All traffic from nearby routers gets sent there
- How can you tell if they are lying?
- Can this happen in real life?
  - It has, several times....

66

## Link State vs. Distance Vector

- Core idea
  - LS: tell all nodes about your immediate neighbors
  - DV: tell your immediate neighbors about (your least cost distance to) all nodes

67

## Link State vs. Distance Vector

- LS: each node learns the complete network map; each node computes shortest paths independently and in parallel
- DV: no node has the complete picture; nodes cooperate to compute shortest paths in a distributed manner

→ LS has higher messaging overhead  
 → LS has higher processing complexity  
 → LS is less vulnerable to looping

68

## Link State vs. Distance Vector

### Message complexity

- LS:  $O(N \times E)$  messages;
  - $N$  is #nodes;  $E$  is #edges
- DV:  $O(\#iterations \times E)$ 
  - where #iterations is ideally  $O(\text{network diameter})$  but varies due to routing loops or the count-to-infinity problem

### Processing complexity

- LS:  $O(N^2)$
- DV:  $O(\#iterations \times N)$

### Robustness: what happens if router malfunctions?

- LS:
  - node can advertise incorrect *link* cost
  - each node computes only its own table
- DV:
  - node can advertise incorrect *path* cost
  - each node's table used by others; error propagates through network

69

## Routing: Just the Beginning

- Link state and distance-vector are the deployed routing paradigms for intra-domain routing
- Inter-domain routing (BGP)
  - more Part II (Principles of Communications)
  - A version of DV

70

## What are desirable goals for a routing solution?

- “Good” paths (least cost)
- Fast convergence after change/failures
  - no/rare loops
- Scalable
  - #messages
  - table size
  - processing complexity
- Secure
- Policy
- Rich metrics (more later)

71

## Delivery models

- What if a node wants to send to more than one destination?
  - broadcast: send to all
  - multicast: send to all members of a group
  - anycast: send to any member of a group
- What if a node wants to send along more than one path?

72

## Metrics

- Propagation delay
- Congestion
- Load balance
- Bandwidth (available, capacity, maximal, bbw)
- Price
- Reliability
- Loss rate
- Combinations of the above

In practice, operators set abstract “weights” (much like our costs); how exactly is a bit of a black art

73

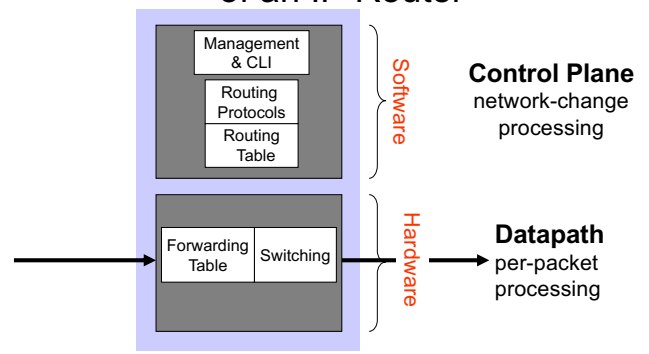


## From Routing back to Forwarding

- Routing: “control plane”
  - Computing paths the packets will follow
  - Routers talking amongst themselves
  - Jointly creating the routing state
- Forwarding: “data plane”
  - Directing a data packet to an outgoing link
  - Individual router using routing state
- Two very different timescales....

74

## Basic Architectural Components of an IP Router



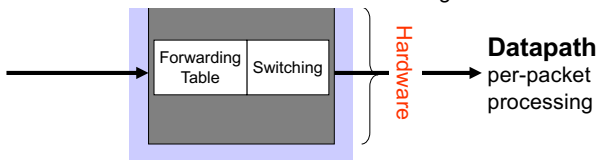
75

## Independent operation!

If the control-plane **fails**....

The data-path is **not affected**...  
like a loyal pet it will keep going using the current (last) table update

This is a feature **not** a bug



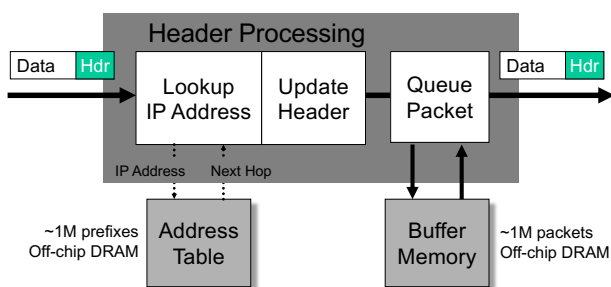
76

## Per-packet processing in an IP Router

1. Accept packet arriving on an incoming link.
2. Lookup packet destination address in the forwarding table, to identify outgoing port(s).
3. Manipulate packet header: e.g., decrement TTL, update header checksum.
4. Send packet to the outgoing port(s).
5. Buffer packet in the queue.
6. Transmit packet onto outgoing link.

77

## Generic Router Architecture



78

## Forwarding tables

IP address } 32 bits wide → ~ 4 billion unique address

### Naïve approach:

One entry per address

Entry	Destination	Port
1	0.0.0.0	1
2	0.0.0.1	2
⋮	⋮	⋮
2 <sup>32</sup>	255.255.255.255	12

~ 4 billion entries

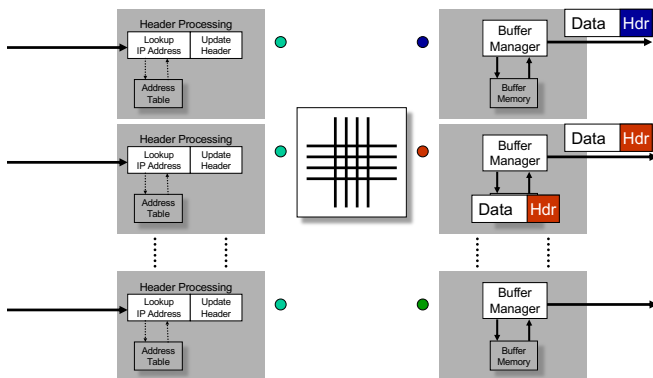
### Improved approach:

Group (and ***SORT***) entries to reduce table size

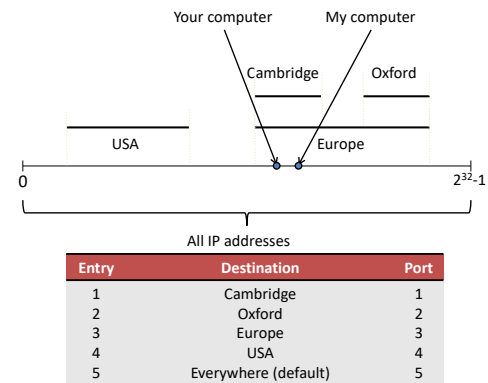
Entry	Destination	Port
1	0.0.0.0 – 127.255.255.255	1
2	128.0.0.1 – 128.255.255.255	2
⋮	⋮	⋮
50	248.0.0.0 – 255.255.255.255	12

79

## Generic Router Architecture

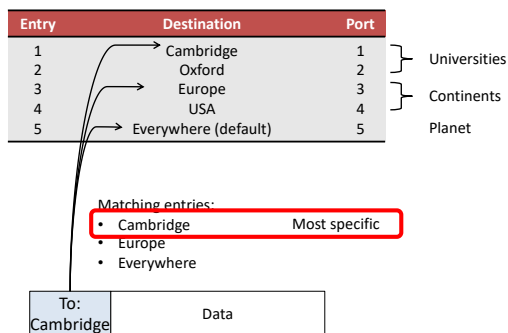


## IP addresses as a line



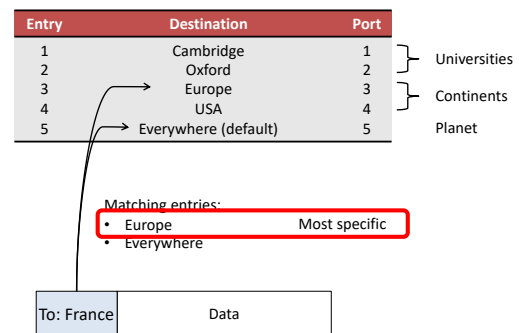
81

## Longest Prefix Match (LPM)



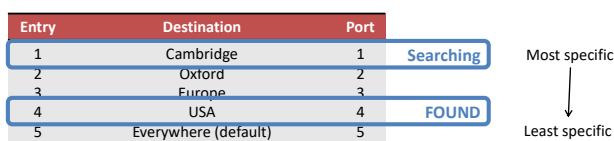
82

## Longest Prefix Match (LPM)



83

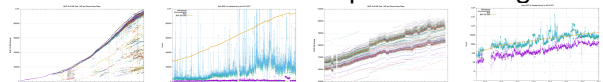
## Implementing Longest Prefix Match



84

## Forwarding table realities

- High Speed: Must be “packet-rate” lookup
  - about 200M lookups / second for 100Gbps
- Large (messy) tables – (BGP Jan 2024 stats)
  - 950,000+ routing prefix entries for IPv4
  - 205,000+ routing prefix entries for IPv6
- Changing and Growing
  - the harsh side of “up and to the right”



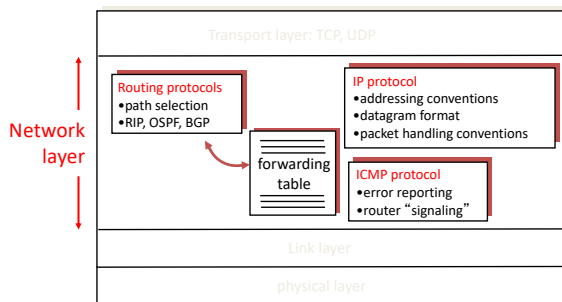
**Open problems :** continual growth is continual demand for innovation opportunities in control, algorithms, & network hardware

Hudson 2023 report <https://blog.apnic.net/2024/01/09/measuring-bgp-in-2023-have-we-reached-peak-ipv4/>

85

## The Internet version of a Network layer

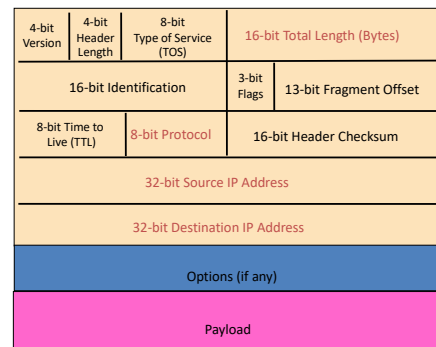
Host, router network layer functions:



86

## IPv4 Packet Structure

### 20 Bytes of Standard Header, then Options



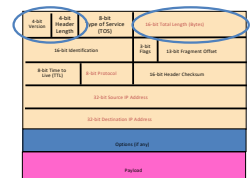
87

## (Packet) Network Tasks One-by-One

- Read packet correctly
- Get packet to the destination
- Get responses to the packet back to source
- Carry data
- Tell host what to do with packet once arrived
- Specify any special network handling of the packet
- Deal with problems that arise along the path

88

## Reading Packet Correctly

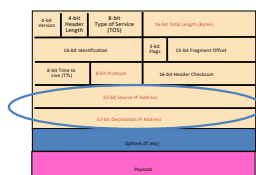


- Version number (4 bits)
  - Indicates the version of the IP protocol
  - Necessary to know what other fields to expect
  - Typically "4" (for IPv4), and sometimes "6" (for IPv6)
- Header length (4 bits)
  - Number of 32-bit words in the header
  - Typically "5" (for a 20-byte IPv4 header)
  - Can be more when IP options are used
- Total length (16 bits)
  - Number of bytes in the packet
  - Maximum size is 65,535 bytes ( $2^{16} - 1$ )
  - ... though underlying links may impose smaller limits

89

## Getting Packet to Destination and Back

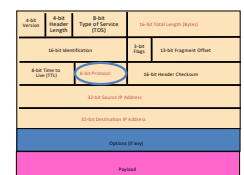
- Two IP addresses
  - Source IP address (32 bits)
  - Destination IP address (32 bits)
- Destination address
  - Unique identifier/locator for the receiving host
  - Allows each node to make forwarding decisions
- Source address
  - Unique identifier/locator for the sending host
  - Recipient can decide whether to accept packet
  - Enables recipient to send a reply back to source



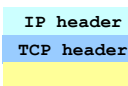
90

## Telling Destination Host How to Handle Packet

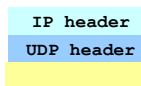
- Protocol (8 bits)
  - Identifies the higher-level protocol
  - Important for demultiplexing at receiving host
- Most common examples
  - E.g., "6" for the Transmission Control Protocol (TCP)
  - E.g., "17" for the User Datagram Protocol (UDP)



protocol=6

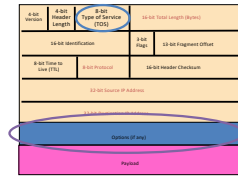


protocol=17



91

## Special Handling



- Type-of-Service (8 bits)
  - Allow packets to be treated differently based on needs
  - E.g., low delay for audio, high bandwidth for bulk transfer
  - Has been redefined several times
  - Rarely actually used and never consistently.....
- Options

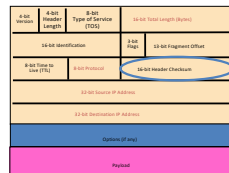
92

## Potential Problems IPv4 solves

- Header Corrupted: **Checksum**
- Loop: **TTL**
- Packet too large: **Fragmentation**

93

## Header Corruption



- Checksum (16 bits)
  - Particular form of checksum over packet header
- If not correct, router discards packets
  - So it doesn't act on bogus information
- Checksum recalculated at every router
  - **Why?**
  - **Why include TTL?**
  - **Why only header?**

94

## Preventing Loops

(aka Internet Zombie plan)

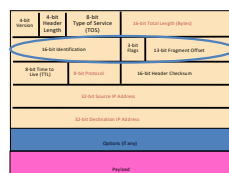


- Forwarding loops cause packets to cycle forever
  - As these accumulate, eventually consume **all** capacity
- Time-to-Live (TTL) Field (8 bits)
  - Decremented at each hop, packet discarded if reaches 0
  - ...and "time exceeded" message is sent to the source
    - Using "ICMP" control message; basis for **traceroute**

95

## Fragmentation

(some assembly required)

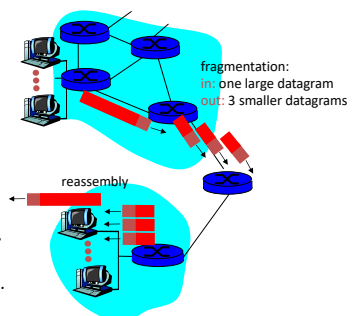


- Fragmentation: when forwarding a packet, an Internet router can **split** it into multiple pieces ("fragments") if too big for next hop link
- Must **reassemble** to recover original packet
  - Need fragmentation information (32 bits)
  - Packet **identifier**, **flags**, and fragment **offset**

96

## 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
- IPv6 does things differently...



97

## IP Fragmentation and Reassembly

### Example

- 4000 byte datagram
- MTU = 1500 bytes

1480 bytes in data field

offset = 1480/8

length	ID	fragflag	offset
=4000	=x	=0	=0

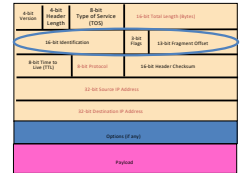
One large datagram becomes several smaller datagrams

length	ID	fragflag	offset
=1500	=x	=1	=0
=1500	=x	=1	=185
=1040	=x	=0	=370

Question: What happens when a fragment is lost?

98

## Fragmentation Details

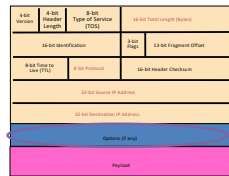


- Identifier (16 bits): used to tell which fragments belong together
- Flags (3 bits):
  - Reserved (**RF**): unused bit
  - Don't Fragment (**DF**): instruct routers to **not** fragment the packet even if it won't fit
    - Instead, they **drop** the packet and send back a "Too Large" ICMP control message
    - Forms the basis for "Path MTU Discovery"
  - More (**MF**): this fragment is not the last one
- Offset (13 bits): what part of datagram this fragment covers in 8-byte units

Pop quiz question: Why do frags use offset and not a frag number?

99

## Options



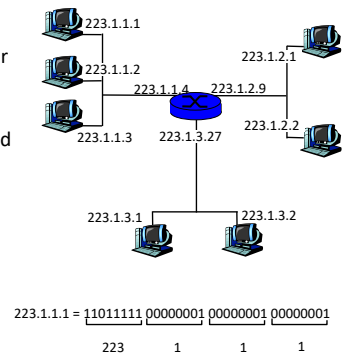
- End of Options List
- No Operation (padding between options)
- Record Route
- Strict Source Route
- Loose Source Route
- Timestamp
- Traceroute
- Router Alert

Few are used as each requires special handling in an IP router.

100

## IP Addressing: introduction

- IP address**: 32-bit identifier for host, router *interface*
- interface**: connection between host/router and physical link
  - routers typically have multiple interfaces
  - host typically has one interface
  - IP addresses associated with each interface



101

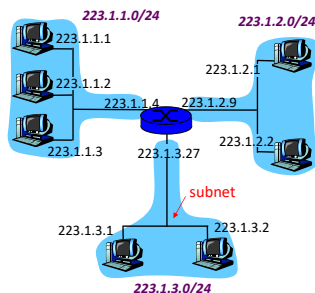
## Subnets

- IP address**:
  - subnet part (high order bits)
  - host part (low order bits)
- What's a subnet?**
  - device interfaces with same subnet part of IP address
  - can physically reach each other without intervening router



**CIDR: Classless InterDomain Routing**

- subnet portion of address of arbitrary length
- address format: **a.b.c.d/x**, where x is # bits in subnet portion of address



Subnet mask: /24

network consisting of 3 subnets

102

## IP addresses: how to get one?

**Q:** How does a *host* get IP address?

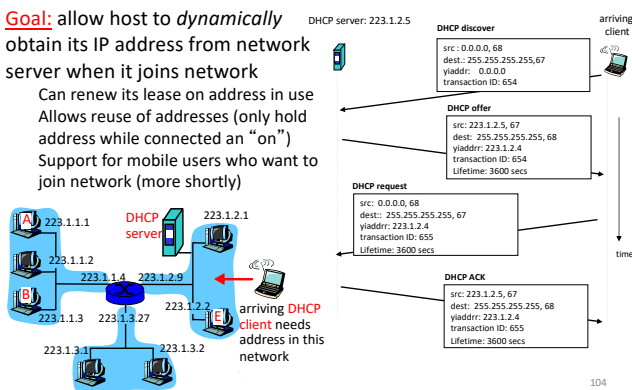
- hard-coded by system admin in a file
  - Windows: control-panel->network->configuration->tcp/ip->properties
  - UNIX: /etc/rc.config (circa 1980's your mileage will vary)
  - R.Pi (last time I checked) /etc/dhcpd.conf
- DHCP: Dynamic Host Configuration Protocol**: dynamically get address from as server
  - "plug-and-play"

103

## DHCP client-server scenario

**Goal:** allow host to *dynamically* obtain its IP address from network server when it joins network

Can renew its lease on address in use  
Allows reuse of addresses (only hold address while connected an "on")  
Support for mobile users who want to join network (more shortly)



104

## IP addresses: how to get one?

**Q:** How does *network* get subnet part of IP addr?

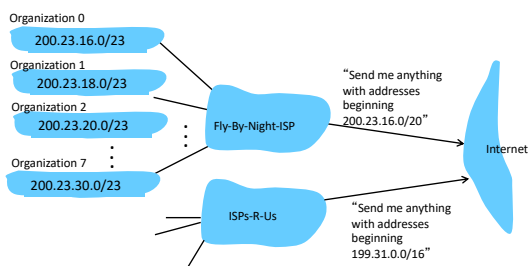
**A:** gets allocated portion of its provider ISP's address space

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

105

## Hierarchical addressing: route aggregation

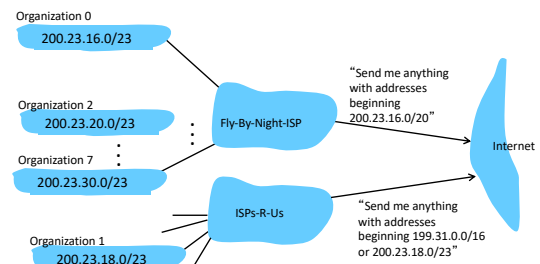
Hierarchical addressing allows efficient advertisement of routing information:



106

## Hierarchical addressing: more specific routes

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



107

## IP addressing: the last word...

**Q:** How does an ISP get a block of addresses?

**A:** ICANN: Internet Corporation for Assigned

Names and Numbers

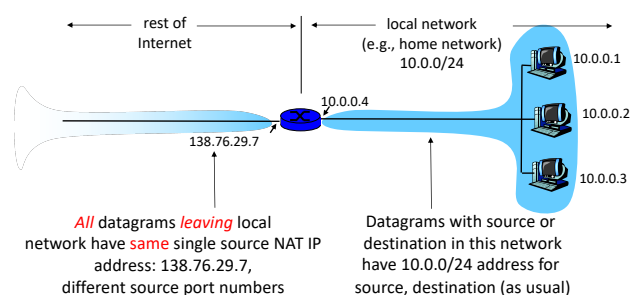
- allocates addresses
- manages DNS
- assigns domain names, resolves disputes

There are regional subordinates but the (US location) of the ICANN dominates proceedings....

108

Cant get more IPv4 addresses? well there is always....

## NAT: Network Address Translation



109

## 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).

110

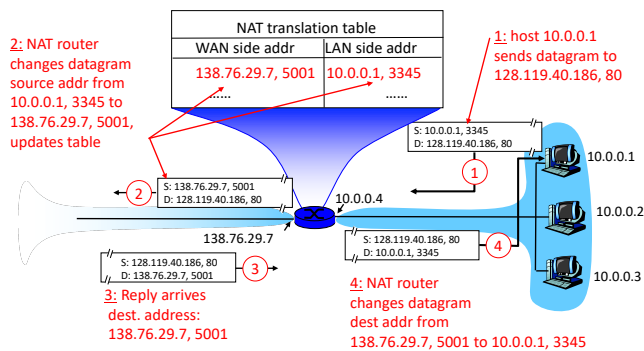
## NAT: Network Address Translation

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

111

## NAT: Network Address Translation



112

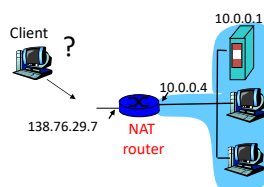
## NAT: Network Address Translation

- 16-bit port-number field:
  - 60,000+ simultaneous connections with a single WAN-side address!
- NAT is controversial:
  - routers should only process up to layer 3
  - violates end-to-end argument (?)
    - NAT possibility must be taken into account by app designers, eg, P2P applications
  - address shortage “should” instead be solved by IPv6

113

## NAT traversal problem

- client wants to connect to server with address 10.0.0.1
  - server address 10.0.0.1 local to LAN (client can't use it as destination addr)
  - only one externally visible NATted address: 138.76.29.7
- solution 1: statically configure NAT to forward incoming connection requests at given port to server
  - e.g., (138.76.29.7, port 2500) always forwarded to 10.0.0.1 port 25000

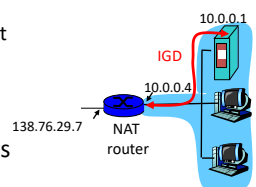


114

## NAT traversal problem

- solution 2: Universal Plug and Play (UPnP) Internet Gateway Device (IGD) Protocol. Allows NATted host to:
  - ❖ learn public IP address (138.76.29.7)
  - ❖ add/remove port mappings (with lease times)

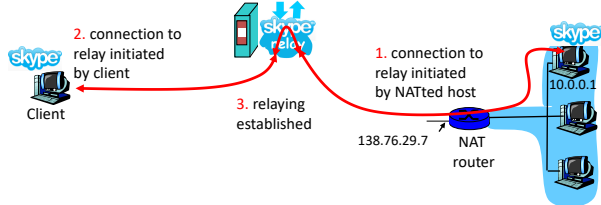
i.e., automate static NAT port map configuration



115

## NAT traversal problem

- solution 3: relaying (was used in (really old) Skype)
  - NATed client establishes connection to relay
  - External client connects to relay
  - relay bridges packets between to connections



116

## Remember this? Traceroute at work...

traceroute: department ssh server to melbourneisp.com (Melbourne)  
(tracepath on winows is similar)

```

awm228svr-ssh-0:-$ traceroute melbourneisp.com
traceroute to melbourneisp.com (116.206.130.24), 30 hops max, 60 byte packets
 1 vian308.gatewick.net.c1.cam.ac.uk (128.232.64.2) 10.209 ms 10.893 ms 11.064 ms
 2 c1-wg0.d-mw.net.cam.ac.uk (193.60.89.5) 0.743 ms 0.789 ms 0.586 ms
 3 d-mw.c-ce.net.cam.ac.uk (131.111.6.63) 0.917 ms 1.123 ms 0.986 ms
 4 c-ce-b-jc.net.cam.ac.uk (131.111.6.82) 0.795 ms 0.771 ms 0.765 ms
 5 ips-out-b-jc.net.cam.ac.uk (131.111.7.217) 1.096 ms 0.986 ms 1.032 ms
 6 ae8.lowdss-bani.ja.net (146.97.41.37) 3.484 ms 3.019 ms 3.076 ms
 7 ae26.lowdss-sbri.ja.net (146.97.35.245) 3.740 ms 3.454 ms 3.376 ms
 8 ae31.lowdss-sbri.ja.net (146.97.33.30) 7.034 ms 6.631 ms 6.962 ms
 9 ae28.lowdss-sbri.ja.net (146.97.33.61) 8.829 ms 16.876 ms 16.954 ms
10 ae8.lowdss-sbri.ja.net (146.97.35.194) 7.328 ms 6.467 ms 6.387 ms
11 ldn-bt1-link.ip.twelve99.net (62.115.175.186) 6.476 ms 6.234 ms 6.585 ms
12 ldn-bt1-link.ip.twelve99.net (62.115.138.168) 7.473 ms * 7.740 ms
13 nyk-bb2-link.ip.twelve99.net (62.115.139.245) 76.127 ms nyk-bb5-link.ip.twelve99.net (62.115.139.244) 75.315 ms *
14 * chl-bb2-link.ip.twelve99.net (62.115.133.135) 148.582 ms 148.856 ms
15 * den-bb1-link.ip.twelve99.net (62.115.115.76) 115.630 ms 114.872 ms
16 den-bb2-link.ip.twelve99.net (62.115.137.114) 113.977 ms * 113.656 ms
17 palp-bb2-link.ip.twelve99.net (62.115.139.112) 143.238 ms 144.527 ms *
18 * tpg-ic-387776.ip.twelve99-cust.net (62.115.188.6) 295.626 ms 296.153 ms
19 tpg-ic-387776.ip.twelve99-cust.net (62.115.188.6) 295.143 ms syd-apt-ros-crt3-be-100.tpgi.com.au (203.29.134.43) 291
20 syd-apt-ros-crt3-be-100.tpgi.com.au (203.29.134.43) 295.545 ms syd-sot-ken-crt2-te-0-0-9.tpgi.com.au (203.26.22.12)
21 syd-sot-ken-crt2-te-0-0-9.tpgi.com.au (203.26.22.12) 380.416 ms AU-VI-1015-IP0-221-Bundle-Ether1.tpgi.com.au (27.
22 AU-VI-4901-IP0-01-Eth-Trunk21.tpgi.com.au (203.220.216.30) 301.396 ms AU-VI-1015-IP0-221-Bundle-Ether2.tpgi.com.au (
23 AU-VI-4901-IP0-01-Eth-Trunk21.tpgi.com.au (203.220.216.30) 300.726 ms 14-202-130-170-static.tpgi.com.au (14.202.130.
24 14-202-130-170-static.tpgi.com.au (14.202.130.170) 303.742 ms * 303.963 ms
25 * * *
26 * * *
27 * * *
28 * * *
29 * * *
30 * * *
* means no response (probe or reply lost, router not replying)
awm228svr-ssh-0:-$
    
```

300ms RTT, 150ms one way Internet, 59.4ms by photon, 42ms by neutron

117

## Traceroute and ICMP

- Source sends series of UDP segments to dest
  - First has TTL=1
  - Second has TTL=2, etc.
  - Unlikely port number
- When nth datagram arrives to nth router:
  - Router discards datagram
  - And sends to source an ICMP message (type 11, code 0)
  - Message includes name of router & IP address
- When ICMP message arrives, source calculates RTT
- Traceroute does this 3 times
- Stopping criterion
- UDP segment eventually arrives at destination host
- Destination returns ICMP “host unreachable” packet (type 3, code 3)
- When source gets this ICMP, stops.

118

## ICMP: Internet Control Message Protocol

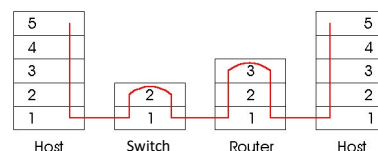
- used by hosts & routers to communicate network-level information
 

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
- 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

119

## Switches vs. Routers Summary

- both (can be implemented as) store-and-forward devices
  - routers: network layer devices (**manipulate** network layer headers eg IP)
  - switches are link layer devices (**examine** Data-Link-Layer headers eg Ethernet)
- Routers: implement routing algorithms, maintain routing tables of the network – create network forwarding tables from routing tables
- Switches: implement learning algorithms, learn switch/DLL forwarding tables



120

121

### Gluing it together:

How does my Network (address) interact with my Data-Link (address) ?



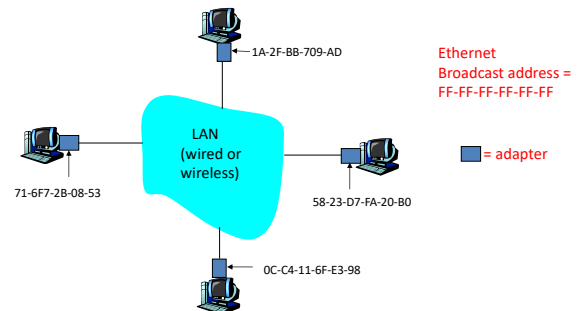
## MAC Addresses (and IPv4 ARP) or How do I glue my network to my data-link?

- 32-bit IP address:
  - *network-layer* address
  - used to get datagram to destination IP subnet
- MAC (or LAN or physical or Ethernet) address:
  - function: *get frame from one interface to another physically-connected interface (same network)*
  - 48 bit MAC address (for most LANs)
    - burned in NIC ROM, firmware, etc.

122

## LAN Addresses and ARP

Each adapter on LAN has unique LAN (MAC) address



123

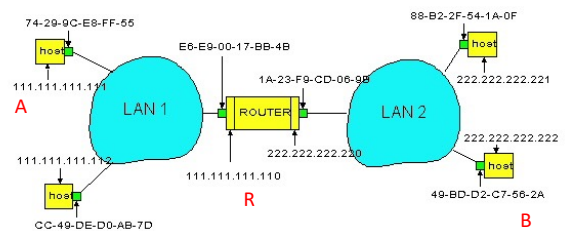
## Address Resolution Protocol

- Every node maintains an **ARP** table
  - <IP address, MAC address> pair
- Consult the table when sending a packet
  - Map destination IP address to destination MAC address
  - Encapsulate and transmit the data packet
- But: what if IP address **not** in the table?
  - Sender **broadcasts**: “Who has IP address 1.2.3.156?”
  - Receiver responds: “MAC address 58-23-D7-FA-20-B0”
  - Sender **caches** result in its ARP table

124

## Example: A Sending a Packet to B

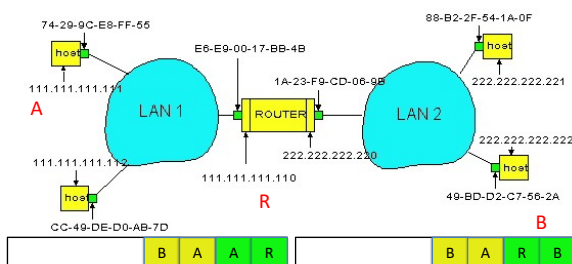
How does host **A** send an IP packet to host **B**?



125

## Example: A Sending a Packet to B

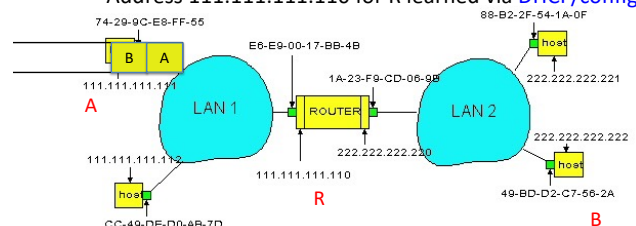
How does host **A** send an IP packet to host **B**?



126

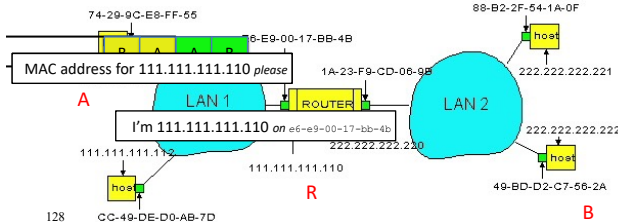
## Host A Decides to Send Through R

- Host **A** constructs an IP packet to send to **B**
  - Source 111.111.111.111, destination 222.222.222.222
- Host **A** has a gateway router **R**
  - Used to reach destinations outside of 111.111.111.0/24
  - Address 111.111.111.110 for R learned via **DHCP/config**



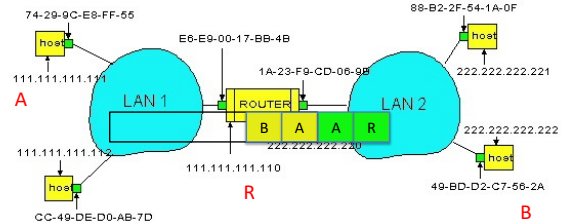
## Host A Sends Packet Through R

- Host **A** learns the MAC address of **R**'s interface
  - ARP request: broadcast request for 111.111.111.110
  - ARP response: **R** responds with E6-E9-00-17-BB-4B
- Host **A** encapsulates the packet and sends to **R**



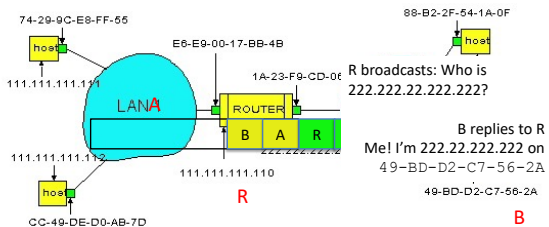
## R Decides how to Forward Packet

- Router **R**'s adaptor receives the packet
  - **R** extracts the IP packet from the Ethernet frame
  - **R** sees the IP packet is destined to 222.222.222.222
- Router **R** consults its forwarding table
  - Packet matches 222.222.222.0/24 via other adaptor



## R Sends Packet to B

- Router **R**'s learns the MAC address of host **B**
  - ARP request: broadcast request for 222.222.222.222
  - ARP response: **B** responds with 49-BD-D2-C7-56-2A
- Router **R** encapsulates the packet and sends to **B**



## Security Analysis of ARP



- Impersonation**
  - Any node that hears request can answer ...
  - ... and can say **whatever** they want
- Actual legit receiver **never sees a problem**
  - Because even though later packets carry its IP address, its NIC doesn't capture them since the (naughty) packets are **not its MAC address**

131

## Key Ideas in Both ARP and DHCP

- Broadcasting**: Can use broadcast to make contact
  - Scalable because of limited size
- Caching**: remember the past for a while
  - Store the information you learn to reduce overhead
  - Remember your own address & other host's addresses
- Soft state**: eventually forget the past
  - Associate a **time-to-live** field with the information
  - ... and either refresh or discard the information
  - Key for **robustness** in the face of unpredictable change

132

## Why Not Use DNS-Like Tables?

- When host arrives:
  - Assign it an IP address that will last as long it is present
  - Add an entry into a table in DNS-server that maps MAC to IP addresses
- Answer:
  - Names: explicit creation, and are plentiful
  - Hosts: come and go without informing network
    - Must do mapping on demand
  - Addresses: not plentiful, need to reuse and remap
    - Soft-state enables dynamic reuse

133

## IPv6



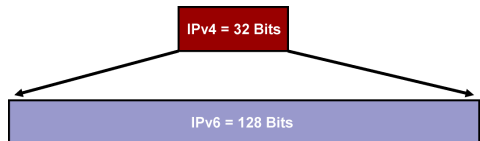
- prematurely
  - Motivated by address exhaustion
    - addresses are larger
    - packet headers are laid out differently
    - address management and configuration are completely different
    - some DNS behavior changes
    - some sockets code changes
    - everybody now has a hard time parsing IP addresses*
- Steve Deering focused on simplifying IP
  - Got rid of all fields that were not absolutely necessary
  - “Spring Cleaning” for IP
- Result is an elegant, if unambitious, protocol



IPv4	IPv6
Addresses are 32 bits (4 bytes) in length.	Addresses are 128 bits (16 bytes) in length
Address (A) resource records in DNS to map host names to IPv4 addresses.	Address (AAAA) resource records in DNS to map host names to IPv6 addresses.
Pointer (PTR) resource records in the IN-ADDR.ARPA DNS domain to map IPv4 addresses to host names.	Pointer (PTR) resource records in the IP6.ARPA DNS domain to map IPv6 addresses to host names.
IPSec is optional and should be supported externally	IPSec support is not optional
Header does not identify packet flow for QoS handling by routers	Header contains Flow Label field, which identifies packet flow for QoS handling by router.
Both routers and the sending host fragment packets.	Routers do not support packet fragmentation. Sending host fragments packets
Header includes a checksum.	Header does not include a checksum.
Header includes options.	Optional data is supported as extension headers.
ARP uses broadcast ARP request to resolve IP to MAC/Hardware address.	Multicast Neighbor Solicitation messages resolve IP addresses to MAC addresses.
Internet Group Management Protocol (IGMP) manages membership in local subnet groups.	Multicast Listener Discovery (MLD) messages manage membership in local subnet groups.
Broadcast addresses are used to send traffic to all nodes on a subnet.	IPv6 uses a link-local scope all-nodes multicast address.
Configured either manually or through DHCP.	Does not require manual configuration or DHCP.
Must support a 576-byte packet size (possibly fragmented).	Must support a 1280-byte packet size (without fragmentation).

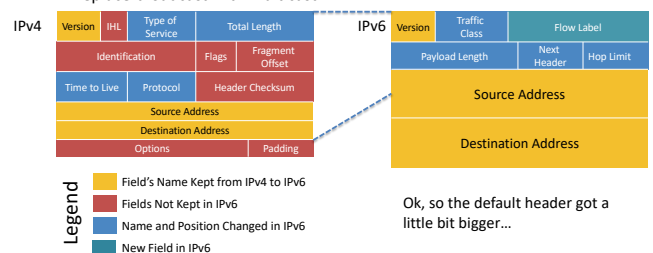
## Larger Address Space

- IPv4 = 4,294,967,295 addresses
- IPv6 = 340,282,366,920,938,463,374,607,432,768,211,456 addresses
- 4x in number of bits translates to huge increase in address space!



## Other Significant Protocol Changes - 1

- Increased minimum MTU from 576 to 1280
- No enroute fragmentation... fragmentation only at source
- Header changes (20bytes to 40bytes)
- Replace broadcast with multicast



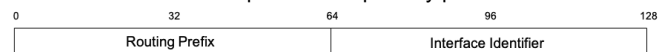
## Other Significant Protocol Changes - 2

operation is intended to be simpler within the network:

- no *in-network* fragmentation
- no checksums in IPv6 header
- UDP checksum required (wasn't in IPv4) rfc6936: **No more zero**
- optional state carried in *extension headers*
  - Extension headers notionally replace IP options
  - Each extension header indicates the type of the *following* header, so they can be chained
  - The final 'next header' either indicates there is no 'next', or escapes into a transport-layer header (e.g., TCP)

## IPv6 Basic Address Structure

IPv6 addresses are split into two primary parts:



- 64 bits is dedicated to an addressable interface (equivalent to the host, if it only has one interface)
- The network prefix allocated to a network by a registry can be up to 64-bits long
- An allocation of a /64 (i.e. a 64-bit network prefix) allows *one* subnet (it cannot be subdivided)
- A /63 allows two subnets; a /62 offers four, etc. /48s are common for older allocations (RFC 3177, obsolete by RFC 6177).
- Longest-prefix matching operates as in IPv4.

## IPv6 Address Representation (quick)

IPv6 addresses represented as eight 16-bit blocks (4 hex chars) separated by colons:

- 2001:4998:000c:0a06:0000:0000:0002:4011

But we can condense the representation by removing leading zeros in each block:

- 2001:4998:c:a06:0:0:2:4011

And by reducing the consecutive block of zeros to a "::" (this double colon rule can only be applied once)

- 2001:4998:c:a06::2:4011

140

## IPv6 Address Families

The address space is carved, like v4, into certain categories<sup>1</sup>:

host-local : localhost; ::1 is equivalent to 127.0.0.1

link-local : not routed: fe80::/10 is equivalent to 169.254.0.0/16

site-local : not routed globally: fc00::/7 is equivalent to 192.168.0.0/16 or 10.0.0.0/8

global unicast : 2000::/3 is basically any v4 address not reserved in some other way

multicast : ff00::/8 is equivalent to 224.0.0.0/4

<sup>1</sup>[http://www.ripe.net/lir-services/new-lir/ipv6\\_reference\\_card.pdf](http://www.ripe.net/lir-services/new-lir/ipv6_reference_card.pdf)

141

## Problem with /64 Subnets

- Scanning a subnet becomes a DoS attack!
  - Creates IPv6 version of 2<sup>64</sup> ARP entries in routers
  - Exhaust address-translation table space

- So now we have:

ping6 ff02::1 All nodes in broadcast domain

ping6 ff02::2 All routers in broadcast domain

- Solutions
  - RFC 6164 recommends use of /127 to protect router-router links
  - RFC 3756 suggest "clever cache management" to address more generally

142

## Neighbour Discovery

- The Neighbour Discovery Protocol<sup>2</sup> specifies a set of ICMPv6 message types that allow hosts to discover other hosts or routing hardware on the network
  - neighbour solicitation
  - neighbour advertisement
  - router solicitation
  - router advertisement
  - redirect
- In short, a host can *solicit* neighbour (host) state to determine the layer-2 address of a host or to check whether an address is in use
- or it can solicit router state to learn more about the network configuration
- In both cases, the solicit message is sent to a well-known multicast address

<sup>2</sup><http://tools.ietf.org/html/rfc4861>

143

## IPv6 Dynamic Address Assignment

We have the two halves of the IPv6 address: the network component and the host component. Those are derived in different ways.

Network (top 64 bits):

- Router Advertisements (RAs)  
Interface

Identifier (bottom 64 bits):

- Stateless, automatic: SLAAC
- Stateful, automatic: DHCPv6

144

## SLAAC: overview

SLAAC is:

- ... intended to make network configuration easy without manual configuration or even a DHCP server
- ... an algorithm for hosts to automatically configure their network interfaces (set up addresses, learn routes) without intervention

145

## SLAAC: overview

- When a host goes live or an interface comes up, the system wants to know more about its environment
- It *can* configure link-local addresses for its interfaces: it uses the interface identifier, the EUI-64
- It uses this to ask (solicit) router advertisements sooner than the next periodic announcements; ask the network for information

146

## SLAAC: overview

The algorithm (assuming one interface):

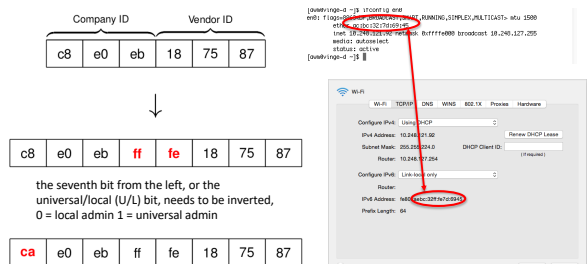
1. Generate potential link-local address
2. Ask the network (multicast<sup>4</sup>) if that address is in use: *neighbour solicitation*
3. Assuming no responses, assign to interface

<sup>4</sup><https://tools.ietf.org/html/rfc2373>

147

## The EUI-64 Interface Identifier

- IEEE 64-bit Extended Unique Identifier (EUI-64)<sup>3</sup>
- There are various techniques to derive a 64-bit value, but often times we derive from the 48-bit MAC address



<sup>3</sup><http://tools.ietf.org/html/rfc2373>

148

## SLAAC: overview; Router Solicitation

Then,

- Once the host has a unique *link-local* address, it can send packets to anything else sharing that link substrate ... but the host doesn't yet know any routers, or public routes ... bootstrap: routers listen to a well-known multicast address
4. host asks the network (multicast) for router information: *router solicitation*
  5. responses from the routers are sent directly (unicast) to the host that sent the router solicitation
  6. the responses *may* indicate that the host should do more (e.g., use DHCP to get DNS information)

149

## Router Advertisement

Without solicitation, regular router advertisements are generated by routing hardware.

Router Advertisements:

- nodes that forward traffic periodically advertise themselves to the network
- periodicity and expiry of the advertisement are configurable

Router Advertisement (RA), among other things, tells a host where to derive its network state with two flags: M(anaged) and O(ther info):

- M: "Managed Address Configuration", which means: use DHCPv6 to find your host address (and ignore option O)
- O: Other information is available via DHCPv6, such as DNS configuration

150

## Uh-oh

What problem(s) arises from totally decentralised address configuration?

Concerns that arise from using an EUI-64:

- Privacy: SLAAC interface identifiers don't change over time, so a host can be identified across networks
- Security: embedding a MAC address into an IPv6 address will carry that vendor's ID(s)<sup>5</sup>, a possible threat vector

<sup>5</sup><http://standards.ieee.org/develop/regauth/oui/public.html>

151

## Address Configuration: SLAAC Privacy Addresses

### Privacy extensions for SLAAC<sup>6</sup>

- temporary addresses for initiating outgoing sessions
- generate one temporary address per prefix
- when they expire, they are not used for new sessions, but can continue to be used for existing sessions
- the addresses should appear random, such that they are difficult to predict
- lifetime is configurable; this OSX machine sets an 86,400s timer (1 day)

<sup>6</sup><https://tools.ietf.org/html/rfc4941>

152

## Address Configuration: SLAAC Privacy Addresses

The algorithm:

- Assume: a stored 64-bit input value from previous iterations, or a pseudo-randomly generated value
  1. take that input value and append it to the EUI-64
  2. compute the MD5 message digest of that value
  3. set bit 6 to zero
  4. compare the leftmost 64-bits against a list of reserved interface identifiers and those already assigned to an address on the local device. If the value is unacceptable, re-run using the rightmost 64 bits of the result instead of the historic input value in step 1
  5. use the leftmost 64-bits as the randomised interface identifier
  6. store the rightmost 64-bits as the history value to be used in the next iteration of the algorithm

153

## IPv6: why has the transition taken so long?

### IPv4 and IPv6 are not compatible:

- different packet formats
- different addressing schemes
- no flag days

as the Internet has grown bigger and accumulated many IPv4-only services, transition has proven ... Tricky

### Incentive issues

e.g. Virgin Media policy in 2010

...When IPv6 is rolled out across the whole of the Internet then a lot of the ISP's will roll out IPv6, ...

Virgin Media are only now (late 2024) "committing" to IPv6

154

## IPv6: why has the transition taken so long?

- IPv4 has/had the momentum
  - ... which led to CIDR
  - ... and encouraged RFC1918 space and NAT
- IPv4 NAT was covered earlier in this topic (reminder)
  - your ISP hands you only one IPv4 address
  - you share that across multiple devices in your household
  - The NAT handles all the translation between internal ("private") and external ("public") space

155

## Transition tech: outline

- Tunnelling
- dual-stacked services, and happy eyeballs
- DNS behaviour

156

## Transition tech: outline

- Tunnelling



Hurricane Electric Free IPv6 Tunnel Broker

### IPv6 Tunnel Broker

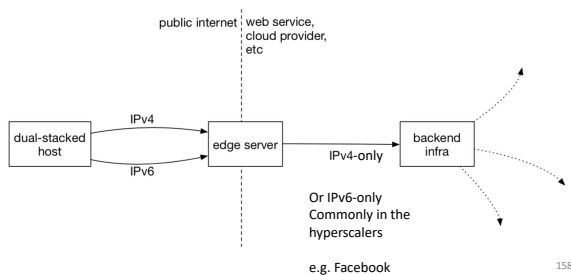
Think of it as an IPv6 VPN service; which is essentially what it is

<sup>8</sup><https://tools.ietf.org/html/rfc6146>

157

## Dual-Stack Services: Common Deployment

It's common for web services to play conservatively: dual-stack your edge services (e.g., load balancers), leaving some legacy infrastructure for later:



158

## Dual-Stack Services: Common Deployment

Aim is to reduce the pain:

- You can dual-stack the edge hosts, and carry state in, say, HTTP headers indicating the user's IP address (common over v4 anyway)
- You can dual-stack the backend opportunistically, over a longer period of time
- You use DNS to enable/disable the v6 side last (if there is no AAAA record in DNS, no real users will connect to the IPv6 infrastructure)

159

## IPV6 sadness and DNS

- The introduction of IPv6 carried with it an obligation that applications attempt to use IPv6 before falling back to IPv4.
- What happens though if you try to connect to a host which doesn't exist?<sup>9</sup>
- But the presence of IPv6 modifies the behaviour of DNS responses and response preference<sup>10</sup>

<sup>9</sup><https://tools.ietf.org/html/rfc5461>  
<sup>10</sup><https://tools.ietf.org/html/rfc3484>

160

## Happy Eyeballs

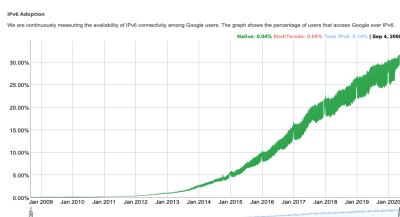
- Happy Eyeballs<sup>11</sup> was the proposed solution
  - the eyeballs in question are yours, or mine, or whoever is sitting in front of their browser getting mad that things are unresponsive
- Modifies application behaviour

<sup>11</sup><https://tools.ietf.org/html/rfc8305>

161

## IPv6: adoption

- Google<sup>1</sup>: ~ 30% of clients access services via IPv6
- NIST: 1/3 of all US government domains are IPv6 capable



<sup>1</sup> <https://www.google.com/intl/en/ipv6/statistics.html>

## Improving on IPv4 and IPv6?

- Why include unverifiable source address?
  - Would like accountability **and** anonymity (now neither)
  - Return address can be communicated at higher layer
- Why packet header used at edge same as core?
  - Edge: host tells network what service it wants
  - Core: packet tells switch how to handle it
    - One is local to host, one is global to network
- Some kind of payment/responsibility field?
  - Who is responsible for paying for packet delivery?
  - Source, destination, other?
- Other ideas?

163

## Summary Network Layer

- understand principles behind network layer services:
  - network layer service models
  - forwarding versus routing (versus switching)
  - how a switch & router works
  - routing (path selection)
  - IPv6
- Algorithms
  - Two routing approaches (LS vs DV)
  - One of these in detail (LS)
  - ARP
- Other Core ideas
  - Caching, soft-state, broadcast
  - Fate-sharing in practice....



## Topic 5 – Transport

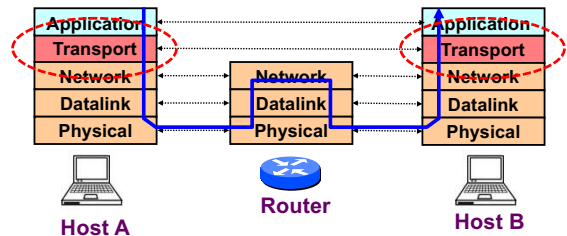
### Our goals:

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

1

## Transport Layer

- Commonly a layer **at end-hosts**, between the application and network layer



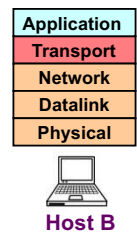
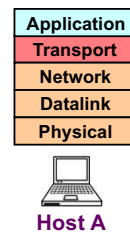
2

## Why a transport layer?

- IP packets are addressed to a host but end-to-end communication is between application/processes/tasks at hosts
  - Need a way to decide which packets go to which applications (*more multiplexing*)

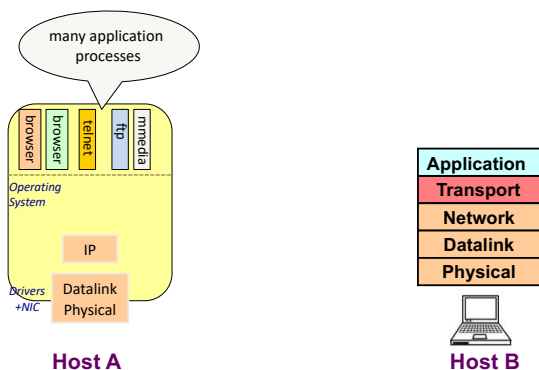
3

## Why a transport layer?



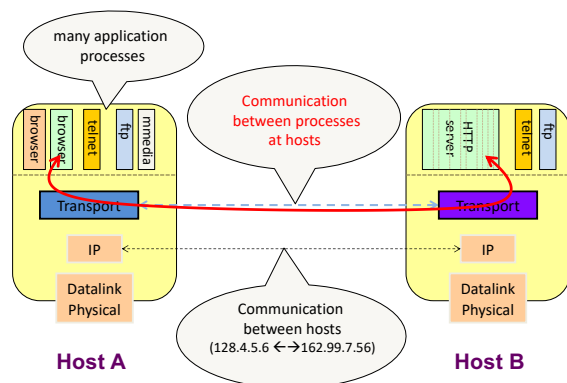
4

## Why a transport layer?



5

## Why a transport layer?



6

## Why a transport layer?

- IP packets are addressed to a host but end-to-end communication is between application processes at hosts
  - Need a way to decide which packets go to which applications (mux/demux)
- IP provides a weak service model (*best-effort*)
  - Packets can be corrupted, delayed, dropped, reordered, duplicated
  - No guidance on how much traffic to send and when
  - Dealing with this is tedious for application developers

7

## Role of the Transport Layer

- Communication between application processes
  - Multiplexing between application processes
  - Implemented using *ports*

8

## Role of the Transport Layer

- Communication between application processes
- Provide common end-to-end services for app layer [optional]
  - Reliable, in-order data delivery
  - Paced data delivery: flow and congestion-control
    - too fast may overwhelm the network
    - too slow is not efficient

(Just Like Computer Networking Lectures....)

9

## Role of the Transport Layer

- Communication between processes
- Provide common end-to-end services for app layer [optional]
- TCP and UDP are the common transport protocols
  - also SCTP, MTCP, SST, RDP, DCCP, ...

10

## Role of the Transport Layer

- Communication between processes
- Provide common end-to-end services for app layer [optional]
- TCP and UDP are the common transport protocols
- UDP is a minimalist, no-frills transport protocol
  - only provides mux/demux capabilities

11

## Role of the Transport Layer

- Communication between processes
- Provide common end-to-end services for app layer [optional]
- TCP and UDP are the common transport protocols
- UDP is a minimalist, no-frills transport protocol
- TCP is the *totus porcus* protocol
  - offers apps a reliable, in-order, byte-stream abstraction
  - with congestion control
  - but **no** performance (delay, bandwidth, ...) guarantees

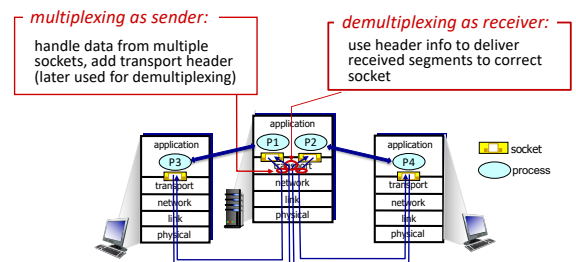
12

## Role of the Transport Layer

- Communication between processes
  - mux/demux from and to application processes
  - implemented using ports

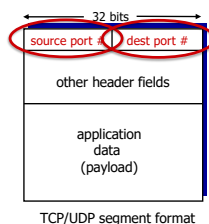
13

## Multiplexing/demultiplexing



## 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



## Connectionless demultiplexing

- when creating socket, must specify **host-local** port #:
 

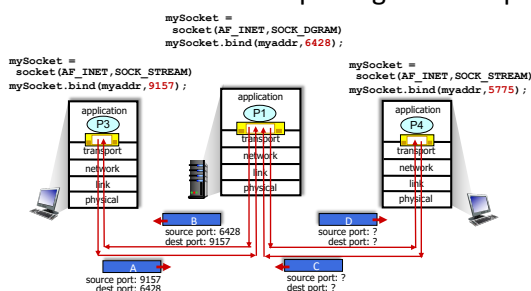
```
DatagramSocket mySocket1 = new DatagramSocket(12534);
```
- when creating datagram to send into UDP socket, must specify
  - destination IP address
  - destination port #

- when receiving host receives **UDP** segment:
- checks destination port # in segment
  - directs UDP segment to socket with that port #

↓

IP/UDP datagrams with **same dest. port #**, but different source IP addresses and/or source port numbers will be directed to **same socket** at receiving host

## Connectionless demultiplexing: an example



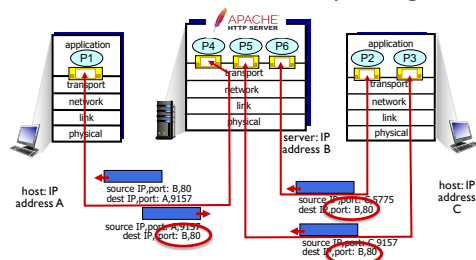
## Connection-oriented demultiplexing

- TCP socket identified by **4-tuple**:
  - source IP address
  - source port number
  - dest IP address
  - dest port number
- demux: receiver uses **all four values (4-tuple)** to direct segment to appropriate socket
- server may support many simultaneous TCP sockets:
  - each socket identified by its own 4-tuple
  - each socket associated with a different connecting client

slight lie alert.... I should say that a common network tuple has FIVE values

- source IP address
- source port number
- dest IP address
- dest port number AND
- protocol e.g. TCP (6) or UDP (17)

## Connection-oriented demultiplexing: example



Three segments, all destined to IP address: B, dest port: 80 are demultiplexed to *different* sockets

## Summary

- Multiplexing, demultiplexing: based on segment, datagram header field values
- **UDP**: demultiplexing using destination port number (only)
- **TCP**: demultiplexing using 4-tuple: source and destination IP addresses, and port numbers
- Multiplexing/demultiplexing can happen at *any* layer

## More on Ports

- Separate 16-bit port address space for UDP and TCP
- “Well known” ports (0-1023): everyone agrees which services run on these ports
  - e.g., ssh:22, http:80, https:443
  - helps client know server’s port
- Ephemeral ports (most 1024-65535): dynamically selected: as the source port for a client process

28

## UDP: User Datagram Protocol

- Lightweight communication between processes
  - Avoid overhead and delays of ordered, reliable delivery
- UDP described in RFC 768 – (1980!)
  - Destination IP address and port to support demultiplexing
  - Optional error checking on the packet contents
    - (checksum field of 0 means “don’t verify checksum”) *not in IPv6!*
    - ((this idea of optional checksum is removed in IPv6))

SRC port	DST port
checksum	length
DATA	

29

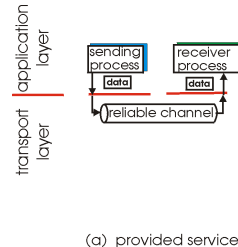
## Why a transport layer?

- IP packets are addressed to a host but end-to-end communication is between application processes at hosts
  - Need a way to decide which packets go to which applications (mux/demux)
- IP provides a weak service model (*best-effort*)
  - Packets can be corrupted, delayed, dropped, reordered, duplicated

30

## Principles of Reliable data transfer

- important in app., transport, link layers
  - top-10 list of important networking topics!
    - In a perfect world, reliable transport is easy
- But the Internet default is *best-effort*
- All the bad things best-effort can do
    - a packet is corrupted (bit errors)
    - a packet is lost
    - a packet is delayed (*why?*)
    - packets are reordered (*why?*)
    - a packet is duplicated (*why?*)

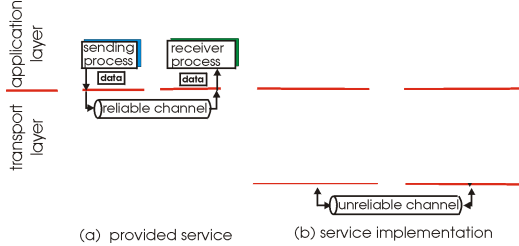


(a) provided service

31

## Principles of Reliable data transfer

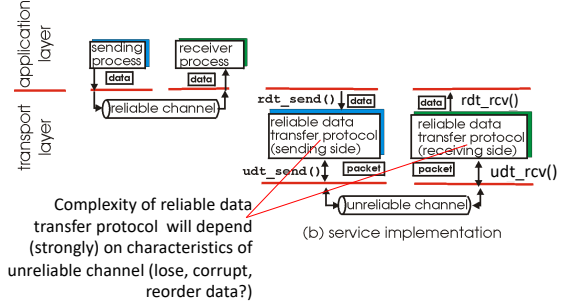
- important in app., transport, link layers
- top-10 list of important networking topics!



32

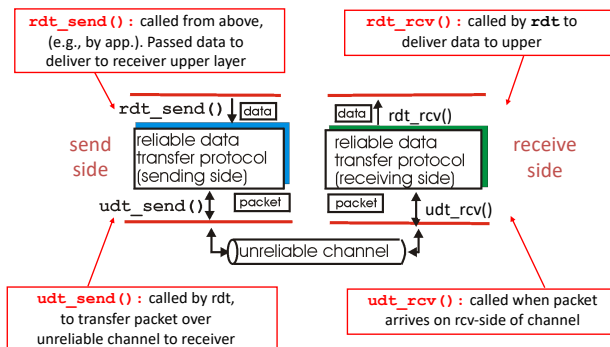
## Principles of Reliable data transfer

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



33

## Reliable data transfer: getting started

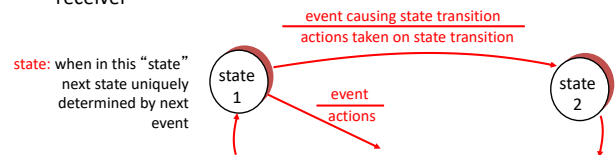


34

## 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



35

## KR state machines – a note.

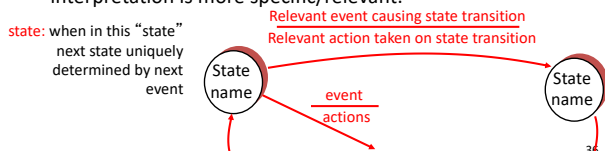
Beware

Kurose and Ross has a confusing/confused attitude to state-machines.

I've attempted to normalise the representation.

UPSHOT: these slides have differing information to the KR book (from which the RDT example is taken.)

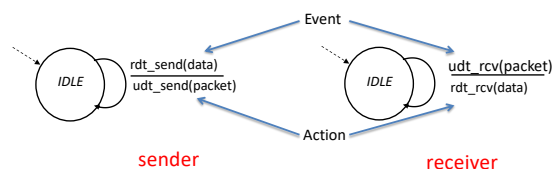
in KR "actions taken" appear wide-ranging, my interpretation is more specific/relevant.



36

## 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 read data from underlying channel



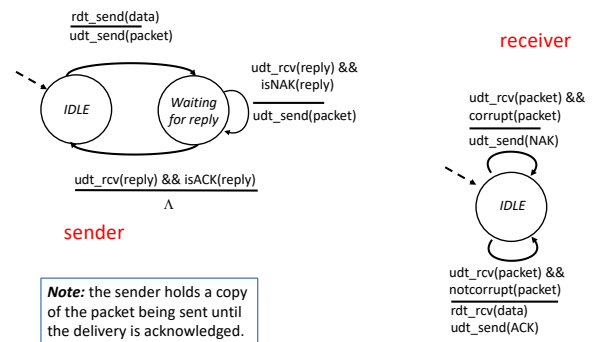
37

## 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 packet received is OK
  - negative acknowledgements (NAKs)**: receiver explicitly tells sender that packet had errors
    - sender retransmits packet on receipt of NAK
- new mechanisms in **rdt2.0** (beyond **rdt1.0**):
  - error detection
  - receiver feedback: control msgs (ACK,NAK) receiver->sender

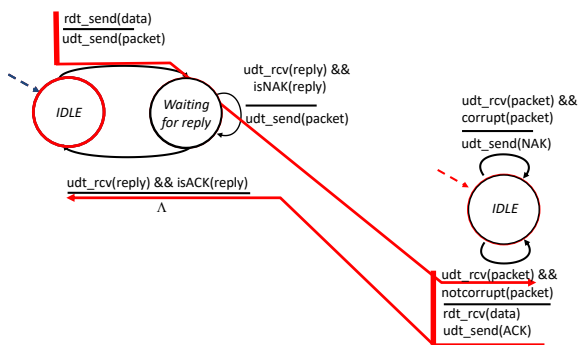
38

## rdt2.0: FSM specification



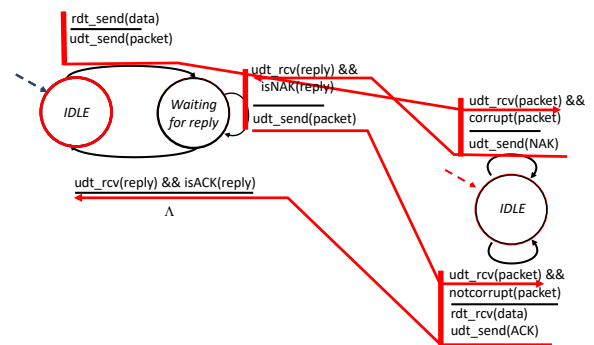
40

## rdt2.0: operation with no errors



41

## rdt2.0: error scenario



42

## 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

### Handling duplicates:

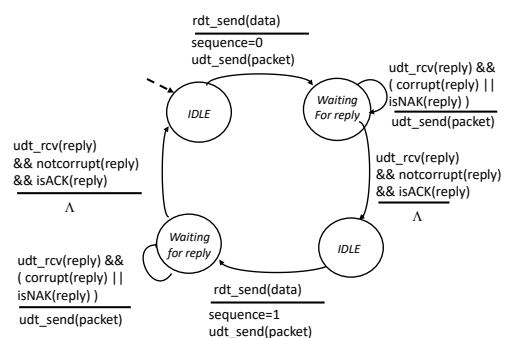
- sender retransmits current packet if ACK/NAK garbled
- sender adds **sequence number** to each packet
- receiver discards (doesn't deliver) duplicate packet

### stop and wait

Sender sends one packet, then waits for receiver response

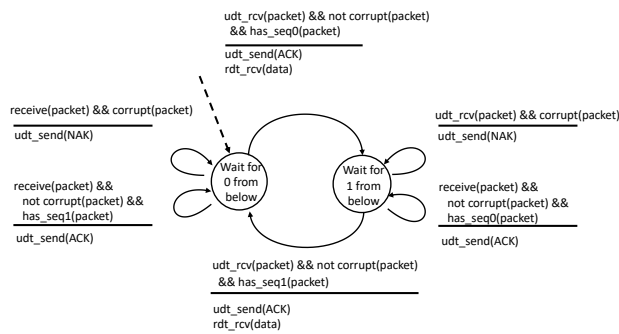
43

## rdt2.1: sender, handles garbled ACK/NAKs



45

## rdt2.1: receiver, handles garbled ACK/NAKs



46

## 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 "current" pkt has a 0 or 1 sequence number

### Receiver:

- 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

47

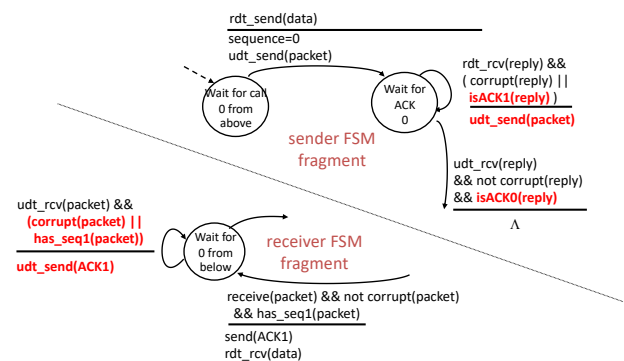
## rdt2.2: a NAK-free protocol

- same functionality as rdt2.1, 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*

As we will see, TCP uses this approach to be NAK-free

48

## rdt2.2: sender, receiver fragments



49

## rdt3.0: channels with errors *and* loss

**New channel assumption:** underlying channel can also *lose* packets (data, ACKs)

- checksum, sequence #'s, ACKs, retransmissions will be of help ... but not quite enough

**Q:** How do *humans* handle lost sender-to-receiver words in conversation?

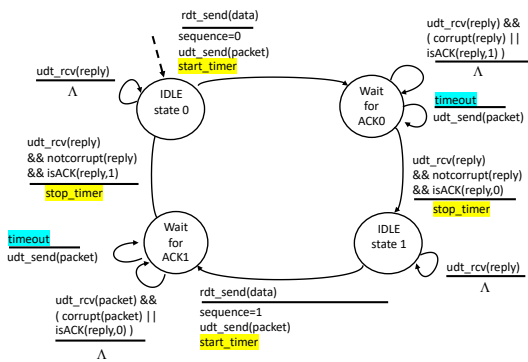
## rdt3.0: channels with errors *and* loss

**Approach:** sender waits "reasonable" amount of time for ACK

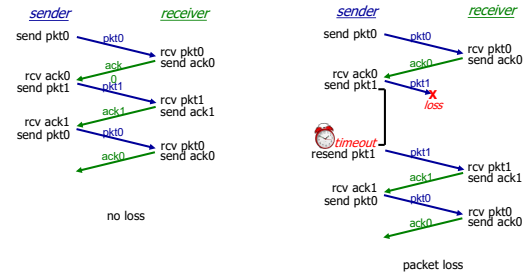
- retransmits if no ACK received in this time
- if pkt (or ACK) just delayed (not lost):
  - retransmission will be duplicate, but seq #'s already handles this!
  - receiver must specify seq # of packet being ACKed
- use countdown timer to interrupt after "reasonable" amount of time



## rdt3.0 sender

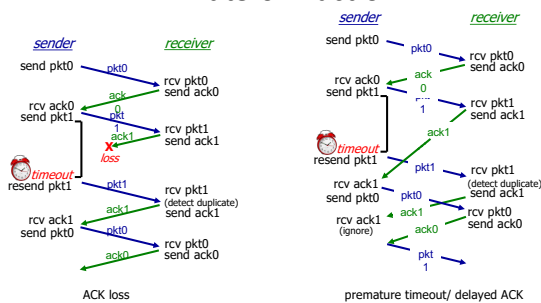


## rdt3.0 in action

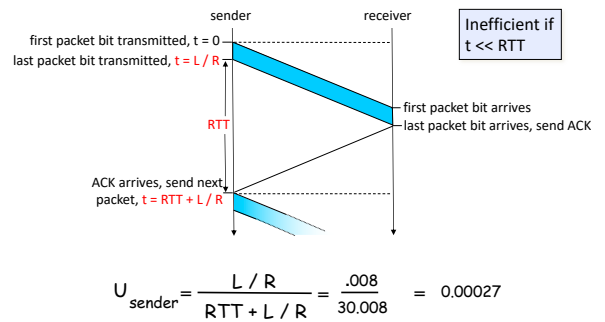


53

## rdt3.0 in action



## rdt3.0: stop-and-wait operation



60

## Performance of rdt3.0 (stop-and-wait)

- rdt3.0 works, but performance stinks
- ex: 1 Gbps link, 15 ms prop. delay, 8000 bit packet:

$$d_{\text{trans}} = \frac{L}{R} = \frac{8000 \text{ bits}}{10^9 \text{ bps}} = 8 \text{ microseconds}$$

- $U_{\text{sender}}$ : utilization – fraction of time sender busy sending

$$U_{\text{sender}} = \frac{L / R}{RTT + L / R} = \frac{.008}{30.008} = 0.00027$$

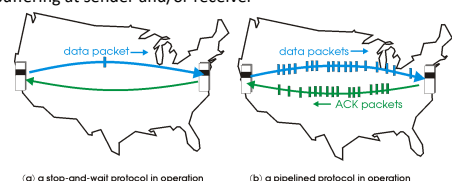
- 1KB pkt every 30 msec -> 33kB/sec throughput over 1 Gbps link
- The network protocol limits use of physical resources!

61

## Pipelined (Packet-Window) protocols

**Pipelining:** sender allows multiple, “in-flight”, yet-to-be-acknowledged pkts

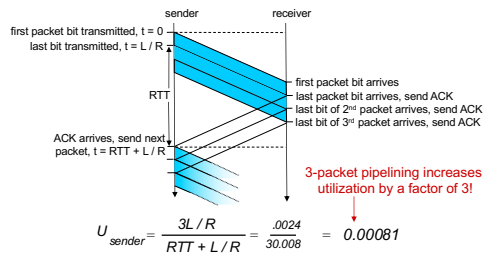
- range of sequence numbers must be increased
- buffering at sender and/or receiver



62



## Pipelining: increased utilization

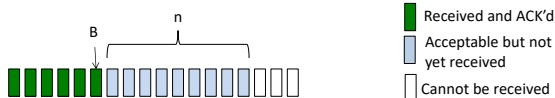


- **window** = set of adjacent sequence numbers
  - The size of the set is the **window size**; assume window size is  $n$
- General idea: send up to  $n$  packets at a time
  - Sender can send packets in its window
  - Receiver can accept packets in its window
  - Window of acceptable packets “slides” on successful reception/acknowledgement

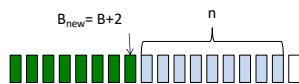
64

## Acknowledgements (1)

- At receiver



- After receiving B+1, B+2

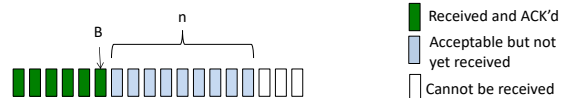


- Receiver sends ACK( $B_{\text{new}}+1$ )

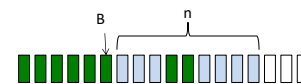
66

## Acknowledgements (2)

- At receiver



- After receiving B+4, B+5



- Receiver sends ACK(B+???)

**Oh....  
how do we  
recover?**

67

## Acknowledgements w/ Sliding Window

## Dealing with loss....

- Two common options
  - Go-Back-N (GBN)
  - Selective Repeat (SR)
    - Also called Selective Acknowledgement (SACK)

68

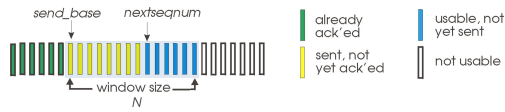
## Go-Back-N (GBN)

- Sender transmits up to  $n$  unacknowledged packets
- Receiver only accepts packets in order
  - discards out-of-order packets (i.e., packets other than  $B+1$ )
- Receiver uses **cumulative acknowledgements**
  - i.e., sequence# in ACK = next expected in-order sequence#
- Sender sets timer for 1<sup>st</sup> outstanding ack ( $A+1$ )
- If timeout, retransmit  $A+1, \dots, A+n$

69

Go-Back-N: sender

- sender: “window” of up to N, consecutive transmitted but unACKed pkts
  - k-bit seq # in pkt header



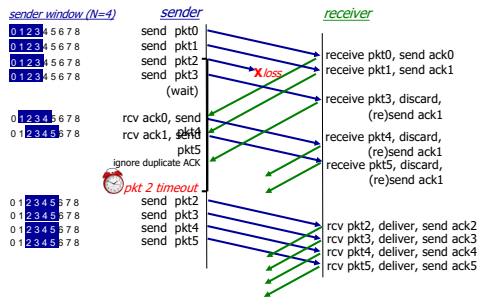
- **cumulative ACK:** ACK( $n$ ): ACKs all packets up to, including seq #  $n$ 
  - on receiving ACK( $n$ ): move window forward to begin at  $n+1$
- timer for oldest in-flight packet
- **timeout( $n$ ):** retransmit packet  $n$  and all higher seq # packets in window

### Go-Back-N: receiver

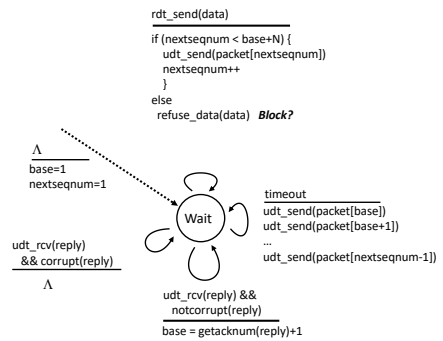
- ACK-only: always send ACK for correctly-received packet so far, with highest *in-order* seq #
  - may generate duplicate ACKs
  - need only remember *rcv\_base*
- on receipt of out-of-order packet:
  - can discard (don't buffer) or buffer: an implementation decision
  - re-ACK pkt with highest in-order seq #



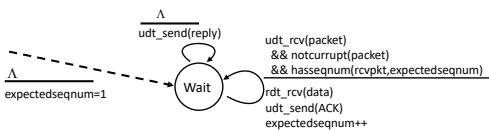
## Go-Back-N in action



## GBN: sender extended FSM



### GBN: receiver extended FSM



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

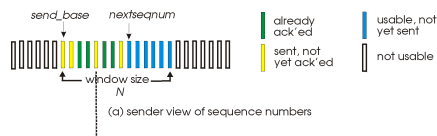
- may generate duplicate ACKs
- need only remember **expectedseqnum**
- out-of-order packet:
  - discard (don't buffer) -> **no receiver buffering!**
  - Re-ACK packet with highest in-order seq #

## Selective repeat

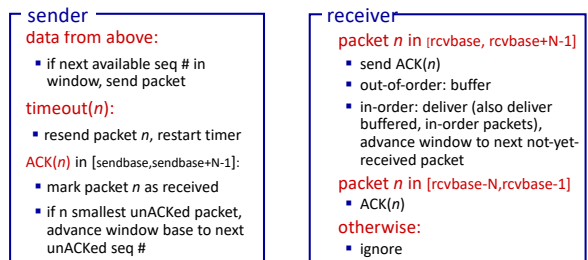
- receiver *individually* acknowledges all correctly received packets
  - buffers packets, as needed, for eventual in-order delivery to upper layer
- sender times-out/retransmits individually for unACKed packets
  - sender maintains timer for each unACKed pkt
- sender window
  - $N$  consecutive seq #s
  - limits seq #s of sent, unACKed packets

This is also known as Selective Acknowledgement or simply SACK

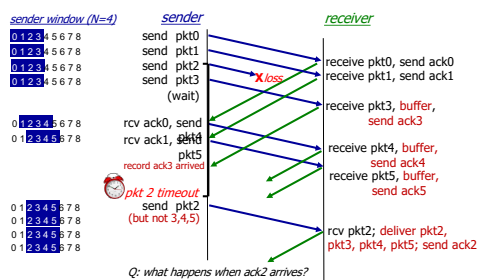
## Selective repeat: sender, receiver windows



## Selective repeat: sender and receiver

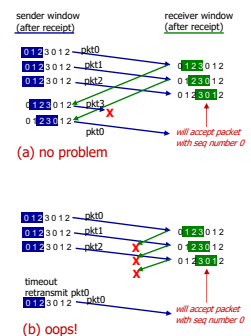


## Selective Repeat in action



## Selective repeat: a dilemma!

- example:
- seq #: 0, 1, 2, 3 (base 4 counting)
  - window size=3



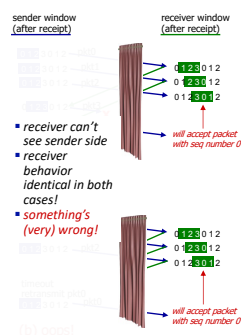
## Selective repeat: a dilemma!

- example:
- seq #: 0, 1, 2, 3 (base 4 counting)
  - window size=3

Q: what relationship is needed between sequence # size and window size to avoid problem in scenario (b)?

Solution:

maximum allowable window size = half the sequence number space.



## Observations

- With sliding windows, it is possible to fully utilize a link, provided the window size ( $n$ ) is large enough. Throughput is  $\sim (n/RTT)$ 
  - Stop & Wait is like  $n = 1$ .
- Sender has to buffer all unacknowledged packets, because they may require retransmission
- Receiver may be able to accept out-of-order packets, but only up to its buffer limits
- Implementation complexity depends on protocol details (GBN vs. SR)

## Recap: components of a solution

- Checksums (for error detection)
- Timers (for loss detection)
- Acknowledgments
  - cumulative
  - selective
- Sequence numbers (duplicates, windows)
- Sliding Windows (for efficiency)
- Reliability protocols use the above to decide when and what to retransmit or acknowledge

91

## What does TCP do?

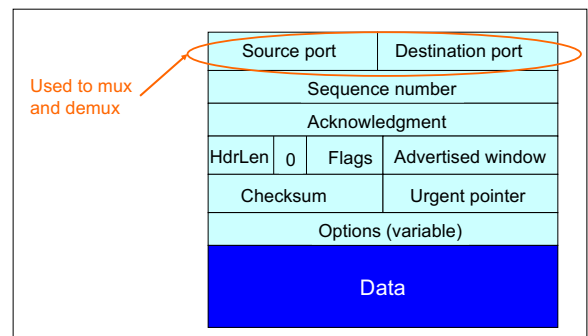
Most of our previous tricks + a few more beside

- Sequence numbers are byte offsets
- Sender and receiver maintain a sliding window
- Receiver sends cumulative acknowledgements (like GBN)
- Sender maintains a single retransmission timer
- Receivers do not drop out-of-sequence packets (like SR)
- Introduces **fast retransmit**: optimization that uses duplicate ACKs to trigger early retransmission
- Introduces timeout estimation algorithms

### TCP: overview RFCs: 793,1122, 2018, 5681, 7323

- **point-to-point:**
  - one sender, one receiver
- **reliable, in-order byte stream:**
  - no "message boundaries"
- **full duplex data:**
  - bi-directional data flow in same connection
  - MSS: maximum segment size
- **cumulative ACKs**
- **pipelining:**
  - TCP congestion and flow control set window size
- **connection-oriented:**
  - handshaking (exchange of control messages) initializes sender, receiver state before data exchange
- **flow controlled:**
  - sender will not overwhelm receiver

## TCP Header



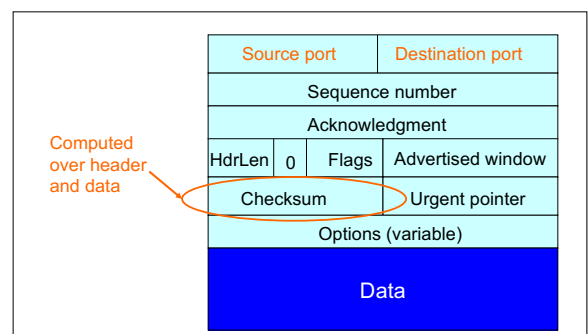
95

## What does TCP do?

Many of our previous ideas, but some key differences

- Checksum

## TCP Header



97

98

## What does TCP do?

Many of our previous ideas, but some key differences

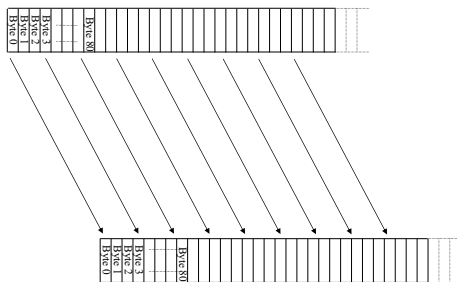
- Checksum
- Sequence numbers are byte offsets

## TCP: Segments and Sequence Numbers

100

### TCP “Stream of Bytes” Service...

Application @ Host A

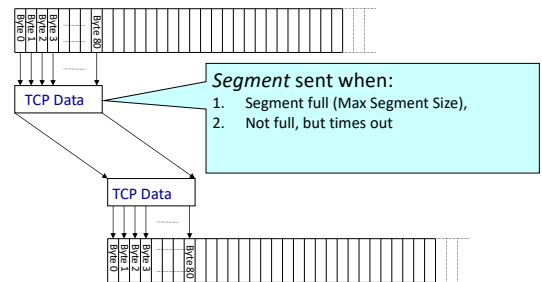


Application @ Host B

101

### ... Provided Using TCP “Segments”

Host A



Host B

102

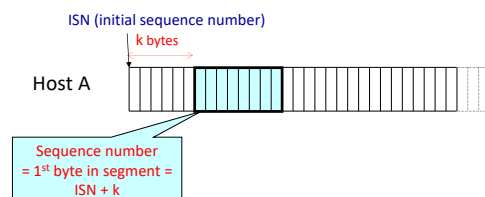
## TCP Segment



- IP packet
  - No bigger than Maximum Transmission Unit (MTU)
  - E.g., up to 1500 bytes with Ethernet
- TCP packet
  - IP packet with a TCP header and data inside
  - TCP header  $\geq 20$  bytes long
- TCP **segment**
  - No more than Maximum Segment Size (MSS) bytes
  - E.g., up to 1460 consecutive bytes from the stream
  - $MSS = MTU - (IP\ header) - (TCP\ header)$

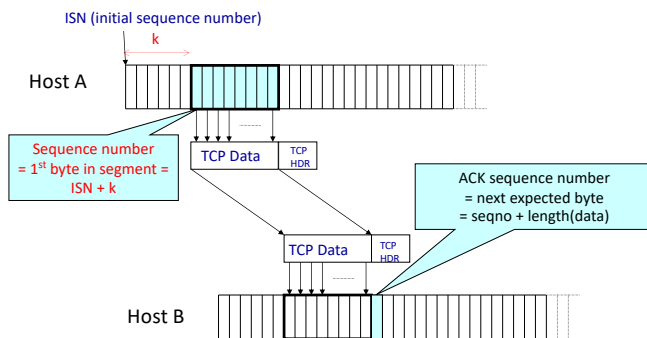
103

## Sequence Numbers



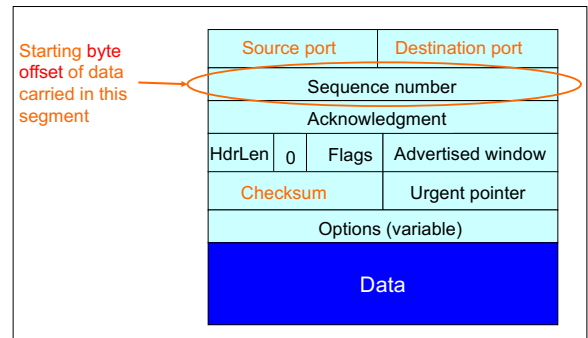
104

## Sequence Numbers



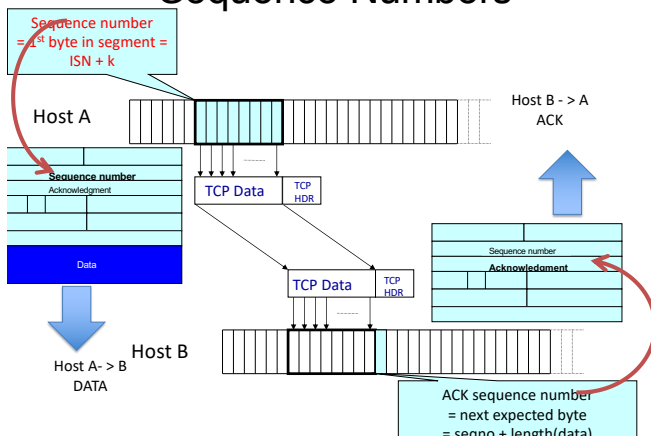
105

## TCP Header



106

## Sequence Numbers



## TCP Sequences and ACKS

TCP is full duplex by default

- two independently flows of sequence numbers

Sequence acknowledgement is given in terms of BYTES (not packets); the window is in terms of bytes.

number of packets = window size (bytes) / Segment Size

Servers and Clients are not Source and Destination

Piggybacking increases efficiency but many flows may only have data moving in one direction

108

## What does TCP do?

Most of our previous tricks, but a few differences

- Checksum
- Sequence numbers are byte offsets
- Receiver sends cumulative acknowledgements (like GBN)

## ACKing and Sequence Numbers

- Sender sends packet
  - Data starts with sequence number X
  - Packet contains B bytes [X, X+1, X+2, ..., X+B-1]
- Upon receipt of packet, receiver sends an ACK
  - If all data prior to X already received:
    - ACK acknowledges X+B (because that is next expected byte)
  - If highest in-order byte received is Y s.t. (Y+1) < X
    - ACK acknowledges Y+1
    - Even if this has been ACKed before

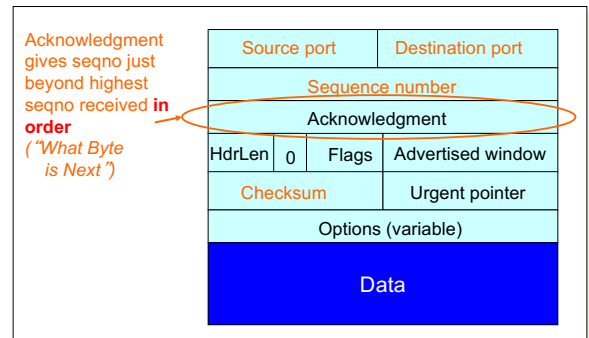
110

## Normal Pattern

- Sender: seqno=X, length=B
- Receiver: ACK=X+B
- Sender: seqno=X+B, length=B
- Receiver: ACK=X+2B
- Sender: seqno=X+2B, length=B
- Seqno of next packet is same as last ACK field

111

## TCP Header



112

## What does TCP do?

Most of our previous tricks, but a few differences

- Checksum
- Sequence numbers are byte offsets
- Receiver sends cumulative acknowledgements (like GBN)
- Receivers **can** buffer out-of-sequence packets (like SR)

113

## Loss with cumulative ACKs

- Sender sends packets with 100B and seqnos.:  
– 100, 200, 300, 400, 500, 600, 700, 800, 900, ...
- Assume the fifth packet (seqno 500) is lost, but no others
- Stream of ACKs will be:  
– 200, 300, 400, 500, 500, 500, 500, ...

114

## What does TCP do?

Most of our previous tricks, but a few differences

- Checksum
- Sequence numbers are byte offsets
- Receiver sends cumulative acknowledgements (like GBN)
- Receivers may not drop out-of-sequence packets (like SR)
- Introduces **fast retransmit**: optimization that uses duplicate ACKs to trigger early retransmission

115

## Loss with cumulative ACKs

- “Duplicate ACKs” are a sign of an isolated loss  
– The lack of ACK progress means 500 hasn’t been delivered  
– Stream of ACKs means some packets are being delivered
- Therefore, could trigger resend upon receiving k duplicate ACKs  
• TCP uses k=3
- But response to loss is trickier....

116

## Loss with cumulative ACKs

- Two choices:
  - Send missing packet and increase W by the number of dup ACKs
  - Send missing packet, and wait for ACK to increase W
- Which should TCP do?

117

## What does TCP do?

Most of our previous tricks, but a few differences

- Checksum
- Sequence numbers are byte offsets
- Receiver sends cumulative acknowledgements (like GBN)
- Receivers do not drop out-of-sequence packets (like SR)
- Introduces fast retransmit: optimization that uses duplicate ACKs to trigger early retransmission
- Sender maintains a single retransmission timer (like GBN) and retransmits on timeout

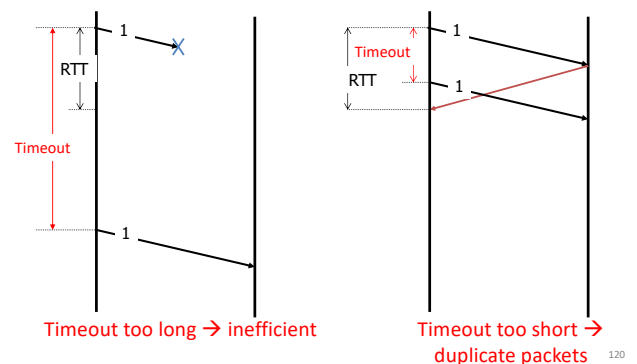
118

## Retransmission Timeout

- If the sender hasn't received an ACK by timeout, retransmit the first packet in the window
- How do we pick a timeout value?

119

## Timing Illustration



120

## Retransmission Timeout

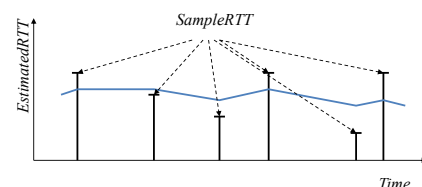
- If haven't received ack by timeout, retransmit the first packet in the window
- How to set timeout?
  - Too long: connection has low throughput
  - Too short: retransmit packet that was just delayed
- Solution: make timeout proportional to RTT
- But how do we measure RTT?

121

## RTT Estimation

- Use exponential averaging of RTT samples

$$\begin{aligned} \text{SampleRTT} &= \text{AckRcvdTime} - \text{SendPacketTime} \\ \text{EstimatedRTT} &= \alpha \times \text{EstimatedRTT} + (1 - \alpha) \times \text{SampleRTT} \\ 0 < \alpha &\leq 1 \end{aligned}$$



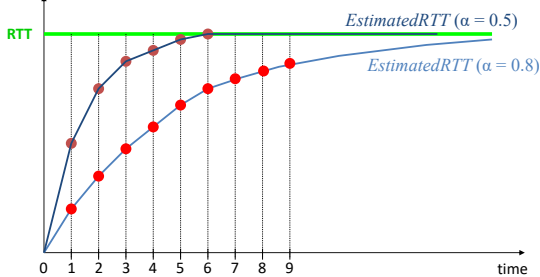
122



## Exponential Averaging Example

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

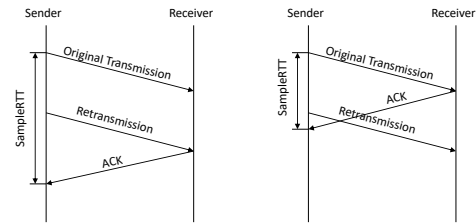
Assume RTT is constant  $\rightarrow \text{SampleRTT} = \text{RTT}$



123

## Problem: Ambiguous Measurements

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



124

## Karn/Partridge Algorithm

Discard junk measures

- Measure *SampleRTT* only for original transmissions
  - Once a segment has been retransmitted, do not use it for any further measurements
- Computes EstimatedRTT using  $\alpha = 0.875$
- Timeout value (RTO) =  $2 \times \text{EstimatedRTT}$
- Employs **exponential backoff**
  - Every time RTO timer expires, set  $\text{RTO} \leftarrow 2 \cdot \text{RTO}$
  - (Up to maximum  $\geq 60$  sec)
  - Every time new measurement comes in (= successful original transmission), collapse RTO back to  $2 \times \text{EstimatedRTT}$

125

## Jacobson/Karels Algorithm

Add a safety margin

- Problem: need to better capture variability in RTT
  - Directly measure **deviation**
- Deviation =  $|\text{SampleRTT} - \text{EstimatedRTT}|$
- EstimatedDeviation: exponential average of Deviation
- $\text{RTO} = \text{EstimatedRTT} + 4 \times \text{EstimatedDeviation}$

126

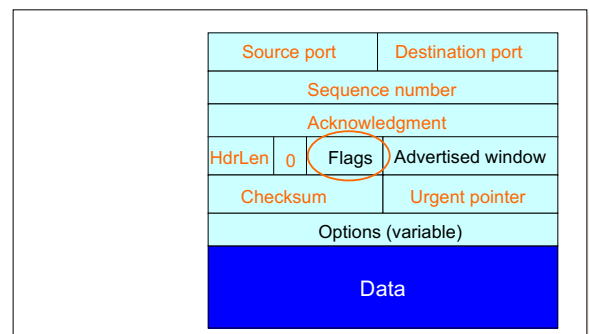
## What does TCP do?

Most of our previous ideas, but some key differences

- Checksum
- Sequence numbers are byte offsets
- Receiver sends cumulative acknowledgements (like GBN)
- Receivers do not drop out-of-sequence packets (like SR)
- Introduces fast retransmit: optimization that uses duplicate ACKs to trigger early retransmission
- Sender maintains a single retransmission timer (like GBN) and retransmits on timeout

127

## TCP Header: What's left?



128

## TCP Connection Establishment and Initial Sequence Numbers

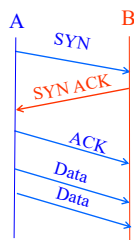
## Initial Sequence Number (ISN)

- Sequence number for the very first byte
- Why not just use ISN = 0?
- Practical issue
  - IP addresses and port #s uniquely identify a connection
  - Eventually, though, these port #s do get **used again**
  - ... small chance an old packet is **still in flight**
- TCP therefore **requires** changing ISN
- Hosts exchange ISNs when they establish a connection

129

130

## Establishing a TCP Connection

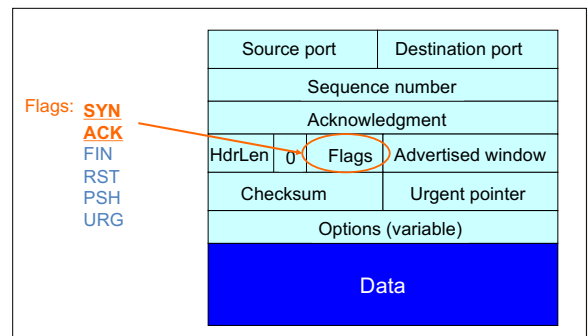


Each host tells its ISN to the other host.

- Three-way handshake to establish connection
  - Host A sends a **SYN** (open; "synchronize sequence numbers") to host B
  - Host B returns a SYN acknowledgment (**SYN ACK**)
  - Host A sends an **ACK** to acknowledge the SYN ACK

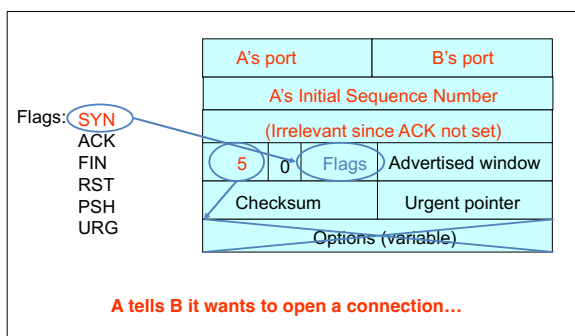
131

## TCP Header



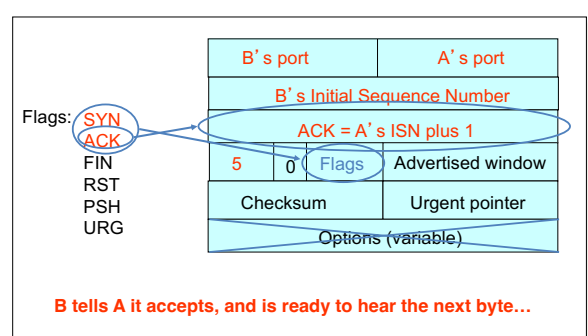
132

## Step 1: A's Initial SYN Packet



133

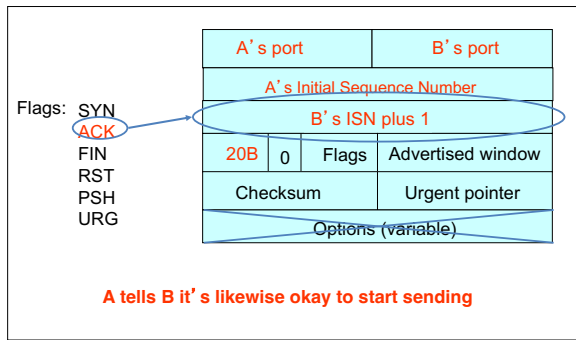
## Step 2: B's SYN-ACK Packet



... upon receiving this packet, A can start sending data

134

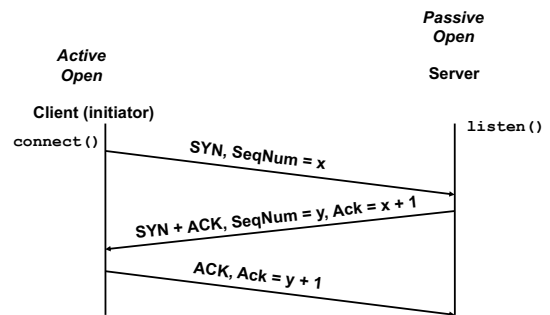
## Step 3: A's ACK of the SYN-ACK



... upon receiving this packet, B can start sending data

135

## Timing Diagram: 3-Way Handshaking



136

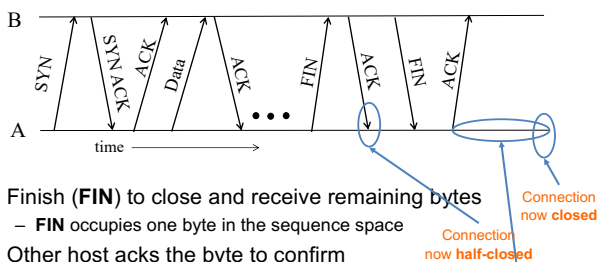
## What if the SYN Packet Gets Lost?

- Suppose the SYN packet gets lost
  - Packet is lost inside the network, or:
  - Server **discards** the packet (e.g., it's too busy)
- Eventually, no SYN-ACK arrives
  - Sender sets a **timer** and **waits** for the SYN-ACK
  - ... and retransmits the SYN if needed
- How should the TCP sender set the timer?
  - Sender has **no idea** how far away the receiver is
  - Hard to guess a reasonable length of time to wait
  - SHOULD** (RFCs 1122 & 2988) use default of **3 seconds**
    - Some implementations instead use 6 seconds

137

## Tearing Down the Connection

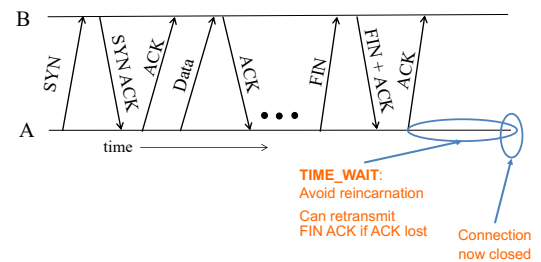
### Normal Termination, One Side At A Time



- Finish (**FIN**) to close and receive remaining bytes
  - FIN** occupies one byte in the sequence space
- Other host acks the byte to confirm
- Closes A's side of the connection, but **not B's**
  - Until B likewise sends a **FIN**
  - Which A then acks

139

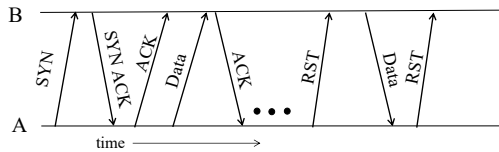
### Normal Termination, Both Together



- Same as before, but B sets **FIN** with their ack of A's **FIN**

140

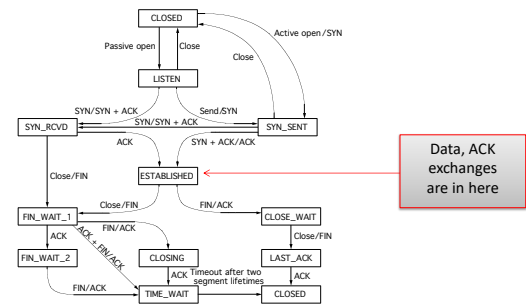
## Abrupt Termination



- A sends a RESET (**RST**) to B
  - E.g., because application process on A **crashed**
- **That's it**
  - B does **not** ack the **RST**
  - Thus, **RST** is **not** delivered **reliably**
  - And: any data in flight is **lost**
  - But: if B sends anything more, will elicit **another RST**

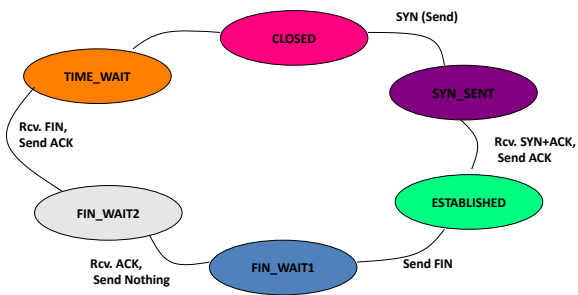
141

## TCP State Transitions



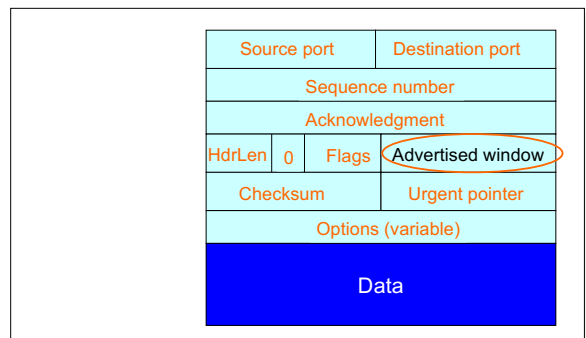
142

## An Simpler View of the Client Side



143

## TCP Header

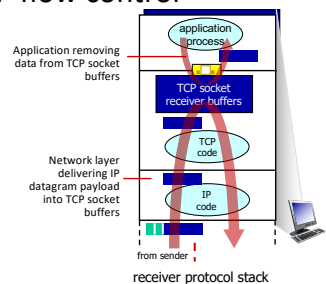


144

- What does TCP do?
  - ARQ windowing, set-up, tear-down
- Flow Control in TCP

## TCP flow control

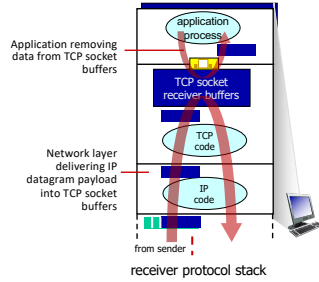
**Q:** What happens if network layer delivers data faster than application layer removes data from socket buffers?



145

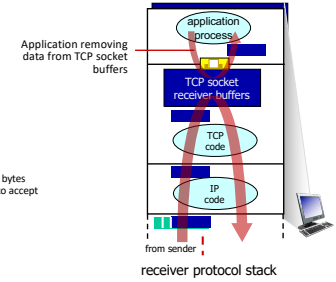
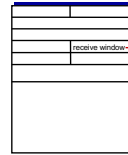
## TCP flow control

Q: What happens if network layer delivers data faster than application layer removes data from socket buffers?



## TCP flow control

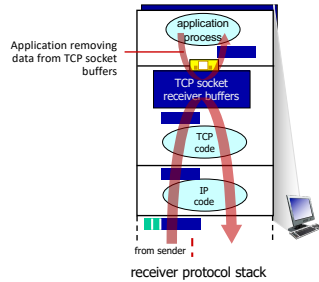
Q: What happens if network layer delivers data faster than application layer removes data from socket buffers?



## TCP flow control

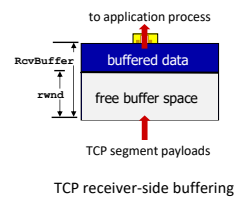
Q: What happens if network layer delivers data faster than application layer removes data from socket buffers?

**flow control**  
receiver controls sender, so sender won't overflow receiver's buffer by transmitting too much, too fast



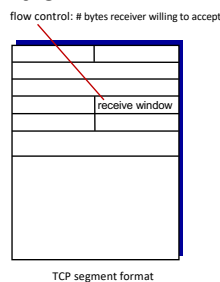
## TCP flow control

- TCP receiver "advertises" free buffer space in **rwnd** field in TCP header
  - RcvBuffer** size set via socket options (typical default is 4096 bytes)
  - many operating systems autoadjust **RcvBuffer**
- sender limits amount of unACKed ("in-flight") data to received **rwnd**
- guarantees receive buffer will not overflow



## TCP flow control

- TCP receiver "advertises" free buffer space in **rwnd** field in TCP header
  - RcvBuffer** size set via socket options (typical default is 4096 bytes)
  - many operating systems autoadjust **RcvBuffer**
- sender limits amount of unACKed ("in-flight") data to received **rwnd**
- guarantees receive buffer will not overflow



## Advertised Window Limits Rate

- Sender can send no faster than  $W/RTT$  bytes/sec
- Receiver only advertises more space when it has consumed old arriving data
- In original TCP design, that was the **sole** protocol mechanism controlling sender's rate
- What's missing?

- What does TCP do?
  - ARQ windowing, set-up, tear-down
- Flow Control in TCP
- Congestion Control in TCP

## We have seen:

- **Flow control**: adjusting the sending rate to keep from overwhelming a slow receiver

## Now lets attend...

- **Congestion control**: adjusting the sending rate to keep from overloading the network

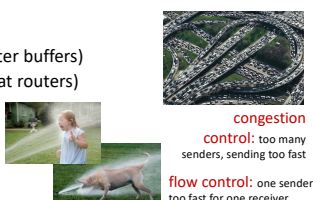
163

164

## Principles of congestion control

### Congestion:

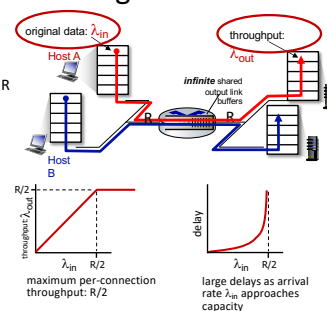
- informally: “too many sources sending too much data too fast for *network* to handle”
- manifestations:
  - long delays (queueing in router buffers)
  - packet loss (buffer overflow at routers)
- different from flow control!
- a top-10 problem!



## Causes/costs of congestion: scenario 1

### Simplest scenario:

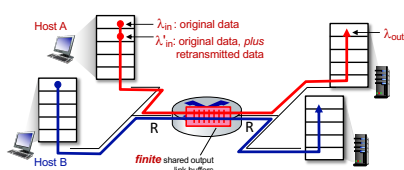
- one router, infinite buffers
- input, output link capacity:  $R$
- two flows
- no retransmissions needed



Q: What happens as arrival rate  $\lambda_{in}$  approaches  $R/2$ ?

## Causes/costs of congestion: scenario 2

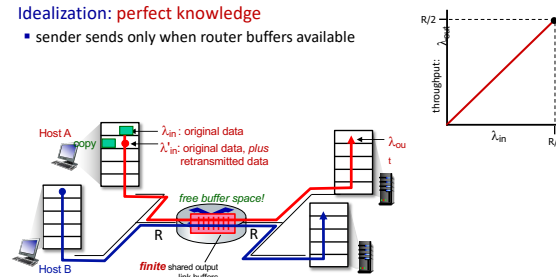
- one router, *finite* buffers
- sender retransmits lost, timed-out packet
  - application-layer input = application-layer output:  $\lambda_{in} = \lambda_{out}$
  - transport-layer input includes *retransmissions*:  $\lambda'_{in} \geq \lambda_{in}$



## Causes/costs of congestion: scenario 2

### Idealization: perfect knowledge

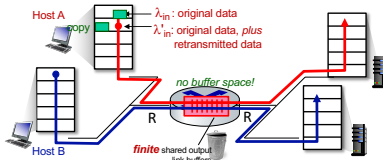
- sender sends only when router buffers available



## Causes/costs of congestion: scenario 2

**Idealization: *some* perfect knowledge**

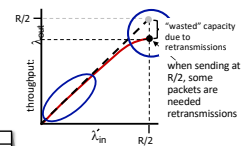
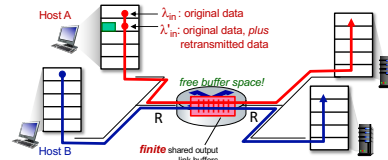
- packets can be lost (dropped at router) due to full buffers
- sender knows when packet has been dropped: only resends if packet *known* to be lost



## Causes/costs of congestion: scenario 2

**Idealization: *some* perfect knowledge**

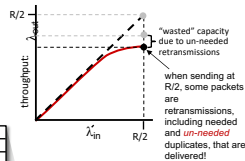
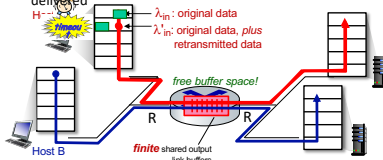
- packets can be lost (dropped at router) due to full buffers
- sender knows when packet has been dropped: only resends if packet *known* to be lost



## Causes/costs of congestion: scenario 2

**Realistic scenario: *un-needed duplicates***

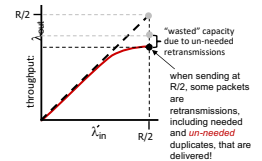
- packets can be lost, dropped at router due to full buffers – requiring retransmissions
- but sender times can time out prematurely, sending *two* copies, *both* of which are delivered



## Causes/costs of congestion: scenario 2

**Realistic scenario: *un-needed duplicates***

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



**“costs” of congestion:**

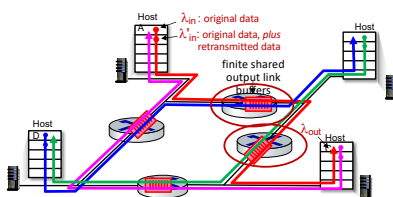
- more work (retransmission) for given receiver throughput
- unnecessary retransmissions: link carries multiple copies of a packet
  - decreasing maximum achievable throughput

## Causes/costs of congestion: scenario 3

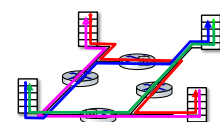
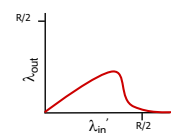
- four senders
- multi-hop paths
- timeout/retransmit

**Q:** what happens as  $\lambda_{in}$  and  $\lambda_{in}'$  increase?

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



## Causes/costs of congestion: scenario 3

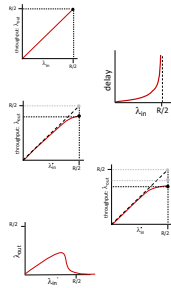


**another “cost” of congestion:**

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

## Causes/costs of congestion: insights

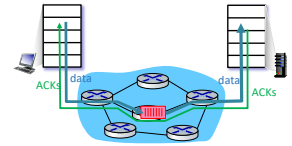
- throughput can never exceed capacity
- delay increases as capacity approached
- loss/retransmission decreases effective throughput
- un-needed duplicates further decreases effective throughput
- upstream transmission capacity / buffering wasted for packets lost downstream



## Approaches towards congestion control

### End-end congestion control:

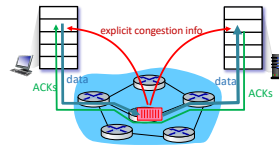
- no explicit feedback from network
- congestion *inferred* from observed loss, delay
- approach taken by TCP



## Approaches towards congestion control

### Network-assisted congestion control:

- routers provide *direct* feedback to sending/receiving hosts with flows passing through congested router
- may indicate congestion level or explicitly set sending rate
- TCP ECN, ATM, DECbit protocols

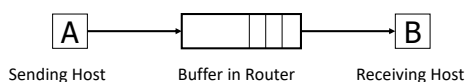


## Three Issues to Consider

- Discovering the available (bottleneck) bandwidth
- Adjusting to variations in bandwidth
- Sharing bandwidth between flows

178

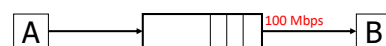
## Abstract View



- Ignore internal structure of router and model it as having a single queue for a particular input-output pair

179

## Discovering available bandwidth

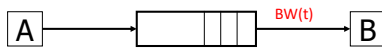


- Pick sending rate to match bottleneck bandwidth
  - Without any *a priori* knowledge
  - Could be gigabit link, could be a modem

180



## Adjusting to variations in bandwidth



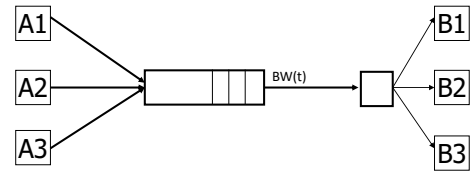
- Adjust rate to match **instantaneous** bandwidth
  - Assuming you have rough idea of bandwidth

181

## Multiple flows and sharing bandwidth

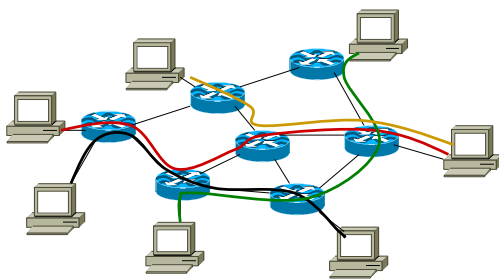
Two Issues:

- Adjust total sending rate to match bandwidth
- Allocation of bandwidth between flows



182

## Reality

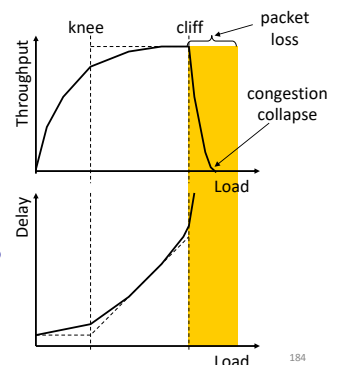


Congestion control is a resource allocation problem involving many flows, many links, and complicated global dynamics

183

## View from a single flow

- Knee – point after which
  - Throughput increases slowly
  - Delay increases fast
- Cliff – point after which
  - Throughput starts to drop to zero (congestion collapse)
  - Delay approaches infinity



184

## General Approaches

- (0) Send without care
  - Many packet drops

185

## General Approaches

- (0) Send without care
- (1) Reservations
  - Pre-arrange bandwidth allocations
  - Requires negotiation before sending packets
  - Low utilization

186

## General Approaches

- (0) Send without care
- (1) Reservations
- (2) Pricing
  - Don't drop packets for the high-bidders
  - Requires payment model

187

## General Approaches

- (0) Send without care
- (1) Reservations
- (2) Pricing
- (3) Dynamic Adjustment
  - Hosts probe network; infer level of congestion; adjust
  - Network reports congestion level to hosts; hosts adjust
  - Combinations of the above
  - Simple to implement but suboptimal, messy dynamics

188

## General Approaches

- (0) Send without care
- (1) Reservations
- (2) Pricing
- (3) Dynamic Adjustment

### All three techniques have their place

- *Generality* of dynamic adjustment has proven powerful
  - Packet drops → senders (repeatedly!) retransmit a full window's worth of packets
- Doesn't presume business model, traffic characteristics, application requirements; does assume good citizenship

189

## Who Takes Care of Congestion?

- Network? End hosts? Both?
- TCP's approach:
  - **End hosts** adjust sending rate
  - Based on **implicit feedback** from network
- Not the only approach
  - A consequence of history rather than planning

190

## Some History: TCP in the 1980s

- Sending rate only limited by flow control
  - Packet drops → senders (repeatedly!) retransmit a full window's worth of packets
- Led to "congestion collapse" starting Oct. 1986
  - Throughput on the NSF network dropped from 32Kbits/s to 40bits/sec
- "Fixed" by Van Jacobson's development of TCP's congestion control (CC) algorithms

191

## Jacobson's Approach

- Extend TCP's existing window-based protocol but adapt the window size in response to congestion
  - **required no upgrades to routers or applications!**
  - patch of a few lines of code to TCP implementations
- A pragmatic and effective solution
  - but many other approaches exist
- Extensively improved on since
  - topic now sees less activity in ISP contexts
  - but is making a comeback in datacenter environments

192

## TCP's Approach in a Nutshell

- TCP connection has window
  - Controls number of packets in flight
- Sending rate:  $\sim \text{Window} / \text{RTT}$
- Vary window size to control sending rate

193

## Windows, Buffers, and TCP



194

## Windows, Buffers, and TCP

- TCP connection has a window
  - Controls number of packets in flight; filling a channel to improve throughput, and vary window size to control sending rate
- Buffers adapt mis-matched channels
  - Buffers smooth bursts
  - Adapt (re-time) arrivals for multiplexing

195

## Windows, Buffers, and TCP

Buffers & TCP can make link utilization 100%

but

Buffers add delay, **variable** delay



196

## Sizing Buffers in Routers



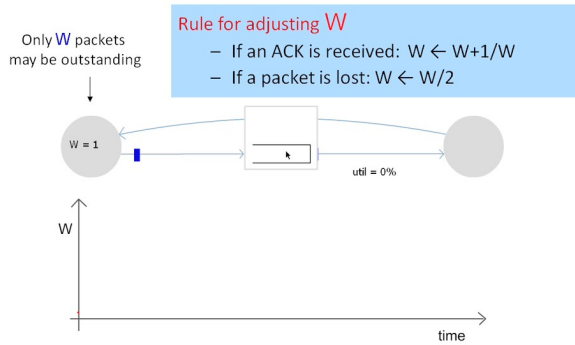
- Packet loss
  - Queue overload, and subsequent packet loss
- End-to-end delay
  - Transmission, propagation, and queueing delay
  - The only variable part is queueing delay
- Router architecture
  - Board space, power consumption, and cost
  - On chip buffers: higher density, higher capacity

197

## Buffer Sizing Story

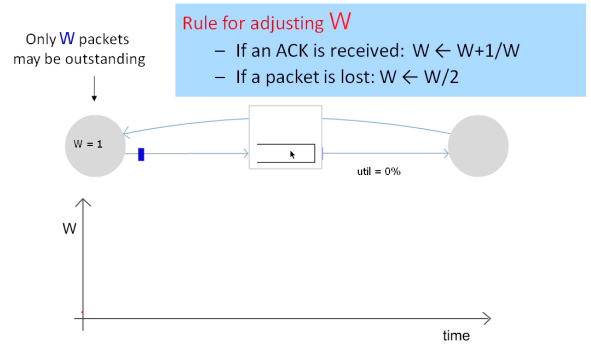
Rule-of-thumb	
# of packets	$2T \times C$
Intuition	1,000,000
Assume	TCP Sawtooth
Evidence	Single TCP Flow, 100% Utilization
	Simulation, Emulation

## Continuous ARQ (TCP) adapting to congestion



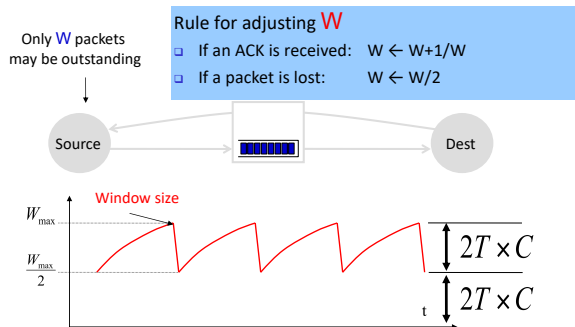
199

## Continuous ARQ (TCP) adapting to congestion



200

## Rule-of-thumb – Intuition



201

## Buffers in Routers

So how large should the buffers be?

### Buffer size matters

- Packet loss
  - Queue overload, and subsequent packet loss
- End-to-end delay
  - Transmission, propagation, and queueing delay
  - The only variable part is queueing delay

202

## Buffer Sizing Story

	Rule-of-thumb	Small Buffers
# of packets	$2T \times C$	$\frac{2T \times C}{\sqrt{n}}$
Intuition	1,000,000	10,000
Assume	TCP Sawtooth	Sawtooth Smoothing
Evidence	Single TCP Flow, 100% Utilization	Many Flows, 100% Utilization
	Simulation, Emulation	Simulations, Test-bed and Real Network Experiments

## Buffers in Routers

So how large should the buffers be?

### Buffer size matters

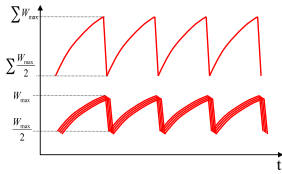
- Packet loss
  - Queue overload, and subsequent packet loss
- End-to-end delay
  - Transmission, propagation, and queueing delay
  - The only variable part is queueing delay
- Router architecture
  - Board space, power consumption, and cost
  - On chip buffers: higher density, higher capacity

204

## Small Buffers – Intuition

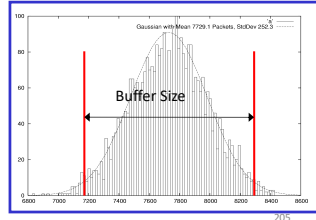
### Synchronized Flows

- Aggregate window has same dynamics
- Therefore buffer occupancy has same dynamics
- Rule-of-thumb still holds.



### Many TCP Flows

- Independent, desynchronized
- Central limit theorem says the aggregate becomes Gaussian
- Variance (buffer size) decreases as N increases



### Buffer Sizing Story

What size do we make the buffer? Well it depends...

Rule-of-thumb:  $2T \times C$

Intuition: One TCP connection? Many Synchronized TCP connections? Just TCP – what about other applications? Small BDP link? Large BDP link? How many devices? W of flows? How many flows?

Assume: Single TCP flow, 100% Utilization

Evidence: Simulation, Emulation, Network Experiments

How much do you know about your traffic? What is best for your traffic?

Equations:  $2T \times C$ ,  $\frac{2T \times C}{\sqrt{n}}$ ,  $O(\log W)$

Other terms: 10,000, 100 - 50, bursty arrivals, 85-90% Utilization, TCP, Simulations, Test-beds, Experiments



## TCP's Approach in a Nutshell

- TCP connection has window
  - Controls number of packets in flight
- Sending rate:  $\sim \text{Window} / \text{RTT}$
- Vary window size to control sending rate

207

## All These Windows...

- Congestion Window: **CWND**
  - How many bytes can be sent without overflowing routers
  - Computed by the sender using congestion control algorithm
- Flow control window: **AdvertisedWindow (RWND)**
  - How many bytes can be sent without overflowing receiver's buffers
  - Determined by the receiver and reported to the sender
- Sender-side window = **minimum{CWND, RWND}**
  - Assume for this material that  $\text{RWND} \gg \text{CWND}$

208

## Note

- This lecture will talk about CWND in units of MSS
  - (Recall MSS: Maximum Segment Size, the amount of payload data in a TCP packet)
  - This is only for pedagogical purposes
- In reality this is a **SIMPLIFICATION**:
  - Real implementations maintain CWND in bytes

209

## Two Basic Questions

- How does the sender detect congestion?
- How does the sender adjust its sending rate?
  - To address three issues
    - Finding available bottleneck bandwidth
    - Adjusting to bandwidth variations
    - Sharing bandwidth

210

## (Recall) Detecting Congestion

- Packet delays
  - Tricky: noisy signal (delay often varies considerably)
- Router tell end-hosts they're congested
- Packet loss
  - Fail-safe signal that TCP already has to detect
  - Complication: non-congestive loss (checksum errors)
- Two indicators of packet loss
  - No ACK after certain time interval: **timeout**
  - Multiple **duplicate ACKs**

211

## Not All Losses the Same

- Duplicate ACKs: isolated loss
  - Still getting ACKs
- Timeout: much more serious
  - Not enough packets in progress to trigger duplicate-acks, OR
  - Suffered several losses
- We will adjust rate differently for each case

212

## Rate Adjustment

- Basic structure:
  - Upon receipt of ACK (of new data): increase rate
  - Upon detection of loss: decrease rate
- How we increase/decrease the rate depends on the phase of congestion control we're in:
  - Discovering available bottleneck bandwidth vs.
  - Adjusting to bandwidth variations

213

## Bandwidth Discovery with Slow Start

- Goal: estimate available bandwidth
  - start slow (for safety)
  - but ramp up quickly (for efficiency)
- Consider
  - RTT = 100ms, MSS=1000bytes
  - Window size to fill 1Mbps of BW = 12.5 packets
  - Window size to fill 1Gbps = 12,500 packets
  - Either is possible!

214

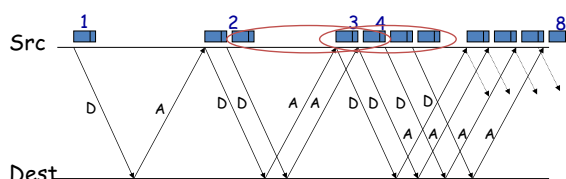
## “Slow Start” Phase

- Sender starts at a slow rate but increases **exponentially** until first loss
- Start with a small congestion window
  - Initially, CWND = 1
  - So, initial sending rate is MSS/RTT
- Double the CWND for each RTT with no loss

215

## Slow Start in Action

- For each RTT: double CWND
- Simpler implementation: for each ACK, CWND += 1



216

## Adjusting to Varying Bandwidth

- Slow start gave an estimate of available bandwidth
- Now, want to track variations in this available bandwidth, oscillating around its current value
  - Repeated probing (rate increase) and backoff (rate decrease)
- TCP uses: “Additive Increase Multiplicative Decrease” (AIMD)
  - We’ll see why shortly...

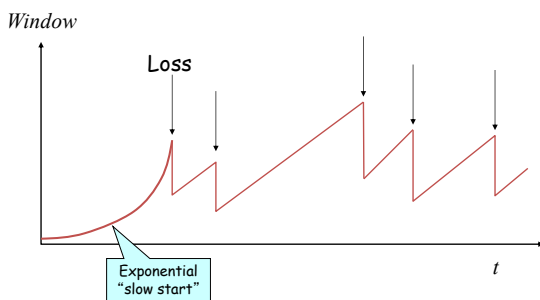
217

## AIMD

- Additive increase
  - Window grows by one MSS for every RTT with no loss
  - For each successful RTT,  $CWND = CWND + 1$
  - Simple implementation:
    - for each ACK,  $CWND = CWND + 1/CWND$
- Multiplicative decrease
  - On loss of packet, divide congestion window in **half**
  - On loss,  $CWND = CWND/2$

218

## Leads to the TCP “Sawtooth”



219

## Slow-Start vs. AIMD

- When does a sender stop Slow-Start and start Additive Increase?
- Introduce a “slow start threshold” (**ssthresh**)
  - Initialized to a large value
  - On timeout,  $ssthresh = CWND/2$
- When  $CWND = ssthresh$ , sender switches from slow-start to AIMD-style increase

220  
220

- What does TCP do?
  - ARQ windowing, set-up, tear-down
- Flow Control in TCP
- Congestion Control in TCP
  - AIMD (slow-start, congestion avoidance)

221

- What does TCP do?
  - ARQ windowing, set-up, tear-down
- Flow Control in TCP
- Congestion Control in TCP
  - AIMD (slow-start, congestion avoidance) and Fast-Recovery

222

## One Final Phase: Fast Recovery

- The problem: congestion avoidance too slow in recovering from an isolated loss

223

## Example (in units of MSS, not bytes)

- Consider a TCP connection with:
  - CWND=10 packets
  - Last ACK was for packet # 101
    - i.e., receiver expecting next packet to have seq. no. 101
- 10 packets [101, 102, 103,..., 110] are in flight
  - Packet 101 is dropped
  - What ACKs do they generate?
  - And how does the sender respond?

224

## The problem – A timeline

- ACK 101 (due to 102) cwnd=10 dupACK#1 (no xmit)
- ACK 101 (due to 103) cwnd=10 dupACK#2 (no xmit)
- ACK 101 (due to 104) cwnd=10 dupACK#3 (no xmit)
- RETRANSMIT 101 ssthresh=5 cwnd= 5
- ACK 101 (due to 105) cwnd=5 + 1/5 (no xmit)
- ACK 101 (due to 106) cwnd=5 + 2/5 (no xmit)
- ACK 101 (due to 107) cwnd=5 + 3/5 (no xmit)
- ACK 101 (due to 108) cwnd=5 + 4/5 (no xmit)
- ACK 101 (due to 109) cwnd=5 + 5/5 (no xmit)
- ACK 101 (due to 110) cwnd=6 + 1/5 (no xmit)
- ACK 111 (due to 101) ← only now can we transmit new packets
- Plus no packets in flight so ACK “clocking” (to increase CWND) stalls for another RTT

225

## Solution: Fast Recovery

Idea: Grant the sender temporary “credit” for each dupACK so as to keep packets in flight

- If dupACKcount = 3
  - ssthresh = cwnd/2
  - cwnd = ssthresh + 3
- While in fast recovery
  - cwnd = cwnd + 1 for each additional duplicate ACK
- Exit fast recovery after receiving new ACK
  - set cwnd = ssthresh

226

## Example

- Consider a TCP connection with:
  - CWND=10 packets
  - Last ACK was for packet # 101
    - i.e., receiver expecting next packet to have seq. no. 101
- 10 packets [101, 102, 103,..., 110] are in flight
  - Packet 101 is dropped

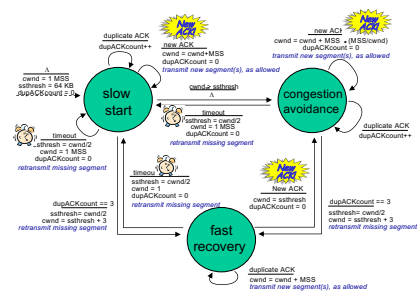
227

## Timeline

- ACK 101 (due to 102) cwnd=10 dup#1
- ACK 101 (due to 103) cwnd=10 dup#2
- ACK 101 (due to 104) cwnd=10 dup#3
- REXMIT 101 ssthresh=5 cwnd= 8 (5+3)
- ACK 101 (due to 105) cwnd= 9 (no xmit)
- ACK 101 (due to 106) cwnd=10 (no xmit)
- ACK 101 (due to 107) cwnd=11 (xmit 111)
- ACK 101 (due to 108) cwnd=12 (xmit 112)
- ACK 101 (due to 109) cwnd=13 (xmit 113)
- ACK 101 (due to 110) cwnd=14 (xmit 114)
- ACK 111 (due to 101) cwnd = 5 (xmit 115) ← exiting fast recovery
- Packets 111-114 already in flight
- ACK 112 (due to 111) cwnd = 5 + 1/5 ← back in congestion avoidance



## Summary: TCP congestion control



- What does TCP do?
  - ARQ windowing, set-up, tear-down
- Flow Control in TCP
- Congestion Control in TCP
  - AIMD (slow-start, congestion avoidance) and Fast-Recovery

Congestion avoidance algorithm has been a fertile field....

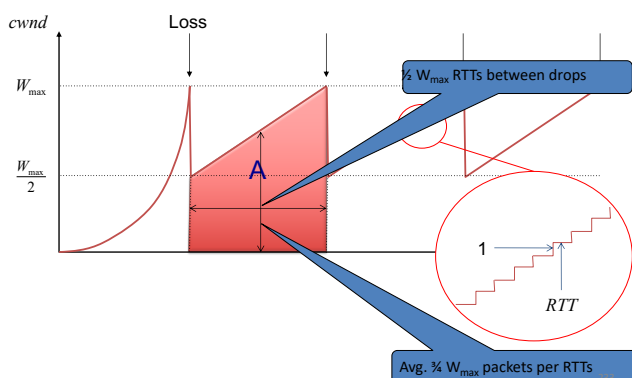
230

Variant	Feedback	Required changes	Benefits	Fairness
(New) Reno	Loss	—	—	Delay
Vegas	Delay	Sender	Less loss	Proportional
High Speed	Loss	Sender	High bandwidth	
BIC	Loss	Sender	High bandwidth	
CUBIC	Loss	Sender	High bandwidth	
• C2TCP <sup>[11][12]</sup>	Loss/Delay	Sender	Ultra-low latency and high bandwidth	
NATCP <sup>[13]</sup>	Multi-bit signal	Sender	Near Optimal Performance	
Elastic-TCP	Loss/Delay	Sender	High bandwidth/short & long-distance	
• Agile-TCP	Loss	Sender	High bandwidth/short-distance	
H-TCP	Loss	Sender	High bandwidth	
FAST	Delay	Sender	High bandwidth	Proportional
Compound TCP	Loss/Delay	Sender	High bandwidth	Proportional
Westwood	Loss/Delay	Sender	L	
• Jersey	Loss/Delay	Sender	L	
BBR <sup>[14]</sup>	Delay	Sender	BLVC, Bufferbloat	
CLAMP	Multi-bit signal	Receiver, Router	V	Max-min
• TFR	Loss	Sender, Receiver	No Retransmission	Minimum delay
XCP	Multi-bit signal	Sender, Receiver, Router	BLFC	Max-min
VCP	2-bit signal	Sender, Receiver, Router	BLF	Proportional
MaxNet	Multi-bit signal	Sender, Receiver, Router	BLFSC	Max-min
JetMax	Multi-bit signal	Sender, Receiver, Router	High bandwidth	Max-min
RED	Loss	Router	Reduced delay	
ECN	Single-bit signal	Sender, Receiver, Router	Reduced loss	

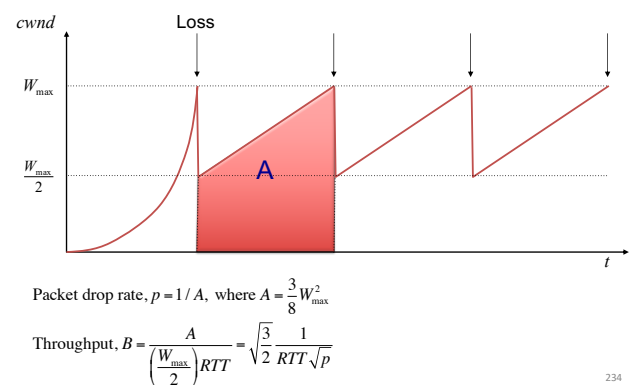
232

## TCP Throughput Equation

### A Simple Model for TCP Throughput



### A Simple Model for TCP Throughput

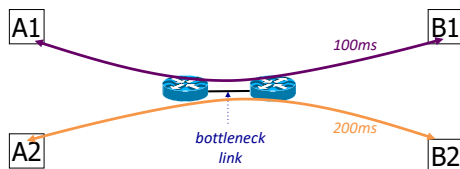


234

## Implications (1): Different RTTs

$$\text{Throughput} = \sqrt{\frac{3}{2}} \frac{1}{RTT \sqrt{p}}$$

- Flows get throughput inversely proportional to RTT
- TCP unfair in the face of heterogeneous RTTs!



235

## Implications (2): High Speed TCP

$$\text{Throughput} = \sqrt{\frac{3}{2}} \frac{1}{RTT \sqrt{p}}$$

- Assume RTT = 100ms, MSS=1500bytes
- What value of  $p$  is required to reach 100Gbps throughput
  - $\sim 2 \times 10^{-12}$
- How long between drops?
  - $\sim 16.6$  hours
- How much data has been sent in this time?
  - $\sim 6$  petabits
- These are not practical numbers!

236

## Adapting TCP to High Speed

- Once past a threshold speed, increase CWND faster
  - A proposed standard [Floyd'03]: once speed is past some threshold, change equation to  $p^{-8}$  rather than  $p^{-5}$
  - Let the additive constant in AIMD depend on CWND
- Other approaches?
  - Multiple simultaneous connections (*hacky* but works today)
  - Router-assisted approaches (will see shortly)

237

## Implications (3): Rate-based CC

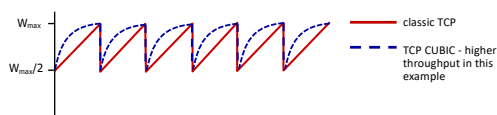
$$\text{Throughput} = \sqrt{\frac{3}{2}} \frac{1}{RTT \sqrt{p}}$$

- TCP throughput is "choppy"
  - repeated swings between  $W/2$  to  $W$
- Some apps would prefer sending at a steady rate
  - e.g., streaming apps
- A solution: "Equation-Based Congestion Control"
  - ditch TCP's increase/decrease rules and just follow the equation
  - measure drop percentage  $p$ , and set rate accordingly
- Following the TCP equation ensures we're "TCP friendly"
  - i.e., use no more than TCP does in similar setting

238

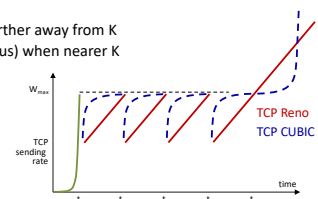
## TCP CUBIC

- Is there a better way than AIMD to "probe" for usable bandwidth?
- Insight/intuition:
  - $W_{\max}$ : sending rate at which congestion loss was detected
  - congestion state of bottleneck link probably (?) hasn't changed much
  - after cutting rate/window in half on loss, initially ramp to to  $W_{\max}$  *faster*, but then approach  $W_{\max}$  more *slowly*



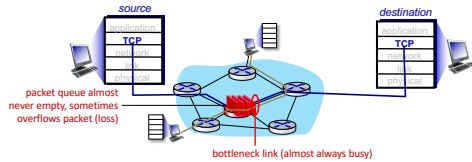
## TCP CUBIC

- K: point in time when TCP window size will reach  $W_{\max}$ 
  - K itself is tuneable
- increase  $W$  as a function of the *cube* of the distance between current time and K
  - larger increases when further away from K
  - smaller increases (cautious) when nearer K
- TCP CUBIC default in Linux, most popular TCP for popular Web servers



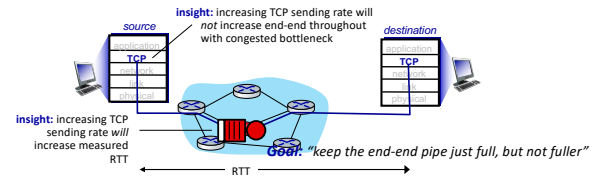
## TCP and the congested “bottleneck link”

- TCP (classic, CUBIC) increase TCP's sending rate until packet loss occurs at some router's output: the **bottleneck link**



## TCP and the congested “bottleneck link”

- TCP (classic, CUBIC) increase TCP's sending rate until packet loss occurs at some router's output: the **bottleneck link**
- understanding congestion: useful to focus on congested bottleneck link



## Delay-based TCP Congestion Control

Keeping sender-to-receiver pipe “just full enough, but no fuller”: keep bottleneck link busy transmitting, but avoid high delays/buffering



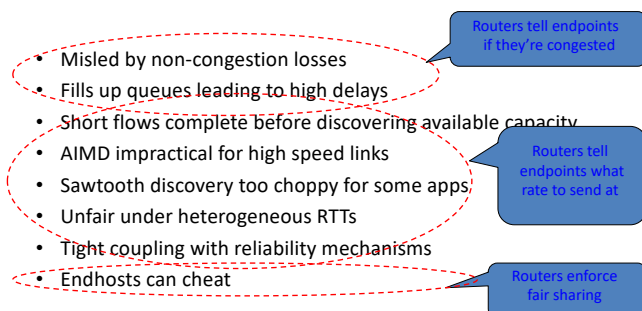
### Delay-based approach:

- $RTT_{min}$  - minimum observed RTT (uncongested path)
- uncongested throughput with congestion window  $cwnd$  is  $cwnd/RTT_{min}$
- if measured throughput “very close” to uncongested throughput  
increase  $cwnd$  linearly /\* since path not congested \*/
- else if measured throughput “far below” uncongested throughput  
decrease  $cwnd$  linearly /\* since path is congested \*/

## Delay-based TCP Congestion Control

- congestion control without inducing/forcing loss
- maximizing throughput (“keeping the just pipe full...”) while keeping delay low (“...but not fuller”)
- a number of deployed TCPs take a delay-based approach
  - BBR deployed on Google's (internal) backbone network

## Recap: TCP problems



Could fix many of these with some help from routers!

## Router-Assisted Congestion Control

- Three tasks for CC:
  - Isolation/fairness
  - Adjustment\*
  - Detecting congestion

\* This may be *automatic* eg loss-response of TCP

## Fairness: General Approach

How can routers ensure each flow gets its “fair share”?

- Routers classify packets into “flows”
  - (For now) flows are packets between same source/destination
- Each flow has its own FIFO queue in router
- Router services flows in a fair fashion
  - When line becomes free, take packet from next flow in a fair order
- What does “fair” mean exactly?

248

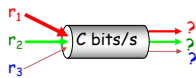
249

## Max-Min Fairness

- Given set of bandwidth demands  $r_i$  and total bandwidth  $C$ , max-min bandwidth allocations are:

$$a_i = \min(f, r_i)$$

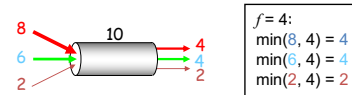
where  $f$  is the unique value such that  $\text{Sum}(a_i) = C$



250

## Example

- $C = 10$ ;  $r_1 = 8$ ,  $r_2 = 6$ ,  $r_3 = 2$ ;  $N = 3$
- $C/3 = 3.33 \rightarrow$ 
  - Can service all of  $r_3$
  - Remove  $r_3$  from the accounting:  $C = C - r_3 = 8$ ;  $N = 2$
- $C/2 = 4 \rightarrow$ 
  - Can't service all of  $r_1$  or  $r_2$
  - So hold them to the remaining fair share:  $f = 4$



251

## Max-Min Fairness

- Given set of bandwidth demands  $r_i$  and total bandwidth  $C$ , max-min bandwidth allocations are:

$$a_i = \min(f, r_i)$$

- where  $f$  is the unique value such that  $\text{Sum}(a_i) = C$
- Property:
  - If you don't get full demand, no one gets more than you
- This is what round-robin service gives if all packets are the same size

252

## How do we deal with packets of different sizes?

- Mental model: Bit-by-bit round robin (“fluid flow”)
- Can you do this in practice?
- No, packets cannot be preempted
- But we can approximate it
  - This is what “fair queuing” routers do

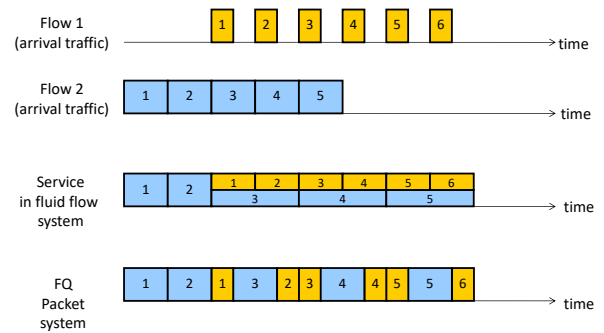
253

## Fair Queuing (FQ)

- For each packet, compute the time at which the last bit of a packet would have left the router *if* flows are served bit-by-bit
- Then serve packets in the increasing order of their deadlines

254

## Example



255

## Fair Queuing (FQ)

- Think of it as an implementation of round-robin generalized to the case where not all packets are equal sized
- **Weighted** fair queuing (WFQ): assign different flows different shares
- Today, some form of WFQ implemented in almost all routers
  - Not the case in the 1980-90s, when CC was being developed
  - Mostly used to isolate traffic at larger granularities (e.g., per-prefix)

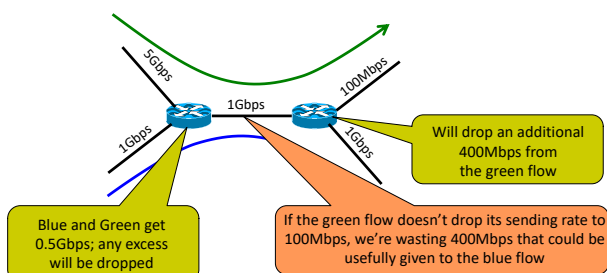
256

## FQ vs. FIFO

- FQ advantages:
  - Isolation: cheating flows don't benefit
  - Bandwidth share does not depend on RTT
  - Flows can pick any rate adjustment scheme they want
- Disadvantages:
  - More complex than FIFO: per flow queue/state, additional per-packet book-keeping

## FQ in the big picture

- FQ does not eliminate congestion → it just manages the congestion



## FQ in the big picture

- FQ does not eliminate congestion → it just manages the congestion
  - robust to cheating, variations in RTT, details of delay, reordering, retransmission, etc.
- But congestion (and packet drops) still occurs
- And we still want end-hosts to discover/adapt to their fair share!
- What would the end-to-end argument say w.r.t. congestion control?

## Fairness is a controversial goal

- What if you have 8 flows, and I have 4?
  - Why should you get twice the bandwidth
- What if your flow goes over 4 congested hops, and mine only goes over 1?
  - Why shouldn't you be penalized for using more scarce bandwidth?
- And what is a flow anyway?
  - TCP connection
  - Source-Destination pair?
  - Source?

## Explicit Congestion Notification (ECN)

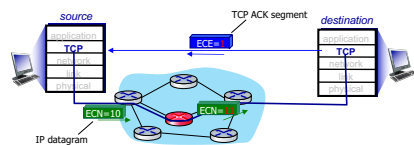
- Single bit in packet header; set by congested routers
  - If data packet has bit set, then ACK has ECN bit set
- Many options for when routers set the bit
  - tradeoff between (link) utilization and (packet) delay
- Congestion semantics can be exactly like that of drop
  - I.e., endhost reacts as though it saw a drop
- Advantages:
  - Doesn't confuse corruption with congestion; recovery w/ rate adjustment
  - Can serve as an early indicator of congestion to avoid delays
  - Easy (easier) to incrementally deploy
    - defined as extension to TCP/IP in RFC 3168 (uses diffserv bits in the IP header)

261

## Explicit congestion notification (ECN)

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

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



## Securing TCP

### Vanilla TCP & UDP sockets:

- no encryption
- cleartext passwords sent into socket traverse Internet in cleartext (!)

### Transport Layer Security (TLS)

- provides encrypted TCP connections
- data integrity
- end-point authentication

### TLS implemented in application layer

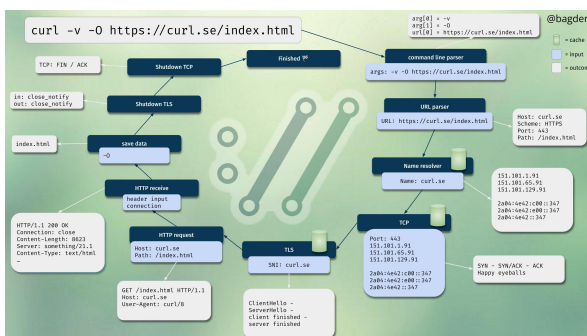
- apps use TLS libraries, that use TCP in turn
- cleartext sent into "socket" traverse Internet **encrypted**

### SSL vs. TLS

Simple: SSL is deprecated

TLS refers to secure socket layers in actual use.

Application Layer: 2.263



264

## Transport Recap

### A "big bag":

Multiplexing, reliability, error-detection, error-recovery, flow and congestion control, ....

- UDP:
  - Minimalist - multiplexing and error detection
- TCP:
  - somewhat hacky
  - but practical/deployable
  - good enough to have raised the bar for the deployment of new approaches
  - though the needs of datacenters change the status quos
- Beyond TCP (discussed in Topic 6):
  - QUIC / application-aware transport layers

265

## Topic 6 – Applications

- Infrastructure Services (DNS)
  - Now with added security...
- Traditional Applications (web)
  - Now with added QUIC
- P2P Networks
  - Every device serves

1

### Some network apps

- social networking
  - Web
  - text messaging
  - e-mail
  - multi-user network games
  - streaming stored video (YouTube, Hulu, Netflix)
  - P2P file sharing
  - voice over IP (e.g., Skype)
  - real-time video conferencing (e.g., Zoom)
  - Internet search
  - remote login
  - ...
- Q: your favorites?

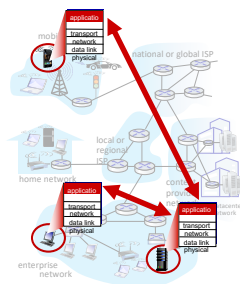
### Creating a network app

#### 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



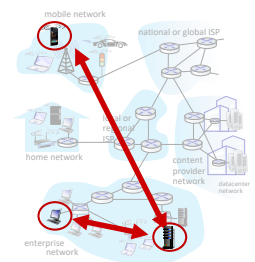
### Client-server paradigm

#### server:

- always-on host
- permanent IP address
- often in data centers, for scaling

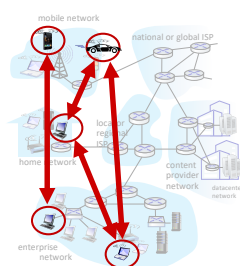
#### clients:

- contact, communicate with server
- may be intermittently connected
- may have dynamic IP addresses
- do *not* communicate directly with each other
- examples: HTTP, IMAP, FTP



### Peer-peer 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
- example: P2P file sharing



### An application-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, everyone has access to protocol definition
- allows for interoperability
- e.g., HTTP, SMTP

#### proprietary protocols:

- e.g., Skype, Zoom



## Relationship Between Names & Addresses

- Addresses can **change** underneath
  - Move `www.bbc.co.uk` to `212.58.246.92`
  - Humans/Apps should be unaffected
- Name could map to **multiple** IP addresses
  - `www.bbc.co.uk` to multiple replicas of the Web site
  - Enables
    - Load-balancing
    - Reducing latency by picking nearby servers
- Multiple names** for the same address
  - E.g., aliases like `www.bbc.co.uk` and `bbc.co.uk`
  - Mnemonic stable name, and dynamic canonical name
    - Canonical name = actual name of host

7

## DNS: Domain Name System

*people*: many identifiers:

- NI #, name, passport #

*Internet hosts, routers*:

- IP address (32 bit or 128bit) - used for addressing datagrams
- "name", e.g., `cam.ac.uk` - used by humans

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

**Domain Name System (DNS)**:

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

## 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

**Q: Why not centralize DNS?**

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

**A: doesn't scale!**

- Comcast DNS servers alone: 770B DNS queries/day
- Akamai DNS servers alone: 2.6T DNS queries/day

## Thinking about the DNS

humongous distributed database:

- ~ billion records, each simple

handles many **trillions of queries/day**:

- many more reads than writes
- performance matters**: almost every Internet transaction interacts with DNS - msec count!

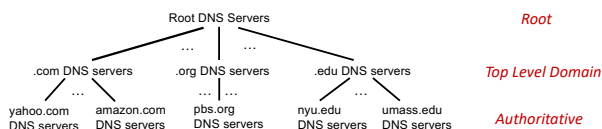
organizationally, physically decentralized:

- millions of different organizations responsible for their records

"bulletproof": reliability, security



## DNS: a distributed, hierarchical database

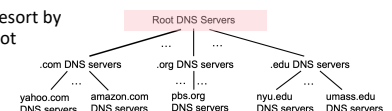


Client wants IP address for `www.amazon.com`; 1<sup>st</sup> approximation:

- 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`

## DNS: root name servers

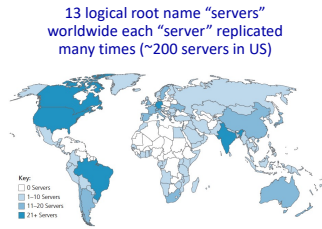
- official, contact-of-last-resort by name servers that can not resolve name





## DNS: root name servers

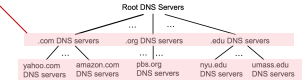
- official, contact-of-last-resort by name servers that can not resolve name
- incredibly important** Internet function
  - Internet couldn't function without it!
  - DNSSEC – provides security (authentication, message integrity)
- ICANN (Internet Corporation for Assigned Names and Numbers) manages root DNS domain



## Top-Level Domain, and authoritative servers

### Top-Level Domain (TLD) servers:

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

## Using DNS

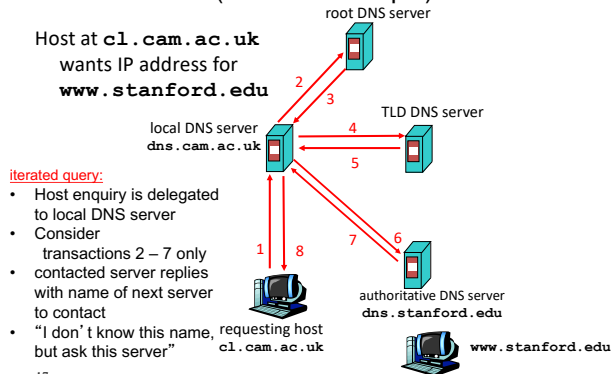
- Two components
  - DNS servers
  - Resolver software on each hosts
- Local DNS server ("default name server")
  - Usually near the endhosts that use it
  - each ISP has local DNS name server; to find yours:
    - MacOS: `% scutil --dns`
    - Windows: `>ipconfig /all`
- Client application
  - Extract server name (e.g., from the URL)
  - Do `gethostbyname()` to trigger resolver code

## Local DNS name Servers

- when host makes DNS query, it is sent to its *local* DNS server
  - Local DNS server returns reply, answering:
    - from its local cache of recent name-to-address translation pairs (possibly out of date!)
    - forwarding request into DNS hierarchy for resolution
  - each ISP has local DNS name server; to find yours:
    - MacOS: `% scutil --dns`
    - Windows: `>ipconfig /all`
- local DNS server doesn't strictly belong to hierarchy, acting as they do on behalf of other hosts.

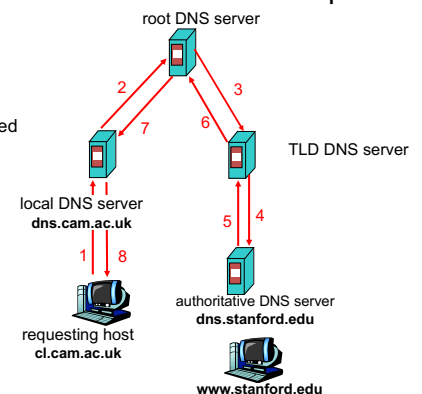
## How Does Resolution Happen?

### (Iterative example)



## DNS name resolution **recursive** example

- recursive query:**
  - puts burden of name resolution on contacted name server
- heavy load?



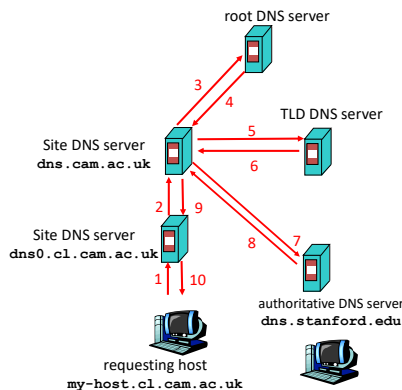
## Recursive and Iterative Queries - Hybrid case

### recursive query:

- Ask server to get answer for you
- E.g., requests 1,2 and responses 9,10

### Iterative query:

- Ask server who to ask next
- E.g., all other request-response pairs



19

## DNS Caching

- Performing all these queries takes time
  - And all this **before** actual communication takes place
  - E.g., 1-second latency before starting Web download
- Caching greatly reduces overhead
  - The top-level servers very rarely change
  - Popular sites (e.g., www.bbc.co.uk) visited often
  - Local DNS servers have regularly used information cached
- How DNS caching works
  - DNS servers will cache responses to queries
  - Responses include a "time to live" (TTL) field
  - Server deletes cached entry after TTL expires
  - Cached entries may be **out-of-date**
    - if named host changes IP address, may not be known Internet-wide until all TTLs expire!
    - best-effort name-to-address translation!**

20

## Reliability

- DNS servers are **replicated**
  - Name service available if at least one replica is up
  - Queries can be load-balanced between replicas
- Anycast provides reliability for ROOT servers
- Usually, UDP is used for queries
  - Need reliability: must implement this on top of UDP
  - DNS spec. supports TCP too, but not always available
- Try alternate servers on timeout
  - Exponential backoff** when retrying same server
- Same identifier for all queries
  - Don't care which server responds

21

## 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=CNAME

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

### type=NS

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

### type=MX

- value is name of SMTP mail server associated with name

## DNS protocol messages

DNS **query** and **reply** messages, both have same **format**:

message header:

- identification**: 16 bit # for query, reply to query uses same #
- flags**:
  - query or reply
  - recursion desired
  - recursion available
  - reply is authoritative

2 bytes	2 bytes
identification	flags
# questions	# answer RRs
# authority RRs	# additional RRs
questions (variable # of questions)	
answers (variable # of RRs)	
authority (variable # of RRs)	
additional info (variable # of RRs)	

## DNS protocol messages

DNS **query** and **reply** messages, both have same **format**:

name, type fields for a query

RRs in response to query

records for authoritative servers

additional "helpful" info that may be used

2 bytes	2 bytes
identification	flags
# questions	# answer RRs
# authority RRs	# additional RRs
questions (variable # of questions)	
answers (variable # of RRs)	
authority (variable # of RRs)	
additional info (variable # of RRs)	

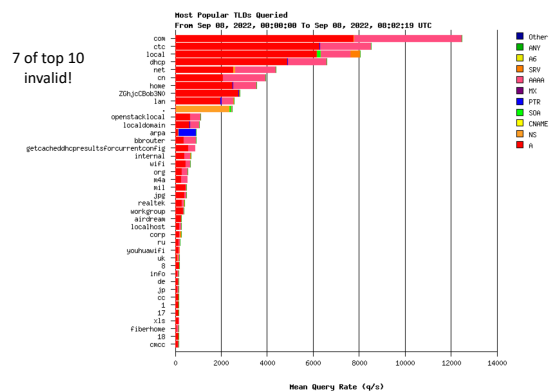
## Getting your info into the DNS

example: new startup "Network Utopia"

- register name `networkuptopia.com` at **DNS registrar** (e.g., Network Solutions)
  - provide names, IP addresses of authoritative name server (primary and secondary)
  - registrar inserts NS, A RRs into `.com` TLD server:  
`(networkuptopia.com, dns1.networkuptopia.com, NS)`  
`(dns1.networkuptopia.com, 212.212.212.1, A)`
- create authoritative server locally with IP address `212.212.212.1`
  - type A record for `www.networkuptopia.com`
  - type MX record for `networkuptopia.com`

## Most popular TLD

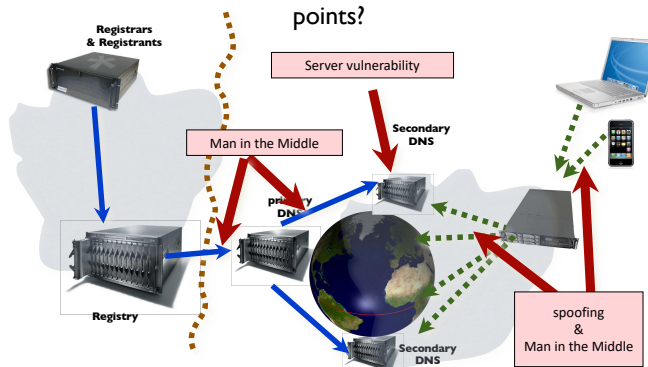
At least WORKGROUP is no longer here!  
It was the top invalid TLD for years...



26

## Data flow through the DNS

### Where are the vulnerable points?



## DNS attack surface

## DDoS attacks

- bombard root servers with traffic
  - not successful to date
  - traffic filtering
  - local DNS servers cache IPs of TLD servers, allowing root server bypass
- bombard TLD servers
  - potentially more dangerous

## Spoofing attacks

- intercept DNS queries, returning bogus replies
  - DNS cache poisoning
  - RFC 4033: DNSSEC authentication services

# DNS Security

- No way to verify answers
  - Opens up DNS to many potential attacks
  - DNSSEC fixes this
- Most obvious vulnerability: recursive resolution
  - Using recursive resolution, host must trust DNS server
  - When at Starbucks, server is under their control
    - And can return whatever values it wants
- More subtle attack: Cache poisoning
  - Those “additional” records can be anything!

DNSSEC protects all these end-to-end

- provides message authentication and integrity verification through cryptographic signatures
  - You know who provided the signature
  - No modifications between signing and validation
- It does **not** provide authorization
- It does **not** provide confidentiality
- It does **not** provide protection against DDOS

## DNSSEC in practice

### Problem: Scaling the key signing and key distribution

### Solution: Using the DNS to Distribute Keys

- Distribute keys through the DNS hierarchy
  - Use one trusted key to establish authenticity of other keys
  - Building chains of trust from the root down
  - Parents need to sign the keys of their children
- Only the root key needed in ideal world
  - Parents always delegate security to child



# Why is the web so successful?

- What do the web, youtube, facebook, twitter, instagram, ..... have in common?
  - The ability to self-publish
- Self-publishing that is easy, independent, *free*
- No interest in collaborative and idealistic endeavors
  - People aren't looking for Nirvana (or even Xanadu)
  - People also aren't looking for technical perfection
- Want to make their mark, and find something neat
  - Two sides of the same coin, creates synergy
  - "Performance" more important than dialogue....

33

On osx "*host -av www.cl.cam.ac.uk*

[illegible]

32

## Web and HTTP

*First, a quick review...*

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

www.university.ac.uk/someDept/pic.gif

host name                      path name

## 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



## HTTP overview (continued)

HTTP 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

- Reminder: Distributed Systems are Hard!

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

## HTTP connections: two types

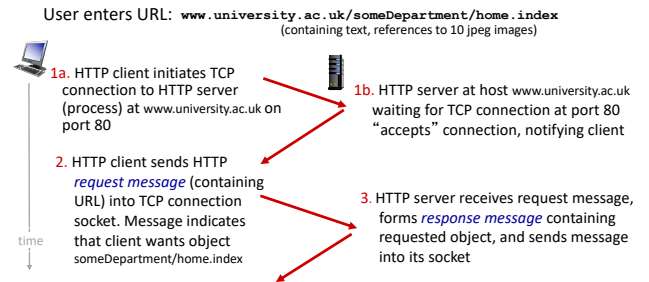
### Non-persistent HTTP

1. TCP connection opened
  2. at most one object sent over TCP connection
  3. TCP connection closed
- downloading multiple objects required multiple connections

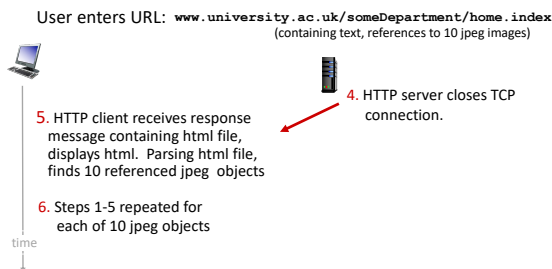
### Persistent HTTP

- TCP connection opened to a server
- multiple objects can be sent over *single* TCP connection between client, and that server
- TCP connection closed

## Non-persistent HTTP: example



## Non-persistent HTTP: example (cont.)

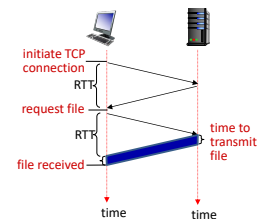


## Non-persistent HTTP: response time

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

### HTTP response time (per object):

- one RTT to initiate TCP connection
- one RTT for HTTP request and first few bytes of HTTP response to return
- object/file transmission time



Non-persistent HTTP response time =  $2RTT + \text{file transmission time}$

## Persistent HTTP (HTTP 1.1)

### Non-persistent HTTP issues:

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

### Persistent HTTP (HTTP1.1):

- 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 (cutting response time in half)

## HTTP request message

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

### HTTP request message:

- ASCII (human-readable format)

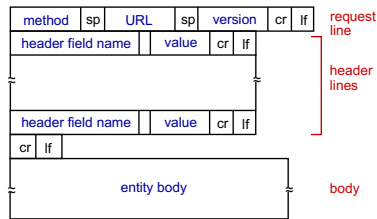
request line (GET, POST, HEAD commands)

carriage return character  
line-feed character

carriage return, line feed at start of line indicates end of header lines

\* Check out the online interactive exercises for more examples: [http://gaia.cs.umass.edu/kurose\\_ross/interactive/](http://gaia.cs.umass.edu/kurose_ross/interactive/)

## HTTP request message: general format



## Other HTTP request messages

### POST method:

- web page often includes form input
- user input sent from client to server in entity body of HTTP POST request message

### HEAD method:

- requests headers (only) that would be returned *if* specified URL were requested with an HTTP GET method.

### PUT method:

- uploads new file (object) to server
- completely replaces file that exists at specified URL with content in entity body of POST HTTP request message

### GET method (for sending data to server):

- include user data in URL field of HTTP GET request message (following a '?'):

`www.somesite.com/animalsearch?monkeys&banana`

## HTTP response message

status line (protocol status code status phrase) → `HTTP/1.1 200 OK`

## 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 message

### 301 Moved Permanently

- requested object moved, new location specified later in this message (in Location: field)

### 400 Bad Request

- request msg not understood by server

### 404 Not Found

- requested document not found on this server

### 505 HTTP Version Not Supported

## Trying out HTTP (client side) for yourself

- Netcat (telnet will also work) to your favorite Web server:

`% nc -c -v www.cl.cam.ac.uk 80` opens TCP connection to port 80 (default HTTP server port) at `www.cl.cam.ac.uk` anything typed in will be sent to port 80 at `www.cl.cam.ac.uk`

- type in a GET HTTP request:

`GET /~awm22/index.php HTTP/1.1`  
`Host: www.cl.cam.ac.uk`

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

- look at response message sent by HTTP server!

(or use Wireshark to look at captured HTTP request/response)

Although in readable ascii – you will notice this is not the webpage but a redirect  
Automatically moving to an https secure connection

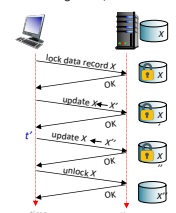
## Maintaining user/server state: cookies

Recall: HTTP GET/response interaction is *stateless*

- no notion of multi-step exchanges of HTTP messages to complete a Web “transaction”

- no need for client/server to track “state” of multi-step exchange
- all HTTP requests are independent of each other
- no need for client/server to “recover” from a partial-but-never-entirely-completed transaction

a stateful protocol: client makes two changes to  $X_0$ , or none at all



Q: what happens if network connection or client crashes at  $t'$ ?

## Maintaining user/server state: cookies

Web sites and client browser use **cookies** to maintain some state between transactions

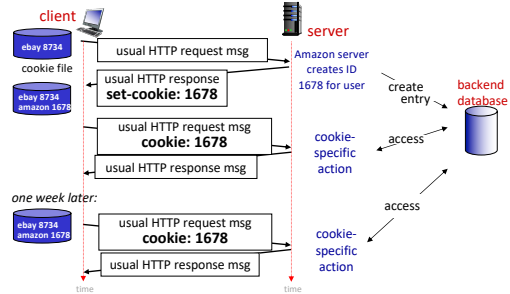
**four components:**

- 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 uses browser on laptop, visits specific e-commerce site for first time
- when initial HTTP requests arrives at site, site creates:
  - unique ID (aka "cookie")
  - entry in backend database for ID
- subsequent HTTP requests from Susan to this site will contain cookie ID value, allowing site to "identify" Susan

## Maintaining user/server state: cookies



## HTTP cookies: comments

**What cookies can be used for:**

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

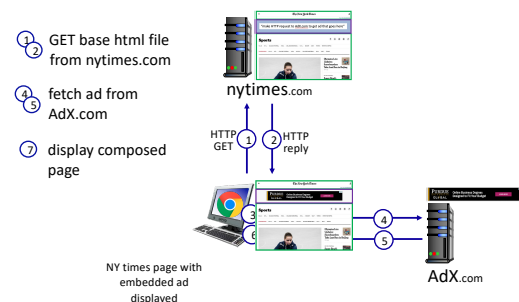
**Challenge: How to keep state?**

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

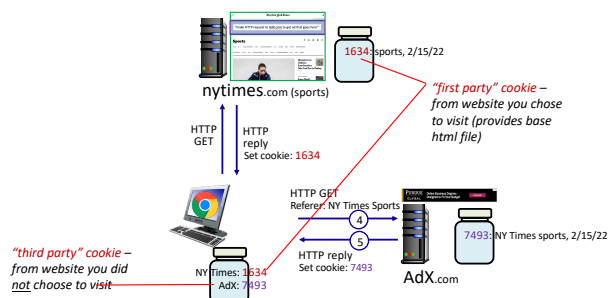
**cookies and privacy:** aside

- cookies permit sites to *learn* a lot about you on their site.
- third party persistent cookies (tracking cookies) allow common identity (cookie value) to be tracked across multiple web sites

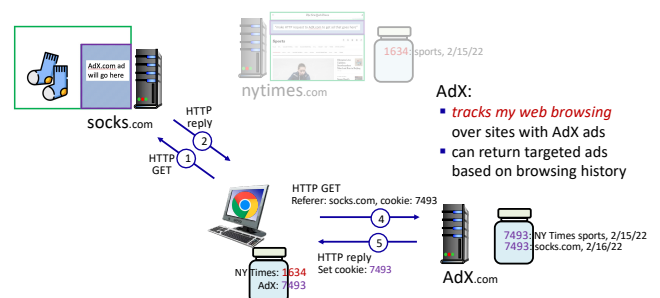
## Example: displaying a NY Times web page



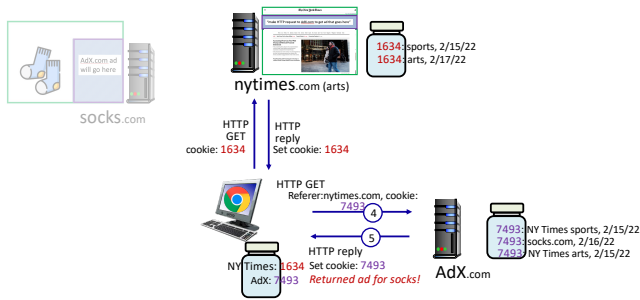
## Cookies: tracking a user's browsing behavior



## Cookies: tracking a user's browsing behavior



## Cookies: tracking a user's browsing behavior (one day later)



## Cookies: tracking a user's browsing behavior

Cookies can be used to:

- track user behavior on a given website (**first party cookies**)
- track user behavior across multiple websites (**third party cookies**) without user ever choosing to visit tracker site (!)
- tracking may be *invisible* to user:
  - rather than displayed ad triggering HTTP GET to tracker, could be an invisible link

third party tracking via cookies:

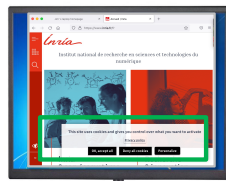
- disabled by default in Firefox, Safari browsers
- to be disabled in Chrome browser in 2023

## GDPR (EU General Data Protection Regulation) and cookies

"Natural persons may be associated with online identifiers [...] such as internet protocol addresses, cookie identifiers or other identifiers [...]. This may leave traces which, in particular when combined with unique identifiers and other information received by the servers, may be used to create profiles of the natural persons and identify them."

GDPR, recital 30 (May 2018)

when cookies can identify an individual, cookies are considered personal data, subject to GDPR personal data regulations

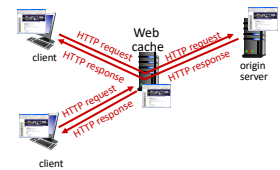


User has explicit control over whether or not cookies are allowed

## Web caches

**Goal:** satisfy client requests without involving origin server

- user configures browser to point to a (local) **Web cache**
- browser sends all HTTP requests to cache
  - if object in cache: cache returns object to client
  - else cache requests object from origin server, caches received object, then returns object to client



## Web caches (aka proxy servers)

- Web cache acts as both client and server
  - server for original requesting client
  - client to origin server
- server tells cache about object's allowable caching in response header:

Cache-Control: max-age=<seconds>

Cache-Control: no-cache

**Why** Web caching?

- reduce response time for client request
  - cache is closer to client
- reduce traffic on an institution's access link
- Internet is dense with caches
  - enables "poor" content providers to more effectively deliver content

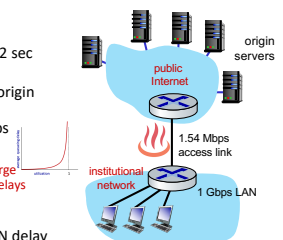
## Caching example

**Scenario:**

- access link rate: 1.54 Mbps
- RTT from institutional router to server: 2 sec
- web object size: 100K bits
- average request rate from browsers to origin servers: 15/sec
  - avg data rate to browsers: 1.50 Mbps

**Performance:**

- access link utilization  $\approx .97$  **problem: large queueing delays at high utilization!**
- LAN utilization: .0015
- end-end delay = Internet delay + access link delay + LAN delay = 2 sec + **minutes** + usecs





## Option 1: buy a faster access link

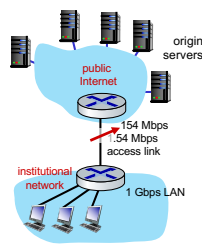
### Scenario:

- access link rate: ~~154~~ 1.54 Mbps
- RTT from institutional router to server: 2 sec
- web object size: 100K bits
- average request rate from browsers to origin servers: 15/sec
  - avg data rate to browsers: 1.50 Mbps

### Performance:

- access link utilization = ~~0.97~~ .0097
- LAN utilization: .0015
- end-end delay = Internet delay + access link delay + LAN delay
  - = 2 sec + ~~minutes~~ + usecs

Cost: faster access link (expensive!)



## Option 2: install a web cache

### Scenario:

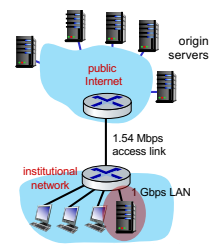
- access link rate: 1.54 Mbps
- RTT from institutional router to server: 2 sec
- web object size: 100K bits
- average request rate from browsers to origin servers: 15/sec
  - avg data rate to browsers: 1.50 Mbps

Cost: web cache (cheap!)

### Performance:

- LAN utilization: .?
- access link utilization = ?
- average end-end delay = ?

How to compute link utilization, delay?

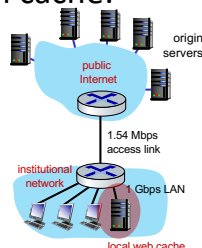


## Calculating access link utilization, end-end delay with cache:

suppose cache hit rate is 0.4:

- 40% requests served by cache, with low (msec) delay
- 60% requests satisfied at origin
  - rate to browsers over access link =  $0.6 * 1.50 \text{ Mbps} = .9 \text{ Mbps}$
  - access link utilization =  $0.9 / 1.54 = .58$  means low (msec) queueing delay at access link
- average end-end delay:
  - =  $0.6 * (\text{delay from origin servers})$
  - +  $0.4 * (\text{delay when satisfied at cache})$
  - =  $0.6 (2.01) + 0.4 (\sim \text{msecs}) \approx \sim 1.2 \text{ secs}$

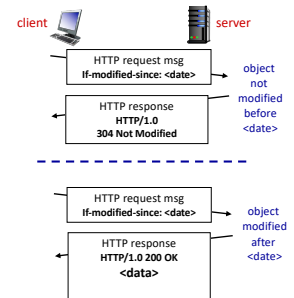
lower average end-end delay than with 154 Mbps link (and cheaper too!)



## Browser caching: Conditional GET

Goal: don't send object if browser has up-to-date cached version

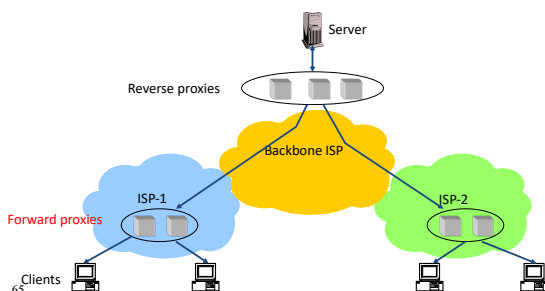
- no object transmission delay (or use of network resources)
- client: specify date of browser-cached copy in HTTP request
  - If-modified-since: <date>
- server: response contains no object if browser-cached copy is up-to-date:
  - HTTP/1.0 304 Not Modified



## Improving HTTP Performance: Caching with Forward Proxies

Cache documents close to **clients**  
→ reduce network traffic and decrease latency

- Typically done by ISPs or corporate LANs to reduce link usage



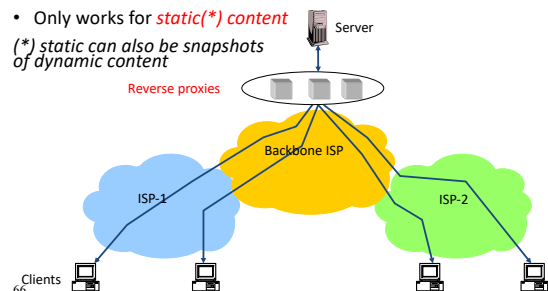
## Improving HTTP Performance: Caching with Reverse Proxies

Cache documents close to **server**  
→ decrease server load

- Typically done by content providers (e.g. scaling capacity for news site)

- Only works for **static(\*) content**

(\*) static can also be snapshots of dynamic content

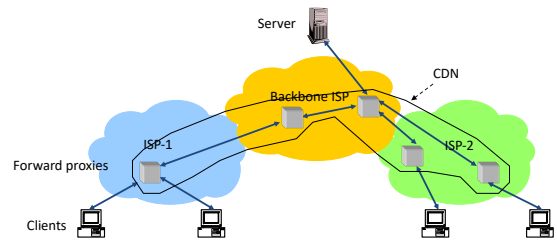


## Improving HTTP Performance: Caching w/ Content Distribution Networks

- Integrate forward and reverse caching functionality
  - One overlay network (usually) administered by one entity
  - *e.g.*, Akamai
- Provide document caching
  - **Pull**: Direct result of clients' requests
  - **Push**: Expectation of high access rate
- Also do some processing
  - Handle *dynamic* web pages
  - *Transcoding*
  - *Maybe do some security function – watermark IP*

67

## Improving HTTP Performance: Caching with CDNs (cont.)



68

## Improving HTTP Performance: CDN Example – Akamai

- Akamai creates new domain names for each client content provider.
  - *e.g.*, [a128.g.akamai.net](http://a128.g.akamai.net)
- The CDN's DNS servers are authoritative for the new domains
- The client content provider modifies its content so that embedded URLs reference the new domains.
  - "Akamaize" content
  - *e.g.*: <http://www.bbc.co.uk/popular-image.jpg> becomes <http://a128.g.akamai.net/popular-image.jpg>
- *Requests now sent to CDN's infrastructure...*

69

## Hosting: Multiple Sites Per Machine

- Multiple Web sites on a single machine
  - Hosting company runs the Web server on behalf of multiple sites (*e.g.*, [www.foo.com](http://www.foo.com) and [www.bar.com](http://www.bar.com))
- Problem: GET /index.html
  - [www.foo.com/index.html](http://www.foo.com/index.html) OR [www.bar.com/index.html](http://www.bar.com/index.html)?
- Solutions:
  - Multiple server processes on the same machine
    - Have a separate IP address (or port) for each server
  - Include site name in HTTP request
    - Single Web server process with a single IP address
    - Client includes "Host" header (*e.g.*, Host: [www.foo.com](http://www.foo.com))
    - *Required header with HTTP/1.1*

70

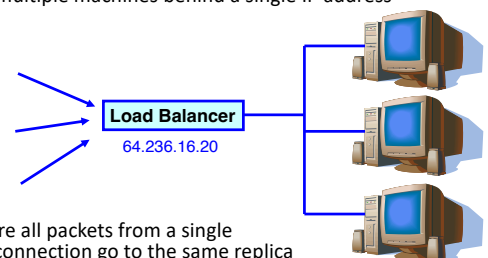
## Hosting: Multiple Machines Per Site

- Replicate popular Web site across many machines
  - Helps to handle the load
  - Places content closer to clients
- Helps when content isn't cacheable
- Problem: Want to direct client to particular replica
  - Balance load across server replicas
  - Pair clients with nearby servers

71

## Multi-Hosting at Single Location

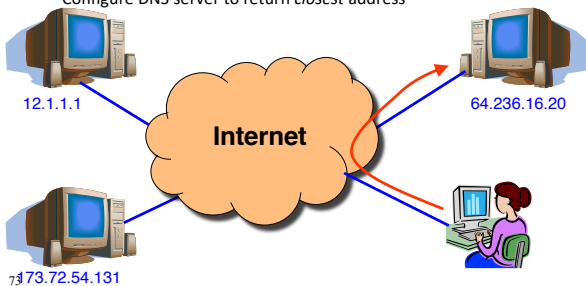
- Single IP address, multiple machines
  - Run multiple machines behind a single IP address



72

## Multi-Hosting at Several Locations

- Multiple addresses, multiple machines
  - Same name but different addresses for all of the replicas
  - Configure DNS server to return *closest* address



## After HTTP/1.1

SPDY (speedy) and its moral successor HTTP/2

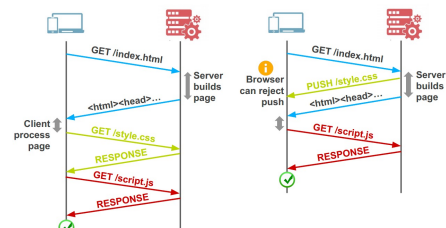
- Binary protocol
- Multiplexing
- Priority control over Frames
- Header Compression
- Server Push

## CDN examples round-up

- CDN using DNS
  - DNS has information on loading/distribution/location (akami uses this one)
- CDN using anycast
  - same address from DNS name but local routes (ROOT DNS servers and 8.8.8.8 use this one)
- CDN based on rewriting HTML URLs
  - (akami example in previous slides)

74

## After HTTP/1.1



- Server Push
  - Proactively push stuff to client that it will need

75

76

## After HTTP/1.1

SPDY (speedy) and its moral successor HTTP/2

- Binary protocol
  - More efficient to parse
  - More compact on the wire
  - Much less error prone as compared to textual protocols

Wireshark decoders for the win

77

## HTTP/2

**Key goal:** decreased delay in multi-object HTTP requests

**HTTP1.1:** introduced **multiple, pipelined GETs** over single TCP connection

- server responds *in-order* (FCFS: first-come-first-served scheduling) to GET requests
- with FCFS, small object may have to wait for transmission (**head-of-line (HOL) blocking**) behind large object(s)
- loss recovery (retransmitting lost TCP segments) stalls object transmission

## HTTP/2

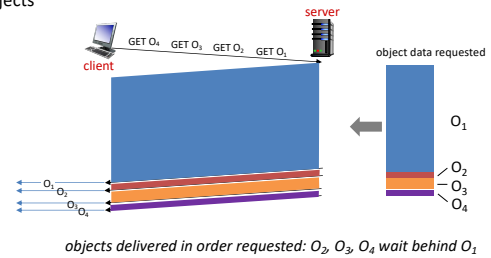
**Key goal:** decreased delay in multi-object HTTP requests

**HTTP/2:** [RFC 7540, 2015] increased flexibility at *server* in sending objects to client:

- methods, status codes, most header fields unchanged from HTTP 1.1
- transmission order of requested objects based on client-specified object priority (not necessarily FCFS)
- *push* unrequested objects to client
- divide objects into frames, schedule frames to mitigate HOL blocking

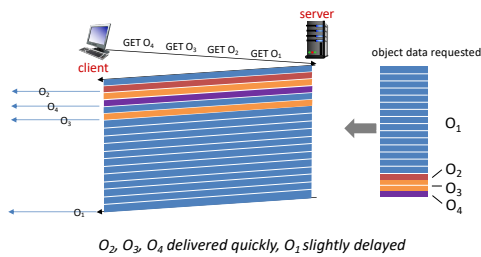
## HTTP/2: mitigating HOL blocking

HTTP 1.1: client requests 1 large object (e.g., video file) and 3 smaller objects



## HTTP/2: mitigating HOL blocking

HTTP/2: objects divided into frames, frame transmission interleaved

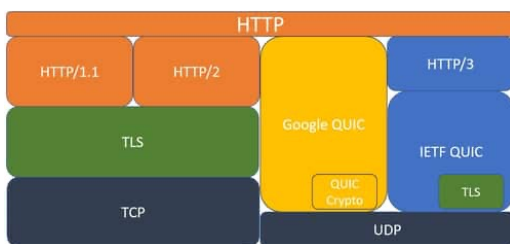


## HTTP/2 to HTTP/3

HTTP/2 over single TCP connection means:

- recovery from packet loss still stalls all object transmissions
  - as in HTTP 1.1, browsers have incentive to open multiple parallel TCP connections to reduce stalling, increase overall throughput
- no security over vanilla TCP connection
- **HTTP/3:** adds security, per object error- and congestion-control (more pipelining) over UDP

As at when I last looked



Other ongoing work includes QUIC for datagrams  
Seriously! It adds QUIC crypto to "UDP" so isn't totally silly.

## Add QUIC and stir...

### Quick UDP Internet Connections

Objective: Combine speed of UDP protocol with TCP's reliability

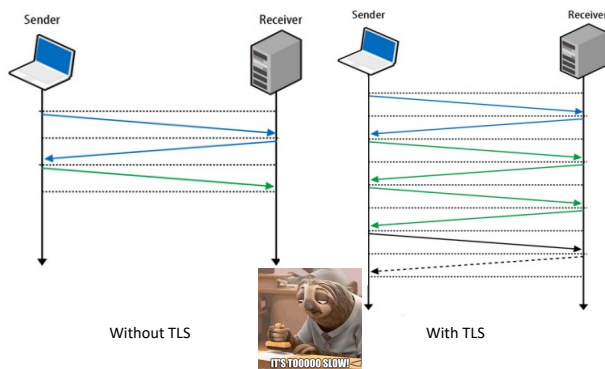
**Problem:** Very hard to make changes to TCP

- *Faster to implement new protocol on top of UDP*
- (Roll out features in TCP if they prove theory)

QUIC (First presented to IETF in ~2013):

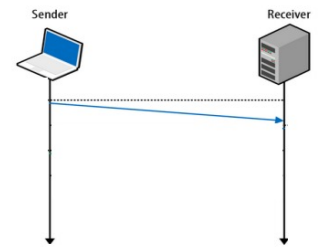
- Reliable transport over UDP
- Uses FEC
- Default crypto
- Restartable connections

## 3-Way Handshake



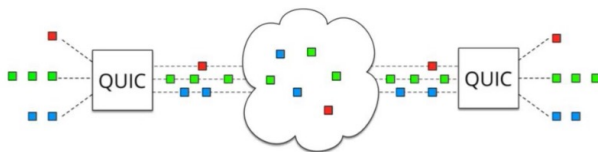
## UDP

- Fire and forget
  - Less time spent to validate packets
  - Downside - no reliability, this has to be added on top of UDP



## QUIC

- UDP does NOT depend on order of arriving packets
- Lost packets will only impact an individual resource, e.g., CSS or JS file.
- QUIC combined the best parts of HTTP/2 over UDP:
  - Multiplexing on top of non-blocking transport protocol

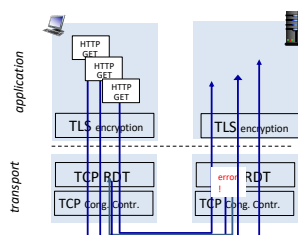


## QUIC: Quick UDP Internet Connections

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

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

QUIC: streams – parallelism  
no HOL blocking in transport or application



(a) HTTP 1.1

## QUIC – more than just UDP

- QUIC outshines TCP under poor network conditions, shaving a full second off the Google Search page load time for the slowest 1% of connections.
- These benefits are even more apparent for video services like YouTube
  - Users report 30% fewer rebuffers with QUIC.

## Why QUIC over UDP and not a new proto

- IP proto value for new transport layer
- Change the protocol – risk the wraith of
  - Legacy code
  - Firewalls
  - Load-balancer
  - NATs (the high-priest of middlebox)
- Same problem faces any significant TCP change

Every host is a server:  
Peer-2-Peer

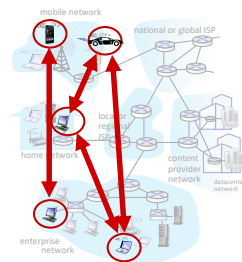
Honda M. et al. "Is it still possible to extend TCP?", IMC'11  
<https://dl.acm.org/doi/abs/10.1145/2068816.2068834>

91

92

### Peer-to-peer (P2P) Architecture

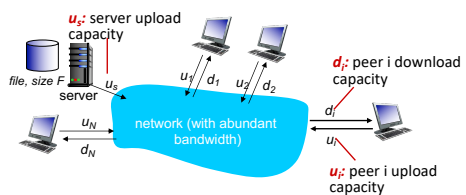
- no always-on server
- arbitrary end systems directly communicate
- peers request service from other peers, provide service in return to other peers
- **self scalability** – new peers bring new service capacity, and new service demands
- peers are intermittently connected and change IP addresses
  - complex management
- examples: P2P file sharing (BitTorrent), streaming (KanKan), VoIP (Skype)



### File distribution: client-server vs P2P

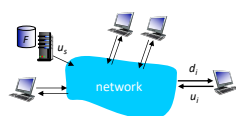
**Q:** how much time to distribute file (size  $F$ ) from one server to  $N$  peers?

- peer upload/download capacity is limited resource



### 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:  $F/d_{\min}$

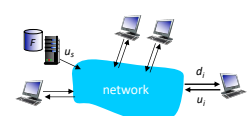


$$\text{time to distribute } F \text{ to } N \text{ clients using client-server approach} \quad D_{c-s} \geq \max\{NF/u_s, F/d_{\min}\}$$

increases linearly in  $N$

### File distribution time: P2P

- **server transmission:** must upload at least one copy:
  - time to send one copy:  $F/u_s$
- **client:** each client must download file copy
  - min client download time:  $F/d_{\min}$
- **clients:** as aggregate must download  $NF$  bits
  - max upload rate (limiting max download rate) is  $u_s + \sum u_i$

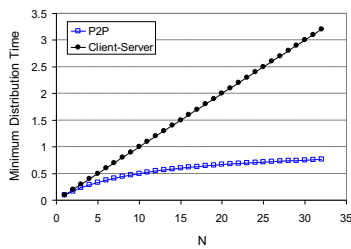


$$\text{time to distribute } F \text{ to } N \text{ clients using P2P approach} \quad D_{p2p} \geq \max\{F/u_s, F/d_{\min}, NF/(u_s + \sum u_i)\}$$

increases linearly in  $N$  ...  
... but so does this, as each peer brings service capacity

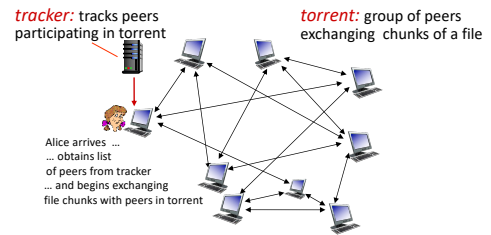
## Client-server vs. P2P: example

client upload rate =  $u$ ,  $F/u = 1$  hour,  $u_s = 10u$ ,  $d_{min} \geq u_s$



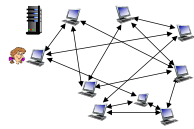
## P2P file distribution: BitTorrent

- file divided into 256Kb chunks
- peers in torrent send/receive file chunks



## 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
- peer exchanges prioritize rarer blocks
- churn**: peers may come and go
- once peer has entire file, it may (selfishly) leave or (altruistically) remain in torrent



## BitTorrent: requesting, sending file chunks

### Requesting chunks:

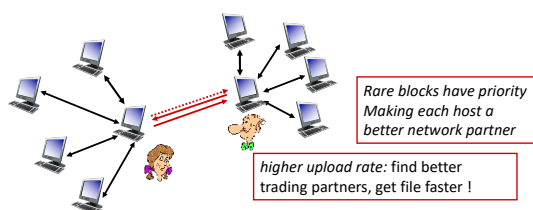
- at any given time, different peers have different subsets of file chunks
- periodically, Alice asks each peer for list of chunks that they have
- Alice requests missing chunks from peers, rarest first

### 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

## BitTorrent: tit-for-tat

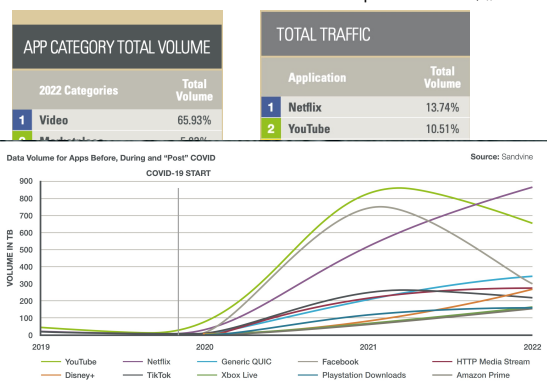
- Alice "optimistically unchokes" Bob
- Alice becomes one of Bob's top-four providers; Bob reciprocates
- Bob becomes one of Alice's top-four providers



## Internet

(current data is \$\$\$ or hard to get)

This info taken from an annual Sandvine report for 2022 <https://www.sandvine.com>



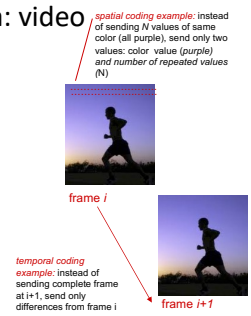
## Video Streaming and CDNs: context

- stream video traffic: major consumer of Internet bandwidth
  - Netflix, YouTube, Amazon Prime: 80% of residential ISP traffic (2020)
- challenge:** scale - how to reach ~1B users?
- challenge:** heterogeneity
  - different users have different capabilities (e.g., wired versus mobile; bandwidth rich versus bandwidth poor)
- solution:** distributed, application-level infrastructure



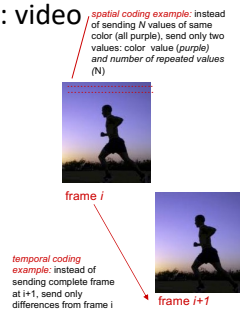
## Multimedia: video

- 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)



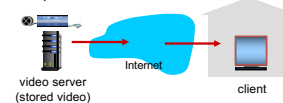
## Multimedia: video

- 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, 64Kbps – 12 Mbps)



## Streaming stored video

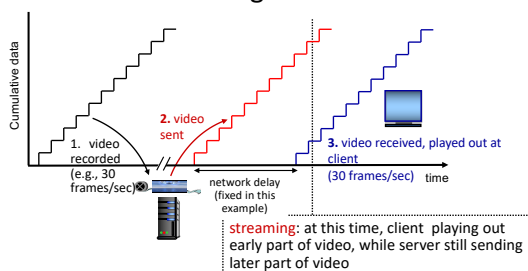
simple scenario:



Main challenges:

- server-to-client bandwidth will *vary* over time, with changing network congestion levels (in house, access network, network core, video server)
- packet loss, delay due to congestion will delay playout, or result in poor video quality

## Streaming stored video



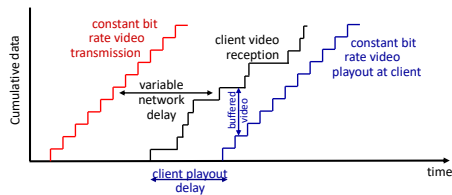
## Streaming stored video: challenges

- continuous playout constraint:** during client video playout, playout timing must match original timing
  - ... but **network delays are variable** (jitter), so will need **client-side buffer** to match continuous playout constraint
- other challenges:**
  - client interactivity: pause, fast-forward, rewind, jump through video
  - video packets may be lost, retransmitted





## Streaming stored video: playout buffering



- **client-side buffering and playout delay:** compensate for network-added delay, delay jitter

## Streaming multimedia: DASH

Dynamic, Adaptive  
Streaming over HTTP

### server:

- divides video file into multiple chunks
- each chunk encoded at multiple different rates
- different rate encodings stored in different files
- files replicated in various CDN nodes
- **manifest file:** provides URLs for different chunks



### client:

- periodically estimates 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), and from different servers

## Streaming multimedia: DASH

- **“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)



Streaming video = encoding + DASH + playout buffering

## Content distribution networks (CDNs)

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

- **option 1:** single, large “mega-server”
  - single point of failure
  - point of network congestion
  - long (and possibly congested) path to distant clients

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

## Content distribution networks (CDNs)

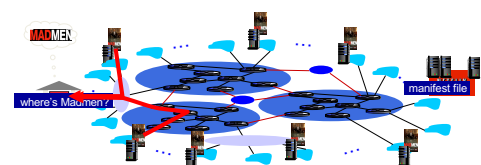
**challenge:** how to stream content (selected from millions of videos) to hundreds of thousands of *simultaneous* users?

- **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
    - Akamai: 240,000 servers deployed in > 120 countries (2015)
  - **bring home:** smaller number (10's) of larger clusters in POPs near access nets
    - used by Limelight



## Content distribution networks (CDNs)

- CDN: stores copies of content (e.g. MADMEN) at CDN nodes
- subscriber requests content, service provider returns manifest
  - using manifest, client retrieves content at highest supportable rate
  - may choose different rate or copy if network path congested



# Summary

## Content distribution networks (CDNs)



**OTT challenges:** coping with a congested Internet from the “edge”

- what content to place in which CDN node?
- from which CDN node to retrieve content? At which rate?

- Applications have protocols too
- We covered examples from
  - Traditional Applications (web)
  - Scaling and Speeding the web (CDN/Cache tricks)
- Infrastructure Services (DNS)
  - Cache and Hierarchy
- P2P Network examples
- Evolving standards (Email)
- Video CDN Stream challenges

116

## Email (as time permits)

Still the best/worst most useful/useless service

Email was the exemplar of the *Electronic Office*

Because every business thought in *memo*

**MEMORANDUM**  
TO: All Employees  
CC: Kevin Smith  
FROM: Jonker Black  
DATE: January 5, 2015  
SUBJECT: Ant Problem in the Office



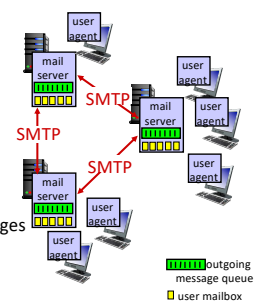
## E-mail

**Three major components:**

- user agents
- mail servers
- simple mail transfer protocol: SMTP

**User Agent**

- a.k.a. “mail reader”
- composing, editing, reading mail messages
- e.g., Outlook, iPhone mail client
- outgoing, incoming messages stored on server



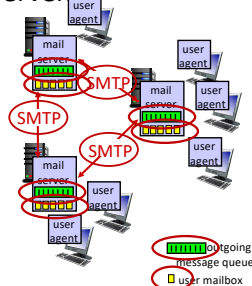
## E-mail: mail servers

**mail servers:**

- **mailbox** contains incoming messages for user
- **message queue** of outgoing (to be sent) mail messages

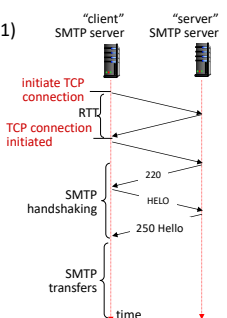
**SMTP protocol** between mail servers to send email messages

- **client:** sending mail server
- **“server”:** receiving mail server



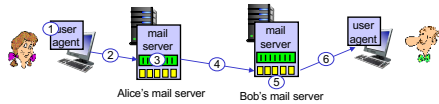
## SMTP RFC (5321)

- uses TCP to reliably transfer email message from client (mail server initiating connection) to server, port 25
  - direct transfer: sending server (acting like client) to receiving server
- three phases of transfer
  - SMTP handshaking (greeting)
  - SMTP transfer of messages
  - SMTP closure
- command/response interaction (like HTTP)
  - **commands:** ASCII text
  - **response:** status code and phrase



## Scenario: Alice sends e-mail to Bob

- 1) Alice uses UA to compose e-mail message "to" bob@someschool.edu
- 2) Alice's UA sends message to her mail server using SMTP; message placed in message queue
- 3) client side of SMTP at mail server opens TCP connection with Bob's mail server
- 4) SMTP client sends Alice's message over the TCP connection
- 5) Bob's mail server places the message in Bob's mailbox
- 6) Bob invokes his user agent to read message



## Sample SMTP interaction

S: 220 hamburger.edu

## SMTP: observations

### comparison with HTTP:

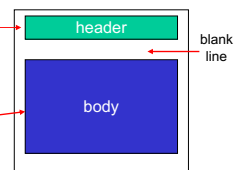
- HTTP: client pull
- SMTP: client push
- both have ASCII command/response interaction, status codes
- HTTP: each object encapsulated in its own response message
- SMTP: multiple objects sent in multipart message
- SMTP uses persistent connections
- SMTP requires message (header & body) to be in 7-bit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

## Mail message format

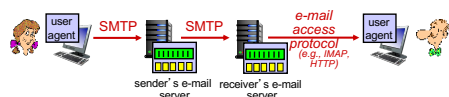
SMTP: protocol for exchanging e-mail messages, defined in RFC 5321 (like RFC 7231 defines HTTP)

RFC 2822 defines *syntax* for e-mail message itself (like HTML defines syntax for web documents)

- header lines, e.g.,
  - To:
  - From:
  - Subject:
 these lines, within the body of the email message area different from SMTP MAIL FROM:, RCPT TO: commands!
- Body: the "message", ASCII characters only



## Retrieving email: mail access protocols



- SMTP: delivery/storage of e-mail messages to receiver's server
- mail access protocol: retrieval from server
  - IMAP: Internet Mail Access Protocol [RFC 3501]: messages stored on server, IMAP provides retrieval, deletion, folders of stored messages on server
- HTTP: gmail, Hotmail, Yahoo!Mail, etc. provides web-based interface on top of SMTP (to send), IMAP (or POP) to retrieve e-mail messages

## Kitkats and questionnaire



Thank you and good luck.