# Networking

- Middleware gives guarantees not provided by networking
- How do you connect computers?
  - LAN
  - WAN

- Let us consider the example of the Internet
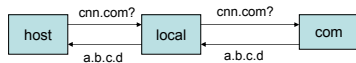
1

# Internet: Example

- Click -> get page
- specifies
  - protocol (http)
  - location
    (www.cnn.com)
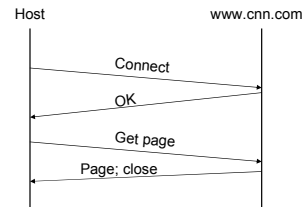


# Internet: Locating Resource

- www.cnn.com
  - name of a computer
  - Implicitly also a file
- Map name to IP address
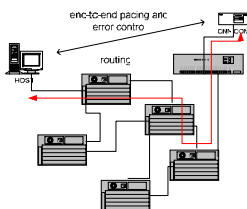  - DNS



3

# Internet: Connection

- Http sets up a connection (tcp)
  - between the host and cnn.com to transfer the page
- The connection transfers page as a byte stream
  - without errors: flow control + error control



4

# Internet: End-to-end

- Byte stream flows end to end across many links/switches:
  - routing (+ addressing)
- That stream is regulated and controlled by both ends:
  - retransmission of erroneous or missing bytes; flow control
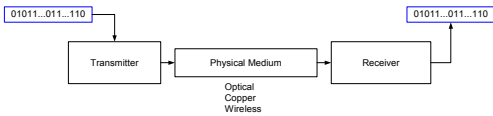


5

# Internet: Packets

- The network transports bytes grouped into packets
- Packets are "self-contained"; routers handle them 1 by 1
- The end hosts worry about errors and pacing
  - Destination sends ACKs; Source checks losses



6

## Internet: Bits

- Equipment in each node sends packets as string of bits
- That equipment is not aware of the meaning of the bits
- Frames (packetizing) vs. streams

01011...011...110 → Transmitter → Physical Medium → Receiver → 01011...011...110

Optical
Copper
Wireless

7

## Internet: Points to remember

- Separation of tasks
  - send bits on a link: transmitter/receiver [clock, modulation,…]
  - send packet on each hop [framing, error detection,…]
  - send packet end to end [addressing, routing]
  - pace transmissions [detect congestion]
  - retransmit erroneous or missing packets [acks, timeout]
  - find destination address from name [DNS]
- Scalability
  - routers don't know full path
  - names and addresses are hierarchical

8

## Internet : Challenges

- Addressing ?
- Routing ?
- Reliable transmission ?
- Interoperability ?
- Resource management ?
- Quality of service ?

9

## Concepts at heart of the Internet

- Protocol
- Layered Architecture
- Packet Switching
- Distributed Control
- Open System

10

## Protocol

- Two communicating entities must agree on:
  - Expected order and meaning of messages they exchange
  - The action to perform on sending/receiving a message

- Asking the time

11

## Layered Architectures

- Human beings can handle lots of complexity in their protocol processing.
  - Ambiguously defined protocols
  - Many protocols all at once
- How computers manage complex protocol processing?
  - Specify well defined protocols to enact.
  - Decompose complicated jobs into layers;
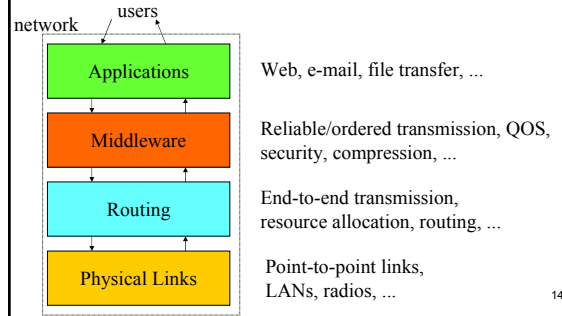    - each has a well defined task

12

# Layered Architectures

- Break-up design problem into smaller problems
  - More manageable
- Modular design: easy to extend/modify.
- Difficult to implement
  - careful with interaction of layers for efficiency

13

# Layered Architecture



network  users

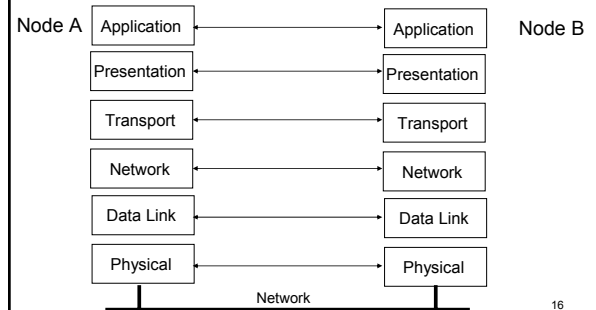| Applications | Web, e-mail, file transfer, ... |
| Middleware | Reliable/ordered transmission, QOS, security, compression, ... |
| Routing | End-to-end transmission, resource allocation, routing, ... |
| Physical Links | Point-to-point links, LANs, radios, ... |

14

# The OSI Model

- Open Systems Interconnect model is a standard way of understanding conceptual layers of network comm.
- This is a model, nobody builds systems like this.
- Each level provides certain functions and guarantees, and communicates with the same level on remote notes.
- A message is generated at the highest level, and is passed down the levels, encapsulated by lower levels, until it is sent over the wire.
- On the destination, it makes its way up the layers,until the high-level msg reaches its high-level destination.

15

# OSI Levels



Node A — Application, Presentation, Transport, Network, Data Link, Physical
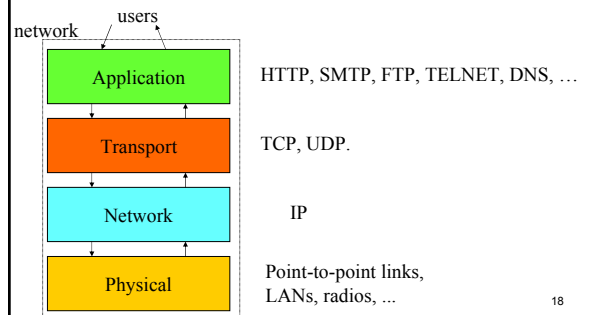Node B — Application, Presentation, Transport, Network, Data Link, Physical

Network

16

# OSI Levels

- Physical Layer: electrical details of bits on the wire
- Data Link: sending "frames" of bits and error detection
- Network Layer:" routing packets to the destination
- Transport Layer: reliable transmission of messages, disassembly/assembly, ordering, retransmission of lost packets
- Session Layer; really part of transport, typ. Not impl.
- Presentation Layer: data representation in the message
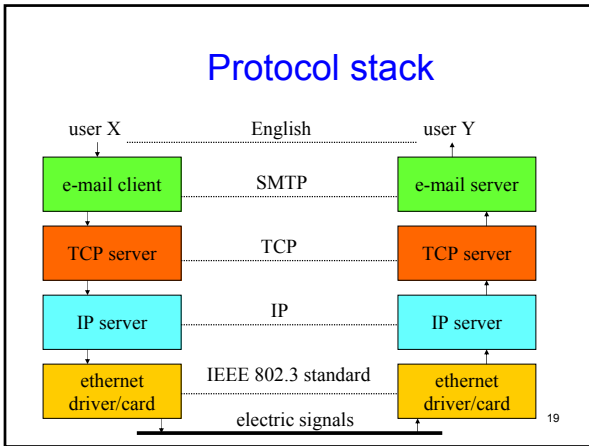- Application: high-level protocols (mail, ftp, etc.)
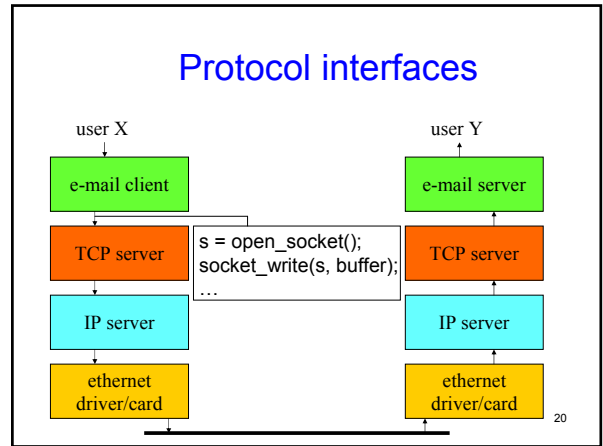
17

# Internet protocol stack



network  users

| Application | HTTP, SMTP, FTP, TELNET, DNS, … |
| Transport | TCP, UDP. |
| Network | IP |
| Physical | Point-to-point links, LANs, radios, ... |

18

## Protocol stack

| user X | English | user Y |
|--------|---------|--------|
| e-mail client | SMTP | e-mail server |
| TCP server | TCP | TCP server |
| IP server | IP | IP server |
| ethernet driver/card | IEEE 802.3 standard | ethernet driver/card |
| | electric signals | |

19

## Protocol interfaces

| user X | | user Y |
|--------|---|--------|
| e-mail client | | e-mail server |
| TCP server | s = open_socket();<br>socket_write(s, buffer);<br>… | TCP server |
| IP server | | IP server |
| ethernet driver/card | | ethernet driver/card |

20

## Addressing

- Each network interface has a hardware address
  - Multiple interfaces ⇒ multiple addresses
- Each application communicates via a *port*
  - Port is a logical connection endpoint
  - Allows multiple local applications to use network resources
  - Up to 65535
    - < 1024 : used by privileged applications
    - 1024 ≤ available for use ≤ 49151
    - 49152 ≤ Dynamic ports/private ports ≤ 65535
  - http ports 80 and 8080
  - telnet 23, ftp 21, etc
- Think of a telephone network …

21

## Addressing and Packet Format

- The ``Data'' segment contains higher level protocol information.
  - Which protocol is this packet destined for?
  - Which process is the packet destined for?
  - Which packet is this in a sequence of packets?
  - What kind of packet is this?
- This is the stuff of the OSI reference model.

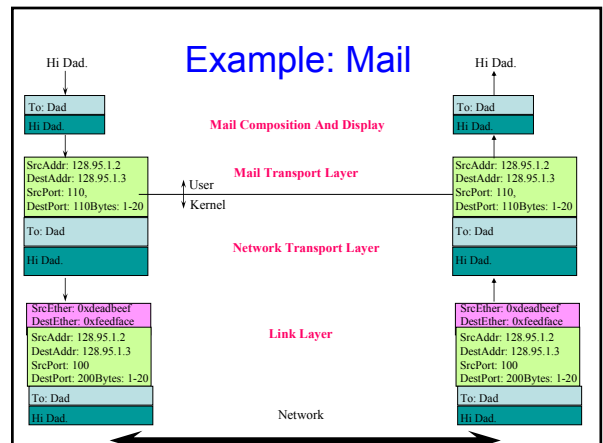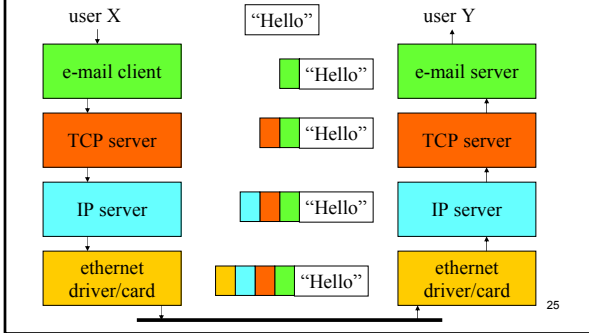| |
|---|
| Start (7 bytes) |
| Destination (6) |
| Source (6) |
| Length (2) |
| Msg Data (1500) |
| Checksum (4) |

22

## Ethernet packet dispatching

- An incoming packet comes into the Ethernet controller.
- The Ethernet controller reads it off the network into a buffer.
- It interrupts the CPU.
- A network interrupt handler reads the packet out of the controller into memory.
- A dispatch routine looks at the Data part and hands it to a higher level protocol
- The higher level protocol copies it out into user space.
- A program manipulates the data.
- The output path is similar.
- Consider what happens when you send mail.

23

## Example: Mail

Hi Dad.

| To: Dad |
| Hi Dad. |

Mail Composition And Display

| SrcAddr: 128.95.1.2 |
| DestAddr: 128.95.1.3 |
| SrcPort: 110, |
| DestPort: 110Bytes: 1-20 |

Mail Transport Layer

User / Kernel

| To: Dad |
| Hi Dad. |

Network Transport Layer

| SrcEther: 0xdeadbeef |
| DestEther: 0xfeedface |
| SrcAddr: 128.95.1.2 |
| DestAddr: 128.95.1.3 |
| SrcPort: 100 |
| DestPort: 200Bytes: 1-20 |

Link Layer

| To: Dad |
| Hi Dad. |

Network

Hi Dad.

| To: Dad |
| Hi Dad. |

| SrcAddr: 128.95.1.2 |
| DestAddr: 128.95.1.3 |
| SrcPort: 110, |
| DestPort: 110Bytes: 1-20 |

| To: Dad |
| Hi Dad. |

| SrcEther: 0xdeadbeef |
| DestEther: 0xfeedface |
| SrcAddr: 128.95.1.2 |
| DestAddr: 128.95.1.3 |
| SrcPort: 100 |
| DestPort: 200Bytes: 1-20 |

| To: Dad |
| Hi Dad. |

## Protocol encapsulation



user X    "Hello"    user Y

e-mail client — "Hello" — e-mail server

TCP server — "Hello" — TCP server

IP server — "Hello" — IP server

ethernet driver/card — "Hello" — ethernet driver/card

25

---

## End-to-End Argument

- What function to implement in each layer?
- Saltzer, Reed, Clarke 1984
  - A function can be correctly and completely implemented only with the knowledge and help of applications standing at the communication endpoints
  - Argues for moving function upward in a layered architecture
- Should the network guarantee packet delivery ?
  - Think about a file transfer program
  - Read file from disk, send it, the receiver reads packets and writes them to the disk

26

---

## End-to-End Argument

- If the network guaranteed packet delivery
  - one might think that the applications would be simpler
    - No need to worry about retransmits
  - But need to check that file was written to the remote disk intact
- A check is necessary if nodes can fail
  - Consequently, applications need to perform their retransmits
- No need to burden the internals of the network with properties that can, and must, be implemented at the periphery

27

---

## End-to-End Argument

- An Occam's razor for Internet design
  - If there is a problem, the simplest explanation is probably the correct one
- Application-specific properties are best provided by the applications, not the network
  - Guaranteed, or ordered, packet delivery, duplicate suppression, security, etc.
- The internet performs the simplest packet routing and delivery service it can
  - Packets are sent on a best-effort basis
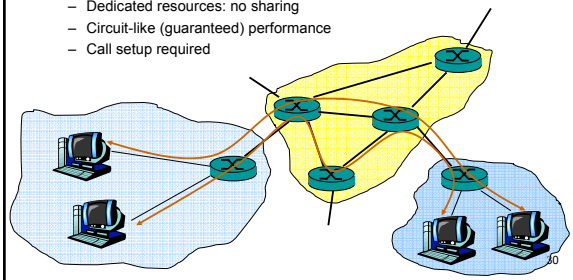  - Higher-level applications do the rest

28

---

## Two ways to handle networking

- Circuit Switching
  - What you get when you make a phone call
  - Dedicated circuit per call
- Packet Switching
  - What you get when you send a bunch of letters
  - Network bandwidth consumed only when sending
  - Packets are routed independently
- Message Switching
  - It's just packet switching, but routers perform store-and-forward

29

---

## Circuit Switching

- End-to-end resources reserved for "call"
  - Link bandwidth, switch capacity
  - Dedicated resources: no sharing
  - Circuit-like (guaranteed) performance
  - Call setup required



30

## Packet Switching

- Each end-to-end data stream divided into packets
  - User's packets *share* network resources
    - Compared to dedicated allocation
  - Each packet uses full link bandwidth
    - Compared to dividing bandwidth into pieces
  - Resources are used as needed
    - Compared to resource reservation
- Resource contention:
  - Aggregate demand can exceed amount available
  - Congestion: packets queue, wait for link use
  - Store and forward: packets move one hop at a time
    - Transmit over link
    - Wait turn at next link

31

## Routing

- Goal: move data among routers from source to dest.
- Datagram packet network:
  - Destination address determines next hop
  - Routes may change during session
  - Analogy: driving, asking directions
  - No notion of call state
- Circuit-switched network:
  - Call allocated time slots of bandwidth at each link
  - Fixed path (for call) determined at call setup
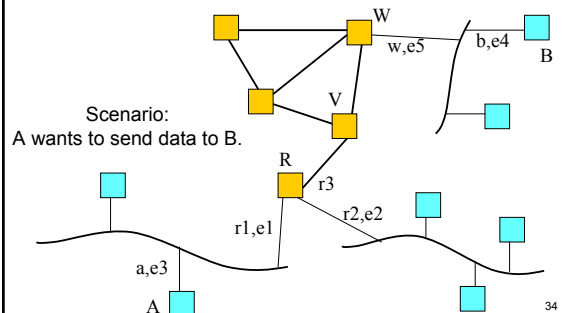  - Switches maintain lots of per call state: resource allocation

32

## Packet vs. Circuit Switching

- Reliability: no congestion, in-order data in circuit-switch
- Packet switching: better bandwidth use
- State, resources: packet switching has less state
  - Good: less control plane processing resources along the way
  - More data plane (address lookup) processing
- Failure modes (routers/links down)
  - Packet switch reconfigures sub-second timescale
  - Circuit switching: more complicated
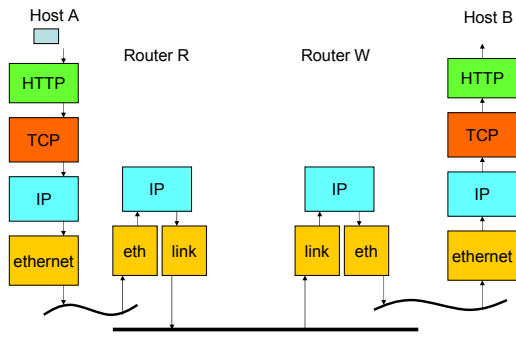    - Involves all switches in the path

33

## A small Internet



Scenario:
A wants to send data to B.

34

## Packet forwarding



35

## The Link Layer

# What is purpose of this layer?

- Physically encode bits on the wire
- Link = pipe to send information
  - E.g. point to point or broadcast



- Can be built out of:
  - Twisted pair, coaxial cable, optical fiber, radio waves, etc
- Links should only be able to send data
  - Could corrupt, lose, reorder, duplicate, (fail in other ways)

37

# How to connect routers/machines?

- WAN/Router Connections
  - Commercial:
    - T1 (1.5 Mbps), T3 (44 Mbps)
    - OC1 (51 Mbps), OC3 (155 Mbps)
    - ISDN (64 Kbps)
    - Frame Relay (1-100 Mbps, usually 1.5 Mbps)
    - ATM (some Gbps)
  - To your home:
    - DSL
    - Cable
- Local Area:
  - Ethernet: IEEE 802.3 (10 Mbps, 100 Mbps, 1 Gbps)
  - Wireless: IEEE 802.11 b/g/a (11 Mbps, 22 Mbps, 54 Mbps)

38

# Link level Issues

- Encoding: map bits to analog signals
- Framing: Group bits into frames (packets)
- Arbitration: multiple senders, one resource
- Addressing: multiple receivers, one wire

39

# Addressing & ARP



128.84.96.89   128.84.96.90

128.84.96.91

"What is the physical address of the host named 128.84.96.89"

"I'm at 1a:34:2c:9a:de:cc"

- ARP is used to discover physical addresses
  - ARP = Address Resolution Protocol

40

# Addressing & RARP



???   128.84.96.90
RARP Server

128.84.96.91

"I just got here. My physical address is 1a:34:2c:9a:de:cc. What's my name ?"
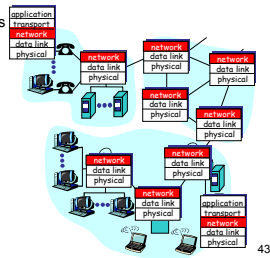
"Your name is 128.84.96.89"

- RARP is used to discover virtual addresses
  - RARP = Reverse Address Resolution Protocol

41

# The Network Layer

## Purpose of Network layer

- Given a packet, send it across the network to destination
- 2 key issues:
  - Portability:
    - connect different technologies
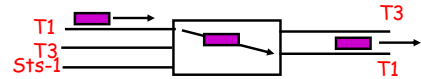  - Scalability
    - To the Internet scale

---

## What does it involve?

Two important functions:
- *routing:* determine path from source to dest.
- *forwarding:* move packets from router's input to output

T1
T3
Sts-1

T3

T1

---

## Network service model

Q: What *service model* for "channel" transporting packets from sender to receiver?

service abstraction

- guaranteed bandwidth?
- preservation of inter-packet timing (no jitter)?
- loss-free delivery?
- in-order delivery?
- congestion feedback to sender?

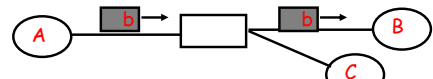The most important abstraction provided by network layer:

virtual circuit
or
datagram?

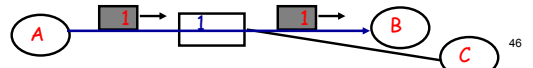Which things can be "faked" at the transport layer?

---

## Two connection models

- Connectionless (or "datagram"):
  - each packet contains enough information that routers can decide how to get it to its final destination
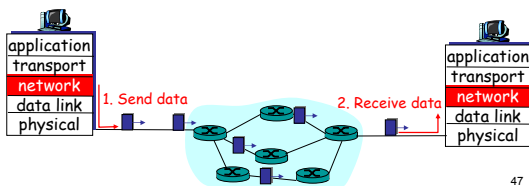
A    b    b    B
          C

- Connection-oriented (or "virtual circuit"):
  - first set up a connection between two nodes
  - label it (called a virtual circuit identifier (VCI))
  - all packets carry label

A    1    1    1    B
                    C

---

## Datagram networks

- no call setup at network layer
- routers: no state about end-to-end connections
  - no network-level concept of "connection"
- packets typically routed using destination host ID
  - packets between same source-dest pair may take different paths
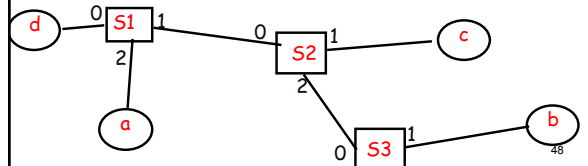- Best effort: data corruption, packet drops, route loops

application
transport
network
data link
physical

1. Send data    2. Receive data

application
transport
network
data link
physical

---

## Datagrams: Forwarding

How does packet get to the destination?
- switch creates a "forwarding table", mapping destinations to output port (ignores input ports)
- when a packet with a destination address in the table arrives, it pushes it out on the appropriate output port
- when a packet with a destination address not in the table arrives, it must find out more routing information (next problem)

d    0  S1  1        0  S2  1        c
        2                  2

     a              0  S3  1    b

# Datagrams

- Plusses:
  - No round trip connection setup time
  - No explicit route teardown
  - No resource reservation $\Rightarrow$ each flow could get max bandwidth
  - Easily handles switch failures; routes around it
- Minuses
  - Difficult to provide resource guarantees
  - Higher per packet overhead

- Internet uses datagrams: IP (Internet Protocol)
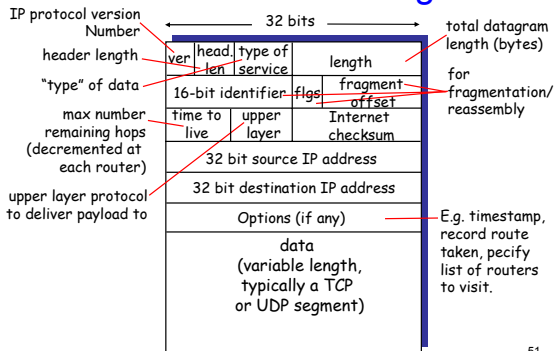
49

# IP addressing: CIDR

- Classless InterDomain Routing
  - network portion of address of arbitrary length
  - address format: a.b.c.d/x, where x is # bits in network portion

```
        network                    host
         part                      part
11001000 00010111 00010000 00000000
```

200.23.16.0/23

  - Examples:
    - Class A: /8
    - Class B: /16
    - Class C: /24

50

# Internet Protocol Datagram

IP protocol version Number
header length
"type" of data
max number remaining hops (decremented at each router)
upper layer protocol to deliver payload to

total datagram length (bytes)
for fragmentation/ reassembly
E.g. timestamp, record route taken, pecify list of routers to visit.

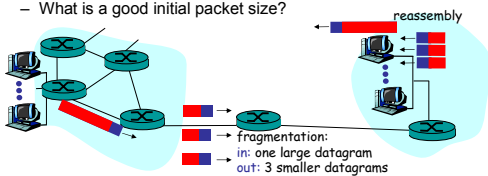| 32 bits | | |
|---|---|---|
| ver | head. len | type of service | length |
| 16-bit identifier | flgs | fragment offset |
| time to live | upper layer | Internet checksum |
| 32 bit source IP address | | |
| 32 bit destination IP address | | |
| Options (if any) | | |
| data (variable length, typically a TCP or UDP segment) | | |

51

# Datagram Portability

- IP Goal: To create one logical network from multiple physical networks
  - All intermediate routers should understand IP
  - IP header information sufficient to carry the packet to destination
  - Goal: Run over anything!
- Problem:
  - Physical networks have different MTUs
  - "max. transmission unit": 1500 for Ethernet, 48 for ATM
- Solution 1:
  - Fit everything in the MTU (!)

52

# IP Fragmentation & Reassembly

- Solution 2: (the one used)
  - If packet size > MTU of network, then fragment into pieces
    - Each fragment is less than MTU size
    - Each has IP headers + frag bit set + frag id + offset
  - Packets may get refragmented on the way to destination
  - Reassembly only done at the destination
  - What is a good initial packet size?



reassembly

fragmentation:
in: one large datagram
out: 3 smaller datagrams

53

# Internet: Names and Addresses

## Naming in the Internet

- What are named? All Internet Resources.
  - Objects: www.cs.cornell.edu/pages/ranveer
  - Services: weather.yahoo.com/forecast
  - Hosts: planetlab1.cs.cornell.edu

- Characteristics of Internet Names
  - human recognizable
  - unique
  - persistent

- Universal Resource Names (URNs)

## Locating the resources

- Internet services and resources are provided by end-hosts
  - ex. www1.cs.cornell.edu and www2.cs.cornell.edu host Ranveer's home page.

- Names are mapped to Locations
  - Universal Resource Locators (URL)
  - Embedded in the name itself: ex. weather.yahoo.com/forecast

- Semantics of Internet naming
  - ✓ human recognizable
  - ✓ uniqueness
  - x persistent

## Locating the Hosts?

- Internet Protocol Addresses (IP Addresses)
  - ex. planetlab1.cs.cornell.edu → 128.84.154.49

- Characteristics of IP Addresses
  - 32 bit fixed-length
  - enables network routers to efficiently handle packets in the Internet

- Locating services on hosts
  - port numbers (16 bit unsigned integer) 65536 ports
  - standard ports: HTTP 80, FTP 20, SSH 22, Telnet 20

## Mapping Not 1 to 1

- One host may map to more than one name
  - One server machine may be the web server (www.foo.com), mail server (mail.foo.com)etc.

- One host may have more than one IP address
  - IP addresses are per network interface

- But IP addresses are generally unique!
  - two globally visible machines should not have the same IP address
  - Anycast is an Exception:
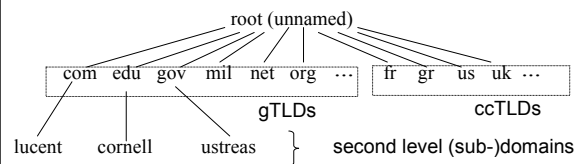    - routers send packets dynamically to the closest host matching an anycast address

## How to get a name?

- Naming in Internet is Hierarchical
  - decreases centralization
  - improves name space management

- First, get a domain name then you are free to assign sub names in that domain
  - How to get a domain name coming up

- Example: weather.yahoo.com belongs to yahoo.com which belongs to .com
  - regulated by global non-profit bodies

## Domain name structure

root (unnamed)

com edu gov mil net org ···    fr gr us uk ···

gTLDs                          ccTLDs

lucent  cornell  ustreas  }  second level (sub-)domains

gTLDs= Generic Top Level Domains
ccTLDs = Country Code Top Level Domains

# Top-level Domains     (TLDs)

- Generic Top Level Domains (gTLDs)
  - .com - commercial organizations
  - .org - not-for-profit organizations
  - .edu - educational organizations
  - .mil - military organizations
  - .gov - governmental organizations
  - .net - network service providers
  - New: .biz, .info, .name, …
- Country code Top Level Domains (ccTLDs)
  - One for each country

# How to get a domain name?

- In 1998, non-profit corporation, Internet Corporation for Assigned Names and Numbers (ICANN), was formed to assume responsibility from the US Government
- ICANN authorizes other companies to register domains in com, org and net and new gTLDs
  - Network Solutions is largest and in transitional period between US Govt and ICANN had sole authority to register domains in com, org and net

# How to get an IP Address?

- Answer 1: Normally, answer is get an IP address from your upstream provider
  - This is essential to maintain efficient routing!
- Answer 2: If you need lots of IP addresses then you can acquire your own block of them.
  - IP address space is a scarce resource - must prove you have fully utilized a small block before can ask for a larger one and pay $$ (Jan 2002 - $2250/year for /20 and $18000/year for a /14)

# How to get lots of IP Addresses? Internet Registries

RIPE NCC (Riseaux IP Europiens Network Coordination Centre) for Europe, Middle-East, Africa

APNIC (Asia Pacific Network Information Centre )for Asia and Pacific

ARIN (American Registry for Internet Numbers) for the Americas, the Caribbean, sub-saharan Africa

Note: Once again regional distribution is important for efficient routing!

Can also get Autonomous System Numnbers (ASNs from these registries

# Are there enough addresses?

- Unfortunately No!
  - 32 bits → 4 billion unique addresses
  - but addresses are assigned in chunks
  - ex. cornell has four chunks of /16 addressed
    - ex. 128.84.0.0 to 128.84.255.255
    - 128.253.0.0, 128.84.0.0, 132.236.0.0, and 140.251.0.0
- Expanding the address space!
  - IPv6 128 bit addresses
  - difficult to deploy (requires cooperation and changes to the core of the Internet)
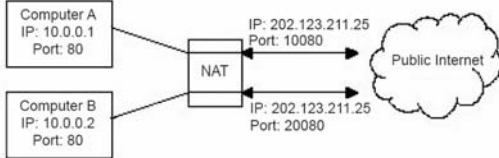
# DHCP and NATs

- Dynamic Host Control Protocol
  - lease IP addresses for short time intervals
  - hosts may refresh addresses periodically
  - ♥ only live hosts need valid IP addresses

- Network Address Translators
  - Hide local IP addresses from rest of the world
  - only a small number of IP addresses are visible outside
  - ♥ solves address shortage for all practical purposes
  - ✢ access is highly restricted
    - ex. peer-to-peer communication is difficult

## NATs in operation

- Translate addresses when packets traverse through NATs
- Use port numbers to increase number of supportable flows



67

## DNS: Domain Name System

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

68

## DNS name servers

How could we provide this service? Why not centralize DNS?
- single point of failure
- traffic volume
- distant centralized database
- maintenance

doesn't *scale!*

- no server has all name-to-IP address mappings

Name server: process running on a host that processes DNS requests
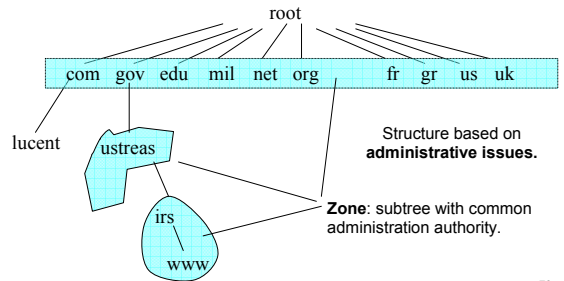
local name servers:
  – each ISP, company has *local (default) name server*
  – host DNS query first goes to local name server

authoritative name server:
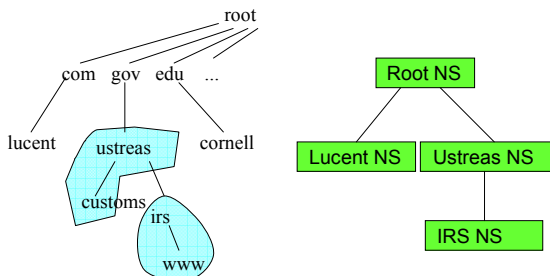  – can perform name/address translation for a specific domain or zone

69

## Name Server Zone Structure



Structure based on **administrative issues.**

**Zone**: subtree with common administration authority.

70

## Name Servers (NS)



71

## Name Servers (NS)

- NSs are **duplicated** for reliability.
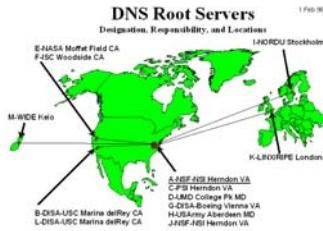- Each domain must have a primary and secondary.
- Anonymous ftp from:
  ftp.rs.internic.net, netinfo/root-server.txt
  gives the current root NSs (about 10).
- Each host knows the IP address of the **local** NS.
- Each NS knows the IP addresses of all root NSs.

72

## DNS: Root name servers

- contacted by local name server that can not resolve name
- root name server:
  - Knows the authoritative name server for main domain
- ~ 60 root name servers worldwide
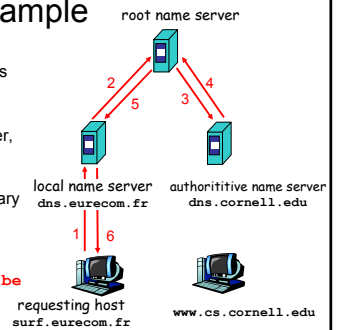  - real-world application of anycast

**DNS Root Servers**
Designation, Responsibility, and Locations

I-NORDU Stockholm
E-NASA Moffet Field CA
F-ISC Woodside CA
M-WIDE Keio
K-LINX/RIPE London
A-NSF-NSI Herndon VA
C-PSI Herndon VA
D-UMD College Pk MD
G-DISA-Boeing Vienna VA
H-USArmy Aberdeen MD
J-NSF-NSI Herndon VA
B-DISA-USC Marina delRey CA
L-DISA-USC Marina delRey CA

73

---

## Simple DNS example

root name server

host `surf.eurecom.fr` wants IP address of `www.cs.cornell.edu`

1. Contacts its local DNS server, `dns.eurecom.fr`
2. `dns.eurecom.fr` contacts root name server, if necessary
3. root name server contacts authoritative name server, `dns.cornell.edu`, if necessary **(what might be wrong with this?)**
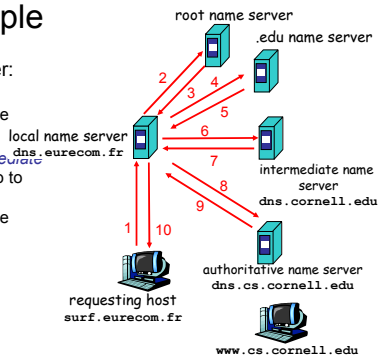
local name server
`dns.eurecom.fr`

authoritive name server
`dns.cornell.edu`

requesting host
`surf.eurecom.fr`

`www.cs.cornell.edu`

74

---

## DNS example

root name server
.edu name server

Root name server:
- may not know authoritative name server
- may know *intermediate name server:* who to contact to find authoritative name server

local name server
`dns.eurecom.fr`

intermediate name server
`dns.cornell.edu`

authoritative name server
`dns.cs.cornell.edu`

requesting host
`surf.eurecom.fr`

`www.cs.cornell.edu`

75

---

## DNS Architecture

- Hierarchical Namespace Management
  - domains and sub-domains
  - distributed and localized authority
- Authoritative Nameservers
  - server mappings for specific sub-domains
  - more than one (at least two for failure resilience)
- Caching to mitigate load on root servers
  - time-to-live (ttl) used to delete expired cached mappings
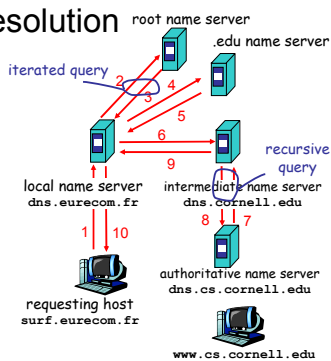
76

---

## DNS: query resolution

root name server
.edu name server

**iterated query:**
- contacted server replies with name of server to contact
- "I don't know this name, but ask this server"
- Takes burden off root servers

iterated query

recursive query

**recursive query:**
- puts burden of name resolution on contacted name server
- reduces latency

local name server
`dns.eurecom.fr`

intermediate name server
`dns.cornell.edu`

authoritative name server
`dns.cs.cornell.edu`

requesting host
`surf.eurecom.fr`

`www.cs.cornell.edu`

77

---

## DNS records: More than Name to IP Address

DNS: distributed db storing resource records (RR)

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

- Type=A
  - **name** is hostname
  - **value** is IP address
  - One we've been discussing; most common
- Type=NS
  - **name** is domain (e.g. foo.com)
  - **value** is IP address of authoritative name server for this domain

- Type=CNAME
  - **name** is an alias name for some "cannonical" (the real) name
  - **value** is cannonical name
- Type=MX
  - **value** is hostname of mailserver associated with **name**

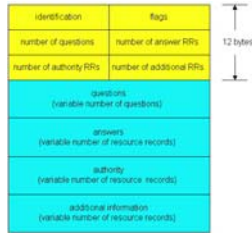78

## DNS protocol, messages

<u>DNS protocol :</u> *query* and *repy* messages, both with same *message format*

msg header
- identification: 16 bit # for query, repy to query uses same #
- flags:
  – query or reply
  – recursion desired
  – recursion available
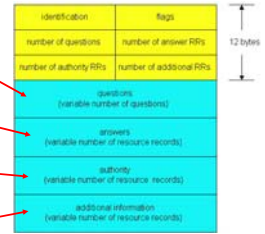  – reply is authoritative
  – reply was truncated



79

## DNS protocol, messages

Name, type fields for a query

RRs in reponse to query

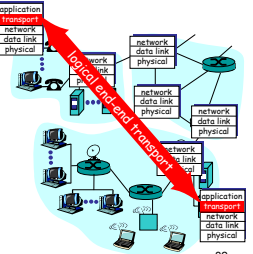records for authoritative servers

additional "helpful" info that may be used



80

# The Transport Layer

# Purpose of this layer

- Interface end-to-end applications and protocols
  – Turn best-effort IP into a usable interface
- Data transfer b/w processes:
  – Compared to end-to-end IP
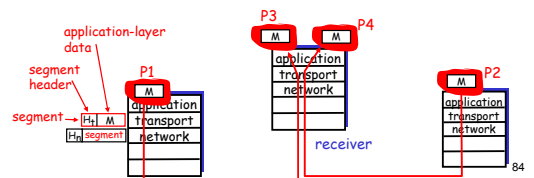- We will look at 2:
  – TCP
  – UDP



82

# UDP

- **U**nreliable **D**atagram **P**rotocol
- Best effort data delivery between processes
  – No frills, bare bones transport protocol
  – Packet may be lost, out of order
- Connectionless protocol:
  – No handshaking between sender and receiver
  – Each UDP datagram handled independently
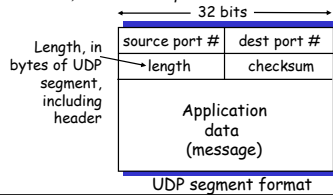
83

# UDP Functionality

- Multiplexing/Demultiplexing
  – Using ports
- Checksums (optional)
  – Check for corruption



84

## Multiplexing/Demultiplexing

- Multiplexing:
  - Gather data from multiple processes, envelope data with header
  - Header has src port, dest port for multiplexing
    - Why not process id?
- Demultiplexing:
  - Separate incoming data in machine to different applications
  - Demux based on *sender addr, src and dest port*



| ← 32 bits → |  |
|---|---|
| source port # | dest port # |
| length | checksum |
| Application data (message) | |

Length, in bytes of UDP segment, including header

UDP segment format

## Implementing Ports

- As a message queue
  - Append incoming message to the end
  - Much like a mailbox file
- If queue full, message can be discarded
- When application reads from socket
  - OS removes some bytes from the head of the queue
- If queue empty, application blocks waiting

86

## UDP Checksum

- Over the headers and data
  - Ensures integrity end-to-end
  - 1's complement sum of segment contents
- Is optional in UDP
- If checksum is non-zero, and receiver computes another value:
  - Silently drop the packet, no error message detected

87

## UDP Discussion

- Why UDP?
  - No delay in connection establishment
  - Simple: no connection state
  - Small header size
  - No congestion control: can blast packets
- Uses:
  - Streaming media, DNS, SNMP
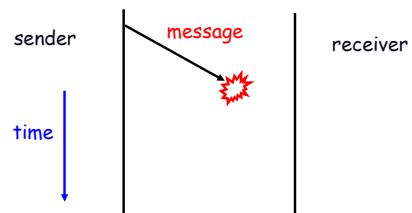  - Could add application specific error recovery

88

## TCP

- Transmission Control Protocol
  - Reliable, in-order, process-to-process, two-way byte stream
- Different from UDP
  - Connection-oriented
  - Error recovery: Packet loss, duplication, corruption, reordering
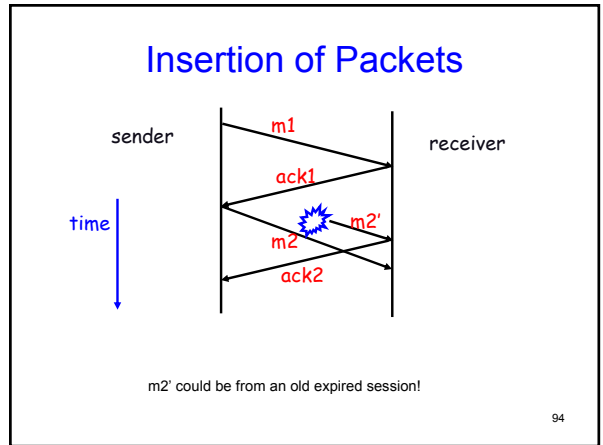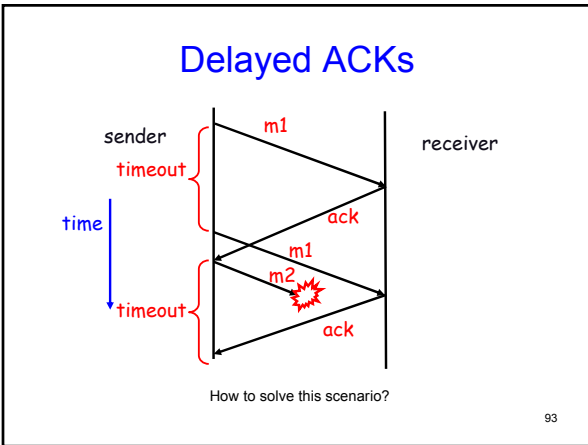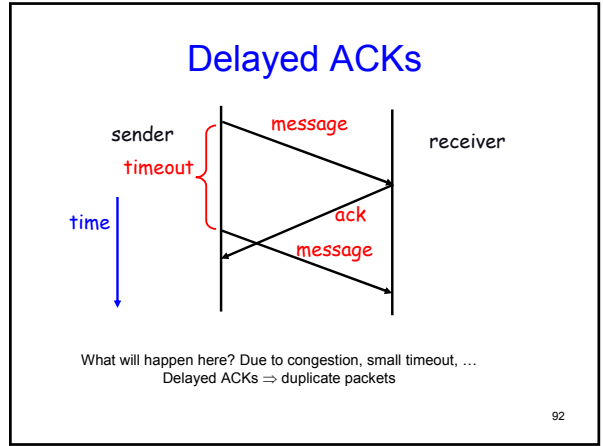- A number of applications require this guarantee
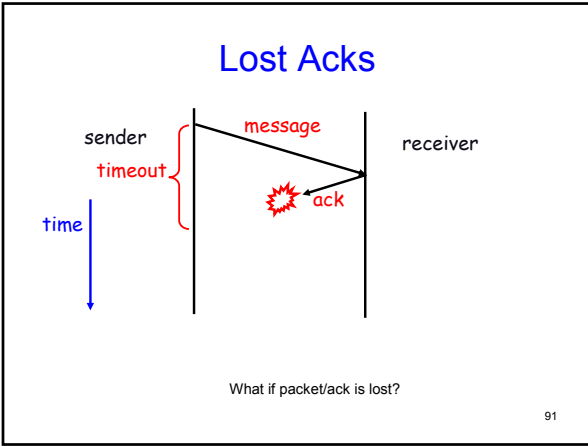  - Web browsers use TCP

89

## Handling Packet Loss



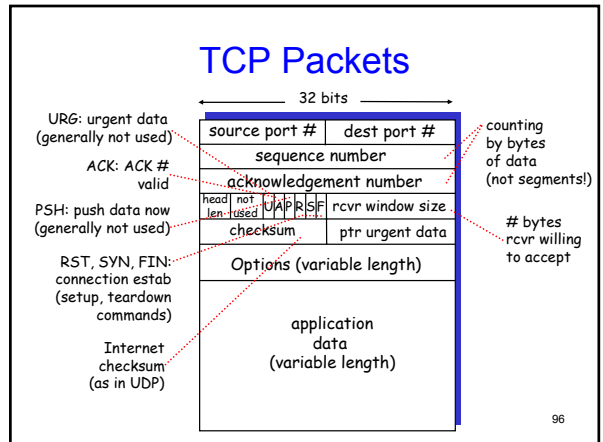sender     message     receiver

time

There are a number of reasons why the packet may get lost:
- router congestion, lossy medium, etc.
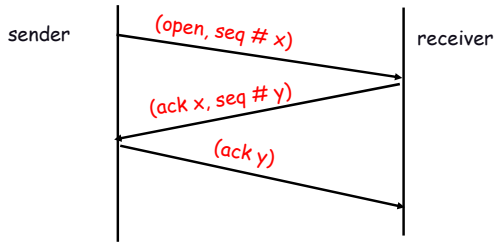
How does sender know of a successful packet send?

90

## Lost Acks

sender
timeout
receiver
time
message
ack

What if packet/ack is lost?

91

## Delayed ACKs

sender
timeout
receiver
time
message
ack
message

What will happen here? Due to congestion, small timeout, …
Delayed ACKs ⇒ duplicate packets

92

## Delayed ACKs

sender
timeout
receiver
time
m1
ack
m1
m2
timeout
ack

How to solve this scenario?

93

## Insertion of Packets

sender
receiver
time
m1
ack1
m2
m2'
ack2

m2' could be from an old expired session!

94

## Message Identifiers

- Each message has <message id, session id>
  - Message id: uniquely identifies message in sender
  - Session id: unique across sessions
- Message ids detect duplication, reordering
- Session ids detect packet from old sessions
- TCP's sequence number has similar functionality:
  - Initial number chosen randomly
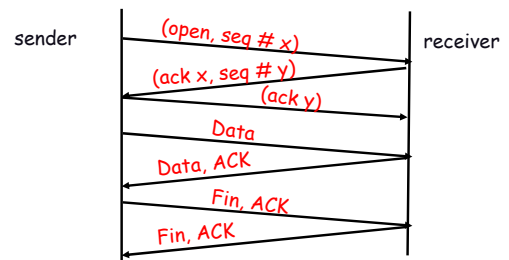  - Unique across packets
  - Incremented by length of data bytes

95

## TCP Packets

32 bits

URG: urgent data
(generally not used)

ACK: ACK #
valid

PSH: push data now
(generally not used)

RST, SYN, FIN:
connection estab
(setup, teardown
commands)

Internet
checksum
(as in UDP)

source port #    dest port #
sequence number
acknowledgement number
head | not | U A P R S F | rcvr window size
len | used |
checksum    ptr urgent data
Options (variable length)

application
data
(variable length)

counting
by bytes
of data
(not segments!)

# bytes
rcvr willing
to accept

96

# TCP Connection Establishment

sender → (open, seq # x) → receiver

(ack x, seq # y)

(ack y)

TCP is connection-oriented. Starts with a 3-way handshake.
Protects against duplicate SYN packets.

97

---

# TCP Usage

sender → (open, seq # x) → receiver

(ack x, seq # y)

(ack y)

Data

Data, ACK

Fin, ACK

Fin, ACK

98

---

# TCP timeouts

• What is a good timeout period ?
– Want to improve throughput without unnecessary transmissions

NewAverageRTT = (1 - $\alpha$) OldAverageRTT + $\alpha$ LatestRTT
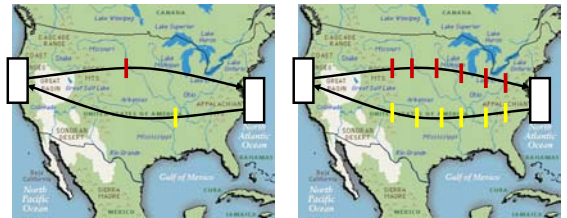NewAverageDev = (1 - $\alpha$) OldAverageDev + $\alpha$ LatestDev
where LatestRTT = (ack_receive_time – send_time),
   LatestDev = |LatestRTT – AverageRTT|,
   $\alpha$ = 1/8, typically.
Timeout = AverageRTT + 4*AverageDev

• Timeout is thus a function of RTT and deviation

99

---

# TCP Windows



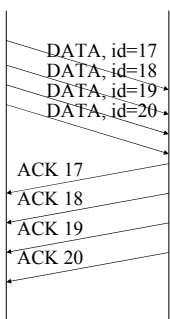• Multiple outstanding packets can increase throughput

100

---

# TCP Windows

DATA, id=17
DATA, id=18
DATA, id=19
DATA, id=20

ACK 17
ACK 18
ACK 19
ACK 20

• Can have more than one packet in transit
• Especially over fat pipes, e.g. satellite connection
• Need to keep track of all packets within the window
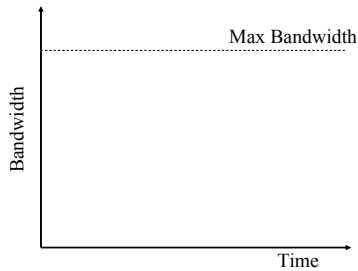• Need to adjust window size

101

---

# TCP Congestion Control

• TCP increases its window size when no packets dropped
• It halves the window size when a packet drop occurs
– A packet drop is evident from the acknowledgements
• Therefore, it slowly builds to the max bandwidth, and hover around the max
– It doesn't achieve the max possible though
– Instead, it shares the bandwidth well with other TCP connections
• This linear-increase, exponential backoff in the face of congestion is termed *TCP-friendliness*
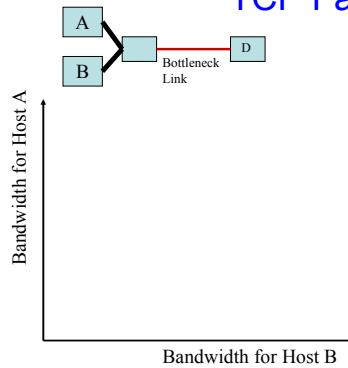
102

## TCP Window Size



Bandwidth vs Time graph with Max Bandwidth line

- Linear increase
- Exponential backoff

- Assuming no other losses in the network except those due to bandwidth

103

## TCP Fairness



Diagram: A and B connected through Bottleneck Link to D. Axes: Bandwidth for Host A (vertical), Bandwidth for Host B (horizontal)

- Want to share the bottleneck link fairly between two flows

104

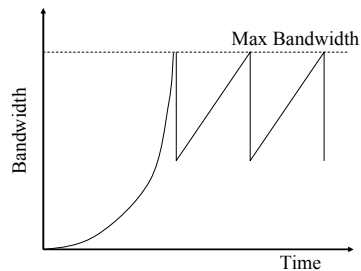## TCP Slow Start

- Linear increase takes a long time to build up a window size that matches the link bandwidth*delay
- Most file transactions are not long enough
- Consequently, TCP can spend a lot of time with small windows, never getting the chance to reach a sufficiently large window size
- Fix: Allow TCP to build up to a large window size initially by doubling the window size until first loss

105

## TCP Slow Start



Bandwidth vs Time graph with Max Bandwidth line

- Initial phase of exponential increase

- Assuming no other losses in the network except those due to bandwidth

106

## TCP Summary

- Reliable ordered message delivery
  - Connection oriented, 3-way handshake
- Transmission window for better throughput
  - Timeouts based on link parameters
- Congestion control
  - Linear increase, exponential backoff
- Fast adaptation
  - Exponential increase in the initial phase

107