# Chapter 2: outline

# DNS: domain name system

*people:* many identifiers:
- SSN, name, passport #

*Internet hosts, routers:*
- IP address (32 bit) - used for addressing datagrams
- "name", e.g., www.yahoo.com - used by humans

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

*Domain Name System:*

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

# DNS: services, structure

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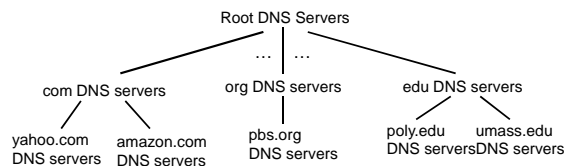  - replicated Web servers: many IP addresses correspond to one name

*why not centralize DNS?*
- single point of failure
- traffic volume
- distant centralized database
- maintenance

*A: doesn't scale!*

# DNS: a distributed, hierarchical database

Root DNS Servers

... ...

com DNS servers        org DNS servers        edu DNS servers

yahoo.com   amazon.com        pbs.org              poly.edu     umass.edu
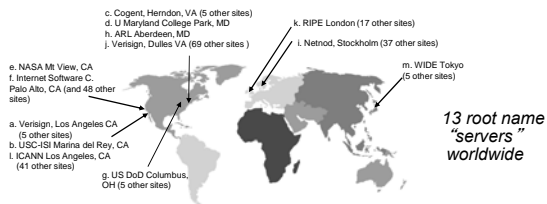DNS servers DNS servers       DNS servers          DNS serversDNS servers

*client wants IP for www.amazon.com; 1st approx:*
- 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

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

c. Cogent, Herndon, VA (5 other sites)
d. U Maryland College Park, MD
h. ARL Aberdeen, MD
j. Verisign, Dulles VA (69 other sites )

k. RIPE London (17 other sites)

i. Netnod, Stockholm (37 other sites)

e. NASA Mt View, CA
f. Internet Software C.
Palo Alto, CA (and 48 other sites)

m. WIDE Tokyo
(5 other sites)

a. Verisign, Los Angeles CA
(5 other sites)
b. USC-ISI Marina del Rey, CA
l. ICANN Los Angeles, CA
(41 other sites)

g. US DoD Columbus,
OH (5 other sites)

*13 root name
"servers"
worldwide*

## TLD, authoritative servers

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

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

## Local DNS name server

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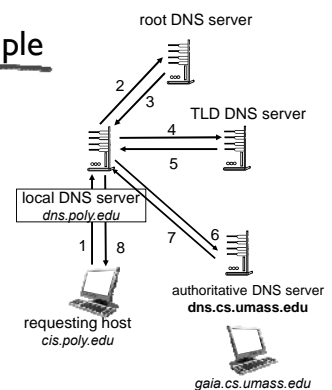  ▪ acts as proxy, forwards query into hierarchy

## DNS name resolution example

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

*iterated query:*
❖ contacted server replies with name of server to contact
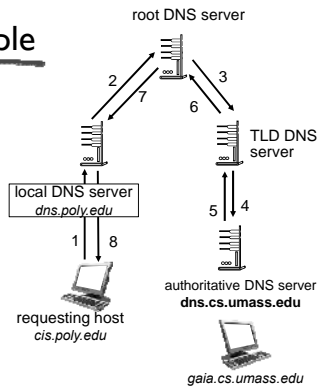❖ "I don't know this name, but ask this server"

root DNS server

2
3

TLD DNS server
4
5

local DNS server
*dns.poly.edu*

1   8        7   6

requesting host
*cis.poly.edu*

authoritative DNS server
**dns.cs.umass.edu**

*gaia.cs.umass.edu*

## DNS name resolution example

*recursive query:*

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

root DNS server

TLD DNS server

local DNS server
*dns.poly.edu*

requesting host
*cis.poly.edu*

authoritative DNS server
**dns.cs.umass.edu**

*gaia.cs.umass.edu*

## DNS: caching, updating records

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

## DNS records

*DNS:* distributed db storing resource records (RR)

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

### type=A
- ▪ **name** is hostname
- ▪ **value** is IP address

### type=NS
- ▪ **name** is domain (e.g., foo.com)
- ▪ **value** is hostname of authoritative name server for this domain

### type=CNAME
- ▪ **name** is alias name for some "canonical" (the real) name
- ▪ **www.ibm.com** is really **servereast.backup2.ibm.com**
- ▪ **value** is canonical name

### type=MX
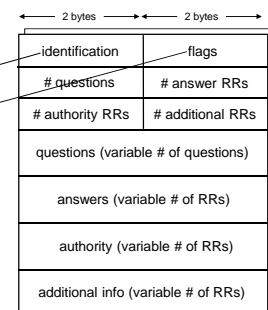- ▪ **value** is name of mailserver associated with **name**

## DNS protocol, messages

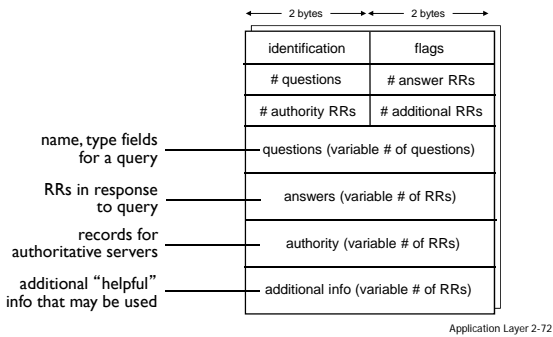- ❖ *query* and *reply* messages, both with same *message format*

msg 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

| 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) | |

name, type fields for a query — questions (variable # of questions)

RRs in response to query — answers (variable # of RRs)

records for authoritative servers — authority (variable # of RRs)

additional "helpful" info that may be used — additional info (variable # of RRs)

## Inserting records into 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 two RRs into .com TLD server:
     `(networkutopia.com, dns1.networkutopia.com, NS)`
     `(dns1.networkutopia.com, 212.212.212.1, A)`
❖ create authoritative server type A record for www.networkuptopia.com; type MX record for networkutopia.com

## Chapter 2: outline

2.1 principles of network applications
   ▪ app architectures
   ▪ app requirements
2.2 Web and HTTP
2.3 FTP
2.4 electronic mail
   ▪ SMTP, POP3, IMAP
2.5 DNS

2.6 P2P applications
2.7 socket programming with UDP and TCP

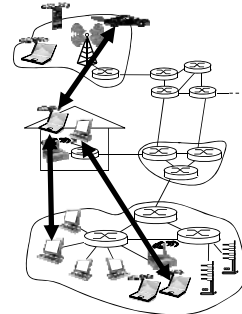## Pure P2P architecture

❖ *no* always-on server
❖ arbitrary end systems directly communicate
❖ peers are intermittently connected and change IP addresses

*examples:*
   ▪ file distribution (BitTorrent)
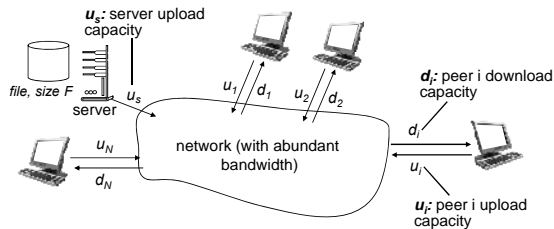   ▪ Streaming (KanKan)
   ▪ VoIP (Skype)

## File distribution: client-server vs P2P

*Question:* how much time to distribute file (size $F$) from one server to $N$ peers?
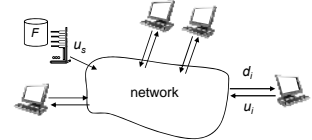- peer upload/download capacity is limited resource



$u_s$: server upload capacity

file, size $F$

server

network (with abundant bandwidth)

$d_i$: peer i download capacity

$u_i$: peer i upload capacity

## 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: $F/d_{min}$



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

increases linearly in N

## File distribution time: P2P

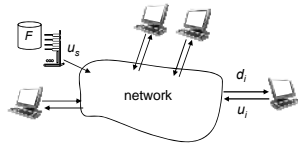- *server transmission:* must upload at least one copy
  - time to send one copy: $F/u_s$
- *client:* each client must download file copy
  - min client download time: $F/d_{min}$
- *clients:* as aggregate must download $NF$ bits
  - max upload rate (limting max download rate) is $u_s + \Sigma u_i$



> time to distribute F to N clients using P2P approach
> $$D_{P2P} \geq max\{F/u_{s}, F/d_{min}, NF/(u_s + \Sigma u_i)\}$$
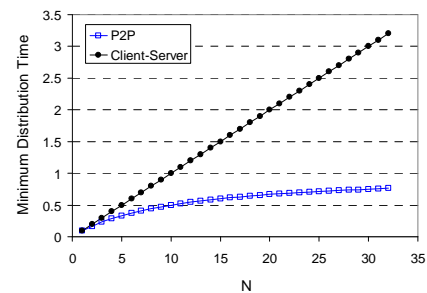
increases linearly in $N$ …
… but so does this, as each peer brings service capacity

## Client-server vs. P2P: example
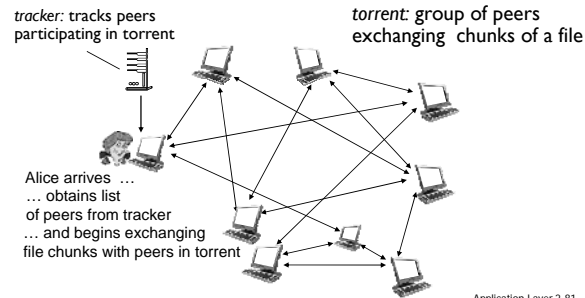
client upload rate = $u$, $F/u$ = 1 hour, $u_s = 10u$, $d_{min} \geq u_s$

5

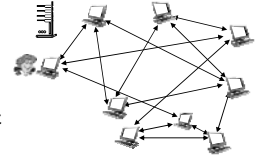## P2P file distribution: BitTorrent

❖ file divided into 256Kb chunks
❖ peers in torrent send/receive file chunks

*tracker:* tracks peers
participating in torrent

*torrent:* group of peers
exchanging chunks of a file

Alice arrives …
… obtains list
of peers from tracker
… and begins exchanging
file chunks with peers in torrent

## P2P file distribution: BitTorrent

❖ peer joining torrent:
  ▪ has no chunks, but will accumulate them over time from other peers
  ▪ registers with tracker to get list of peers, connects to subset of peers ("neighbors")
❖ while downloading, peer uploads chunks to other peers
❖ peer may change peers with whom it exchanges chunks
❖ *churn:* peers may come and go
❖ once peer has entire file, it may (selfishly) leave or (altruistically) remain in torrent

## 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

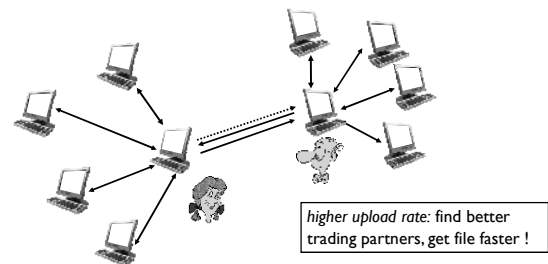*sending chunks: tit-for-tat*
❖ Alice sends chunks to those four peers currently sending her chunks *at highest rate*
  ▪ other peers are choked by Alice (do not receive chunks from her)
  ▪ re-evaluate top 4 every10 secs
❖ every 30 secs: randomly select another peer, starts sending chunks
  ▪ "optimistically unchoke" this peer
  ▪ newly chosen peer may join top 4

## BitTorrent: tit-for-tat

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

*higher upload rate:* find better trading partners, get file faster !

## Distributed Hash Table (DHT)

❖ DHT: a *distributed P2P database*
❖ database has (key, value) pairs; examples:
  ▪ key: ss number; value: human name
  ▪ key: movie title; value: IP address
❖ Distribute the (key, value) pairs over the (millions of peers)
❖ a peer queries DHT with key
  ▪ DHT returns values that match the key
❖ peers can also insert (key, value) pairs

## Q: how to assign keys to peers?

❖ central issue:
  ▪ assigning (key, value) pairs to peers.
❖ basic idea:
  ▪ convert each key to an integer
  ▪ Assign integer to each peer
  ▪ put (key,value) pair in the peer that is closest to the key

## DHT identifiers

❖ assign integer identifier to each peer in range $[0,2^n-1]$ for some $n$.
  ▪ each identifier represented by $n$ bits.

❖ require each key to be an integer in same range
❖ to get integer key, hash original key
  ▪ *e.g.,* key = hash("Led Zeppelin IV")
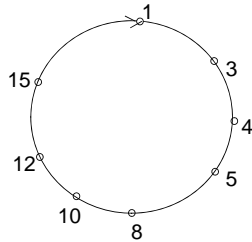  ▪ this is why its is referred to as a *distributed "hash" table*

## Assign keys to peers

❖ rule: assign key to the peer that has the *closest* ID.
❖ convention in lecture: closest is the *immediate successor* of the key.
❖ e.g., $n$=4; peers: 1,3,4,5,8,10,12,14;
  ▪ key = 13, then successor peer = 14
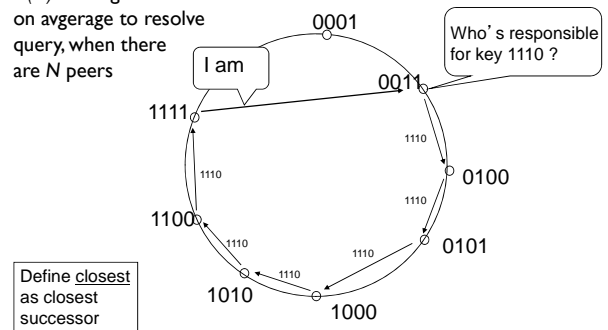  ▪ key = 15, then successor peer = 1

## Circular DHT (1)



- each peer *only* aware of immediate successor and predecessor.
- "overlay network"

## Circular DHT (1)

*O(N)* messages on avgerage to resolve query, when there are *N* peers

Who's responsible for key 1110 ?

I am

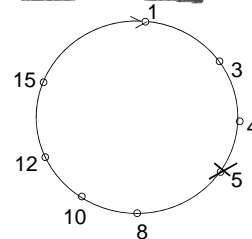Define <u>closest</u> as closest successor

## Circular DHT with shortcuts



Who's responsible for key 1110?

- each peer keeps track of IP addresses of predecessor, successor, short cuts.
- reduced from 6 to 2 messages.
- possible to design shortcuts so *O(log N)* neighbors, *O(log N)* messages in query

## Peer churn



handling peer churn:
- peers may come and go (churn)
- each peer knows address of its two successors
- each peer periodically pings its two successors to check aliveness
- if immediate successor leaves, choose next successor as new immediate successor

*example: peer 5 abruptly leaves*

- peer 4 detects peer 5 departure; makes 8 its immediate successor; asks 8 who its immediate successor is; makes 8's immediate successor its second successor.
- what if peer 13 wants to join?

# Chapter 2: outline
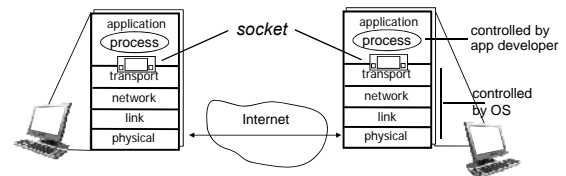
2.1 principles of network applications
- app architectures
- app requirements

2.2 Web and HTTP

2.3 FTP

2.4 electronic mail
- SMTP, POP3, IMAP

2.5 DNS

2.6 P2P applications

2.7 socket programming with UDP and TCP

# Socket programming

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

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

## Socket programming

*Two socket types for two transport services:*
- *UDP:* unreliable datagram
- *TCP:* reliable, byte stream-oriented

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

# Socket programming *with UDP*

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

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

Application viewpoint:
- UDP provides *unreliable* transfer of groups of bytes ("datagrams") between client and server

## Client/server socket interaction: UDP

**server** (running on *serverIP*)

create socket, port= x:
serverSocket =
socket(AF_INET,SOCK_DGRAM)

↓

read datagram from
serverSocket

↓

write reply to
serverSocket
specifying
client address,
port number

**client**

create socket:
clientSocket =
socket(AF_INET,SOCK_DGRAM)

Create datagram with server IP and
port=x; send datagram via
clientSocket

↓

read datagram from
clientSocket

close
clientSocket

---

## Example app: UDP client

*Python UDPClient*

include Python's socket library → from socket import *

serverName = 'hostname'

serverPort = 12000

create UDP socket for server → clientSocket = socket(socket.AF_INET,
socket.SOCK_DGRAM)

get user keyboard input → message = raw_input('Input lowercase sentence:')

Attach server name, port to message; send into socket → clientSocket.sendto(message,(serverName, serverPort))

read reply characters from socket into string → modifiedMessage, serverAddress =
clientSocket.recvfrom(2048)

print out received string and close socket → print modifiedMessage

clientSocket.close()

---

## Example app: UDP server

*Python UDPServer*

from socket import *

serverPort = 12000

create UDP socket → serverSocket = socket(AF_INET, SOCK_DGRAM)

bind socket to local port number 12000 → serverSocket.bind(('', serverPort))

print *"The server is ready to receive"*

loop forever → while 1:

Read from UDP socket into message, getting client's address (client IP and port) → message, clientAddress = serverSocket.recvfrom(2048)

modifiedMessage = message.upper()

send upper case string back to this client → serverSocket.sendto(modifiedMessage, clientAddress)

---

## Socket programming *with TCP*

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

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

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

application viewpoint:
TCP provides reliable, in-order byte-stream transfer ("pipe") between client and server

## Client/server socket interaction: TCP

**server** (running on **hostid**)        **client**

create socket,
port=**x**, for incoming
request:
serverSocket = socket()

wait for incoming
connection request
connectionSocket =
serverSocket.accept()

┄┄┄ TCP ┄┄┄
connection setup

create socket,
connect to **hostid**, port=**x**
clientSocket = socket()

send request using
clientSocket

read request from
connectionSocket

write reply to
connectionSocket

read reply from
clientSocket

close
connectionSocket

close
clientSocket

## Example app: TCP client

*Python TCPClient*

create TCP socket for
server, remote port 12000

```
from socket import *
serverName = 'servername'
serverPort = 12000
clientSocket = socket(AF_INET,SOCK_STREAM)
clientSocket.connect((serverName,serverPort))
sentence = raw_input('Input lowercase sentence:')
clientSocket.send(sentence)
modifiedSentence = clientSocket.recv(1024)
print 'From Server:', modifiedSentence
clientSocket.close()
```

No need to attach server
name, port

## Example app: TCP server

*Python TCPServer*

create TCP welcoming
socket

server begins listening for
incoming TCP requests

loop forever

server waits on accept()
for incoming requests, new
socket created on return

read bytes from socket (but
not address as in UDP)

close connection to this
client (but *not* welcoming
socket)

```
 from socket import *
serverPort = 12000
serverSocket = socket(AF_INET,SOCK_STREAM)
serverSocket.bind(('',serverPort))
serverSocket.listen(1)
print 'The server is ready to receive'
while 1:
    connectionSocket, addr = serverSocket.accept()

    sentence = connectionSocket.recv(1024)
    capitalizedSentence = sentence.upper()
    connectionSocket.send(capitalizedSentence)
    connectionSocket.close()
```

## Chapter 2: summary

*our study of network apps now complete!*

- ❖ application architectures
  - ▪ client-server
  - ▪ P2P
- ❖ application service requirements:
  - ▪ reliability, bandwidth, delay
- ❖ Internet transport service model
  - ▪ connection-oriented, reliable: TCP
  - ▪ unreliable, datagrams: UDP

- ❖ specific protocols:
  - ▪ HTTP
  - ▪ FTP
  - ▪ SMTP, POP, IMAP
  - ▪ DNS
  - ▪ P2P: BitTorrent, DHT
- ❖ socket programming: TCP, UDP sockets

# Chapter 2: summary

*most importantly: learned about protocols!*

* typical request/reply message exchange:
  * client requests info or service
  * server responds with data, status code
* message formats:
  * headers: fields giving info about data
  * data: info being communicated

*important themes:*

* control vs. data msgs
  * in-band, out-of-band
* centralized vs. decentralized
* stateless vs. stateful
* reliable vs. unreliable msg transfer
* "complexity at network edge"

Application Layer 2-105