# CSC358 *Intro. to Computer Networks*

**Lecture 4:** *FTP App, DNS App, P2P App, Introduction to Transport Layer*

Amir H. Chinaei, Winter 2016

ahchinaei@cs.toronto.edu
*http://www.cs.toronto.edu/~ahchinaei/*

Many slides are (inspired/adapted) from the above source
© all material copyright; all rights reserved for the authors

Office Hours: T 17:00–18:00 R 9:00–10:00 BA4222

TA Office Hours: W 16:00-17:00 BA3201 R 10:00-11:00 BA7172
csc358ta@cdf.toronto.edu
*http://www.cs.toronto.edu/~ahchinaei/teaching/2016jan/csc358/*

---

## Review

❖ Many applications run on Internet application layer
  1. some have open protocols, such as HTTP, DNS, etc., many others have proprietary protocols.
  2. use some underlying protocols as a black-box
  3. architecture: C/S, P2P, Hybrid

❖ HTTP
❖ Cookies provide user-server state
❖ Non-persistent vs persistent HTTP connections

---

## Architecture examples

| myApp | |
|---|---|
| HTTP | |
| TCP | |
| .... | |

| myApp | |
|---|---|
| HTTP | DNS |
| TCP | TCP,UDP |
| .... | |

| myApp | |
|---|---|
| HTTP | DNS |
| TLS | |
| TCP | TCP,UDP |
| .... | |

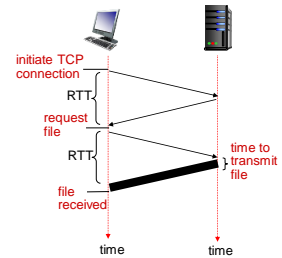| myApp | |
|---|---|
| HTTP | Skype |
| TCP | TCP,UDP |
| .... | |

| myApp |
|---|
| TCP or UDP or (TCP,UDP) or DCCP |
| .... |

---

## Non-persistent HTTP: response time

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

HTTP response time:
❖ one RTT to initiate TCP connection
❖ one RTT for HTTP request and first few bytes of HTTP response to return
❖ file transmission time
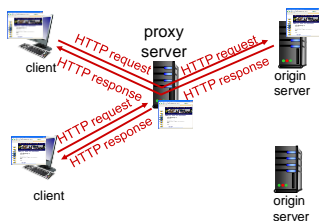❖ non-persistent HTTP response time =
  $2RTT + t_f$

---

## Web caches (proxy server)

*goal:* satisfy client request without involving origin server

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

---

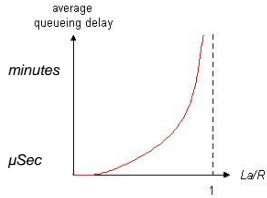## More about Web caching

❖ cache acts as both client and server
  ▪ server for original requesting client
  ▪ client to origin server
❖ typically cache is installed by ISP (university, company, residential ISP)

*why Web caching?*
❖ reduce response time for client request
❖ reduce traffic on an institution's access link
❖ Internet dense with caches: enables "poor" content providers to effectively deliver content (so too does P2P file sharing)

1

## Utilization

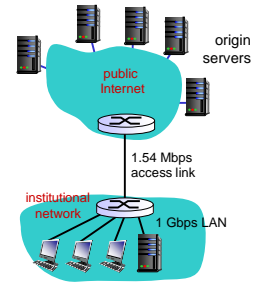average queueing delay

*minutes*

*μSec*

$La/R$

1

---

## Caching example:

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

*consequences:*
- LAN utilization: ?    *problem!*
- access link utilization = ?
- total delay = LAN$_{outbound}$ delay+ access$_{outbound}$ delay + Internet delay + access$_{inbound}$ delay + LAN$_{inbound}$ delay
  = ?

origin servers

public Internet

1.54 Mbps access link

institutional network

1 Gbps LAN

---

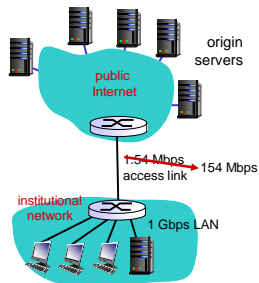## Caching example: fatter access link

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

*consequences:*
- LAN utilization: ?
- access link utilization = ? ?
- total delay =
  = 2 sec + minutes + μsecs
  msecs

*Cost:* increased access link speed (not cheap!)

origin servers

public Internet

~~1.54 Mbps~~ 154 Mbps access link

institutional network

1 Gbps LAN

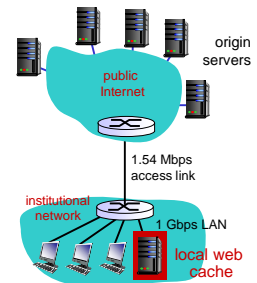---

## Caching example: install local cache

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

*consequences:*
- LAN utilization: 15%
- access link utilization = ?
- total delay = ?

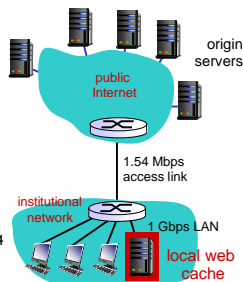  *How to compute link utilization, delay?*

*Cost:* web cache (cheap!)

origin servers

public Internet

1.54 Mbps access link

institutional network

1 Gbps LAN

local web cache

---

## Caching example: install local cache

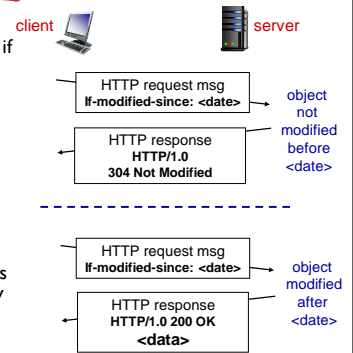*Calculating access link utilization, delay with cache:*
- suppose cache hit rate is 0.4
  - 40% requests satisfied at cache, 60% requests satisfied at origin
- access link utilization:
  - 60% of requests use access link
- data rate to browsers over access link = 0.6*1.50 Mbps = .9 Mbps
  - utilization = 0.9/1.54 = .58
- total delay
  - = 0.6 * (delay from origin servers) +0.4 * (delay when satisfied at cache)
  - = 0.6 * (2.01) + 0.4 * (~μsecs)
  - = ~ 1.2 secs
  - less than with 154 Mbps link (and cheaper too!)

origin servers

public Internet

1.54 Mbps access link

institutional network
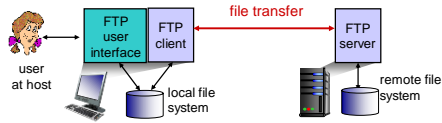
1 Gbps LAN

local web cache

---

## Conditional GET

- *Goal:* don't send object if cache has up-to-date cached version
  - no object transmission delay
  - lower link utilization
- *cache:* specify date of cached copy in HTTP request
  `If-modified-since: <date>`
- *server:* response contains no object if cached copy is up-to-date:
  `HTTP/1.0 304 Not Modified`

client

server

HTTP request msg
**If-modified-since: <date>**

object not modified before <date>

HTTP response
**HTTP/1.0 304 Not Modified**

- - - - - - - - - - - - - - - - -

HTTP request msg
**If-modified-since: <date>**

object modified after <date>

HTTP response
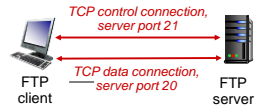**HTTP/1.0 200 OK <data>**

## FTP: the file transfer protocol



* transfer file to/from remote host
* client/server model
  * *client:* side that initiates transfer (either to/from remote)
  * *server:* remote host
* ftp: RFC 959
* ftp server: port 21

## FTP: separate control, data connections

* FTP client contacts FTP server at port 21, using TCP
* client authorized over control connection
* client browses remote directory, sends commands over control connection
* when server receives file transfer command, *server* opens *2nd* TCP data connection (for file) *to* client
* after transferring one file, server closes data connection



* server opens another TCP data connection to transfer another file
* control connection: *"out of band"*
* FTP server maintains "state": current directory, earlier authentication

## FTP commands, responses

### *sample commands:*
* sent as ASCII text over control channel
* **USER *username***
* **PASS *password***
* **LIST** return list of file in current directory
* **RETR filename** retrieves (gets) file
* **STOR filename** stores (puts) file onto remote host

### *sample return codes*
* status code and phrase (as in HTTP)
* **331 Username OK, password required**
* **125 data connection already open; transfer starting**
* **425 Can't open data connection**
* **452 Error writing file**

## 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
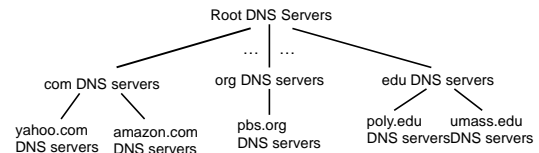  * replicated Web servers: many IP addresses correspond to one name

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

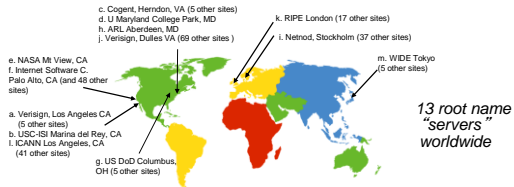*A: doesn't scale!*

## DNS: a distributed, hierarchical database



*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!)
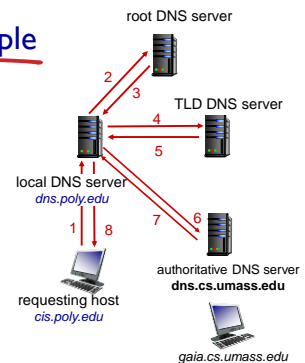  - acts as proxy, forwards query into hierarchy

---

# 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
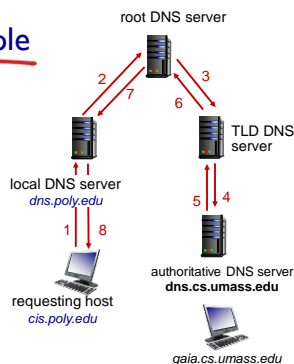- ❖ "I don't know this name, but ask this server"

root DNS server

TLD DNS server

local DNS server
*dns.poly.edu*

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

requesting host
*cis.poly.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*

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

requesting host
*cis.poly.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**

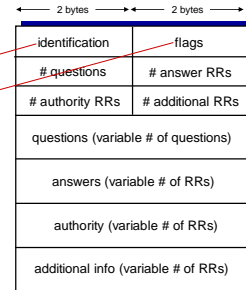Application Layer 2-25

---

## DNS protocol, messages

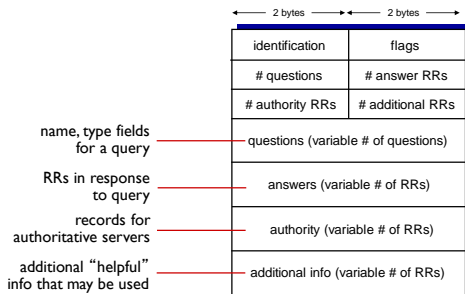- *query* and *reply* messages, both with same *message format*

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) ||

Application Layer 2-26

---

## 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)

Application Layer 2-27

---

## 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)**
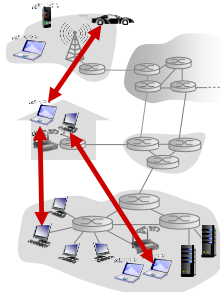- create authoritative server type A record for www.networkuptopia.com; type MX record for networkutopia.com

Application Layer 2-28

---

## Pure P2P architecture

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

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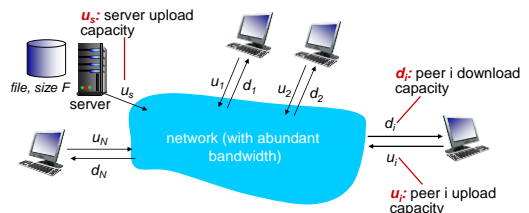- file distribution (BitTorrent)
- Streaming (KanKan)
- VoIP (Skype)

Application Layer 2-29

---

## 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
$u_s$
$u_1$ / $d_1$   $u_2$ / $d_2$
$d_i$: peer i download capacity
$u_N$
$d_N$
network (with abundant bandwidth)
$d_i$
$u_i$
$u_i$: peer i upload capacity

Application Layer 2-30

5

## File distribution time: client-server

- *server transmission:* must sequentially send (upload) *N* file copies:
    - time to send one copy: $F/u_s$
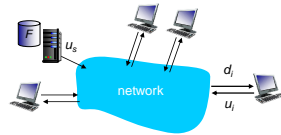    - time to send N copies: $NF/u_s$
- *client:* each client must download file copy
    - $d_{min}$ = min client download rate
    - min client download time: $F/d_{min}$

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

increases linearly in N

Application Layer 2-31

## 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 $u_s + \Sigma u_i$

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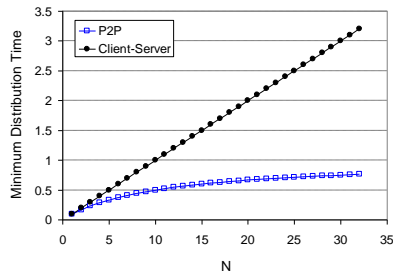
increases linearly in *N* …
… but so does this, as each peer brings service capacity

Application Layer 2-32

## Client-server vs. P2P: example

client upload rate = *u*,  *F/u* = 1 hour,  $u_s = 10u$,  $d_{min} \geq u_s$
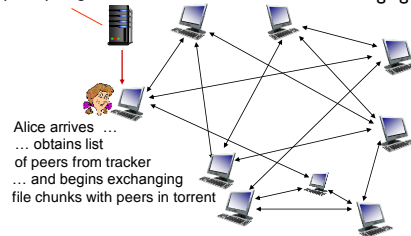
Application Layer 2-33

## 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
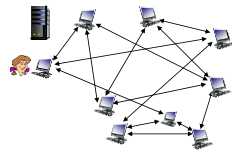… and begins exchanging file chunks with peers in torrent

Application Layer 2-34

## P2P file distribution: BitTorrent

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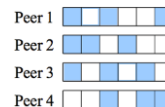- once peer has entire file, it may (selfishly) leave or (altruistically) remain in torrent

Application Layer 2-35

## 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
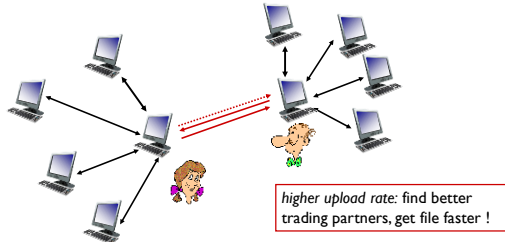- Alice requests missing chunks from peers, rarest first

*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

Peer 1
Peer 2
Peer 3
Peer 4

Application Layer 2-36

6

# 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)

❖ Hash table

❖ DHT paradigm

❖ Circular DHT and overlay networks

❖ Peer churn

---

# Simple Database

Simple database with (key, value) pairs:

- key: human name; value: social security #

| Key | Value |
|-----|-------|
| John Washington | 132-54-3570 |
| Diana Louise Jones | 761-55-3791 |
| Xiaoming Liu | 385-41-0902 |
| Rakesh Gopal | 441-89-1956 |
| Linda Cohen | 217-66-5609 |
| ....... | ......... |
| Lisa Kobayashi | 177-23-0199 |

- key: movie title; value: IP address

---

# Hash Table

- More convenient to store and search on numerical representation of key
- key = hash(original key)

| Original Key | Key | Value |
|--------------|-----|-------|
| John Washington | 8962458 | 132-54-3570 |
| Diana Louise Jones | 7800356 | 761-55-3791 |
| Xiaoming Liu | 1567109 | 385-41-0902 |
| Rakesh Gopal | 2360012 | 441-89-1956 |
| Linda Cohen | 5430938 | 217-66-5609 |
| ....... | | ......... |
| Lisa Kobayashi | 9290124 | 177-23-0199 |

---

# Distributed Hash Table (DHT)

❖ Distribute (key, value) pairs over millions of peers
  ▪ pairs are evenly distributed over peers
❖ Any peer can query database with a key
  ▪ database returns value for the key
  ▪ To resolve query, small number of messages exchanged among peers
❖ Each peer only knows about a small number of other peers
❖ Robust to peers coming and going (churn)

---

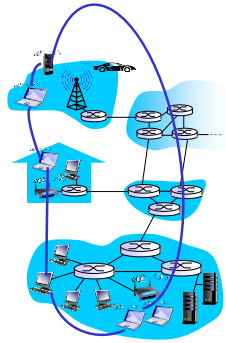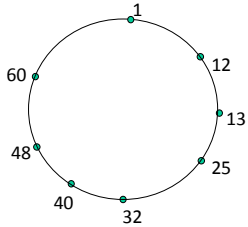# Assign key-value pairs to peers

❖ rule: assign key-value pair to the peer that has the *closest* ID.
❖ convention: closest is the *immediate successor* of the key.
❖ e.g., ID space {0,1,2,3,...,63}
❖ suppose 8 peers: 1,12,13,25,32,40,48,60
  ▪ If key = 51, then assigned to peer 60
  ▪ If key = 60, then assigned to peer 60
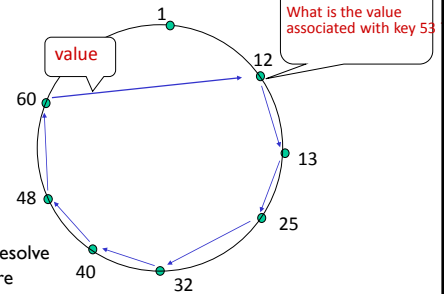  ▪ If key = 61, then assigned to peer 1

# Circular DHT

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



"overlay network"

# Resolving a query
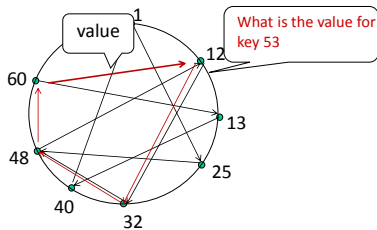


What is the value associated with key 53

value

*O(N)* messages on avgerage to resolve query, when there are *N* peers
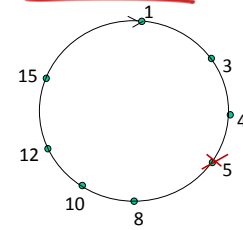
# Circular DHT with shortcuts



value

What is the value for key 53

- each peer keeps track of IP addresses of predecessor, successor, short cuts.
- reduced from 6 to 3 messages.
- possible to design shortcuts with *O(log N)* neighbors, *O(log N)* messages in query
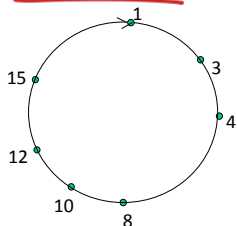
# Peer churn



*example: peer 5 abruptly leaves*

### 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

# Peer churn



*example: peer 5 abruptly leaves*
❖peer 4 detects peer 5's departure; makes 8 its immediate successor

❖ 4 asks 8 who its immediate successor is; makes 8's immediate successor its second successor.

### 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
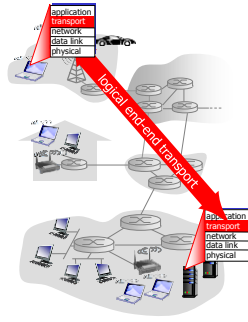
# Let's move on to **Transport Layer**

- TCP, UDP
- principles, services
- multiplexing, demultiplexing
- reliable data transfer
- flow control
- congestion control

## Transport services and protocols

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

## Transport vs. network layer

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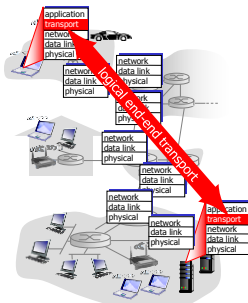  - relies on, enhances, network layer services

*household analogy:*
*12 kids in Ann's house sending letters to 12 kids in Bill's house:*
- hosts = houses
- processes = kids
- app messages = (long) letters
- segments= letters in envelopes
- transport protocol = Ann and Bill who demux to in-house siblings
- network-layer protocol = postal service

## Internet transport-layer protocols

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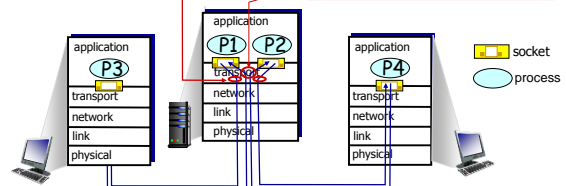- services not available:
  - delay guarantees
  - bandwidth guarantees

## Multiplexing/demultiplexing

*multiplexing at sender:*
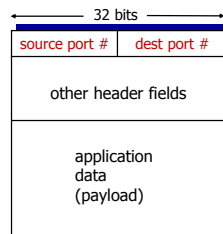handle data from multiple sockets, add transport header (later used for demultiplexing)

*demultiplexing at receiver:*
use header info to deliver received segments to correct socket

## 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

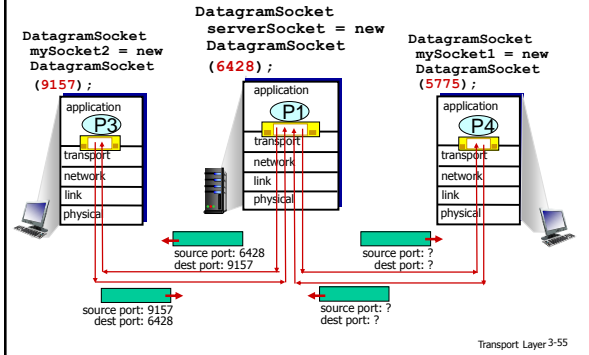| 32 bits | |
|---|---|
| source port # | dest port # |
| other header fields | |
| application data (payload) | |

TCP/UDP segment format

## Connectionless demultiplexing

- *recall:* created socket has host-local port #:
  ```
  DatagramSocket mySocket1
  = new DatagramSocket(12534);
  ```
- when host receives UDP segment:
  - checks destination port # in segment
  - directs UDP segment to socket with that port #

- *recall:* when creating datagram to send into UDP socket, must specify
  - destination IP address
  - destination port #

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

## Connectionless demux: example

```
DatagramSocket
mySocket2 = new
DatagramSocket
(9157);
```

```
DatagramSocket
serverSocket = new
DatagramSocket
(6428);
```

```
DatagramSocket
mySocket1 = new
DatagramSocket
(5775);
```



source port: 6428
dest port: 9157

source port: ?
dest port: ?

source port: 9157
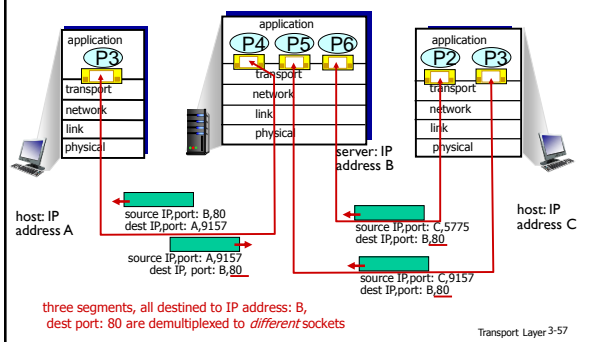dest port: 6428

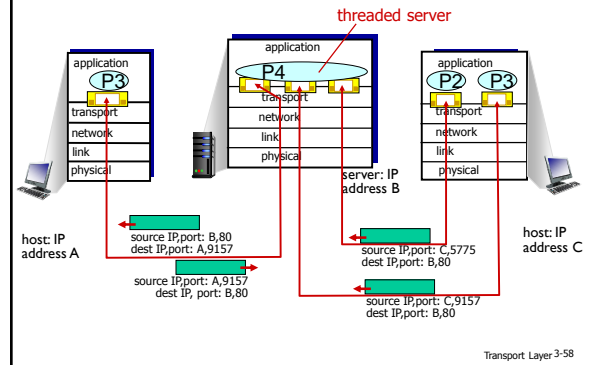source port: ?
dest port: ?

## Connection-oriented demux

❖ TCP socket identified by 4-tuple:
  ▪ source IP address
  ▪ source port number
  ▪ dest IP address
  ▪ dest port number
❖ demux: receiver uses all four values to direct segment to appropriate socket

❖ server host may support many simultaneous TCP sockets:
  ▪ each socket identified by its own 4-tuple
❖ web servers have different sockets for each connecting client
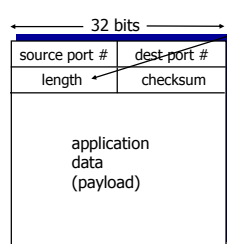  ▪ non-persistent HTTP will have different socket for each request

## Connection-oriented demux: example



source IP,port: B,80
dest IP,port: A,9157

source IP,port: A,9157
dest IP, port: B,80

source IP,port: C,5775
dest IP,port: B,80

source IP,port: C,9157
dest IP,port: B,80

host: IP address A

server: IP address B

host: IP address C

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

## Connection-oriented demux: example

threaded server



source IP,port: B,80
dest IP,port: A,9157

source IP,port: A,9157
dest IP, port: B,80

source IP,port: C,5775
dest IP,port: B,80

source IP,port: C,9157
dest IP,port: B,80

host: IP address A

server: IP address B

host: IP address C

## UDP: segment header

length, in bytes of UDP segment, including header



← 32 bits →

| source port # | dest port # |
|---|---|
| length | checksum |

application
data
(payload)

UDP segment format

why is there a UDP?
❖ no connection establishment (which can add delay)
❖ simple: no connection state at sender, receiver
❖ small header size
❖ no congestion control: UDP can blast away as fast as desired

## UDP checksum

*Goal:* detect "errors" (e.g., flipped bits) in transmitted segment

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

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

## Internet checksum: example

example: add two 16-bit integers

```
            1 1 1 0 0 1 1 0 0 1 1 0 0 1 1 0
            1 1 0 1 0 1 0 1 0 1 0 1 0 1 0 1
wraparound (1) 1 0 1 1 1 0 1 1 1 0 1 1 1 0 1 1

     sum      1 0 1 1 1 0 1 1 1 0 1 1 1 1 0 0
checksum      0 1 0 0 0 1 0 0 0 1 0 0 0 0 1 1
```

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

## Next weeks

❖ Reliable data transfer
❖ Connection-oriented transport
❖ TCP