# Content Distribution

Professor Keith W. Ross
Institut Eurécom
Sophia Antipolis, France
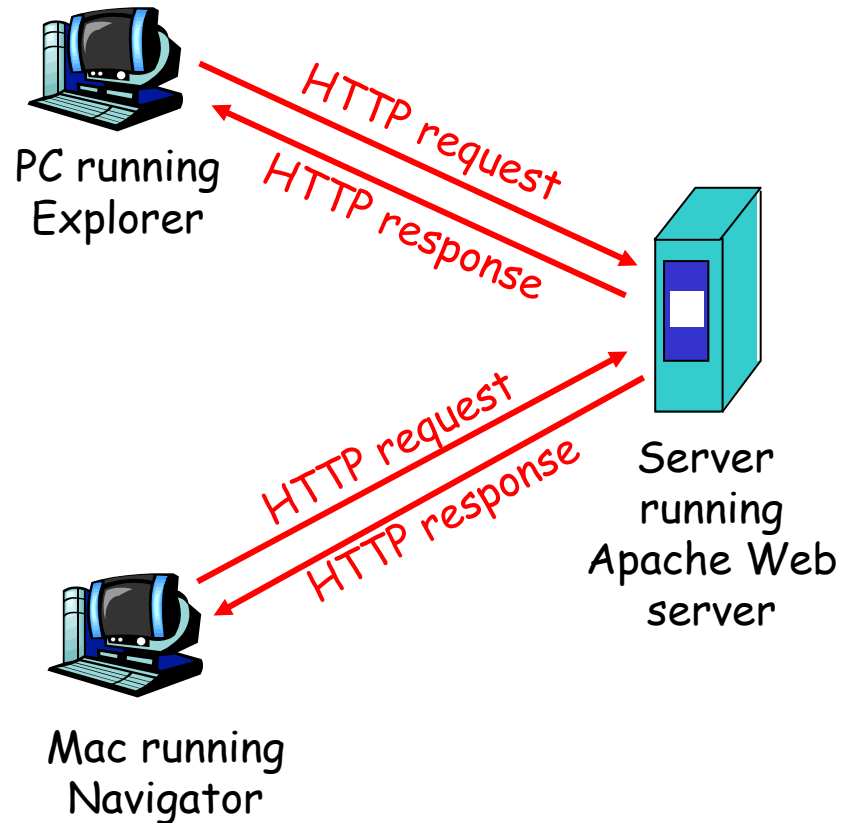
**http://www.eurecom.fr/~ross**

# Content Distribution

r Review of HTTP, DNS, TCP

r Server Farms

r Proxy Web Caches

r Content Distribution Networks (CDNs)
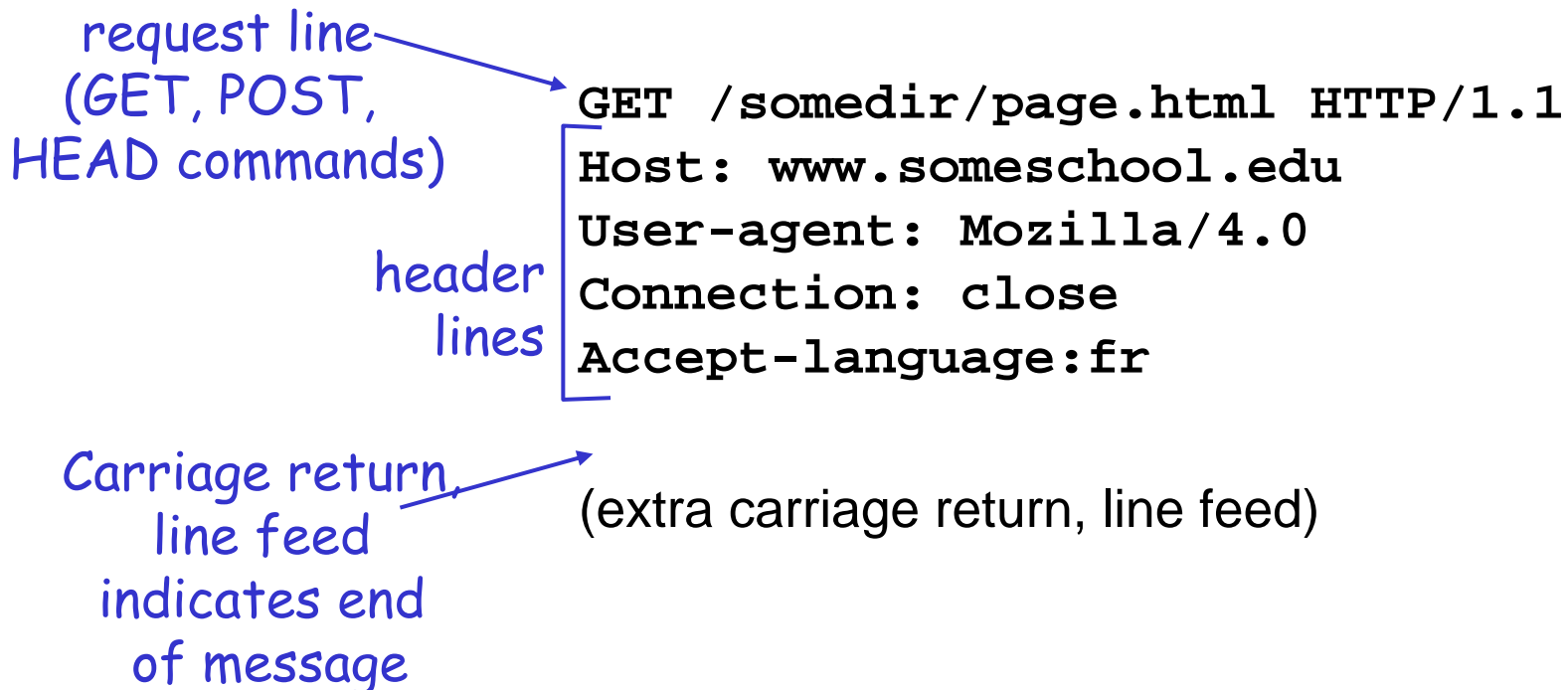
r Peer-to-peer file sharing (P2P)

# HTTP overview

## HTTP: hypertext transfer protocol

r  Web's application layer protocol

r  client/server model

  m  *client:* browser that requests, receives, "displays" Web objects

  m  *server:* Web server sends objects in response to requests
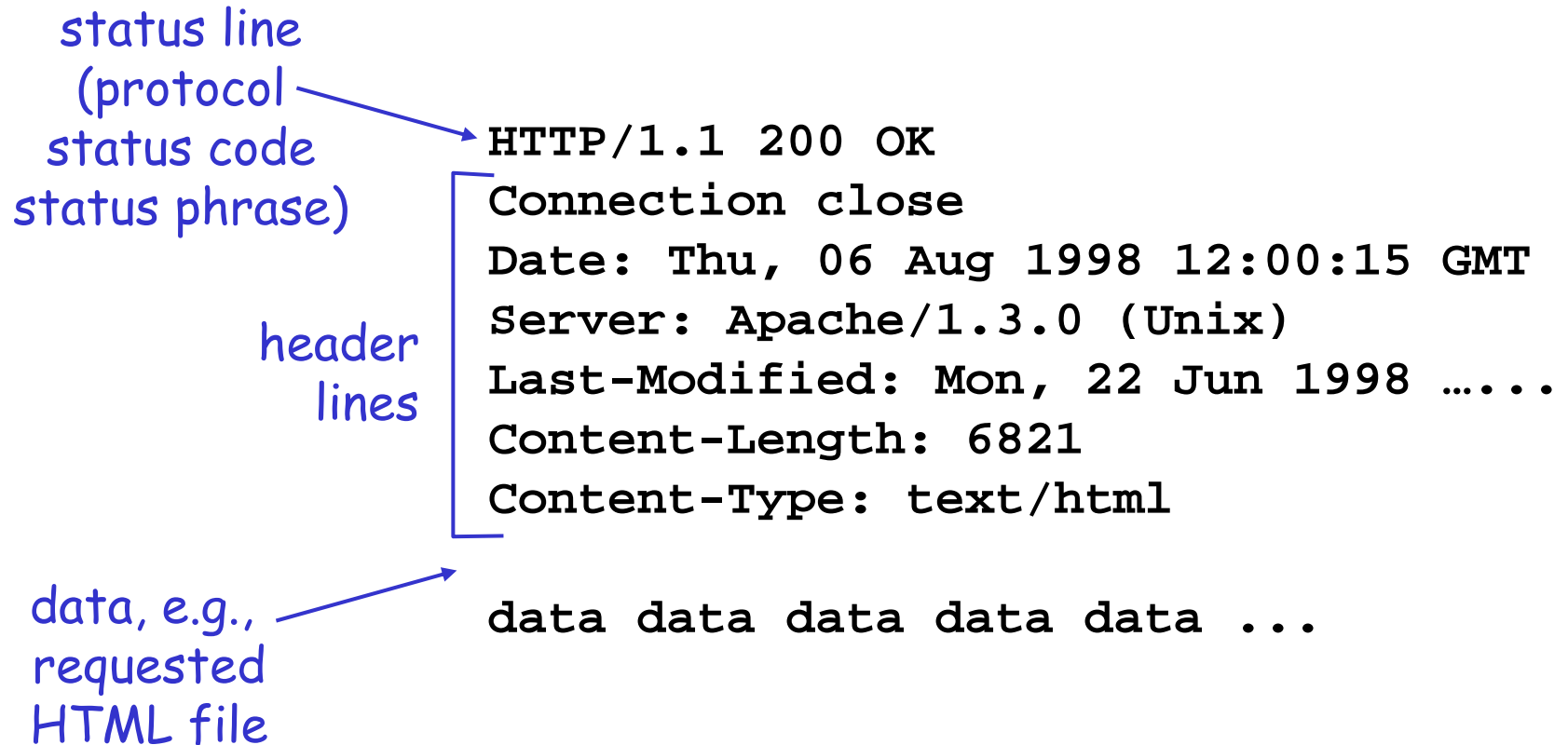
r  HTTP 1.0: RFC 1945

r  HTTP 1.1: RFC 2068

PC running Explorer

HTTP request

HTTP response

Server running Apache Web server

HTTP request

HTTP response

Mac running Navigator

# HTTP request message

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

r  HTTP request message:

m  ASCII (human-readable format)

request line
(GET, POST,
HEAD commands)

header
lines

```
GET /somedir/page.html HTTP/1.1
Host: www.someschool.edu
User-agent: Mozilla/4.0
Connection: close
Accept-language:fr
```

Carriage return,
line feed
indicates end
of message

(extra carriage return, line feed)

# HTTP response message

status line
(protocol
status code
status phrase)

header
lines

data, e.g.,
requested
HTML file

```
HTTP/1.1 200 OK
Connection close
Date: Thu, 06 Aug 1998 12:00:15 GMT
Server: Apache/1.3.0 (Unix)
Last-Modified: Mon, 22 Jun 1998 …...
Content-Length: 6821
Content-Type: text/html

data data data data data ...
```

# DNS: Domain Name System

**People:** many identifiers:
- m   SSN, name, passport #

**Internet hosts, routers:**
- m   IP address (32 bit) - used for addressing datagrams
- m   "name", e.g., gaia.cs.umass.edu - used by humans
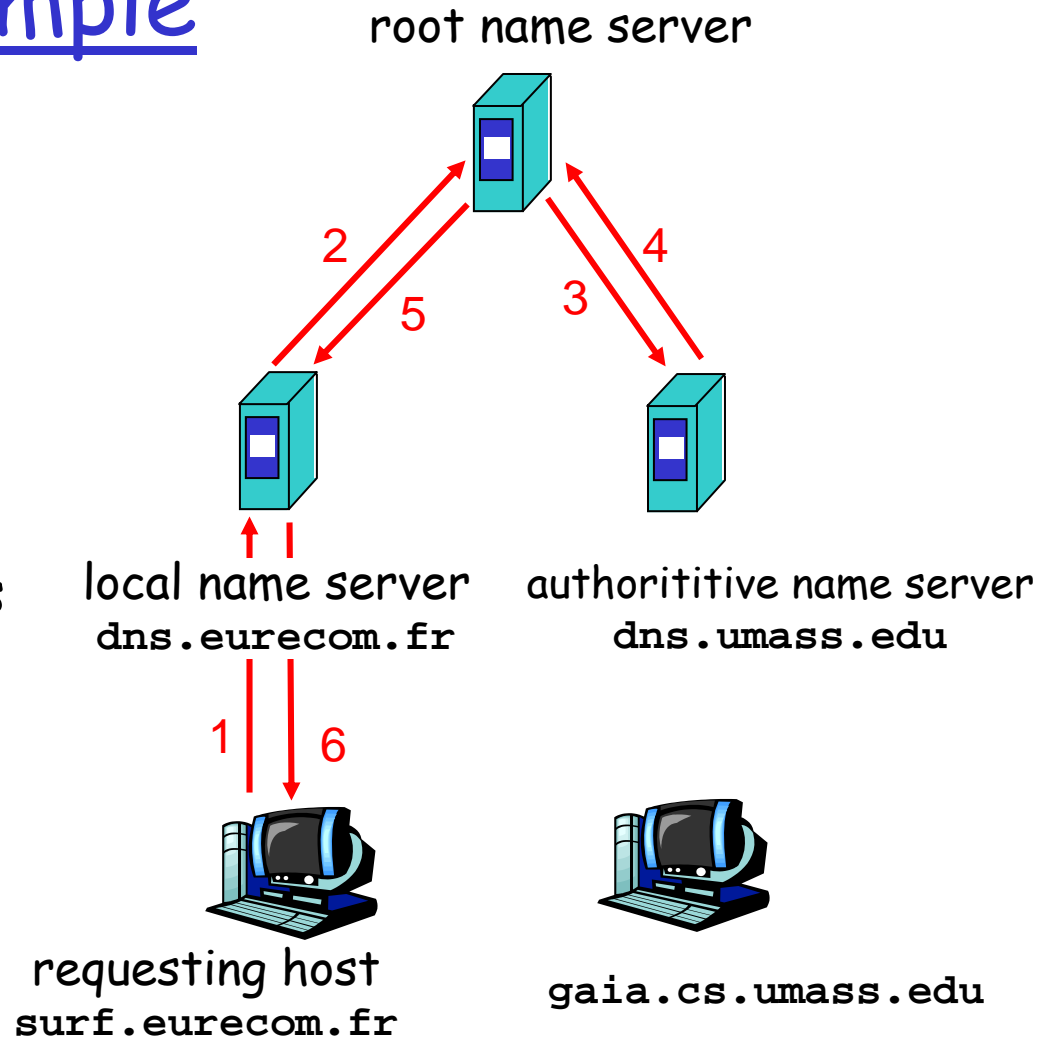
**Q:** map between IP addresses and name ?

**Domain Name System:**
- r   *distributed database* implemented in hierarchy of many *name servers*

# Simple DNS example

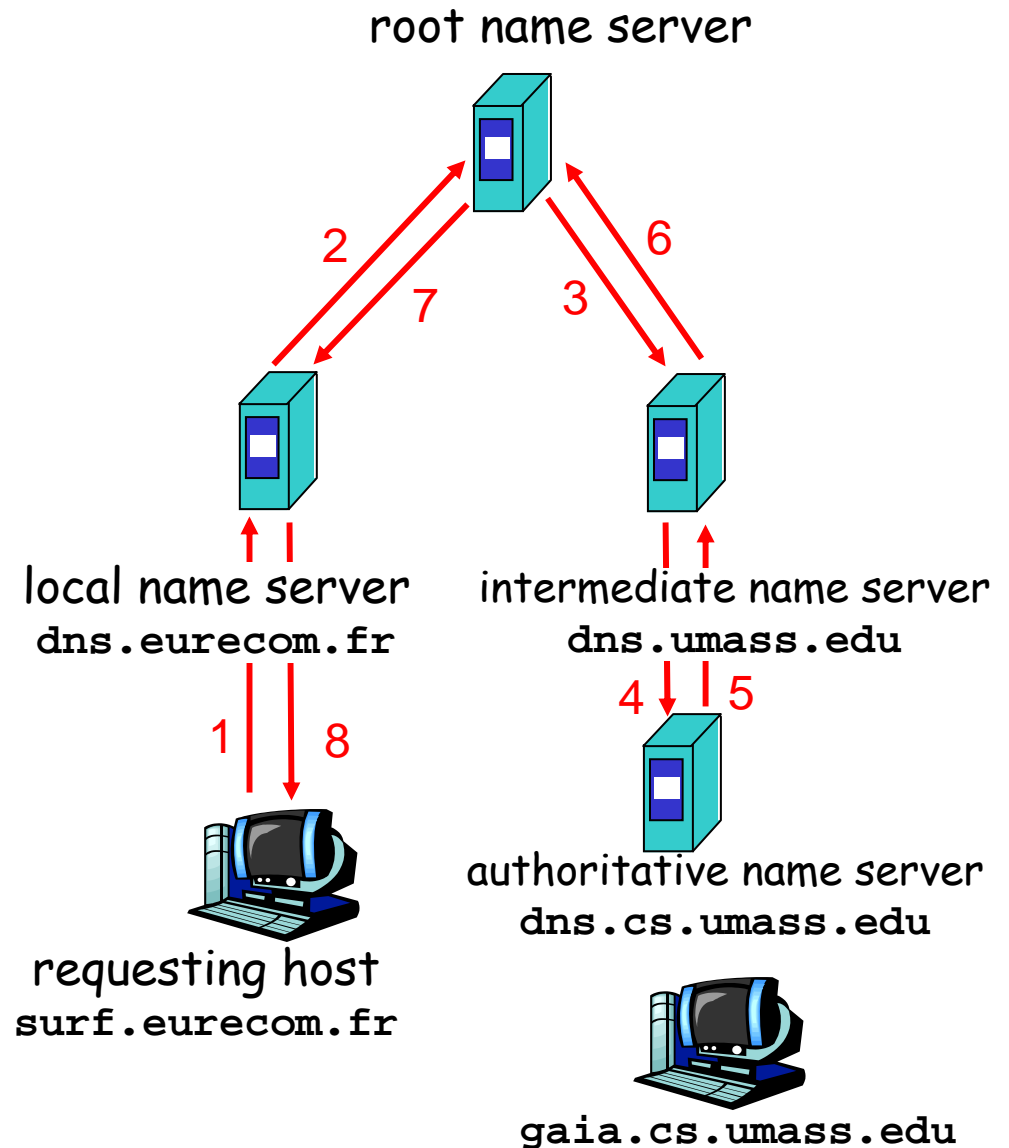host **surf.eurecom.fr** wants IP address of **gaia.cs.umass.edu**

1. contacts its local DNS server, **dns.eurecom.fr**
2. **dns.eurecom.fr** contacts root name server
3. root name server contacts authoritative name server, **dns.umass.edu**

root name server

2

4

3

5

local name server
**dns.eurecom.fr**

authorititive name server
**dns.umass.edu**

1

6

requesting host
**surf.eurecom.fr**

**gaia.cs.umass.edu**

# DNS example

Root name server:

r may not know authoritative name server

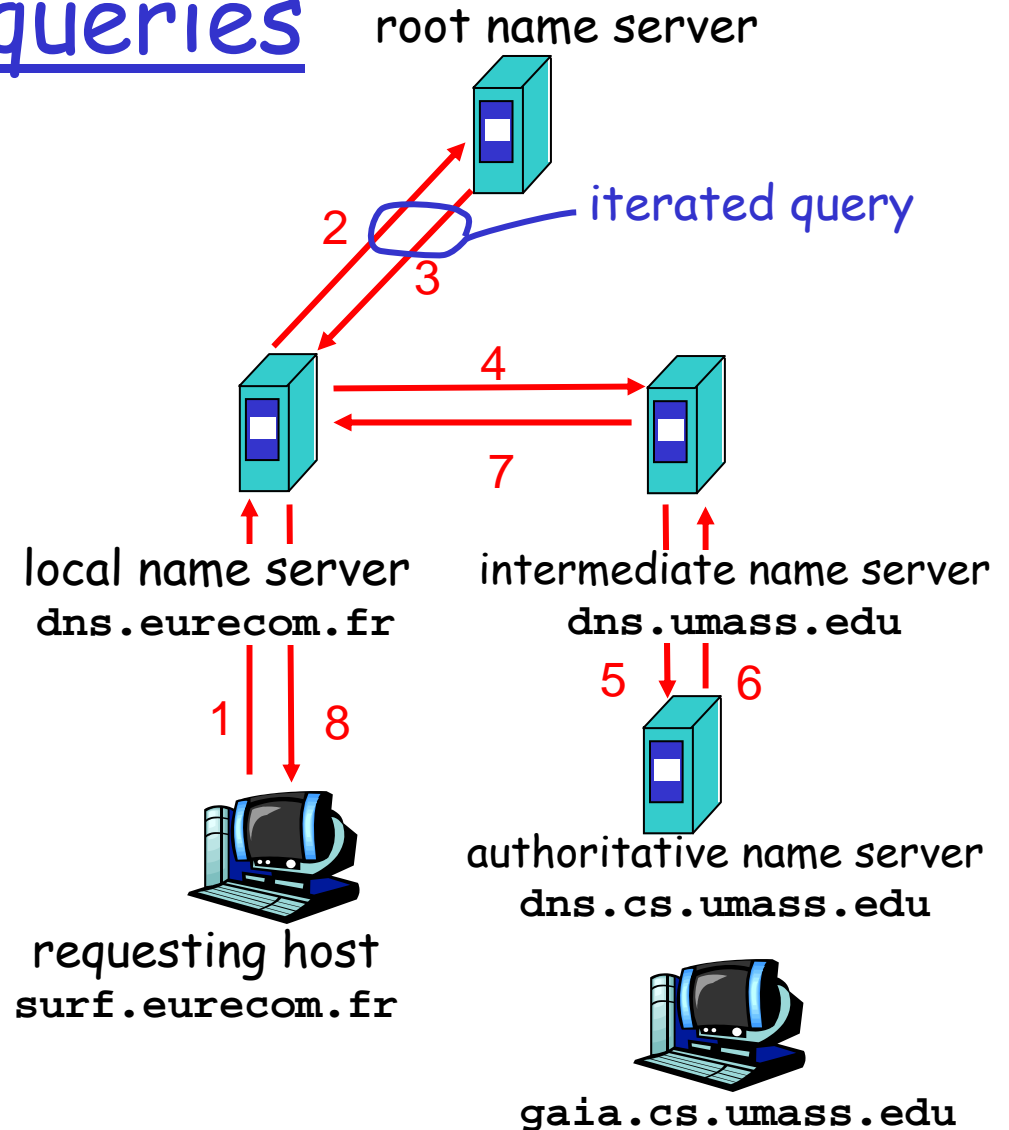r may know *intermediate name server:* who to contact to find authoritative name server

root name server

local name server
dns.eurecom.fr

intermediate name server
dns.umass.edu

authoritative name server
dns.cs.umass.edu

requesting host
surf.eurecom.fr

gaia.cs.umass.edu

1  2  3  4  5  6  7  8

# DNS: iterated queries

root name server

**recursive query:**

r    puts burden of name resolution on contacted name server

r    heavy load?

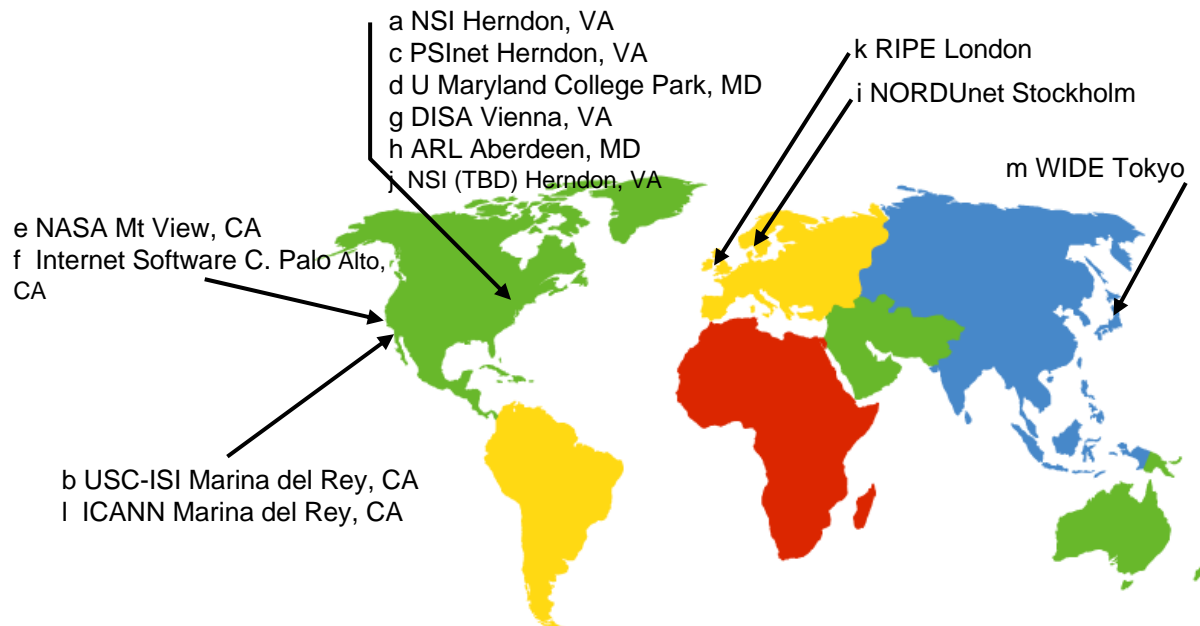**iterated query:**

r    contacted server replies with name of server to contact

r    "I don't know this name, but ask this server"

iterated query

2
3
4
7

local name server
`dns.eurecom.fr`

intermediate name server
`dns.umass.edu`

5    6

1    8

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

requesting host
`surf.eurecom.fr`

`gaia.cs.umass.edu`

# DNS: Root name servers

r  contacted by local name server that can not resolve name

r  root name server:

    m  contacts authoritative name server if name mapping not known

    m  gets mapping

    m  returns mapping to local name server

a NSI Herndon, VA
c PSInet Herndon, VA
d U Maryland College Park, MD
g DISA Vienna, VA
h ARL Aberdeen, MD
j NSI (TBD) Herndon, VA

k RIPE London
i NORDUnet Stockholm

m WIDE Tokyo

e NASA Mt View, CA
f  Internet Software C. Palo Alto, CA

b USC-ISI Marina del Rey, CA
l  ICANN Marina del Rey, CA

13 root name servers worldwide

# DNS: caching and updating records

r  once (any) name server learns mapping, it *caches* mapping

   m  cache entries timeout (disappear) after some time

r  update/notify mechanisms under design by IETF

   m  RFC 2136

   m  http://www.ietf.org/html.charters/dnsind-charter.html

# TCP connection establishment

TCP 3-way handshake:

r   client sends to server TCP segment with SYN bit is set.

r   server responds with segment that has both SYN and ACK bits set.

r   client responds with another segment with ACK bit set.

# Content Distribution

r  Review of HTTP, DNS, TCP

r  Server Farms

r  Proxy Web Caches

r  Content Distribution Networks (CDNs)

r  Peer-to-peer file sharing (P2P)

# Architectures for server farms

**Issue:** In a Web application, multiple server machines may be needed to handle traffic.

**Goal:** Architect server farm so that it appears to be running on a single machine to the client.

client

Internet

124.0.0.1
www.foo1.com

124.0.0.2
www.foo2.com

124.0.0.3
www.foo3.com

**www.foo.com**

Server farm

# Architectures for server farms

Running example:

Site:

r   http://www.foo.com

Three machines:

r   foo1, foo2, foo3

Company address space:

r   124.x.x.x

For simplicity:

r   assume static Web pages

3 architectures

r   DNS rotation

r   Surrogate server

r   NAT

# DNS rotation

r DNS translates host.names to IP addresses

r eg, www.foo.com to 124.0.0.7

r Authoritative DNS server can provide multiple IP addresses:
  124.0.0.1, 124.0.0.2, 124.0.0.3

r Rotates addresses

Issues

r DNS caching in name servers

r Name server directs its clients to the same Web server

r Hard to control how much traffic is directed at each of the servers

# Surrogate server (1)

124.0.0.1

www.foo1.com

TCP connection

TCP connection

124.0.0.7    124.0.0.0

Internet                LAN            124.0.0.2

www.foo.com
**surrogate server**

www.foo2.com

124.0.0.3

- DNS provides client with IP address 124.0.0.7 .
- Client establishes TCP connection with load balancer.
- Client sends HTTP request message to load balancer.

www.foo3.com

- Surrogate (load balancer) chooses one of the
  three Web servers.
- Surrogate establishes TCP connection with chosen server
   and forwards  request.
- Surrogate forwards data it receives from server to client.

# Surrogate server (2)

Flexibility in determining how to direct requests to servers:

r load on servers

r requested content

r cookie in request

r application-layer switch: L7 switch

Inefficient:

r surrogate has to serve at the combined rate of servers

Implementation

r Develop on top of Linux, Solaris, Windows 2000.

r Performance bottleneck: lots of TCP processing, application-layer systems calls

r Can develop customized operating system for surrogate server, resulting in efficient load balancer

# NAT (1)

124.0.0.1

www.foo1.com

TCP connection

124.0.0.207    124.0.0.0

Internet

LAN

www.foo.com

**NAT**

124.0.0.2

www.foo2.com

124.0.0.3

www.foo3.com

- NAT detects new connection by observing SYN bit set.
- NAT makes load balancing decision.
- Forwards IP datagram to chosen server :
    - Needs to modify destination IP address
    - Needs to forward subsequent packets in same connection to the same server; modify destination IP addresses
    - Needs to modify datagrams arriving from servers
- TCP connection is established end-to-end.
    - But client thinks it has TCP connection with 124.0.0.7

19

# NAT (2): Table

r   Suppose client has IP address 227.68.68.19 and source port number 9876

r   NAT table might have for this connection:

Internet                    to          LAN

SA: 227.68.68.19                        SA: 227.68.68.19

DA: 124.0.0.7                           DA: 124.0.0.1

SP: 9876                                SP: 9876

DP: 80                                  DP: 80

LAN                         to          Internet

SA: 124.0.0.1                           SA: 124.0.0.7

DA: 227.68.68.19                        DA: 227.68.68.19

SP: 9876                                SP: 9876

DP: 80                                  DP: 80

# NAT (3): Benefits

r **No TCP processing by NAT**

  m no TCP stack handling TCP buffers, congestion windows, connection establishment, etc.

r **All the packet manipulation takes place by special purpose device**

# Stickiness in load balancing

r **Cookies are often used to maintain session state, e.g., shopping cart**

   m Cookie header line identifies user to server

   m Server maintains a file about user

r **Desirable to send same user to same server in server farm**

**NAT problem:**

r Load balancing taken when SYN segment arrives

r SYN segment does not include HTTP header; thus no cookie header

**Surrogate solution:**

r Surrogate terminates TCP connection and receives HTTP request before choosing server

# Content Distribution

r  Server Farms

r  Proxy Web Caches

r  Content Distribution Networks (CDNs)

r  Peer-to-peer file sharing (P2P)

# Web caches (proxy server)

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

r   user sets browser: Web accesses via  cache

r   browser sends all HTTP requests to  cache

  m   object in cache: cache returns object

  m   else cache requests object from origin server, then returns object to client

# Configuring a Web Browser

In Netscape, go to Edit, Preferences, Advanced, Proxies:

# More about Web caching

r  Cache acts as both client and server

r  Cache can do up-to-date check using `If-modified-since` HTTP header

   m  Issue: should cache take risk and deliver cached object without checking?

   m  Heuristics are used.

r  Typically cache is installed by ISP (university, company, residential ISP)

## Why Web caching?

r  Reduce response time for client request.

r  Reduce traffic on an institution's access link.

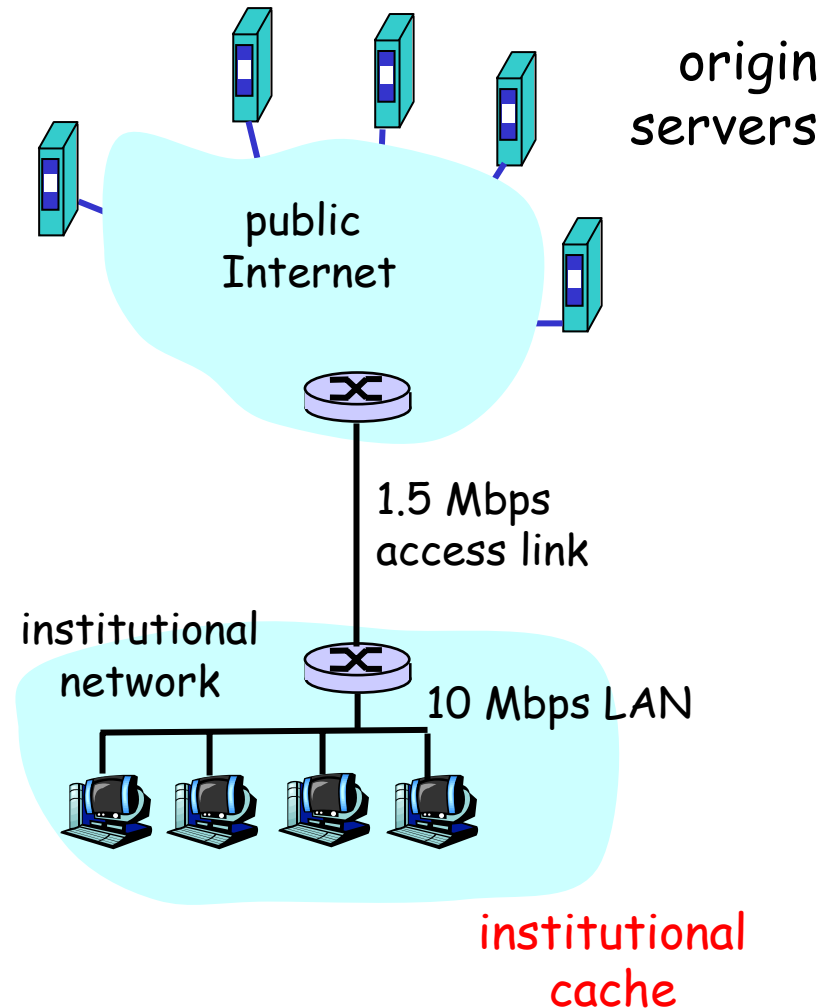r  Internet dense with caches enables "poor" content providers to effectively deliver content

# Caching example (1)

## Assumptions

r   average object size = 100,000 bits

r   avg. request rate from institution's browser to origin serves = 15/sec

r   delay from institutional router to any origin server and back to router  = 2 sec

## Consequences

r   utilization on LAN = 15%

r   utilization on access link = 100%

r   total delay   = Internet delay + access delay + LAN delay
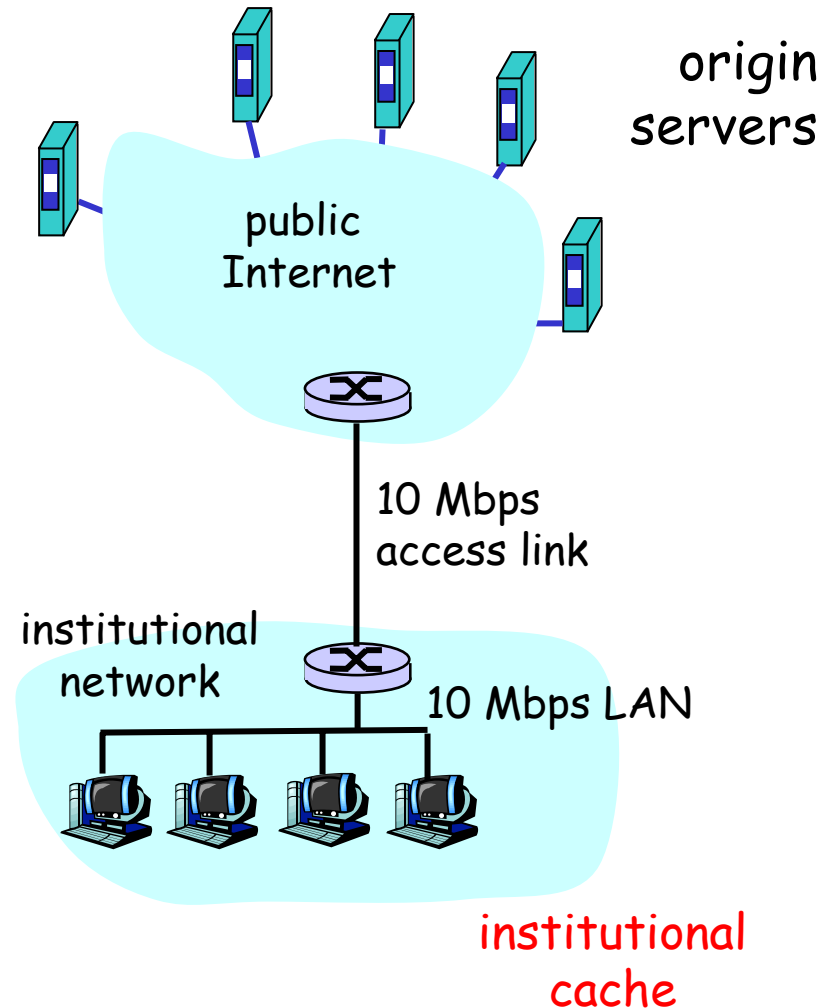
  =  2 sec + minutes + milliseconds



origin servers

public Internet

1.5 Mbps access link

institutional network

10 Mbps LAN

institutional cache

# Caching example (2)

## Possible solution

r   increase bandwidth of access link to, say, 10 Mbps

## Consequences

r   utilization on LAN = 15%

r   utilization on access link = 15%

r   Total delay = Internet delay + access delay + LAN delay

  = 2 sec + msecs + msecs

r   often a costly upgrade
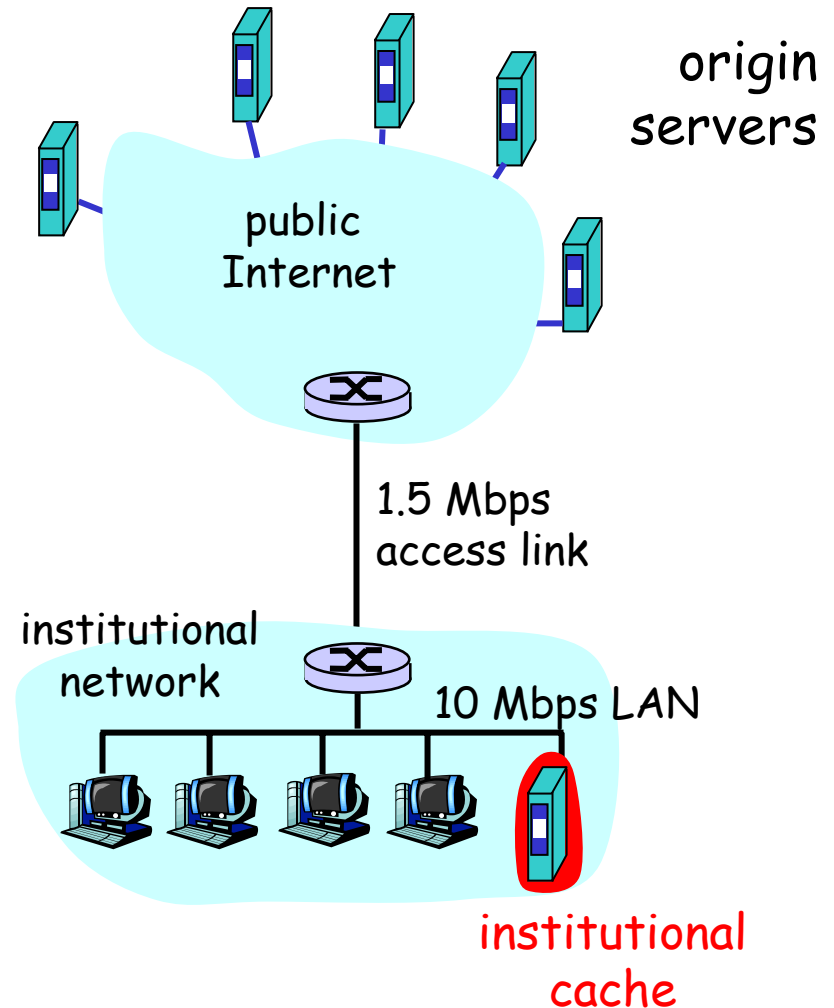
origin servers

public Internet

10 Mbps access link

institutional network

10 Mbps LAN

institutional cache

# Caching example (3)

**Install cache**

r    suppose hit rate is .4

**Consequence**

r    40% requests will be satisfied almost immediately

r    60% requests satisfied by origin server

r    utilization of access link reduced to 60%, resulting in negligible  delays (say 10 msec)

r    total delay   = Internet delay + access delay + LAN delay

   =  .6*2 sec + .6*.01 secs + milliseconds < 1.3 secs

public Internet

1.5 Mbps access link

institutional network

10 Mbps LAN

institutional cache

29

# Caching Challenges

r Cache consistency:

- m Caches often must guess whether a stored object is stale or fresh.

r Dynamic content:

- m Caches shouldn't cache outputs of CGI scripts.

r Hit counts and personalization:

- m Caches can cause hit count calculations and cookie transactions to fail.

r Less-savvy users and privacy-concerned users:

- m How do you get a user to point his browser to a cache?

r Enormous multimedia files:

- m Fortunately, disk storage is increasing at a rate of 60% a year!

# Cache Storage Management

- r LRU (Least Recently Used): Remove objects that have not been accessed for a long time. Example:
  - m $value_{LRU}$ = (50 - days_since_accessed)/50
  - m purge objects which have small $value_{LRU}$ when disk space starts to fill.
- r Weighted by retrieval time:
  - m record transfer time of object: retrieval_time
  - m value = $value_{LRU}$ * log ( retrieval_time +1)
- r Weighted by object size:
  - m value = $value_{LRU}$ * log (size +1)
- r Purge documents that have expired:
  - m But only if space is tight, as an up-to-date check can make an expired document fresh.

# Transparent Caching

## Non-transparent caching

r   Each browser in an ISP is manually configured to point to the correct cache.

  m   All AOL browsers

  m   System admin can configure campus browsers

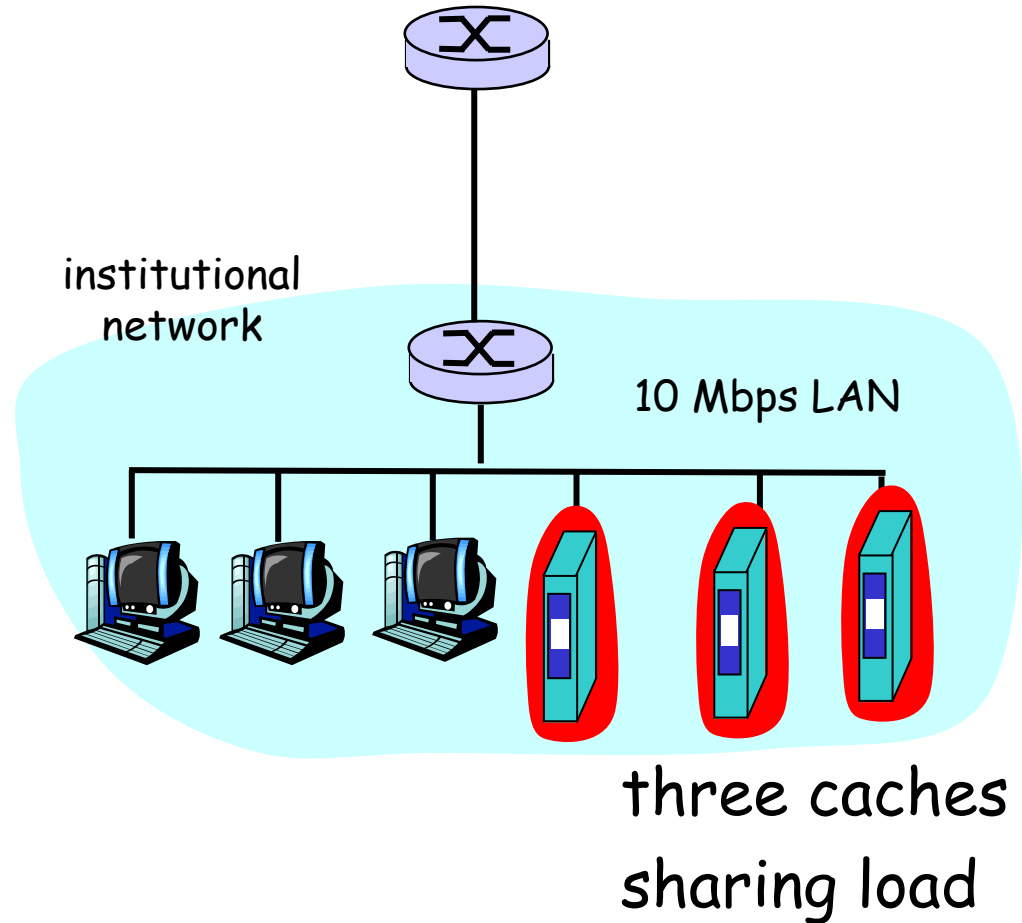r   But not when system administrator has no control over configuration parameters of browsers.

## Transparent caching

r   ISP redirects all datagrams with port number 80 (HTTP) to cache: "layer-4 switching".

r   Cache either retreives object from its cache or fetches document from Internet.

r   In the IP datagrams that carry the HTTP response message, cache must insert origin server's IP address!

# Hash Routing

r   Load balancing issue

r   Overview of hash routing

r   Robust hashing

r   Heterogeneous caches

r   CARP

r   Proxy Automatic Configuration

# Load balancing to multiple caches

institutional
network

10 Mbps LAN

three caches
sharing load

# Hash Routing Overview

r  Choose a hash function h() which maps URLs to a hash space.

  m Example:
  - hash space is {1,…,60}.
  - h() is the sum of the ASCII representation of the characters in the URL, modulo 60.

r  Partition hash space: one set for each cache.

  1) Client hashes URL, determines set to which hashed URL belongs, and sends request to corresponding cache.
  - Example: N=2 caches, set for cache 1 = {1,..,30}, set for cache 2 = {31,..,60}. If h(URLa) = 35, then client sends HTTP request to cache 2.

  2) If cache does not have object, it obtains object from origin server, stores a copy, and forwards a copy to the client.

r  Each object resides in at most one cache!

r  Client is immediately directed to the correct cache.
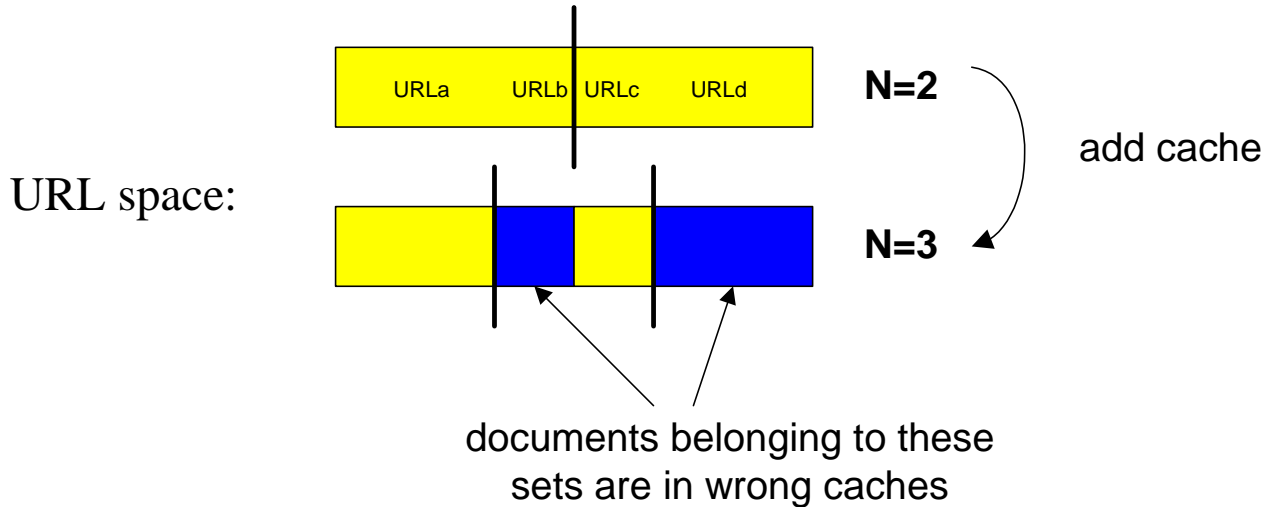
r  Do not pass through surrogate or NAT

# Hash Routing: Robustness Problem

When a cache is added or removed, a cached object can reside in the wrong cache.

**Example:** hash space = {1,2,..,60}
h(URLa) = 10, h(URLb) = 25, h(URLc) = 35, h(URLd) = 50

When there are N=2, URL space partitioned into {1,…,30}, {31,…,60}.

URL space:



documents belonging to these
sets are in wrong caches

- When we add a third cache, URL space is partitioned into {1,..,20}, {21,..,40}, {41,..,60}.
- Request for URLb is directed to cache 2; but URLb is in cache 1 so there is a MISS.
- Cache 2 will get URLb from origin server, so both caches 1 and 2 will contain URLb.

# Disruption Coefficient

D= Disruption Coefficient = fraction of objects in incorrect cache after adding or deleting a cache.

Without loss of generality, suppose hash space is set continuous interval [0,1] . Suppose there are N caches. The set [0,1] is partitioned:

$$[0,\frac{1}{N}], [\frac{1}{N},\frac{2}{N}],...,[\frac{N-1}{N},1]$$

Now suppose we add sibling N+1. The new partition is:

$$[0,\frac{1}{N+1}], [\frac{1}{N+1},\frac{2}{N+1}],...,[\frac{N}{N+1},1]$$

URLs in $[0,\frac{1}{N+1}]$ are in cache 1, which is correct. But URLs in $[\frac{1}{N+1},\frac{1}{N}]$ are in cache 1 when they should be in cache 2. Good intervals:

$$[0,\frac{1}{N+1}], [\frac{1}{N},\frac{2}{N+1}], [\frac{2}{N},\frac{3}{N+1}],...,[\frac{N-1}{N},\frac{N}{N+1}]$$

Sum of intervals is D=.5. Same result for deleting a cache.

**Thus hit rate is cut in half after a disruption: Not good!**

# Robust Hashing

r   Assign name to each cache: Joseph, Richard, Mary, Jane.

r   Choose hash function h(u,s) which is a function of both the URL u and the cache name s.

   m  for example, use h(u+s), where h() is a standard hash function and u+s is the string concatenation of u and s.

r   When a client wants URL u, client calculates the "scores" $h(u,s_1),…,h(u,s_N)$ for each of the N caches.

r   Client directs request for URL u to the cache s that has the highest score.

➔ If a cache fails, all remaining objects are still were they are supposed to be. Disruption coefficient is 1/N, which is typically small.

➔ If a cache is added, roughly the fraction $D = 1/(N+1)$ objects are in the wrong place.

# Heterogeneous Siblings

r   Processing power and storage capacity can vary greatly among siblings.

r   Target probabilities: $p_1, p_2, \ldots, p_N$

r   Introduce multipliers: $x_1, x_2, \ldots, x_N$

  (1) Calculate hash $h(u, s_n)$

  (2) Route URL u to sibling with highest weighted score

$$Z_n = x_n h(u, s_n)$$

r   How do we pick multipliers $x_1, x_2, \ldots, x_N$ to achieve target probabilities:

$$p_1, p_2, \ldots, p_N$$

EURECOM

# Multipliers

$$x_1 = (Np_1)^{1/N}$$

$$x_n = \left[ \frac{(N-n+1)(p_n - p_{n-1})}{\prod_{i=1}^{n-1} x_i} + x_{n-1}^{N-n+1} \right]^{\frac{1}{N-n+1}}$$

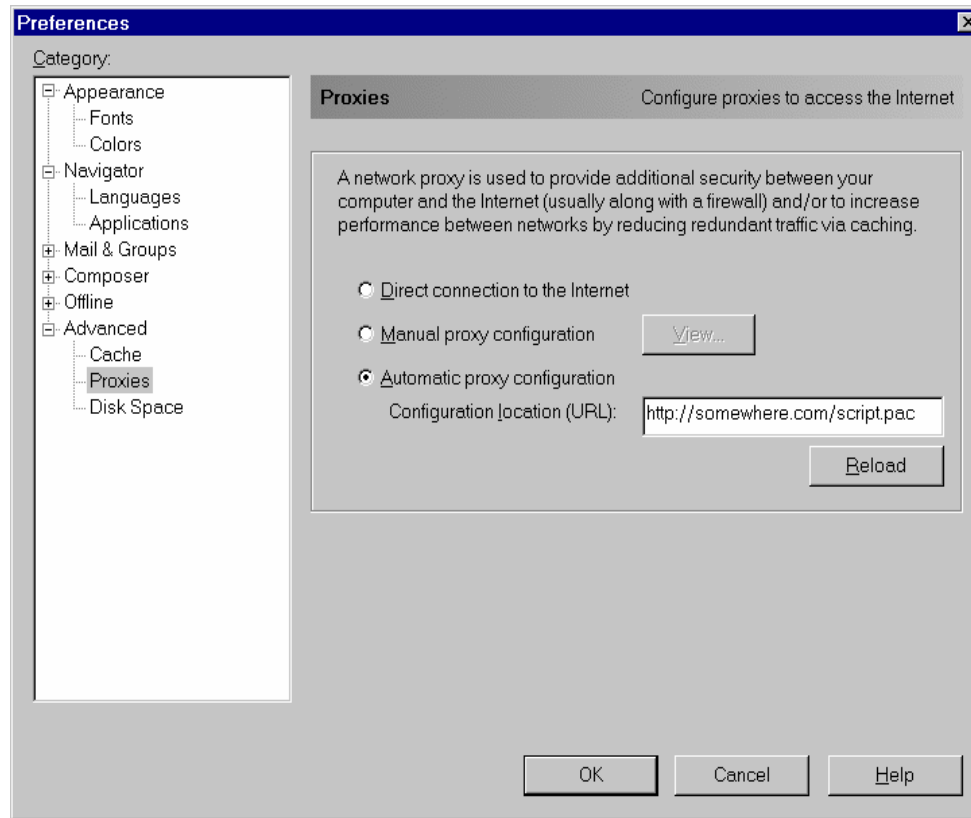Example:   $N = 3,\ p_1 = p_2 = \dfrac{1}{81},\ p_3 = \dfrac{79}{81}$

$$x_1 = x_2 = \frac{1}{3},\quad x_3 = 9$$

# CARP (Cache Array Routing Protocol)

r   Uses robust hash routing with multipliers.

r   All queries done over HTTP

    m   no new application-layer protocol such as ICP

    m   can take advantage of HTTP/1.1's rich set of headers

r   Internet draft (Valloppillil and Ross)

r   Implemented in Microsoft and Netscape cache server products
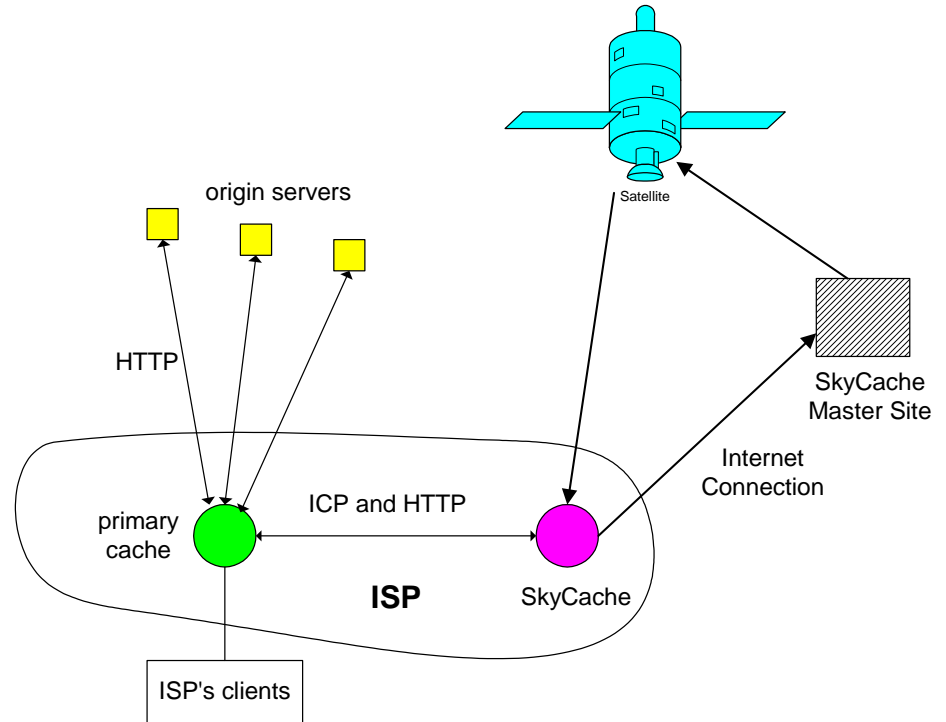
# Proxy Automatic Configuration (PAC)

r  Allows client browser to dynamically choose among a set of caches.
- m  Create an auto-configuration JavaScript file.
- m  Put file on nearby Web server (e.g., http://somewhere.com/script.pac).
- m  Have clients configure their browsers with the URL of the JavaScript.



Each time browser is initiated, browser obtains script.pac from nearby server. Script is run for each URL request.
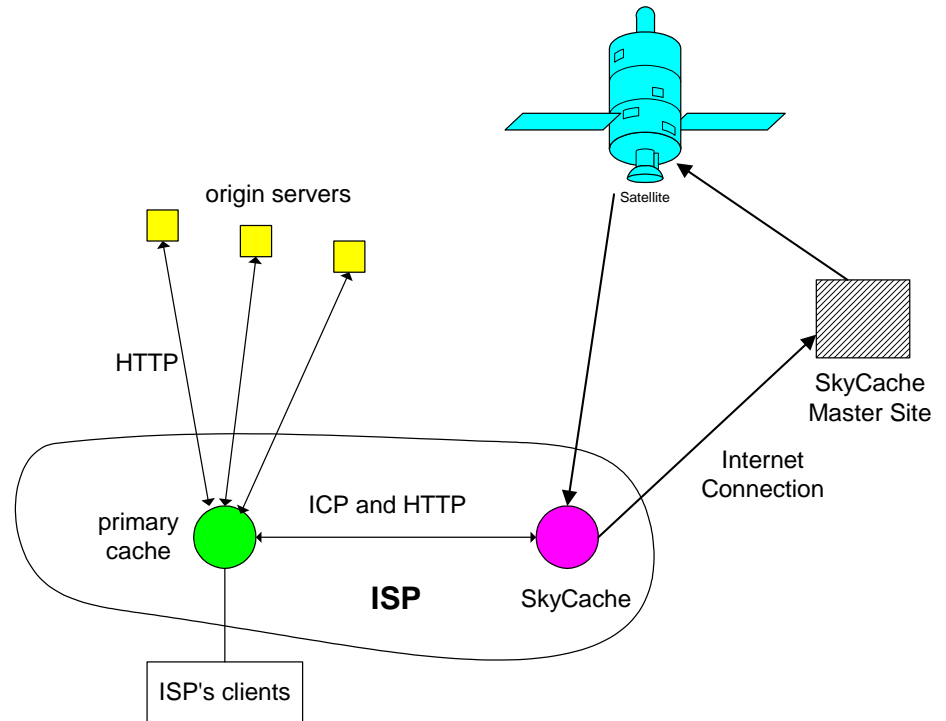
# Satellite Technology

r   Each local ISP has a cache with:
   m   Internet connection
   m   Huge storage capacity
   m   Satellite dish for receiving
r   Master site has:
   m   Internet connection
   m   Satellite transmitter

origin servers

Satellite

HTTP

SkyCache
Master Site

Internet
Connection

ICP and HTTP

primary
cache

ISP

SkyCache

ISP's clients

# Satellite Technology

r   How it works: When there is a miss at some local cache:

  m   that local cache obtains document from origin server using HTTP.

  m   local cache sends URL to master site.

  m   master site obtains document from origin server using HTTP.

  m   master site transmits document into satellite channel.

  m   *all* local caches receive document and cache it.

origin servers

HTTP

primary cache

ICP and HTTP

**ISP**

ISP's clients

Satellite

SkyCache Master Site

Internet Connection

SkyCache

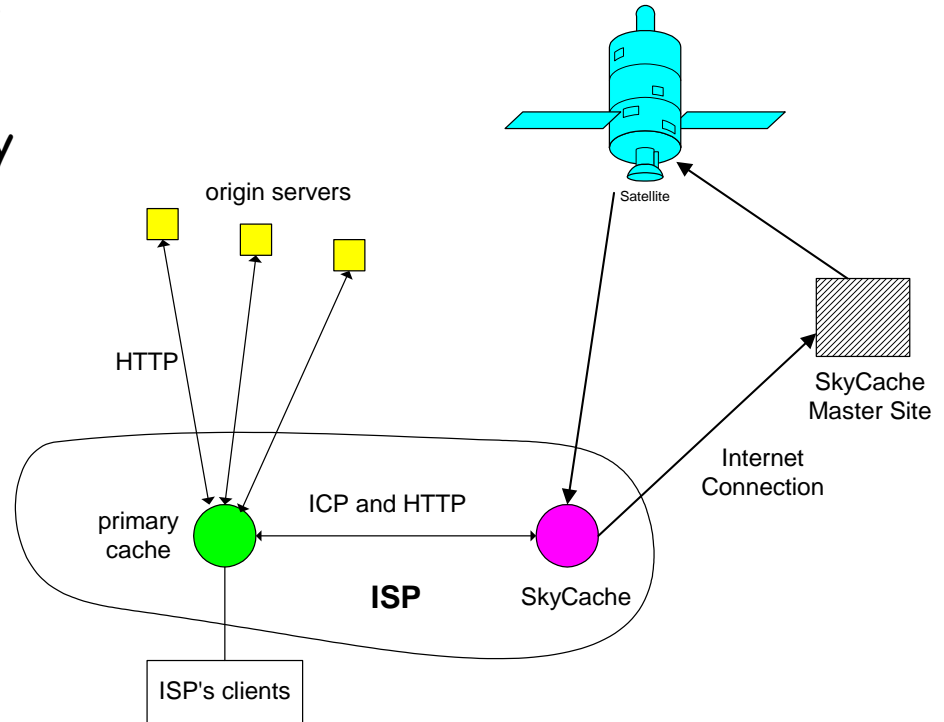# Satellite Technology: The Result

r   The user populations at each of the local ISPs are aggregated together to form one huge user population.

  m   The greater the user population, the greater the likelihood of repeated requests, the greater the hit rate.

r   Brings the Web to the edge of the network.

# A Product: SkyCache

r **SkyCache**: founded in 1997.

r Performance enhancer for an ISP's primary cache.

r Leases all equipment: fixed monthly payment to ISPs

r 4 Mbps one-way satellite link

r When there is a miss at primary cache:

  m primary cache queries SkyCache using ICP.

  m If SkyCache has the object, primary cache obtains object from SkyCache using HTTP.

  m If SkyCache doesn't have the object, primary cache directly obtains object from origin server.

  m SkyCache remembers the URLs for the misses and reports misses to Master Site.

origin servers

Satellite

HTTP

SkyCache Master Site

Internet Connection

ICP and HTTP

primary cache

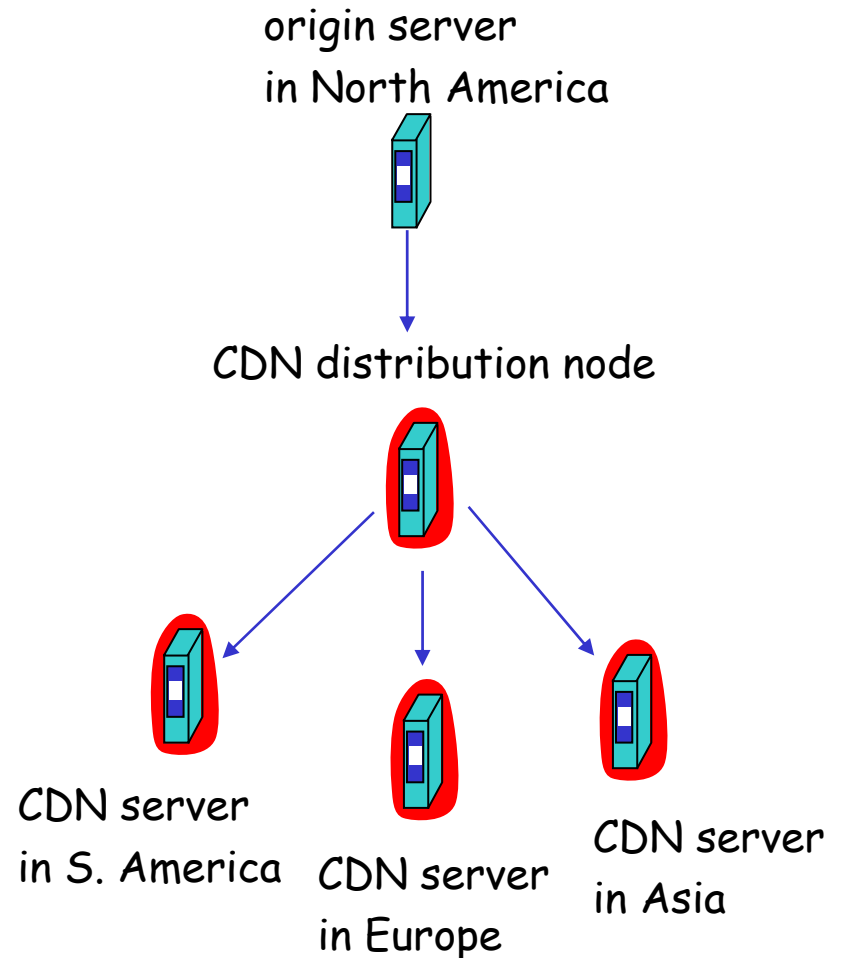**ISP**

SkyCache

ISP's clients

# Content Distribution

r  Server Farms

r  Proxy Web Caches

r  Content Distribution Networks (CDNs)

r  Peer-to-peer file sharing (P2P)
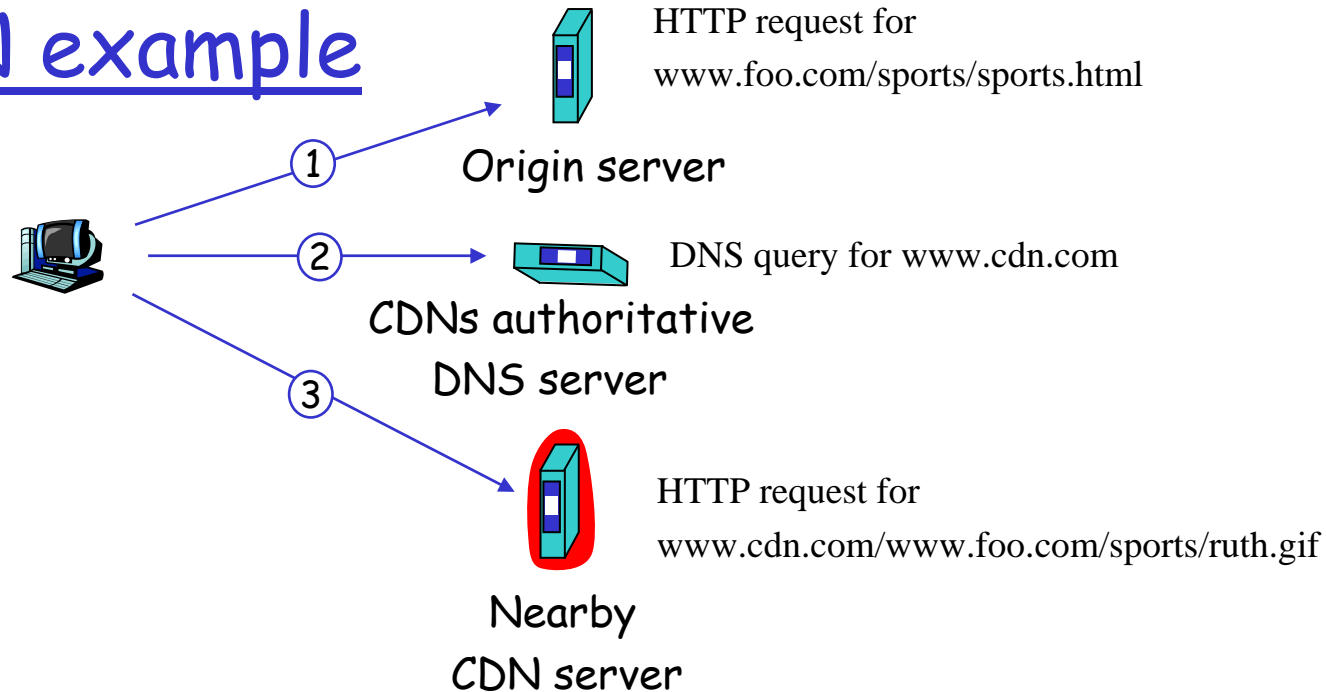
# Content distribution networks (CDNs)

r  The content providers are the CDN customers.

Content replication

r  CDN company installs hundreds of CDN servers throughout Internet
  m  in lower-tier ISPs, close to users

r  CDN replicates its customers' content in CDN servers. When provider updates content, CDN updates servers

origin server
in North America

CDN distribution node

CDN server
in S. America

CDN server
in Europe

CDN server
in Asia

# CDN example



1 → Origin server — HTTP request for www.foo.com/sports/sports.html

2 → CDNs authoritative DNS server — DNS query for www.cdn.com

3 → Nearby CDN server — HTTP request for www.cdn.com/www.foo.com/sports/ruth.gif

## origin server

r  www.foo.com

r  distributes HTML

r  Replaces:

http://www.foo.com/sports.ruth.gif

with

http://www.cdn.com/www.foo.com/sports/ruth.gif

## CDN company

r  cdn.com

r  distributes gif files

r  uses its authoritative DNS server to route redirect requests

# More about CDNs

## routing requests

r  CDN creates a "map", indicating distances from leaf ISPs and CDN nodes

r  when query arrives at authoritative DNS server:

  m  server determines ISP from which query originates

  m  uses "map" to determine best CDN server

## not just Web pages

r  streaming stored audio/video

r  streaming real-time audio/video

  m  CDN nodes create application-layer overlay network

# Content Distribution

r   Server farms

r   Proxy Web Caches

r   Content Distribution Networks (CDNs)

r   Peer-to-peer file sharing (P2P)

# P2P file sharing

r   Alice runs P2P client application on her notebook computer

r   Intermittently connects to Internet; gets new IP address for each connection

r   Asks for "Hey Jude"

r   Application displays other peers that have copy of Hey Jude.

r   Alice chooses one of the peers, Bob.

r   File is copied from Bob's PC to Alice's notebook: HTTP

r   While Alice downloads, other users uploading from Alice.

r   Alice's peer is both a Web client and a transient Web server.

All peers are servers = highly scalable!

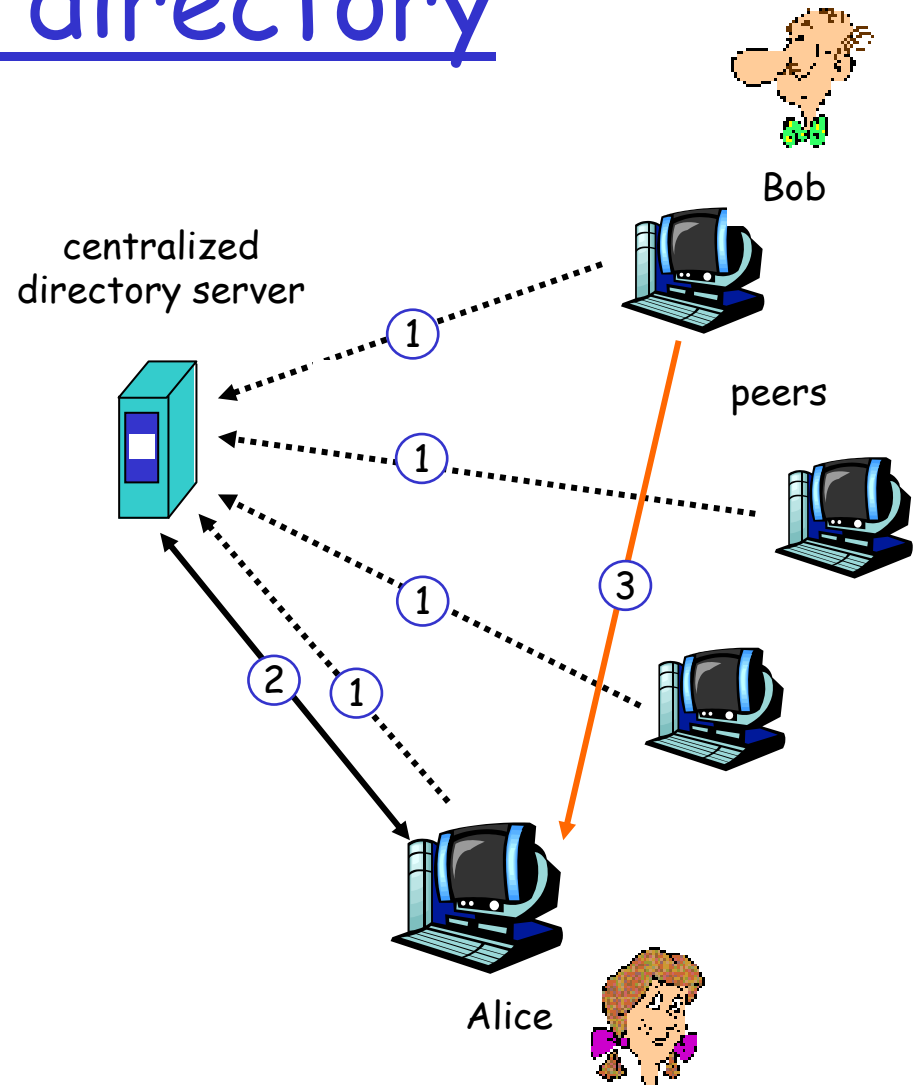# P2P: centralized directory

original "Napster" design
1) when peer connects, it informs central server:
   - m  IP address
   - m  content
2) Alice queries for "Hey Jude"
3) Alice requests file from Bob

Bob
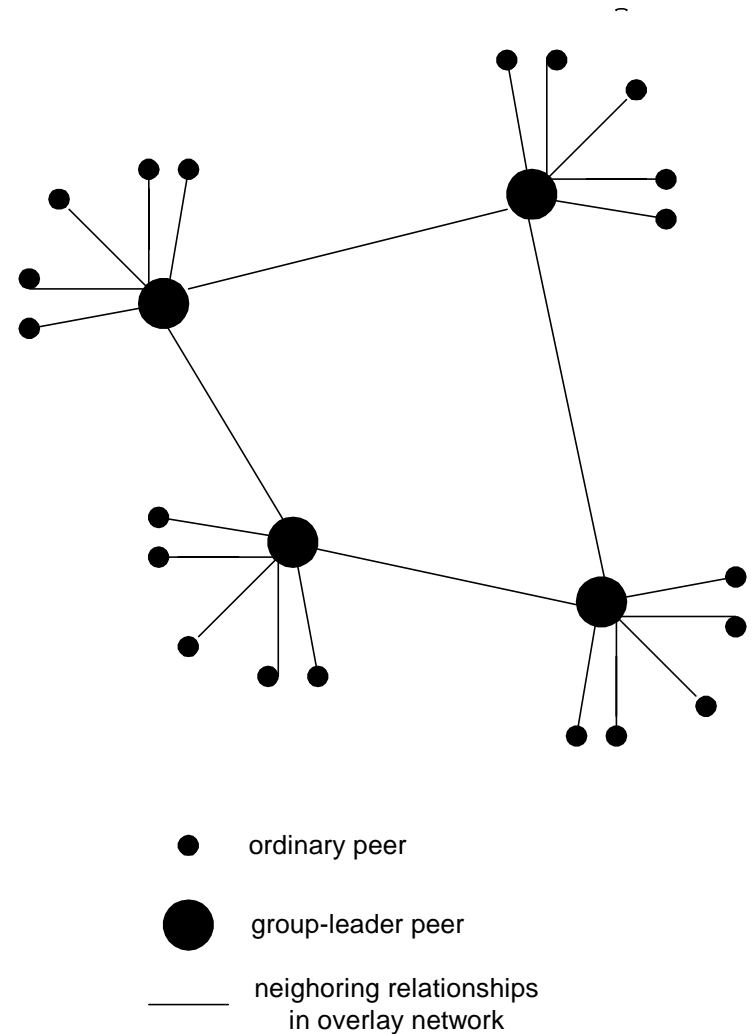
centralized directory server

peers

Alice

① ① ① ① ② ③

# P2P: problems with centralized directory

r  Single point of failure
r  Performance bottleneck
r  Copyright infringement

file transfer is decentralized, but locating content is highly  decentralized

# P2P: decentralized directory

r Each peer is either a group leader or assigned to a group leader.

r Group leader tracks the content in all its children.

r Peer queries group leader; group leader may query other group leaders.



● ordinary peer

● group-leader peer

_____ neighoring relationships
in overlay network

# More about decentralized directory

## overlay network

r   peers are nodes

r   edges between peers and their group leaders

r   edges between some pairs of group leaders

r   virtual neighbors

## bootstrap node

r   connecting peer is either assigned to a group leader or designated as leader
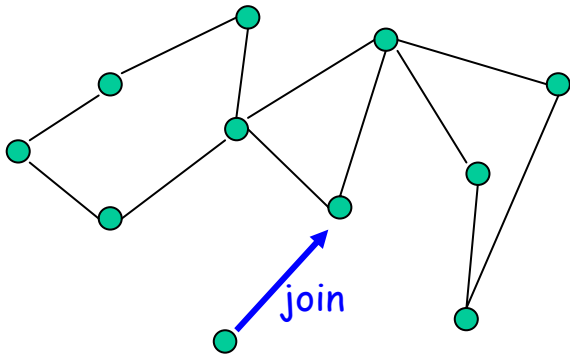
## advantages of approach

r   no centralized directory server

- m   location service distributed over peers
- m   more difficult to  shut down
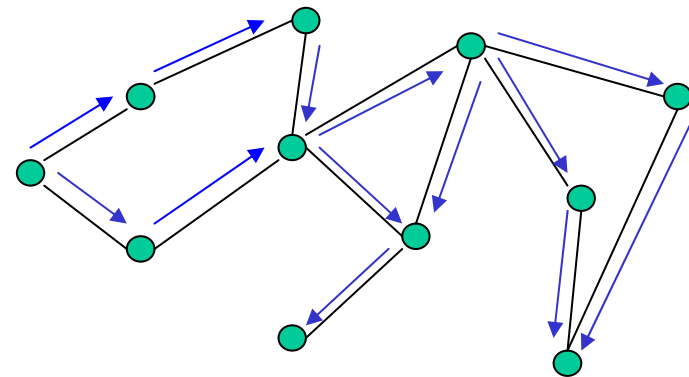
## disadvantages of approach

r   bootstrap node  needed

r   group leaders can get overloaded

# P2P: Query flooding

r  Gnutella

r  no hierarchy

r  use bootstrap node to learn about others

r  join message

r  Send query to neighbors

r  Neighbors forward query

r  If queried peer has object, it sends message back to querying peer

join

# P2P: more on query flooding

## Pros

r  peers have similar responsibilities: no group leaders

r  highly decentralized

r  no peer maintains directory info

## Cons

r  excessive query traffic

r  query radius: may not have content when present

r   bootstrap node

r  maintenance of overlay network