# Content Distribution Networks

COS 418: *Distributed Systems*
Lecture 24

Kyle Jamieson

---

## Today

**1. Domain Name System (DNS) primer**

2. The Web: HTTP, hosting, and caching

3. Content distribution networks (CDNs)

---

## DNS hostname versus IP address

- **DNS host name** (*e.g.* www.cs.princeton.edu)
  – **Mnemonic** name appreciated by humans
  – **Variable length**, full alphabet of characters
  – Provides **little** (if any) information about **location**

- **IP address** (*e.g.* 128.112.136.35)
  – Numerical address appreciated by **routers**
  – **Fixed length**, decimal number
  – **Hierarchical** address space, related to host **location**

---

## Many uses of DNS

- Hostname to IP address translation
  – IP address to hostname translation (*reverse lookup*)

- Host name *aliasing*: other DNS names for a host
  – *Alias* host names point to *canonical* hostname

- **Email:** Lookup domain's mail server by domain name

## Original design of the DNS

- Per-host file named /etc/hosts
  - Flat namespace: each **line = IP address & DNS name**
  - SRI (Menlo Park, California) kept the master copy
  - Everyone else downloads regularly

- **But, a single server doesn't scale**
  - Traffic implosion (lookups and updates)
  - Single point of failure

- Need a distributed, hierarchical **collection** of servers

5

## DNS: Goals and non-goals

- A wide-area **distributed database**

- Goals:
  - **Scalability**; decentralized maintenance
  - **Robustness**
  - Global scope
    - Names mean the same thing everywhere
  - Distributed updates/queries
  - Good **performance**
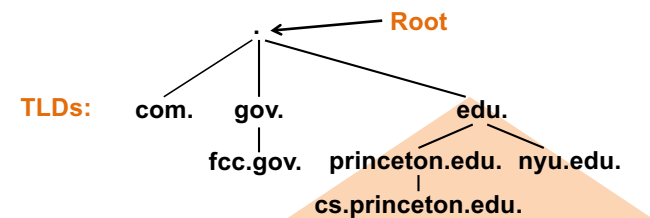
- But **don't need** strong **consistency** properties

6

## Domain Name System (DNS)

- **Hierarchical name space** divided into contiguous sections called *zones*
  - Zones are distributed over a collection of DNS servers

- **Hierarchy** of DNS **servers**:
  - *Root* servers (identity hardwired into other servers)
  - *Top-level domain (TLD)* servers
  - *Authoritative* DNS servers

- Performing the translations:
  - *Local DNS servers* located near clients
  - *Resolver* software running on clients

7

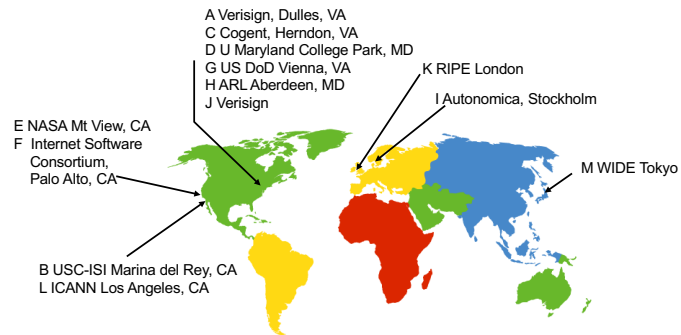## The DNS namespace is hierarchical



- **Hierarchy of namespace matches hierarchy of servers**

- Set of nameservers answers queries for names within zone

- Nameservers store names and links to other servers in tree

8

2

## DNS root nameservers

- 13 root servers.  *Does this scale?*

A Verisign, Dulles, VA
C Cogent, Herndon, VA
D U Maryland College Park, MD
G US DoD Vienna, VA
H ARL Aberdeen, MD
J Verisign

K RIPE London

I Autonomica, Stockholm

E NASA Mt View, CA
F  Internet Software
   Consortium,
   Palo Alto, CA

M WIDE Tokyo

B USC-ISI Marina del Rey, CA
L ICANN Los Angeles, CA

9

## DNS root nameservers

- 13 root servers.  *Does this scale?*
- Each server is really a **cluster** of servers (some geographically distributed), replicated via **IP anycast**

A Verisign, Dulles, VA
C Cogent, Herndon, VA (also Los Angeles, NY, Chicago)
D U Maryland College Park, MD
G US DoD Vienna, VA
H ARL Aberdeen, MD
J Verisign (21 locations)

K RIPE London (plus 16 other locations)

I Autonomica, Stockholm
(plus 29 other locations)

E NASA Mt View, CA
F  Internet Software
   Consortium,
   Palo Alto, CA
   (and 37 other locations)

M WIDE Tokyo
plus Seoul, Paris,
San Francisco

B USC-ISI Marina del Rey, CA
L ICANN Los Angeles, CA

10

## TLD and Authoritative Servers

- *Top-level domain (TLD)* servers
  - Responsible for com, org, net, edu, etc, and all top-level country domains: uk, fr, ca, jp
  - *Network Solutions* maintains servers for com TLD
  - *Educause* non-profit for edu TLD

- *Authoritative* DNS servers
  - An organization's DNS servers, providing authoritative information for that organization
  - May be maintained by organization itself, or ISP

11

## Local name servers

- Do not strictly belong to hierarchy

- Each ISP (or company, or university) has one
  - Also called *default* or *caching* name server

- When host makes DNS query, query is sent to its local DNS server
  - Acts as proxy, forwards query into hierarchy
  - Does work for the client

12

## DNS resource records

- DNS is a distributed database storing **resource records**
- Resource record includes: (**name**, type, **value**, time-to-live)

Type = **A** (address)
- **name** = hostname
- **value** is IP address

Type = **CNAME**
- **name** = alias for some "canonical" (real) name
- **value** is canonical name

Type = **NS** (name server)
- **name** = domain (e.g. princeton.edu)
- **value** is hostname of authoritative name server for this domain
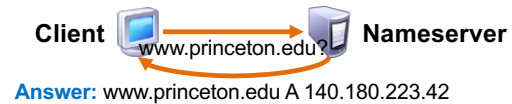
Type = **MX** (mail exchange)
- **name** = domain
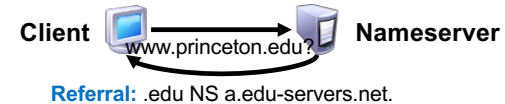- **value** is name of mail server for that domain

13

## DNS in operation

- Most queries and responses are UDP datagrams
  – Two types of queries:

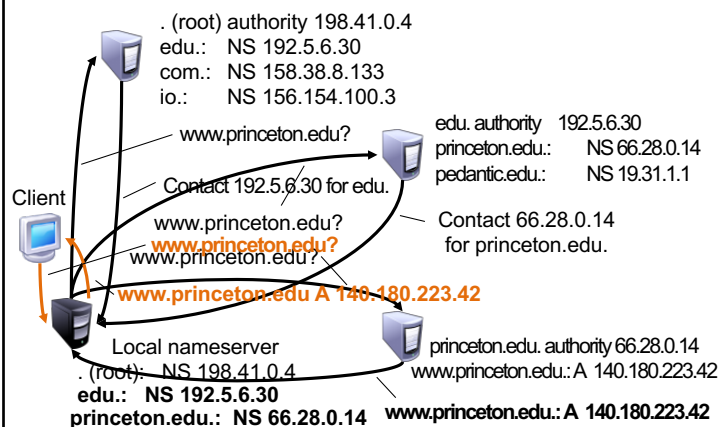- *Recursive*: Nameserver responds with answer or error

  **Client**  www.princeton.edu? **Nameserver**

  **Answer:** www.princeton.edu A 140.180.223.42

- *Iterative*: Nameserver may respond with a referral

  **Client**  www.princeton.edu? **Nameserver**

  **Referral:** .edu NS a.edu-servers.net.

14

## A recursive DNS lookup

. (root) authority 198.41.0.4
edu.:   NS 192.5.6.30
com.:   NS 158.38.8.133
io.:    NS 156.154.100.3

edu. authority    192.5.6.30
princeton.edu.:     NS 66.28.0.14
pedantic.edu.:      NS 19.31.1.1

www.princeton.edu?

Contact 192.5.6.30 for edu.

www.princeton.edu?
**www.princeton.edu?**
www.princeton.edu?

Client

Contact 66.28.0.14 for princeton.edu.

**www.princeton.edu A 140.180.223.42**

Local nameserver
. (root):   NS 198.41.0.4
**edu.:   NS 192.5.6.30**
**princeton.edu.:   NS 66.28.0.14**

princeton.edu. authority 66.28.0.14
www.princeton.edu.: A 140.180.223.42

**www.princeton.edu.: A 140.180.223.42**

15

## Recursive versus iterative queries

**Recursive query**

- Less burden on entity initiating the query

- **More burden on nameserver** (has to return an answer to the query)

- Most root and TLD servers won't answer (shed load)
  – Local name server answers recursive query

**Iterative query**

- **More burden on query initiator**

- Less burden on nameserver (simply refers the query to another server)

16

4

## Slide 17

```
$ dig @a.root-servers.net www.freebsd.org +norecurse
;; Got answer:
;; ->>HEADER<<- opcode: QUERY, status: NOERROR, id: 57494
;; QUERY: 1, ANSWER: 0, AUTHORITY: 2, ADDITIONAL: 2

;; QUESTION SECTION:
;www.freebsd.org.        IN A

;; AUTHORITY SECTION:
org.            172800 IN NS b0.org.afilias-nst.org.
org.            172800 IN NS d0.org.afilias-nst.org.

;; ADDITIONAL SECTION:
b0.org.afilias-nst.org.   172800 IN A   199.19.54.1
d0.org.afilias-nst.org.   172800 IN A   199.19.57.1
                                        Glue records
```

[Output edited for clarity]                                    17

## Slide 18

(authoritative for org.)

```
$ dig @199.19.54.1 www.freebsd.org +norecurse
;; Got answer:
;; ->>HEADER<<- opcode: QUERY, status: NOERROR, id: 39912
;; QUERY: 1, ANSWER: 0, AUTHORITY: 3, ADDITIONAL: 0

;; QUESTION SECTION:
;www.freebsd.org.        IN A

;; AUTHORITY SECTION:
freebsd.org.        86400   IN NS ns1.isc-sns.net.
freebsd.org.        86400   IN NS ns2.isc-sns.com.
freebsd.org.        86400   IN NS ns3.isc-sns.info.
```

[Output edited for clarity]                                    18

## Slide 19

(authoritative for freebsd.org.)

```
$ dig @ns1.isc-sns.net www.freebsd.org +norecurse
;; Got answer:
;; ->>HEADER<<- opcode: QUERY, status: NOERROR, id: 17037
;; QUERY: 1, ANSWER: 1, AUTHORITY: 3, ADDITIONAL: 3

;; QUESTION SECTION:
;www.freebsd.org.        IN A

;; ANSWER SECTION:
www.freebsd.org.    3600   IN A   69.147.83.33

;; AUTHORITY SECTION:
freebsd.org.    3600    IN NS ns2.isc-sns.com.
freebsd.org.    3600    IN NS ns1.isc-sns.net.
freebsd.org.    3600    IN NS ns3.isc-sns.info.

;; ADDITIONAL SECTION:
ns1.isc-sns.net.    3600    IN A    72.52.71.1
ns2.isc-sns.com.    3600    IN A    38.103.2.1
ns3.isc-sns.info.   3600    IN A    63.243.194.1
```

[Output edited for clarity]                                    19

## Slide 20

# DNS caching

- Performing all these queries takes time
  - And all this **before actual communication** takes place

- Caching can **greatly reduce overhead**
  - The top-level servers very rarely change
    - Popular sites visited often
  - Local DNS server often has the information cached

- How DNS caching works
  - All DNS servers **cache responses to queries**
  - Responses include a time-to-live (TTL) field
    - Server deletes cached entry after TTL expires

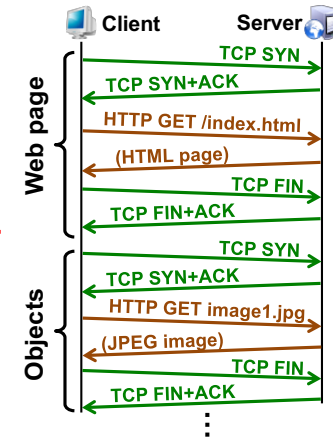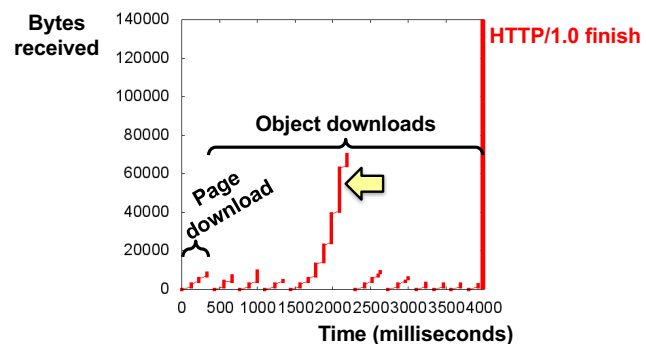Plays a key role in **CDN** (Akamai) **load balancing**

20

5

## Today

1. Domain Name System (DNS) primer

2. **The Web: HTTP, hosting, and caching**

3. Content distribution networks (CDNs)

## Anatomy of an HTTP/1.0 web page fetch

- Web page = HTML file + embedded images/objects

- *Stop-and-wait* at the granularity of objects:
  - Close then open new TCP connection for **each object**
    - Incurs a **TCP round-trip-time delay** each time
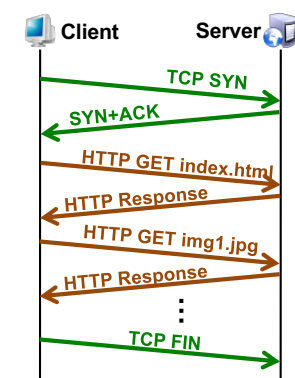
  - Each TCP connection may stay in "slow start"

**Client**      **Server**

**Web page**
- TCP SYN
- TCP SYN+ACK
- HTTP GET /index.html
- (HTML page)
- TCP FIN
- TCP FIN+ACK

**Objects**
- TCP SYN
- TCP SYN+ACK
- HTTP GET image1.jpg
- (JPEG image)
- TCP FIN
- TCP FIN+ACK

## HTTP/1.0 webpage fetch: Timeline

**Bytes received**



Object downloads

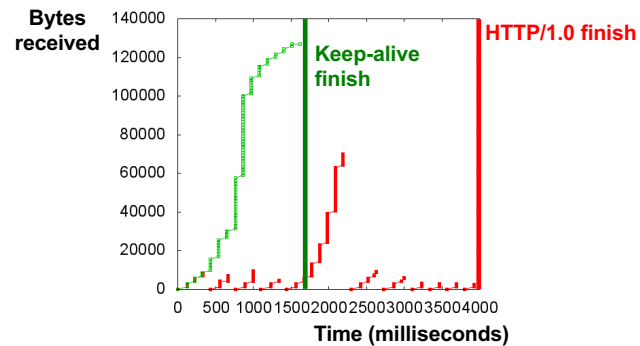HTTP/1.0 finish

Page download

**Time (milliseconds)**

- Fetch 8.5 Kbyte page with 10 objects, most < 10 Kbyte

## Letting the TCP connection persist

- Known as *HTTP keepalive*

- **Still stop-and-wait** at the granularity of objects, at the application layer

  - HTTP response fully received before next HTTP GET dispatched
    - ≥ 1 RTT per object

**Client**      **Server**
- TCP SYN
- SYN+ACK
- HTTP GET index.html
- HTTP Response
- HTTP GET img1.jpg
- HTTP Response
- TCP FIN

## HTTP Keepalive avoids TCP slow starts



**Bytes received** (y-axis: 0 to 140000)
**HTTP/1.0 finish**
**Keep-alive finish**
**Time (milliseconds)** (x-axis: 0 500 1000 1500 2000 2500 3000 3500 4000)

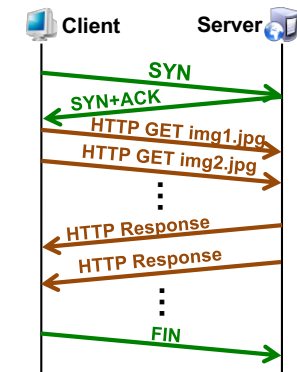Incur **one slow start**, **but stop-and-wait** to issue next request
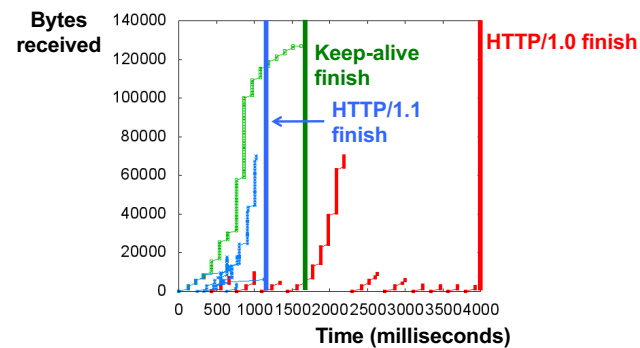
---

## Pipelining within HTTP

- **Idea: Pipeline** HTTP GETs and their responses

- Main benefits:
1. **Amortizes the RTT** across multiple objects retrieved

2. **Reduces overhead** of HTTP requests, packing multiple requests into one packet

- Implemented in HTTP/1.1



**Client**     **Server**
SYN
SYN+ACK
HTTP GET img1.jpg
HTTP GET img2.jpg
...
HTTP Response
HTTP Response
...
FIN

---

## Pipelined HTTP requests overlap RTTs



**Bytes received** (y-axis: 0 to 140000)
**HTTP/1.0 finish**
**Keep-alive finish**
**HTTP/1.1 finish**
**Time (milliseconds)** (x-axis: 0 500 1000 1500 2000 2500 3000 3500 4000)

- Many **HTTP requests** and **TCP connections** at once
- **Overlaps RTTs of all requests**

---

## Today

1. Domain Name System (DNS) primer

2. The Web: HTTP, **hosting, and caching**
   – **Handling heavy loads**

3. Content distribution networks (CDNs)
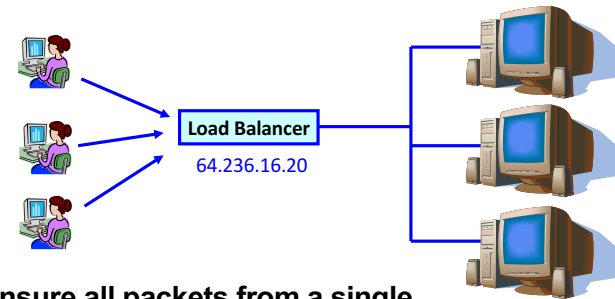
## Hosting: Multiple machines per site

- **Problem: Overloaded** popular web site
  - **Replicate** the site across multiple machines
    - Helps to handle the load

- Want to direct client to a particular replica. Why?
  - **Balance load** across server replicas

- **Solution #1:** Manual selection by clients
  - Each replica has its own site name
  - Some Web page lists replicas (*e.g.*, by name, location), asks clients to click link to pick

## Hosting: Load-balancer approach

- **Solution #2:** Single IP address, multiple machines
  - Run multiple machines behind a single IP address
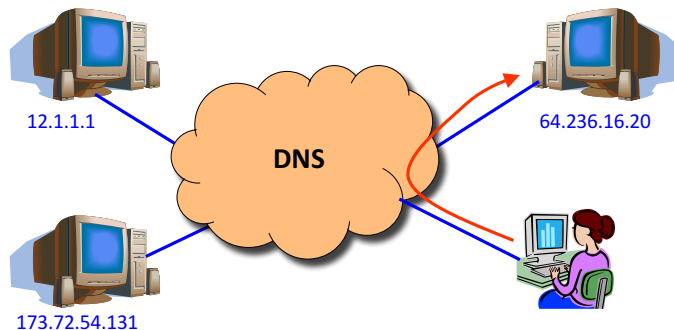


Load Balancer
64.236.16.20

  - **Ensure all packets from a single TCP connection go to the same replica**

## Hosting: DNS redirection approach

- **Solution #3:** Multiple IP addresses, multiple machines
  - Same DNS name but different IP for each replica
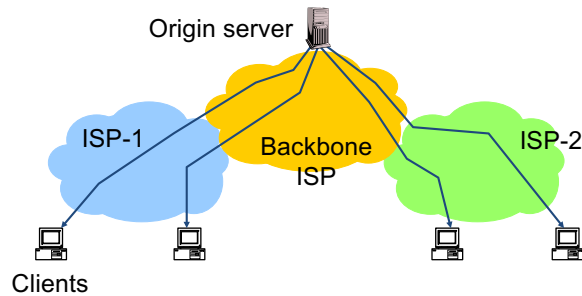    - DNS server returns IP addresses "round robin"



12.1.1.1

64.236.16.20

DNS

173.72.54.131

## Hosting: Summary

- Load-balancer approach
  - No geographical diversity ✘
  - TCP connection issue ✘
  - Does not reduce network traffic ✘

- DNS redirection
  - No TCP connection issues ✔
  - Simple round-robin server selection
    - May be less responsive ✘
  - Does not reduce network traffic ✘

## Web caching

- Many clients transfer the **same information**
  - Generates **redundant** server and network load
  - Also, clients may experience high **latency**

Origin server

ISP-1

Backbone
ISP
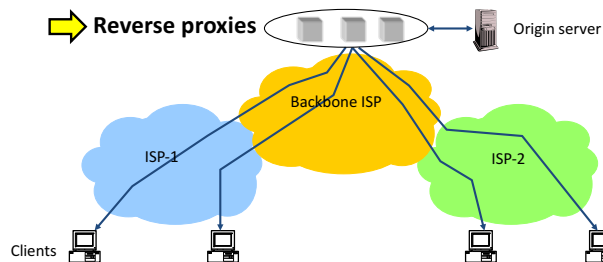
ISP-2

Clients

33

## Why web caching?

- Motivation for **placing content closer to client**:
  - User gets **better response time**
    - Content providers get happier users
  - Network gets **reduced load**

- Why does caching work?  Exploits locality of reference

- How well does caching work?
  - Very well, **up to a limit**
  - Large overlap in content
  - But many unique requests
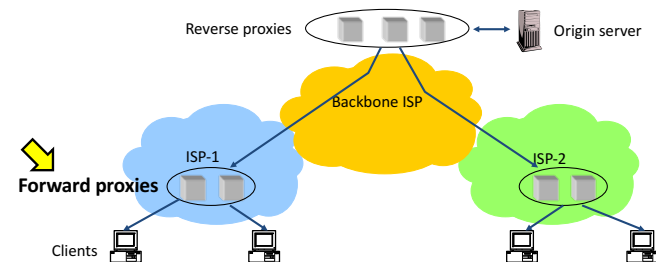
34

## Caching with Reverse Proxies

- Cache data close to origin server → decrease server load
  - Typically done by content providers
  - Client thinks it is talking to the origin server (the server with content)
- Does not work for **dynamic content**

**Reverse proxies**    Origin server

Backbone ISP

ISP-1    ISP-2

Clients

35

## Caching with Forward Proxies

- Cache close to clients → less network traffic, less latency
  - Typically done by ISPs or corporate LANs
  - **Client configured** to send HTTP requests to forward proxy

- Reduces traffic on ISP-1's access link, origin server, and backbone ISP

Reverse proxies    Origin server

Backbone ISP

ISP-1    ISP-2

**Forward proxies**

Clients

36

9

## Caching & Load-Balancing: Outstanding problems

- Problem *ca.* 2002: *How to reliably deliver large amounts of content to users worldwide?*

  - Popular event: **"Flash crowds" overwhelm** (replicated) web server, access link, or back-end database infrastructure

  - More rich content: audio, video, photos

- Web caching: Diversity causes **low cache hit rates (25−40%)**
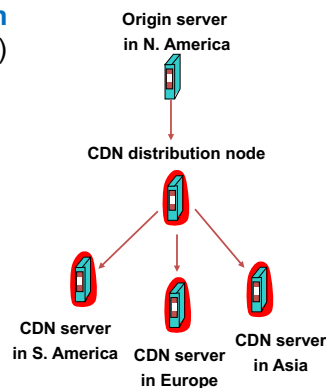
## Today

1. Domain Name System (DNS) primer

2. The Web: HTTP, hosting, and caching

3. **Content distribution networks (CDNs)**
   - Akamai case study

## Content Distribution Networks

- **Proactive content replication**
  - Content provider (*e.g.* CNN) pushes content out from its own *origin server*

- CDN **replicates** the content
  - On many servers spread throughout the Internet

- Updating the replicas
  - Updates **pushed to replicas** when the content changes



Origin server in N. America

CDN distribution node

CDN server in S. America

CDN server in Europe

CDN server in Asia

## Replica selection: Goals

- **Live** server
  - For availability

  **Requires continuous monitoring of liveness, load, and performance**
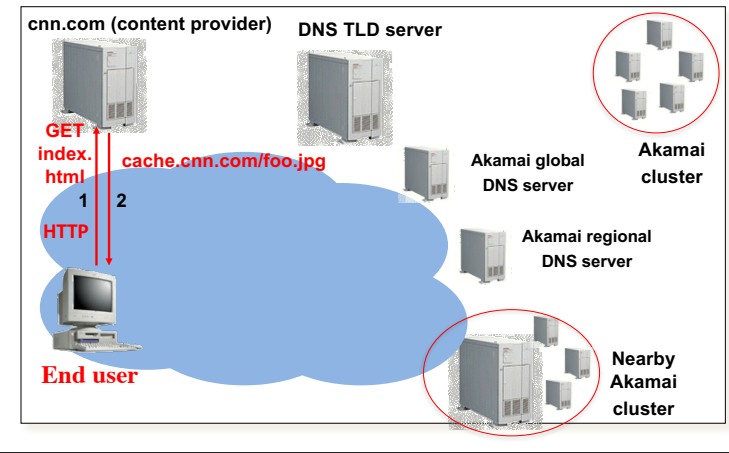
- Lowest **load**
  - To balance load across the servers

- **Closest**
  - Nearest geographically, or in round-trip time

- Best **performance**
  - Throughput, latency, reliability…

## Akamai statistics

- Distributed servers
  - Servers: ~100,000
  - Networks: ~1,000
  - Countries: ~70

- Many customers
  - Apple, BBC, FOX, GM IBM, MTV, NASA, NBC, NFL, NPR, Puma, Red Bull, Rutgers, SAP, …

- Client requests
  - 20+M per second
  - Half in the top 45 networks
  - 20% of all Web traffic worldwide

## How Akamai Uses DNS



**cnn.com (content provider)**  **DNS TLD server**

**GET index. html**  **cache.cnn.com/foo.jpg**

**Akamai global DNS server**

**Akamai cluster**

**1  2**

**HTTP**

**Akamai regional DNS server**

**End user**

**Nearby Akamai cluster**

## How Akamai Uses DNS



**cnn.com (content provider)**  **DNS TLD server**

**DNS lookup cache.cnn.com**

**Akamai global DNS server**

**Akamai cluster**

**1  2       3**

**4   ALIAS: g.akamai.net**

**Akamai regional DNS server**

**End user**

**Nearby Akamai cluster**

## How Akamai Uses DNS



**cnn.com (content provider)**  **DNS TLD server**

**DNS lookup g.akamai.net**

**Akamai global DNS server**

**Akamai cluster**

**1  2     3        5**

**4        6**

**ALIAS a73.g.akamai.net**

**Akamai regional DNS server**

**End user**

**Nearby Akamai cluster**

**11**

# How Akamai Uses DNS

cnn.com (content provider)　　DNS TLD server

Akamai global
DNS server

Akamai
cluster

1　2　　3　　　　5

4　　　　　6

7

DNS a73.g.akamai.net

Akamai regional
DNS server

8

Address
1.2.3.4

End user

Nearby
Akamai
cluster

---

# How Akamai Uses DNS

cnn.com (content provider)　　DNS TLD server

Akamai global
DNS server

Akamai
cluster

1　2　　3　　　　5

4　　　　　6

7

Akamai regional
DNS server

8

9

End user

GET /foo.jpg
Host: cache.cnn.com

Nearby
Akamai
cluster

---

# How Akamai Uses DNS

cnn.com (content provider)　　DNS TLD server

GET foo.jpg

11

CNN　12

Akamai global
DNS server

Akamai
cluster

1　2　　3　　　　5

4　　　　　6

7

Akamai regional
DNS server

8

9

End user

GET /foo.jpg
Host: cache.cnn.com

Nearby
Akamai
cluster
CNN

---

# How Akamai Uses DNS

cnn.com (content provider)　　DNS TLD server

11

12

Akamai global
DNS server

Akamai
cluster

1　2　　3　　　　5

4　　　　　6

7

Akamai regional
DNS server

8

9

End user

CNN　10

Nearby
Akamai
cluster
CNN

## How Akamai Works: Cache Hit



cnn.com (content provider)  DNS TLD server

Akamai global DNS server

Akamai cluster

1  2

Akamai regional DNS server

3

4

5

End user

6

Nearby Akamai cluster

CNN  CNN

## Mapping System

- Equivalence classes of IP addresses
  - IP addresses experiencing similar performance
  - Quantify how well they connect to each other

- **Collect and combine** measurements
  - Ping, traceroute, BGP routes, server logs
    - *e.g.*, over 100 TB of logs per days
  - Network latency, loss, throughput, and connectivity

50

## Routing client requests with the map

- Map each **IP class** to a preferred **server cluster**
  - Based on performance, cluster health, etc.
  - Updated roughly every minute

    - **Short, 60-sec DNS TTLs** in Akamai regional DNS accomplish this

- Map client request to a server in the cluster
  - **Load balancer** selects a specific server
  - *e.g.,* to **maximize** the **cache hit rate**

51

## Adapting to failures

- Failing **hard drive** on a server
  - Suspends after finishing "in progress" requests

- Failed **server**
  - Another server takes over for the IP address
  - Low-level map updated **quickly** (load balancer)

- Failed **cluster**, or **network path**
  - High-level map updated **quickly** (ping/traceroute)

52

## Take-away points: CDNs

- Content distribution is hard
  - Many, diverse, changing objects
  - Clients distributed all over the world

- **Moving content to the client** is key

  - Reduces latency, improves throughput, reliability

- Content distribution solutions **evolved:**
  - Load balancing, reactive caching, to
  - Proactive content distribution networks

53

**Friday precept:**
Extended office hours

54