

14-736

Distributed Systems

Lecture 25 * Spring 2019 * Kesden

Why Talk About DNS?

- ...this isn't a networks class, right?
- I argue that DNS is the most successful distributed system ever deployed:
 - Global scale
 - Nearly transparent, taken for granted
 - Highly available and reliable
 - Distributed Management
 - Etc, etc, etc.

Naming

- How do we efficiently locate resources?
 - DNS: name → IP address
- Challenge
 - How do we scale these to the wide area?

Obvious Solutions (1)

Why not centralize DNS?

- Single point of failure
- Traffic volume
- Distant centralized database
- Single point of update
- Doesn't *scale!*

Obvious Solutions (2)

Why not use /etc/hosts?

- Original Name to Address Mapping
 - Flat namespace
 - /etc/hosts
 - SRI kept main copy
 - Downloaded regularly
- Count of hosts was increasing: machine per domain → machine per user
 - Many more downloads
 - Many more updates

Domain Name System Goals

- Basically a wide-area distributed directory service
- Scalability
- Decentralized maintenance
- Robustness
- Global scope
 - Names mean the same thing everywhere
- Don't need
 - Atomicity
 - Strong consistency

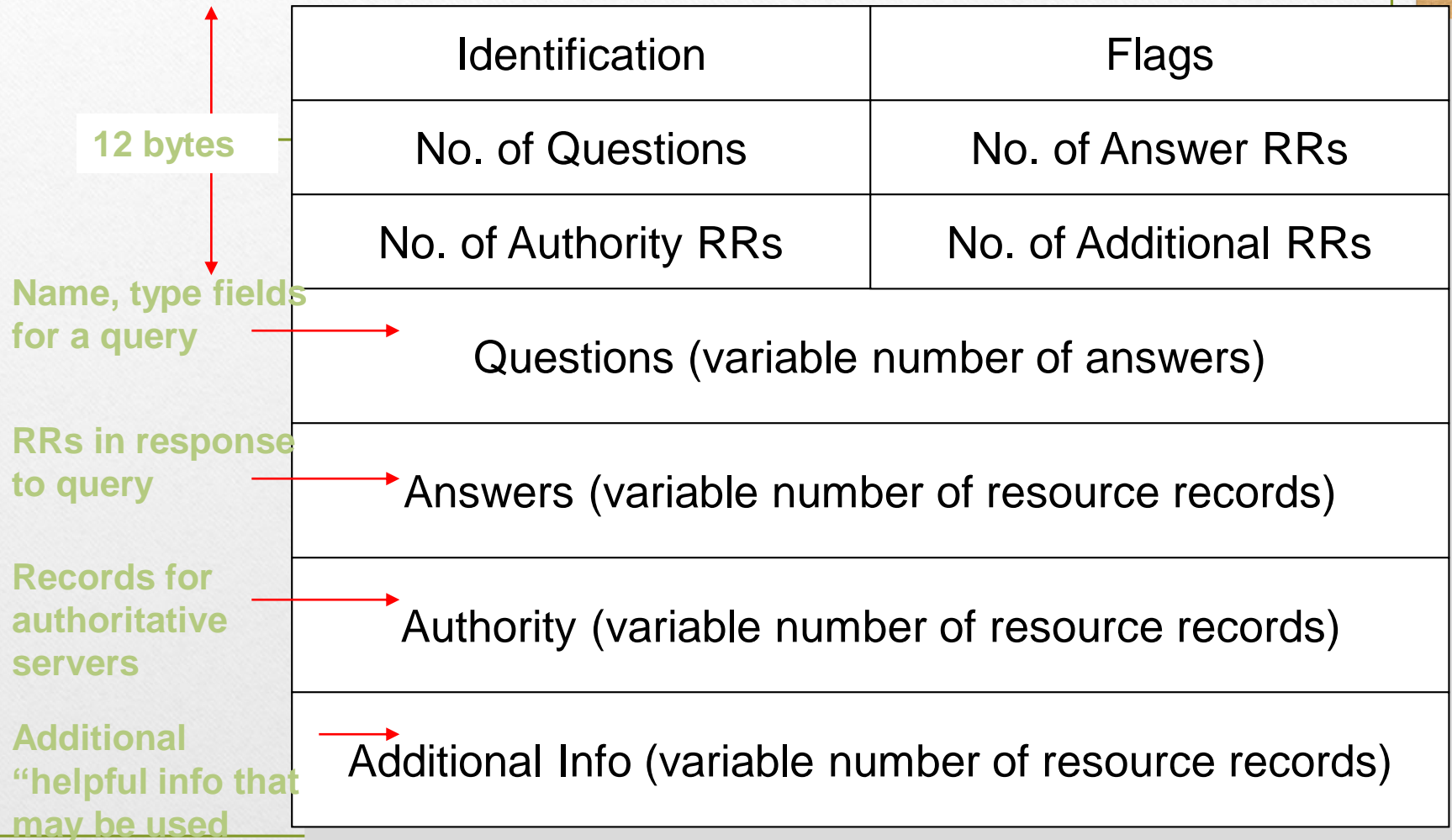
Programmer's View of DNS

- Conceptually, programmers can view the DNS database as a collection of millions of *host entry structures*:

```
/* DNS host entry structure */  
struct addrinfo {  
    int    ai_family; /* host address type (AF_INET) */  
    size_t ai_addrlen; /* length of an address, in bytes */  
    struct sockaddr *ai_addr; /* address! */  
    char  *ai_canonname; /* official domain name of host */  
    struct addrinfo *ai_next; /* other entries for host */  
};
```

- Functions for retrieving host entries from DNS:
 - `getaddrinfo`: query key is a DNS host name.
 - `getnameinfo`: query key is an IP address.

DNS Message Format



DNS Header Fields

- Identification
 - Used to match up request/response
- Flags
 - 1-bit to mark query or response
 - 1-bit to mark authoritative or not
 - 1-bit to request recursive resolution
 - 1-bit to indicate support for recursive resolution

DNS Records

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

- DB contains tuples called resource records (RRs)
 - Classes = Internet (IN), Chaosnet (CH), etc.
 - Each class defines value associated with type

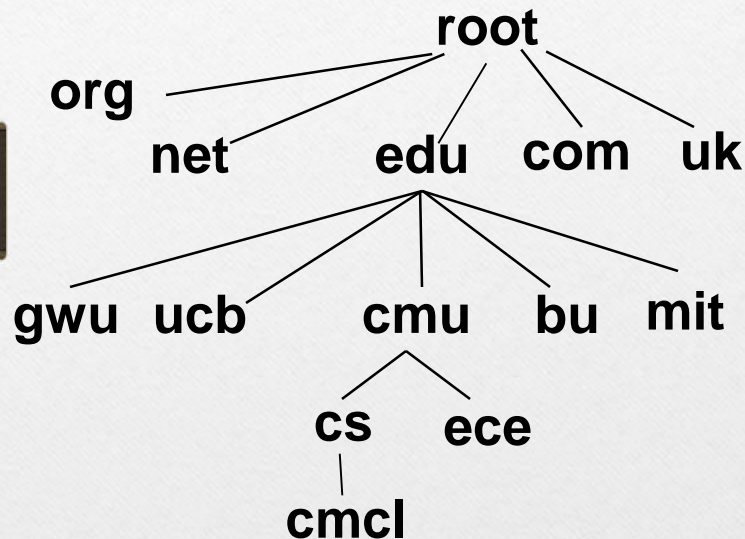
FOR IN class:

- Type=A
 - **name** is hostname
 - **value** is IP address
- Type=NS
 - **name** is domain (e.g. foo.com)
 - **value** is name of authoritative name server for this domain
- Type=CNAME
 - **name** is an alias name for some “canonical” (the real) name
 - **value** is canonical name
- Type=MX
 - **value** is hostname of mailserver associated with **name**

Properties of DNS Host Entries

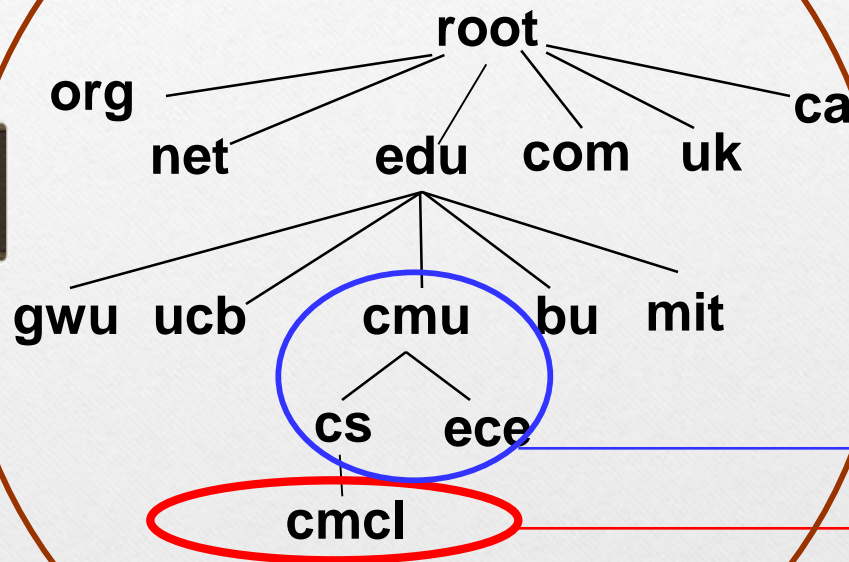
- Different kinds of mappings are possible:
 - Simple case: 1-1 mapping between domain name and IP addr:
 - `kittyhawk.cmcl.cs.cmu.edu` maps to `128.2.194.242`
 - Multiple domain names maps to the same IP address:
 - `eecs.mit.edu` and `cs.mit.edu` both map to `18.62.1.6`
 - Single domain name maps to multiple IP addresses:
 - `aol.com` and `www.aol.com` map to multiple IP addrs.
 - Some valid domain names don't map to any IP address:
 - for example: `cmcl.cs.cmu.edu`

DNS Design: Hierarchy Definitions



- Each node in hierarchy stores a list of names that end with same suffix
 - Suffix = path up tree
- E.g., given this tree, where would following be stored:
 - Fred.com
 - Fred.edu
 - Fred.cmu.edu
 - Fred.cmcl.cs.cmu.edu
 - Fred.cs.mit.edu

DNS Design: Zone Definitions



- Zone = contiguous section of name space
 - E.g., Complete tree, single node or subtree
- A zone has an associated set of name servers
 - Must store list of names and tree links

Subtree

Single node

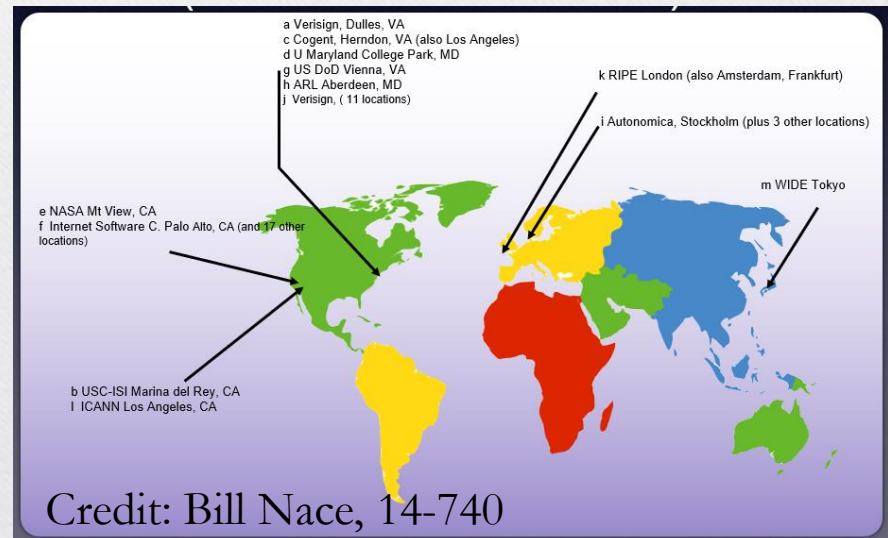
Complete Tree

DNS Design: Cont.

- Zones are created by convincing owner node to create/delegate a subzone
 - Records within zone stored multiple redundant name servers
 - Primary/master name server updated manually
 - Secondary/redundant servers updated by zone transfer of name space
 - Zone transfer is a bulk transfer of the “configuration” of a DNS server – uses TCP to ensure reliability
- Example:
 - CS.CMU.EDU created by CMU.EDU administrators
 - Who creates CMU.EDU or .EDU?

DNS: Root Name Servers

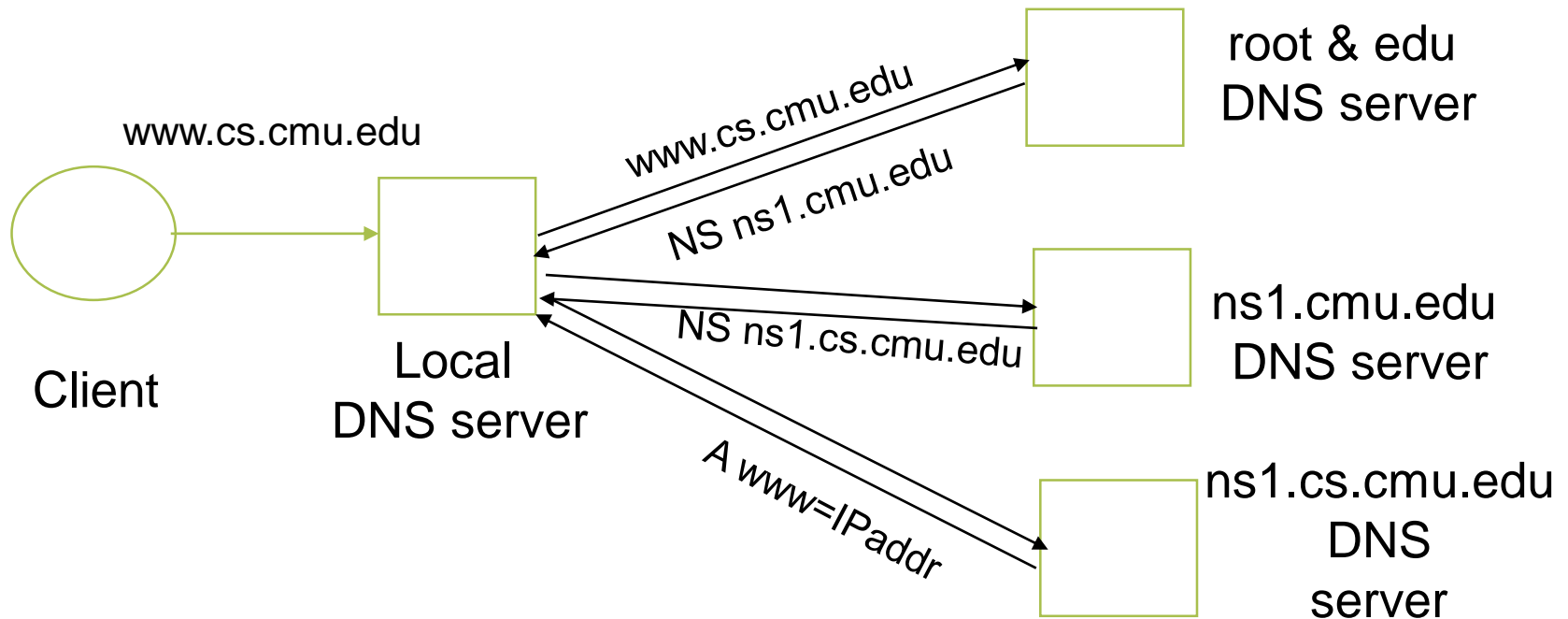
- Responsible for “root” zone
- Approx. 13 root name servers worldwide (well, clusters, thereof)
 - Currently {a-m}.root-servers.net
- Local name servers contact root servers when they cannot resolve a name
 - Configured with well-known root servers
 - www.root-servers.org



Servers/Resolvers

- Each host has a resolver
 - Typically a library that applications can link to
 - Local name servers hand-configured (e.g. /etc/resolv.conf)
- Name servers
 - Either responsible for some zone or...
 - Local servers
 - Do lookup of distant host names for local hosts
 - Typically answer queries about local zone

Typical Resolution



Typical Resolution

- Steps for resolving `www.cmu.edu`
 - Application calls `gethostbyname()` (RESOLVER)
 - Resolver contacts local name server (S_1)
 - S_1 queries root server (S_2) for (www.cmu.edu)
 - S_2 returns NS record for `cmu.edu` (S_3)
 - What about A record for S_3 ?
 - This is what the additional information section is for (PREFETCHING)
 - S_1 queries S_3 for www.cmu.edu
 - S_3 returns A record for www.cmu.edu
- Can return multiple A records → what does this mean?

Lookup Methods

Recursive query:

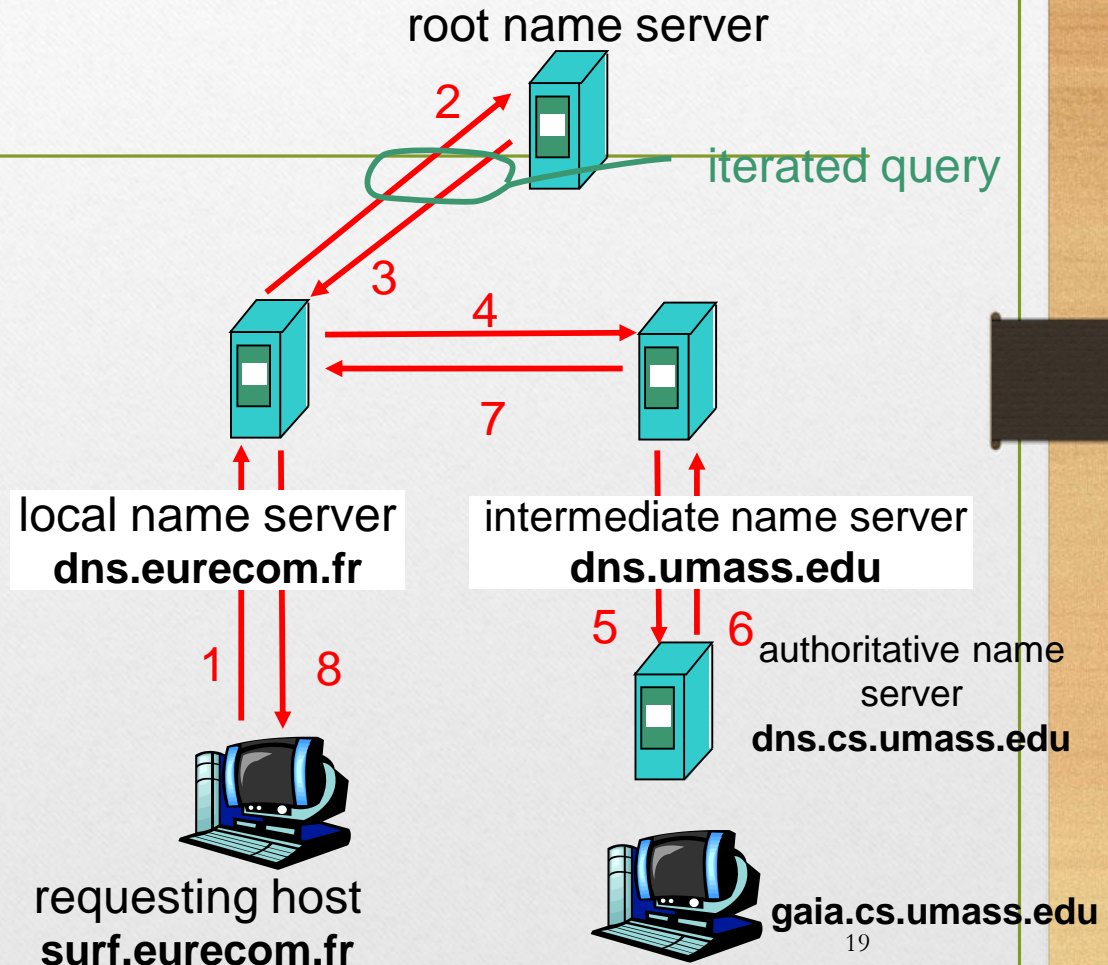
- Server goes out and searches for more info (recursive)
- Only returns final answer or “not found”

Iterative query:

- Server responds with as much as it knows (iterative)
- “I don’t know this name, but ask this server”

Workload impact on choice?

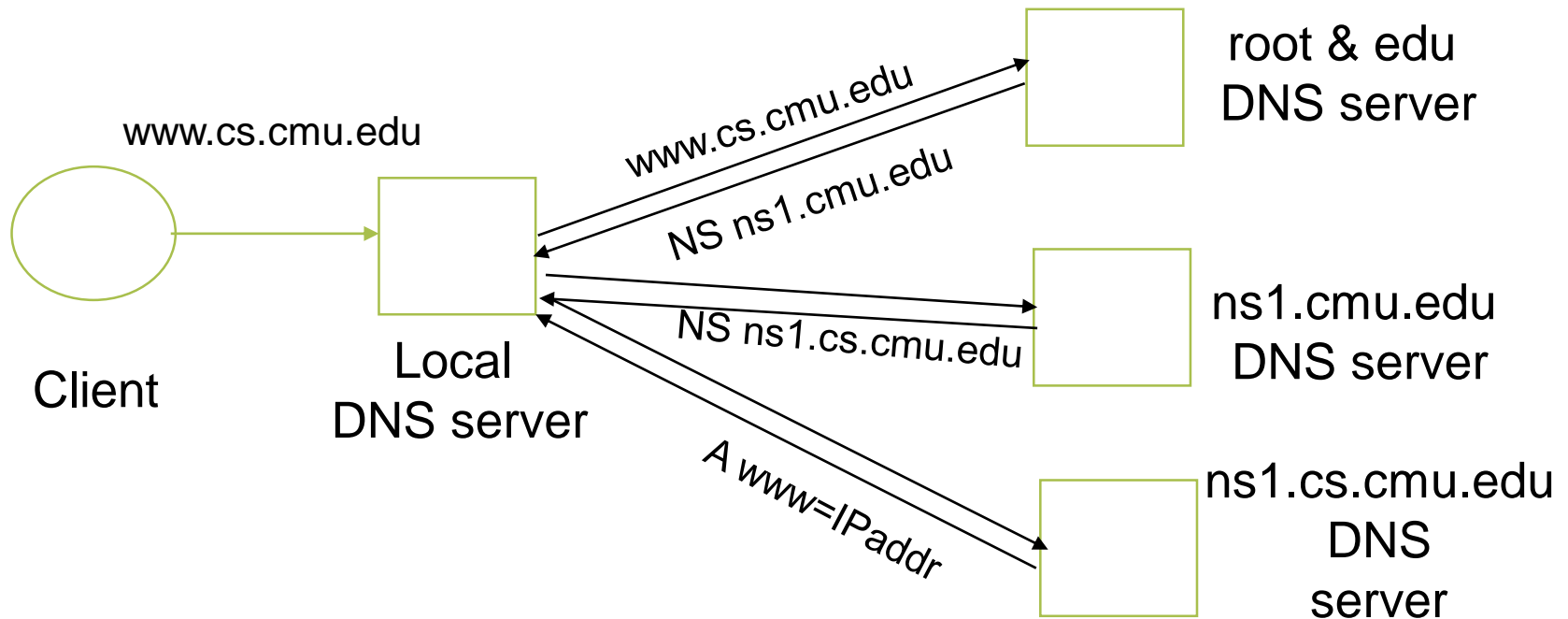
- Local server typically does recursive
- Root/distant server does iterative



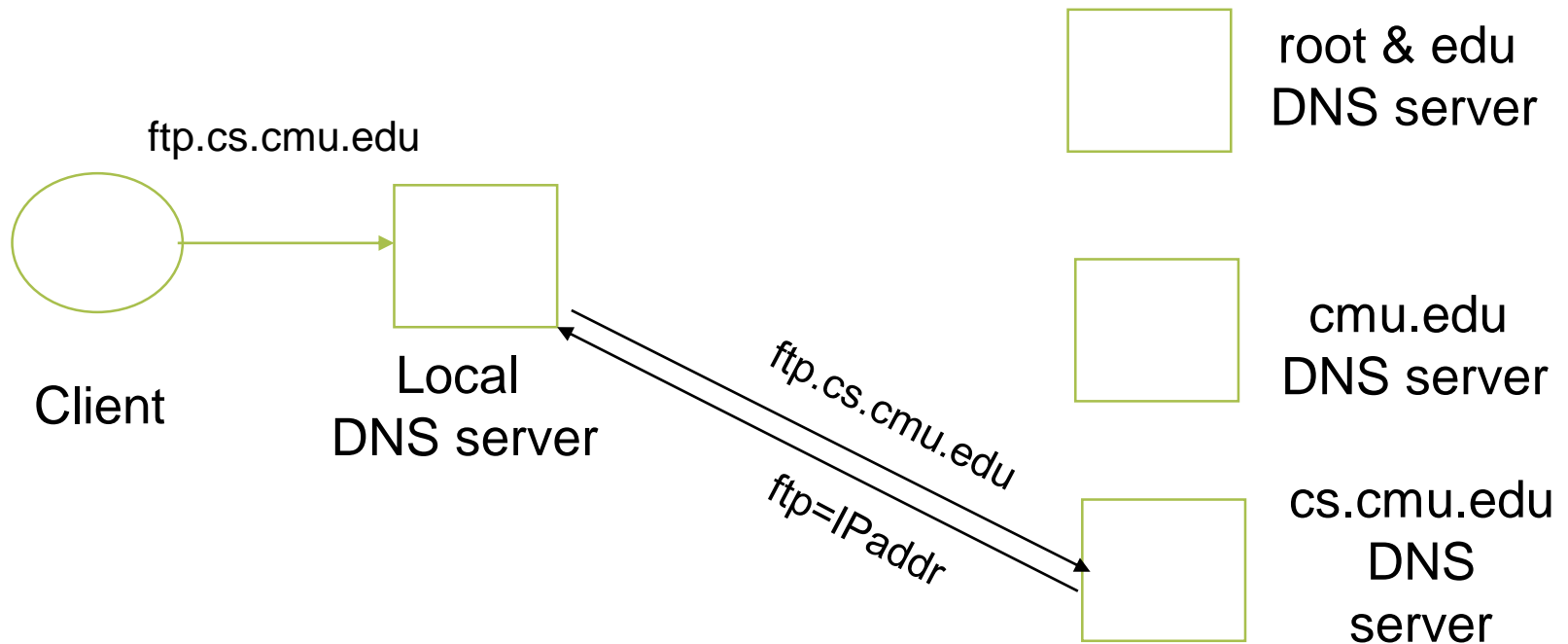
Workload and Caching

- Are all servers/names likely to be equally popular?
 - Why might this be a problem? How can we solve this problem?
- DNS responses are cached
 - Quick response for repeated translations
 - Other queries may reuse some parts of lookup
 - NS records for domains
- DNS negative queries are cached
 - Don't have to repeat past mistakes
 - E.g. misspellings, search strings in resolv.conf
- Cached data periodically times out
 - Lifetime (TTL) of data controlled by owner of data
 - TTL passed with every record

Typical Resolution



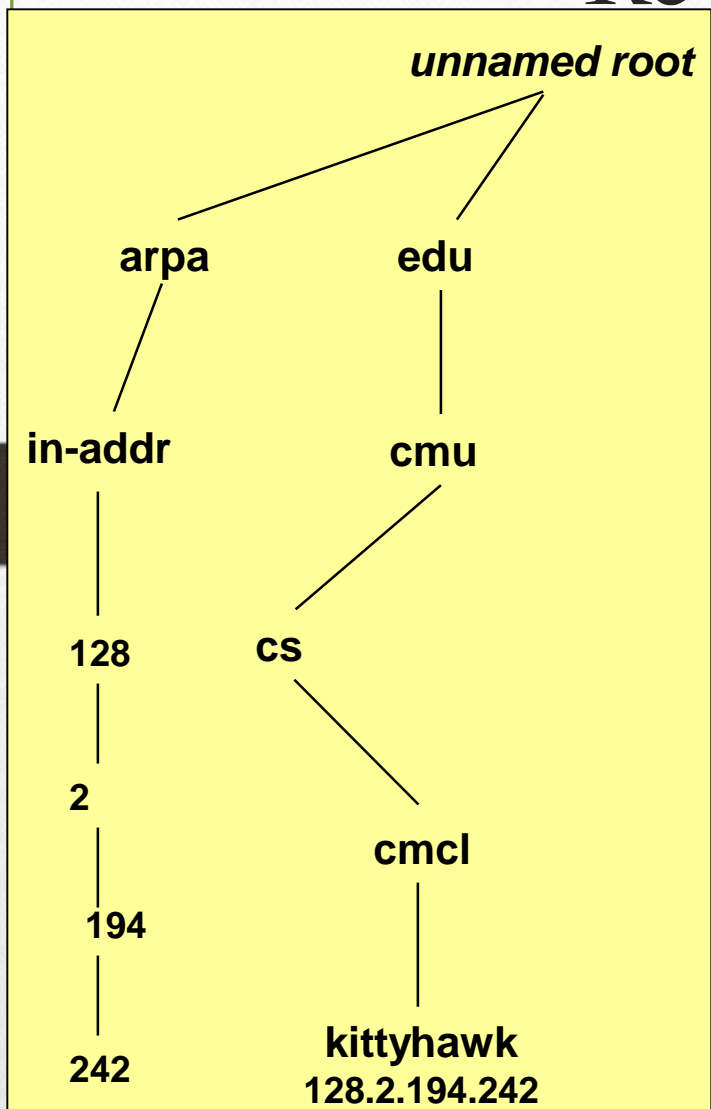
Subsequent Lookup Example



Reliability

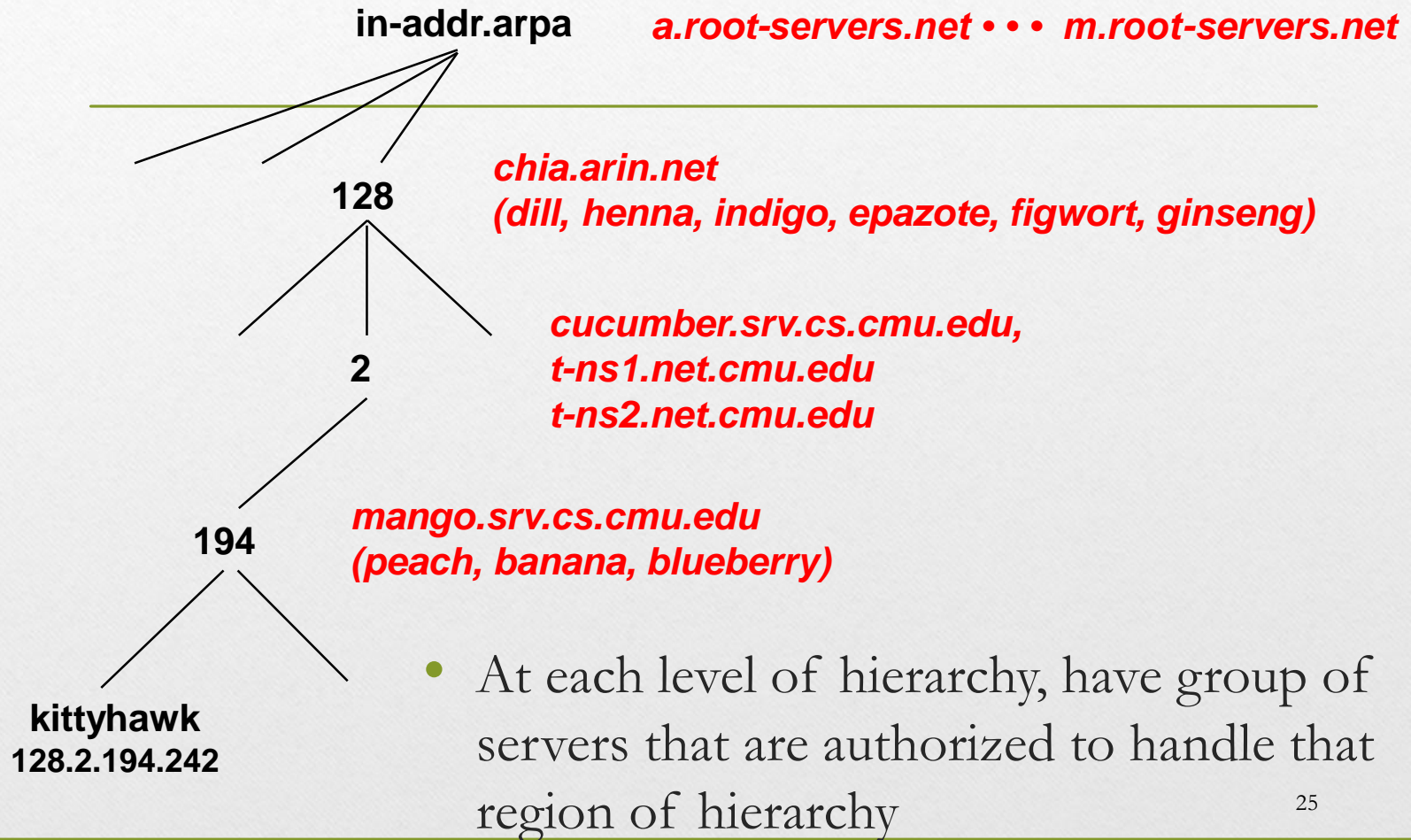
- DNS servers are replicated
 - Name service available if \geq one replica is up
 - Queries can be load balanced between replicas
- UDP used for queries
 - Need reliability \rightarrow must implement this on top of UDP!
 - Why not just use TCP?
- Try alternate servers on timeout
 - Exponential backoff when retrying same server
- Same identifier for all queries
 - Don't care which server responds

Reverse DNS



- Task
 - Given IP address, find its name
- Method
 - Maintain separate hierarchy based on IP names
 - Write 128.2.194.242 as 242.194.128.2.in-addr.arpa
 - Why is the address reversed?
- Managing
 - Authority manages IP addresses assigned to it
 - E.g., CMU manages name space 128.2.in-addr.arpa

.arpa Name Server Hierarchy



Prefetching

- Name servers can add additional data to response
- Typically used for prefetching
 - CNAME/MX/NS typically point to another host name
 - Responses include address of host referred to in “additional section”

Mail Addresses

- MX records point to mail exchanger for a name
 - E.g. mail.acm.org is MX for acm.org
- Addition of MX record type proved to be a challenge
 - How to get mail programs to lookup MX record for mail delivery?
 - Needed critical mass of such mailers

Root Zone

- Generic Top Level Domains (gTLD) = .com, .net, .org, etc...
- Country Code Top Level Domain (ccTLD) = .us, .ca, .fi, .uk, etc...
- Root server ({a-m}.root-servers.net) also used to cover gTLD domains
 - Wow, how times have changed!

gTLDs

- Unsponsored
 - .com, .edu, .gov, .mil, .net, .org
 - .biz → businesses
 - .info → general info
 - .name → individuals
- Sponsored (controlled by a particular association)
 - .aero → air-transport industry
 - .cat → catalan related
 - .coop → business cooperatives
 - .jobs → job announcements
 - .museum → museums
 - .pro → accountants, lawyers, and physicians
 - .travel → travel industry
- Starting up
 - .mobi → mobile phone targeted domains
 - .post → postal
 - .tel → telephone related
- Etc.

Measurements of DNS

- No centralized caching per site
 - Each machine runs own caching local server
 - Why is this a problem?
 - How many hosts do we need to share cache? → recent studies suggest 10-20 hosts
- “Hit rate for DNS = 60-80% → $1 - (\#DNS / \#connections)$ ”
 - Depends upon number of users of cache.
 - Larger community means greater hit rate, to a point.
- Lower TTLs for A records does not affect performance
- DNS performance really relies more on NS-record caching

DNS (Summary)

- Motivations → large distributed database
 - Scalability
 - Independent update
 - Robustness
- Hierarchical database structure
 - Zones
 - How is a lookup done
- Caching/prefetching and TTLs
- Reverse name lookup
- What are the steps to creating your own domain?