



Distributed Systems Naming

Olaf Landsiedel

Previous Session?

- Election
 - Bully algorithm
 - Ring algorithm
- Mutual Exclusion

Today...

- Naming
 - What is it?
 - Why do we need it?
 - How does it work?

What is a name?

- Definition (general, non computer science):
 - a word or set of words by which a person or thing is known, addressed, or referred to.



Naming things

- User names
 - Login, email
- Machine names
 - rlogin, email, web
- Files
- Devices
- Variables in programs
- Network services

Naming Service

Allows you to look up names

• Often returns an address as a response

Might be implemented as

- Search through file
- Client-server program
- Database query
- ...

What's a name?

Name: identifies what you want

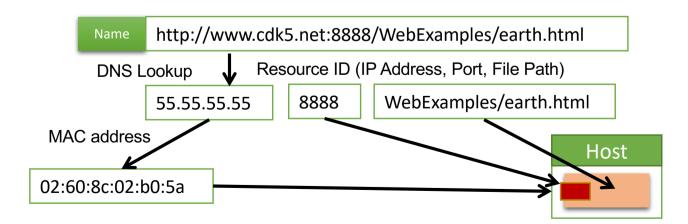
Address: identifies where it is

Route: identifies how to get there

Binding: associates a name with an address

Naming

- Names are used to uniquely identify entities in Distributed Systems
 - Entities may be processes, remote objects, newsgroups, ...
- Names are mapped to an entity's location using a name resolution
- An example of name resolution



Names, Addresses and Identifiers

An entity can be identified by three types of references

1. Name

- A name is a set of bits or characters that references an entity
- Names can be human-friendly (or not)

2. Address

- Every entity resides on an access point, and access point has an address
- Addresses may be location-dependent (or not)
- e.g., IP Address + port

3. Identifier

- Identifiers are names that *uniquely* identify an entity
- A true identifier is a name with following properties:
 - a. An identifier refers to at-most one entity
 - b. Each entity is referred to by at-most one identifier
 - c. An identifier always refers to the same entity (i.e. it is never reused)

Naming Systems

- A naming system is simply a middleware that assists in name resolution
- Naming systems are classified into three classes based on the type of names used:
 - a. Flat naming
 - b. Structured naming
 - c. Attribute-based naming

Classes of Naming

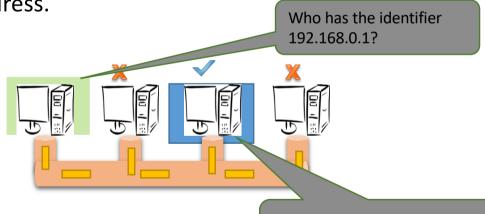
- Flat naming
- Structured naming
- Attribute-based naming

Flat Naming

- In Flat Naming, identifiers are simply random bits of strings (known as unstructured or flat names)
- Flat name does not contain any information on how to locate an entity
- We will study four types of name resolution mechanisms for flat names:
 - 1. Broadcasting
 - 2. Forwarding pointers
 - 3. Home-based approaches
 - 4. Distributed Hash Tables

1. Broadcasting

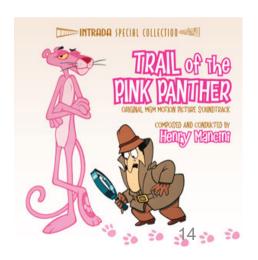
- Approach: Broadcast the identifier to the complete network. The entity associated with the identifier responds with its current address
- Example: Address Resolution Protocol (ARP)
 - Resolve an IP address to a MAC address.
 - In this application,
 - IP address is the identifier of the entity
 - MAC address is the address of the access point
- Challenges:
 - Not scalable in large networks
 - This technique leads to flooding the network with broad
 - Requires all entities to listen to all requests



I am 192.168.0.1. My address is 02:AB:4A:3C:59:85

2. Forwarding Pointers

- Forwarding Pointers enables locating mobile entities
 - Mobile entities move from one access point to another
- When an entity moves from location A to location B, it leaves behind (in A) a reference to its new location at B
- Name resolution mechanism
 - Follow the chain of pointers to reach the entity
 - Update the entity's reference when the present location is found
- Challenges:
 - Reference to at-least one pointer is necessary
 - Long chains lead to longer resolution delays
 - Long chains are prone to failure due to broken links



Forwarding Pointers – An Example

- Stub-Scion Pair (SSP) chains implement remote invocation for mobile entities using forwarding pointers
 - Server stub is referred to as Scion in the original paper
- Each forwarding pointer is implemented as a pair:

```
(client stub, server stub)
```

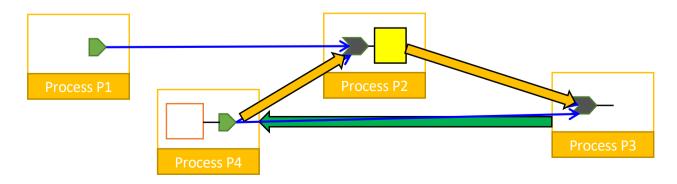
 The server stub contains a local reference to the actual object or a local reference to the remote client stub

= Remote Object; = Caller Object; = Client stub; = Server stub

- When object moves from A to B,
 - It leaves client stub in its place

= Process n;

It installs server stub that refers to the new remote client stub on B

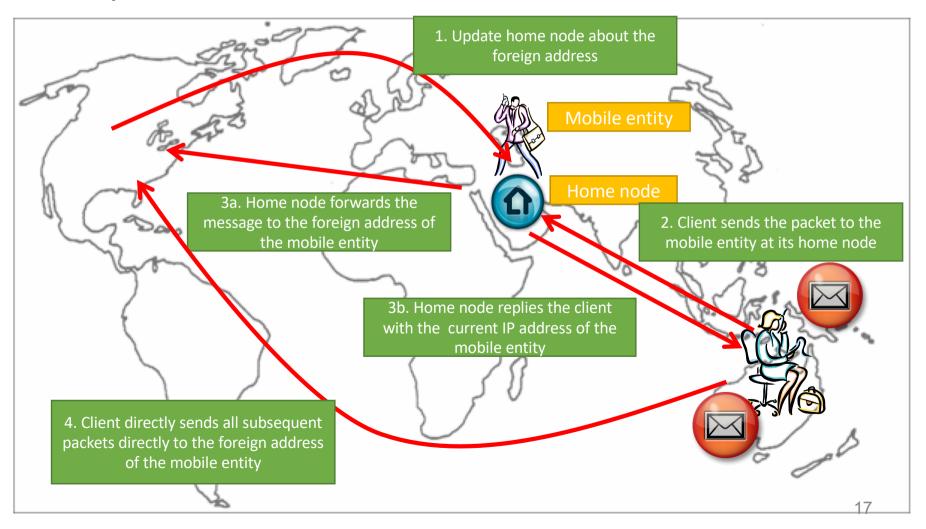


3. Home-based approaches

- Each entity is assigned a home node
 - Home node is typically static (has fixed access point and address)
 - Home node keeps track of current address of the entity
- Entity-home interaction:
 - Entity's home address is registered at a naming service
 - Entity updates the home about its current address (foreign address) whenever it moves
- Name resolution
 - Client contacts the home to obtain the foreign address
 - Client then contacts the entity at the foreign location

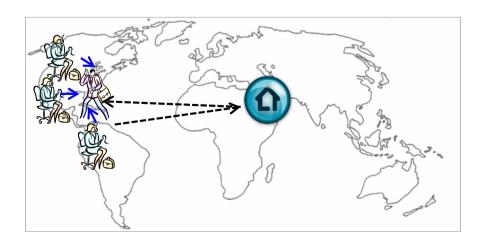
3. Home-based approaches – An example

Example: Mobile-IP



3. Home-based approaches – Challenges

- Home address is permanent for an entity's lifetime
 - If the entity permanently moves, then a simple home-based approach incurs higher communication overhead
- Connection set-up overheads due to communication between the client and home can be excessive
 - Consider the scenario where the clients are nearer to the mobile entity than the home entity

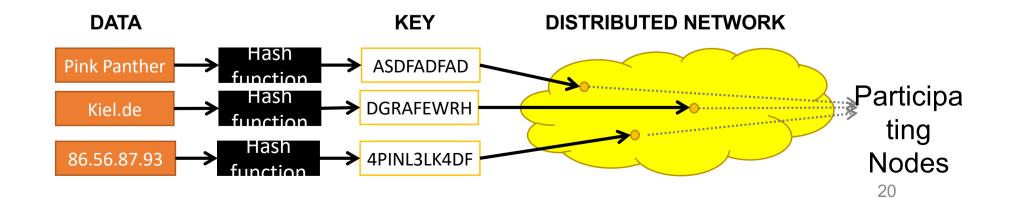


4. Distributed Hash Table (DHT)

- What is a hash table?
 - Efficient storage and lookup system for (key, value) pairs
 - What is a distributed hash table?
- DHT: Like a hash-table, just distributed
- DHT supports two operations:
 - Put (k, v) stores v (the value) according to the key k, the name of the value
 - Get (k) retrieves whatever value is stored associated with key k

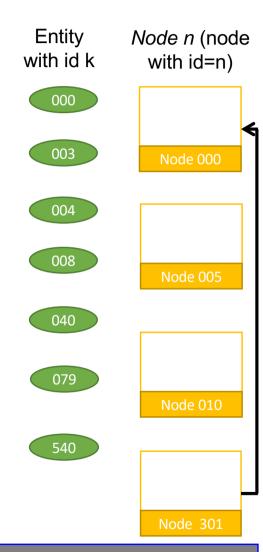
4. Distributed Hash Table (DHT)

- DHT is a class of decentralized distributed system that provides a lookup service similar to a hash table
 - (key, value) pair is stored in the nodes participating in the DHT
 - The responsibility for maintaining the mapping from keys to values is distributed among the nodes
 - Any participating node can retrieve the value for a given key
- We will study a representative DHT known as Chord



Chord

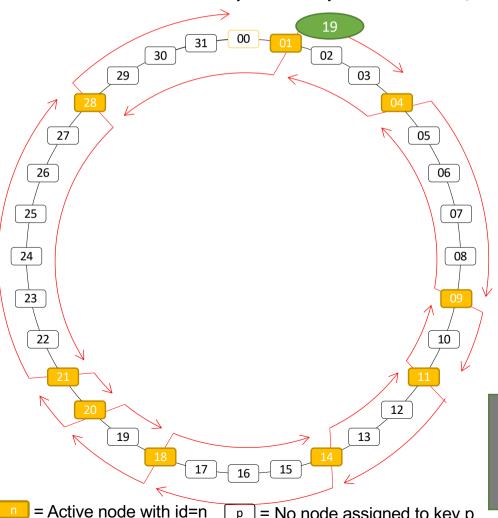
- Chord assigns an m-bit identifier key (randomly chosen) to each node
 - Each node can be contacted through its network address
- Chord also maps each entity to an m-bit identifier key
 - Entities can be processes, files, etc.
- Mapping of entities to nodes
 - Each node is responsible for a set of entities
 - An entity with key k falls under the jurisdiction of the node with smallest identifier id >= k. This node is known as the successor of k, and is denoted by succ(k)



Match each entity with key k with node succ(k)

A Naïve Key Resolution Algorithm

- The main issue in DHT-based solution is to efficiently resolve a key k to the network location of succ(k)
 - Given an entity with key k on node n, how to find the node succ(k)?



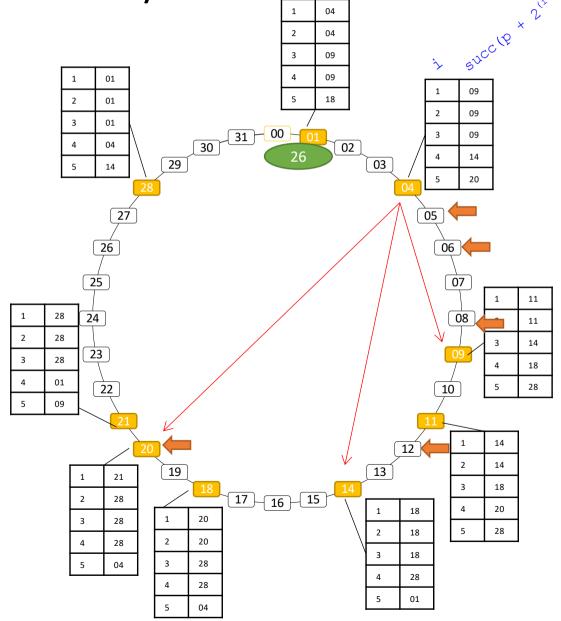
p = No node assigned to key p

- 1. All nodes are arranged in a logical ring according to their keys
- 2. Each node 'p' keeps track of its immediate neighbors: succ(p) and pred(p)
- 3. If node 'n' receives a request to resolve key 'k':
 - **If** pred(p) < k <=p, node will handle it
 - Else it will simply forward it to succ(n) or pred(n)

Problems? Solution is not scalable:

- As the network grows, forwarding delays increase
 - Key resolution has a time complexity of O (n)

Key Resolution in Chard



Chord improves key resolution by reducing the time complexity to O(log n)

- 1. All nodes are arranged in a logical ring according to their keys
- 2. Each node 'p' keeps a table FT_p of atmost m entries. This table is called Finger Table

$$FT_p[i] = succ(p + 2^{(i-1)})$$

NOTE: $FT_p[i]$ increases logarithmically

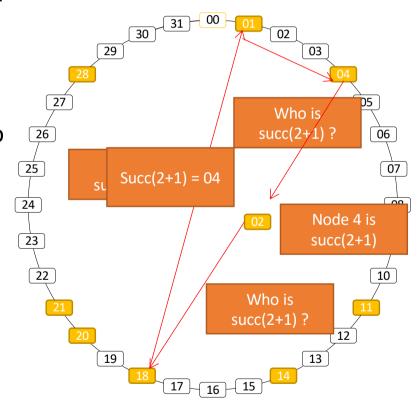
- 3. If node 'n' receives a request to resolve key 'k':
 - Node p will forward it to node q with index j in F_p where

$$q = FT_p[\dot{j}] \le k \le FT_p[\dot{j}+1]$$

- If $k > FT_p[m]$, then node p will forward it to $FT_p[m]$
- 4. Node replies to source

Chord – Join and Leave Protocol

- In large Distributed Systems, nodes dynamically join and leave (voluntarily or due to failure)
- If a node p that wants to join:
 - Node p contacts arbitrary node, looks up for succ (p), and inserts itself into the ring
- If node p wants to leave
 - Node p contacts pred (p), and updates it
- What about the data?
 - One join: split it
 - One leave: copy to pred (p)



Chord – Finger Table Update Protocol

- For any node q, $FT_{\alpha}[1]$ should be up-to-date
 - It refers to the next node in the ring
 - Protocol:
 - Periodically, request succ (q) to return pred (succ (q))
 - If q = pred(succ(q)), then information is up-to-date
 - Otherwise, a new node p has been added to the ring such that q
 - $FT_q[1] = p$
 - Request p to update pred (p) = q
 - Similarly, node q updates each entry i by finding $succ(p + 2^{(i-1)})$

Exploiting Network Proximity in Chord

- The logical organization of nodes in the overlay network may lead to inefficient message transfers in the underlying Internet
 - Node k and node succ(k +1) may be far apart
- Chord can be optimized by considering the network location of nodes
 - 1. Topology-aware Node Assignment
 - Two nearby nodes have identifiers that are close to each other
 - 2. Proximity Routing
 - Each node q maintains 'r' successors for ith entry in the finger table
 - FT_q[i] now refers to successors first r nodes in the range $[p + 2^{(i-1)}, p + 2^{i} 1]$
 - To forward the lookup request, pick one of the r successors closest to the node \triangleleft

Classes of Naming

- Flat naming
- Structured naming
- Attribute-based naming

Structured Naming

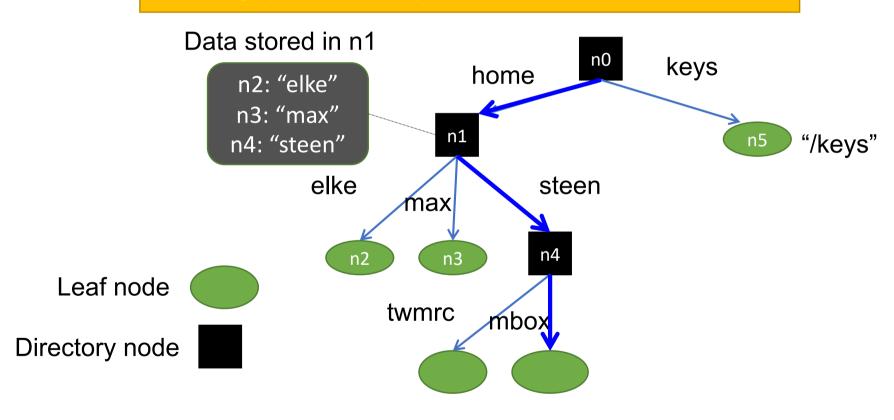
- Structured Names are composed of simple human-readable names
 - Names are arranged in a specific structure
- Examples
 - File-systems utilize structured names to identify files
 - /home/userid/work/dist-systems/naming.txt
 - Websites can be accessed through structured names
 - www.cs.qatar.cmu.edu

Name Spaces

- Structured Names are organized into name spaces
- Name-spaces is a directed graph consisting of:
 - Leaf nodes
 - Each leaf node represents an entity
 - Leaf node generally stores the address of an entity (e.g., in DNS), or the state of an entity (e.g., in file system)
 - Directory nodes
 - Directory node refers to other leaf or directory nodes
 - Each outgoing edge is represented by (edge label, node identifier)
- Each node can store any type of data
 - e.g., type of the entity, address of the entity

Example Name Space

Looking up for the entity with name "/home/steen/mbox"



Name Resolution

 The process of looking up a name is called Name Resolution

- Closure mechanism
 - Name resolution cannot be accomplished without an initial directory node
 - Closure mechanism selects the implicit context from which to start name resolution
 - Examples
 - www.inf.uni-kiel.de: start at the DNS Server
 - /home/steen/mbox: start at the root of the file-system

Name Linking

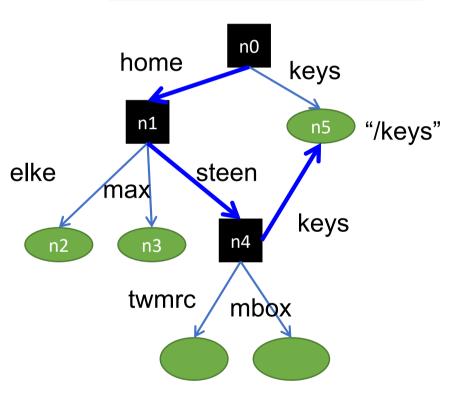
- Name space can be effectively used to link two different entities
- Two types of links can exist between the nodes
 - 1. Hard Links
 - 2. Symbolic Links

1. Hard Links

 There is a directed link from the hard link to the actual node

- Name Resolution
 - Similar to the general name resolution
- Constraint:
 - There should be no cycles in the graph

"/home/steen/keys" is a hard link to "/keys"

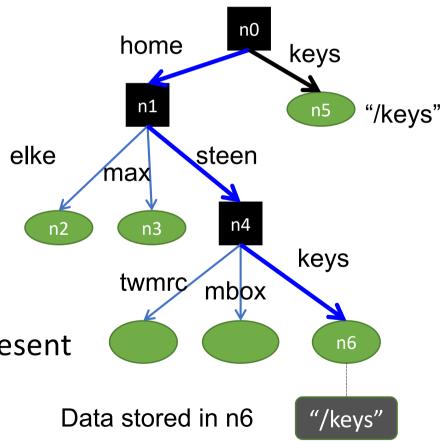


2. Symbolic Links

• Symbolic link stores the name of the original node as data

"/home/steen/keys" is a symbolic link to "/keys"

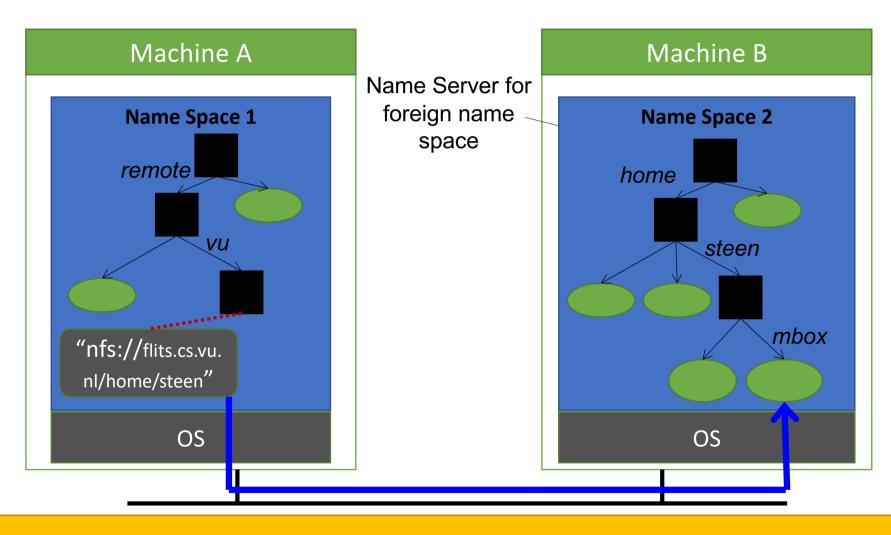
- Name Resolution for a symbolic link SL
 - First resolve SL's name
 - Read the content of SL
 - Name resolution continues with content of SL
- Constraint:
 - No cyclic references should be present



Mounting of Name Spaces

- Two or more name spaces can be merged transparently by a technique known as mounting
- In mounting, a directory node in one name space will store the identifier of the directory node of another name space
- Network File System (NFS) is an example where different name spaces are mounted
 - NFS enables transparent access to remote files

Example of Mounting Name Spaces in NFS



Name resolution for "/remote/vu/home/steen/mbox" in a distributed file system

Distributed Name Spaces

- In large Distributed Systems, it is essential to distribute name spaces over multiple name servers
 - Distribute nodes of the naming graph
 - Distribute name space management
 - Distribute name resolution mechanisms

Layers in Distributed Name Spaces

Distributed Name Spaces can be divided into three layers

Global Layer

- Consists of high-level directory nodes
- Directory nodes are jointly managed by different administrations

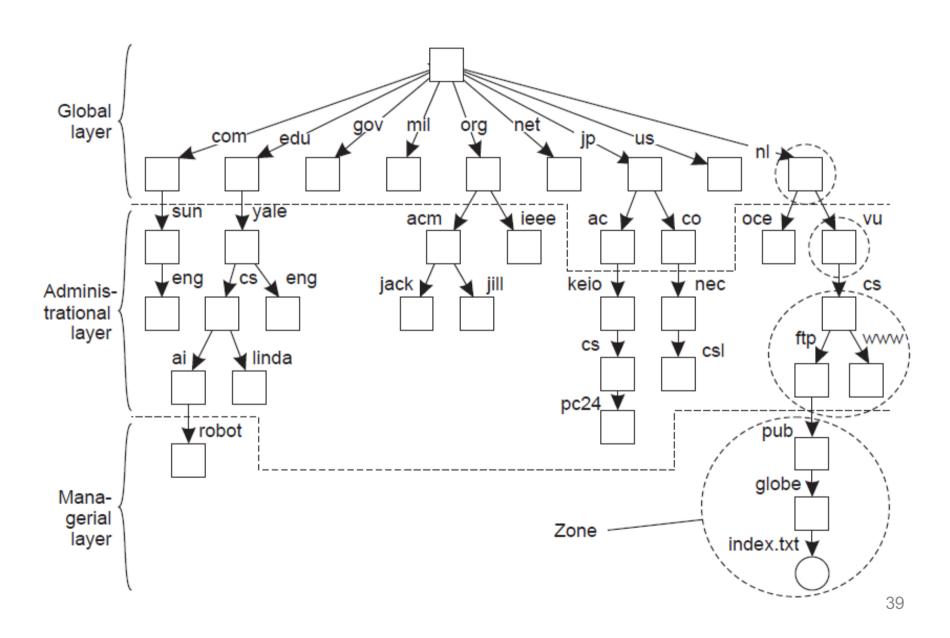
Administrational Layer

- Contains mid-level directory nodes
- Directory nodes grouped together in such a way that each group is managed by an administration

Managerial Layer

- Contains low-level directory nodes within a single administration
- The main issue is to efficiently map directory nodes to local name servers

Distributed Name Spaces – An Example



Comparison of Name Servers at Different Layers

	Global	Administrational	Managerial
Geographical scale of the network	Worldwide	Organization	Department
Total number of nodes	Few	Many	Vast numbers
Number of replicas	Many	None or few	None
Update propagation	Lazy	Immediate	Immediate
Is client side caching applied?	Yes	Yes	Sometimes
Responsiveness to lookups	Seconds	Milliseconds	Immediate

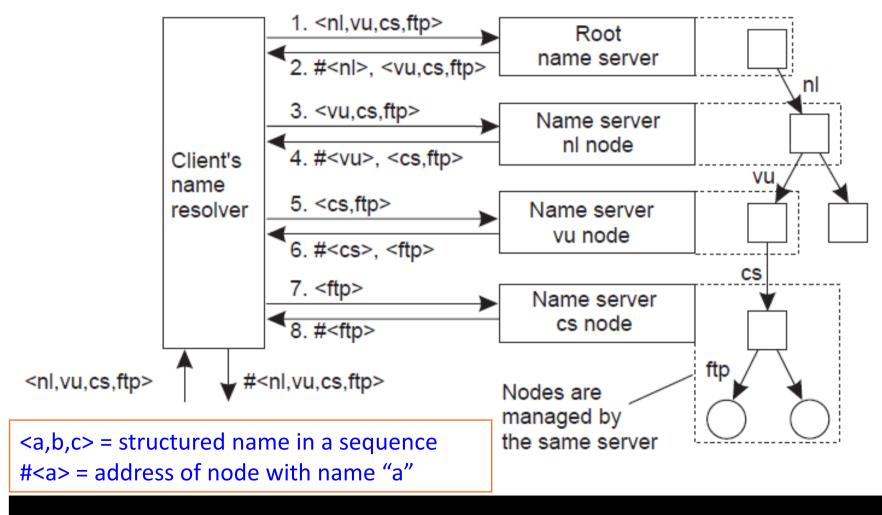
Distributed Name Resolution

- Distributed Name Resolution is responsible for mapping names to address in a system where:
 - Name servers are distributed among participating nodes
 - Each name server has a local name resolver
- We will study two distributed name resolution algorithms:
 - 1. Iterative Name Resolution
 - 2. Recursive Name Resolution

1. Iterative Name Resolution

- 1. Client hands over the complete name to root name server
- Root name server resolves the name as far as it can, and returns the result to the client
 - The root name server returns the address of the next-level name server (say, NLNS) if address is not completely resolved
- 3. Client passes the unresolved part of the name to the NLNS
- 4. NLNS resolves the name as far as it can, and returns the result to the client (and probably its next-level name server)
- 5. The process continues till the full name is resolved

1. Iterative Name Resolution — An Example



2. Recursive Name Resolution

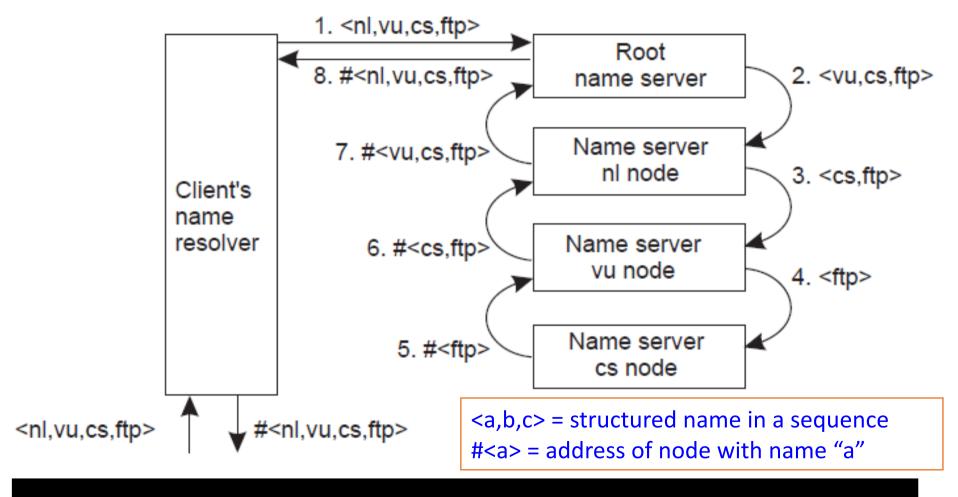
Approach

- Client provides the name to the root name server
- The root name server passes the result to the next name server it finds
- The process continues till the name is fully resolved

• Drawback:

 Large overhead at name servers (especially, at the highlevel name servers)

2. Recursive Name Resolution – An Example



Classes of Naming

- Flat naming
- Structured naming
- Attribute-based naming

Attribute-based Naming

- In many cases, it is much more convenient to name, and look up entities by means of their attributes
 - Similar to traditional directory services (e.g., yellow pages)
- However, the lookup operations can be extremely expensive
 - They require to match requested attribute values, against actual attribute values, which needs to inspect all entities
- Solution: Implement basic directory service as database, and combine with traditional structured naming system
- We will study Light-weight Directory Access Protocol (LDAP); an example system that uses attribute-based naming

Light-weight Directory Access Protocol (LDAP)

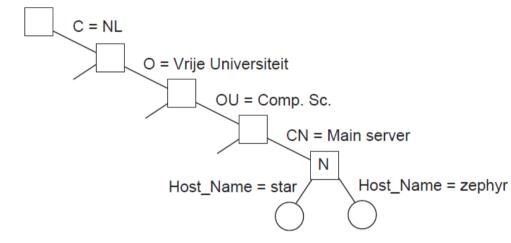
- LDAP Directory Service consists of a number of records called "directory entries"
 - Each record is made of (attribute, value) pair
 - LDAP Standard specifies five attributes for each record
- Directory Information Base (DIB) is a collection of all directory entries
 - Each record in a DIB is unique
 - Each record is represented by a distinguished name

```
e.g., /C=NL/O=Vrije Universiteit/OU=Comp. Sc.
```

Attribute	Value
Country	NL
Locality	Amsterdam
Organization	Vrije Universiteit
OrganizationalUnit	Comp. Sc.
CommonName	Main server
Host_Name	star
Host_Address	192.31.231.42

Directory Information Tree in LDAP

- All the records in the DIB can be organized into a hierarchical trecalled Directory Information Tree (DIT
- LDAP provides advanced search mechanisms based or attributes by traversing the DIT
- Example syntax for searching all Main_Servers in Vrije Universiteit:
- search("&(C = NL) ((= Vrije Universiteit) (OU = *) (CN = Main server)")



Attribute	Value
Country	NL
Locality	Amsterdam
Organization	Vrije Universiteit
OrganizationalUnit	Comp. Sc.
CommonName	Main server
Host_Name	star
Host_Address	192.31.231.42

Attribute	Value
Country	NL
Locality	Amsterdam
Organization	Vrije Universiteit
OrganizationalUnit	Comp. Sc.
CommonName	Main server
Host_Name	zephyr
Host_Address	137.37.20.10

Summary

- Naming and name resolutions enable accessing entities in a Distributed System
- Three types of naming
 - Flat Naming
 - Home-based approaches, Distributed Hash Table
 - Structured Naming
 - Organizes names into Name Spaces
 - Distributed Name Spaces
 - Attribute-based Naming
 - Entities are looked up using their attributes

Next Class

• Time and Clocks

Questions?

In part, inspired from / based on slides from Majd F. Sakr, Vinay Kolar, Mohammad Hammoud and many others