

Chapter 04

(version 21st September 2001)

Maarten van Steen

Vrije Universiteit Amsterdam, Faculty of Science

Dept. Mathematics and Computer Science

Room R4.20. Tel: (020) 444 7784

E-mail: steen@cs.vu.nl, URL: www.cs.vu.nl/~steen

01	Introduction
02	Communication
03	Processes
04	Naming
05	Synchronization
06	Consistency and Replication
07	Fault Tolerance
08	Security
09	Distributed Object-Based Systems
10	Distributed File Systems
11	Distributed Document-Based Systems
12	Distributed Coordination-Based Systems

00 – 1

/

Naming Entities

- Names, identifiers, and addresses
- Name resolution
- Name space implementation

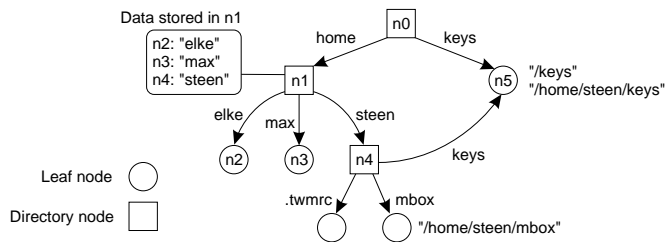
This image shows a single sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

[illegible]

This image shows a single sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

Name Space (1/2)

Essence: a graph in which a **leaf node** represents a (named) entity. A **directory node** is an entity that refers to other nodes.



Note: A directory node contains a (directory) table of (edge label, node identifier) pairs.

Name Space (2/2)

Observation: We can easily store all kinds of **attributes** in a node, describing aspects of the entity the node represents:

- Type of the entity
- An identifier for that entity
- Address of the entity's location
- Nicknames
- ...

Observation: Directory nodes can also have attributes, besides just storing a directory table with (edge label, node identifier) pairs.

Name Resolution

Problem: To resolve a name we need a directory node. How do we actually find that (initial) node?

Closure mechanism: The mechanism to select the implicit context from which to start name resolution:

- `www.cs.vu.nl`: start at a DNS name server
- `/home/steen/mbox`: start at the local NFS file server (possible recursive search)
- `0031204447784`: dial a phone number
- `130.37.24.8`: route to the VU's Web server

Question: Why are closure mechanisms always *implicit*?

Observation: A closure mechanism may also determine how name resolution should proceed

Name Linking (1/2)

Hard link: What we have described so far as a **path name**: a name that is resolved by following a specific path in a naming graph from one node to another.

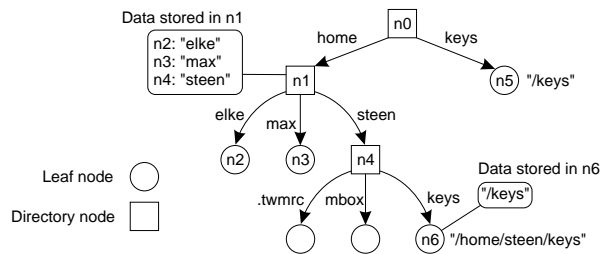
Soft link: Allow a node *O* to contain a *name* of another node:

- First resolve *O*'s name (leading to *O*)
- Read the content of *O*, yielding name
- Name resolution continues with name

Observations:

- The name resolution process determines that we read the *content* of a node, in particular, the name in the other node that we need to go to.
- One way or the other, we know where and how to start name resolution given name

Name Linking (2/2)



Observation: Node **n5** has only one name

Merging Name Spaces (1/3)

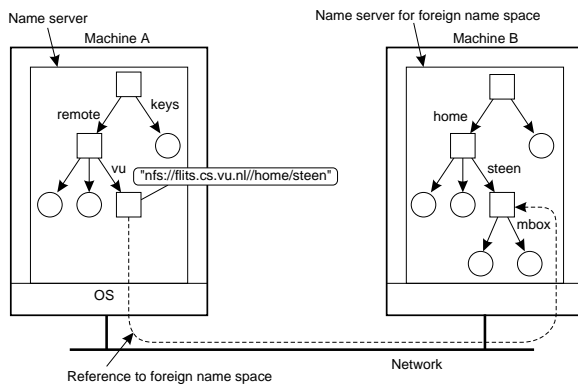
Problem: We have different name spaces that we wish to access from any given name space.

Solution 1: Introduce a naming scheme by which pathnames of different name spaces are simply concatenated (URLs).

ftp://ftp.cs.vu.nl/pub/steen/	
ftp	Name of protocol used to talk with server
://	Name space delimiter
ftp.cs.vu.nl	Name of a node representing an FTP server, and containing the IP address of that server
/	Name space delimiter
pub/steen/	Name of a (context) node in the name space rooted at the context node mapped to the FTP server

Merging Name Spaces (2/3)

Solution 2: Introduce nodes that contain the name of a node in a “foreign” name space, along with the information how to select the initial context in that foreign name space (Jade).

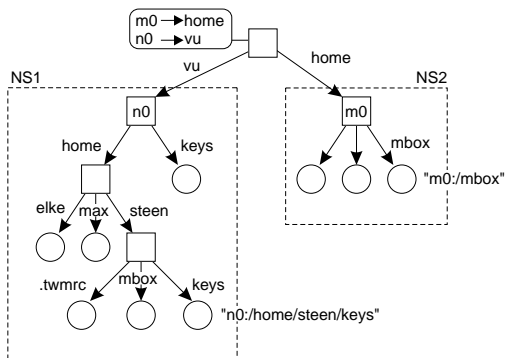


Mount point: (Directory) node in naming graph that refers to other naming graph

Mounting point: (Directory) node in other naming graph that is referred to.

Merging Name Spaces (3/3)

Solution 3: Use only *full pathnames*, in which the starting context is explicitly identified, and merge by adding a new root node (DEC's Global Name Space).



Note: In principle, you *a/ways* have to start in the new root

Name Space Implementation (1/2)

Basic issue: Distribute the name resolution process as well as name space management across multiple machines, by distributing nodes of the naming graph.

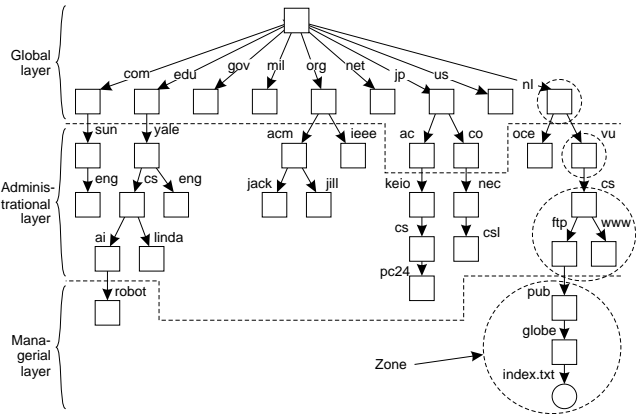
Consider a hierarchical naming graph and distinguish three levels:

Global level: Consists of the high-level directory nodes.
Main aspect is that these directory nodes have to be jointly managed by different administrations

Administrational level: Contains mid-level directory nodes that can be grouped in such a way that each group can be assigned to a separate administration.

Managerial level: Consists of low-level directory nodes within a single administration. Main issue is effectively mapping directory nodes to local name servers.

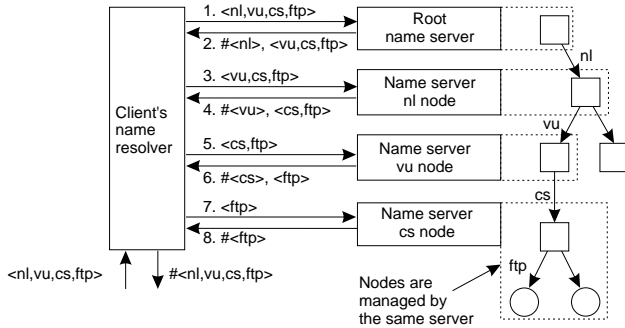
Name Space Implementation (2/2)



Item	Global	Administrational	Managerial
1	Worldwide	Organization	Department
2	Few	Many	Vast numbers
3	Seconds	Milliseconds	Immediate
4	Lazy	Immediate	Immediate
5	Many	None or few	None
6	Yes	Yes	Sometimes
1: Geographical scale		4: Update propagation	
2: # Nodes		5: # Replicas	
3: Responsiveness		6: Client-side caching?	

Iterative Name Resolution

- $\text{resolve}(\text{dir}, [\text{name1}, \dots, \text{nameK}])$ is sent to Server0 responsible for dir
- Server0 resolves $\text{resolve}(\text{dir}, \text{name1}) \rightarrow \text{dir1}$, returning the identification (address) of Server1, which stores dir1.
- Client sends $\text{resolve}(\text{dir1}, [\text{name2}, \dots, \text{nameK}])$ to Server1
- etc.

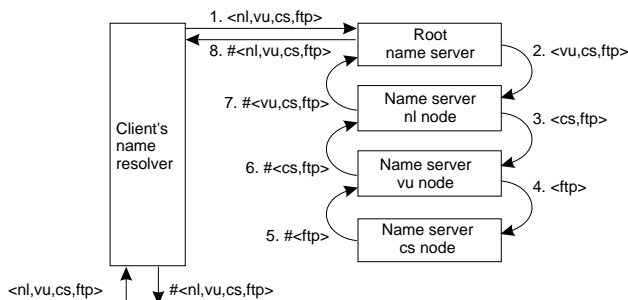


04 – 14

Naming/4.1 Naming Entities

Recursive Name Resolution

- $\text{resolve}(\text{dir}, [\text{name1}, \dots, \text{nameK}])$ is sent to Server0 responsible for dir
- Server0 resolves $\text{resolve}(\text{dir}, \text{name1}) \rightarrow \text{dir1}$, and sends $\text{resolve}(\text{dir1}, [\text{name2}, \dots, \text{nameK}])$ to Server1, which stores dir1.
- Server0 waits for the result from Server1, and returns it to the client.



04 – 15

Naming/4.1 Naming Entities

This image shows a single sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

This image shows a single sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

This image shows a single sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

This image shows a single sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

This image shows a single sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

This image shows a single sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

[illegible][illegible]

1000

This image shows a single sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

- 1000

Naming & Locating Objects (1/2)

Location service: Solely aimed at providing the addresses of the *current* locations of entities.

Assumption: Entities are mobile, so that their current address may change frequently.

Naming service: Aimed at providing the content of nodes in a name space, given a (compound) name. Content consists of different (attribute,value) pairs.

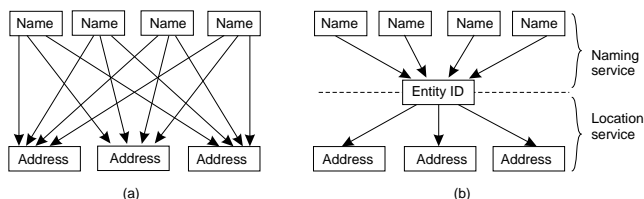
Assumption: Node contents at global and administrative level is relatively stable for scalability reasons.

Observation: If a traditional naming service is used to locate entities, we also have to assume that node contents at the managerial level is stable, as we can use only names as identifiers (think of Web pages).

Naming & Locating Objects (2/2)

Problem: It is not realistic to assume stable node contents down to the local naming level

Solution: Decouple naming from locating entities



Name: Any name in a traditional naming space

Entity ID: A true identifier

Address: Provides *all* information necessary to contact an entity

Observation: An entity's name is now completely independent from its location.

Question: What may be a typical address?

Simple Solutions for Locating Entities

Broadcasting: Simply broadcast the ID, requesting the entity to return its current address.

- Can never scale beyond local-area networks (think of ARP/RARP)
- Requires all processes to listen to incoming location requests

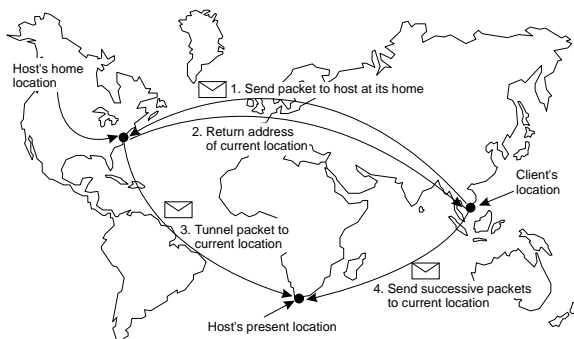
Forwarding pointers: Each time an entity moves, it leaves behind a pointer telling where it has gone to.

- Dereferencing can be made entirely transparent to clients by simply following the chain of pointers
 - Update a client's reference as soon as present location has been found
 - Geographical scalability problems:
 - Long chains are not fault tolerant
 - Increased network latency at dereferencing
- Essential to have separate chain reduction mechanisms

Home-Based Approaches (1/2)

Single-tiered scheme: Let a **home** keep track of where the entity is:

- An entity's **home address** is registered at a naming service
- The home registers the **foreign address** of the entity
- Clients always contact the home first, and then continues with the foreign location



Home-Based Approaches (2/2)

Two-tiered scheme: Keep track of **visiting** entities:

- Check local visitor register first
- Fall back to home location if local lookup fails

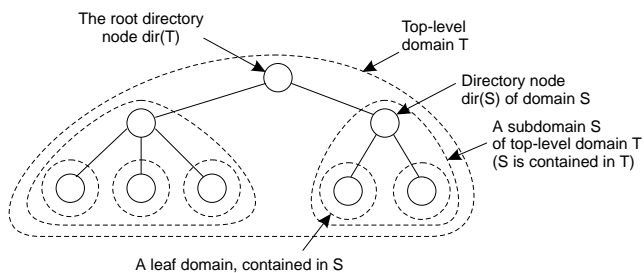
Problems with home-based approaches:

- The home address has to be supported as long as the entity lives.
- The home address is fixed, which means an unnecessary burden when the entity permanently moves to another location
- Poor geographical scalability (the entity may be next to the client)

Question: How can we solve the “permanent move” problem?

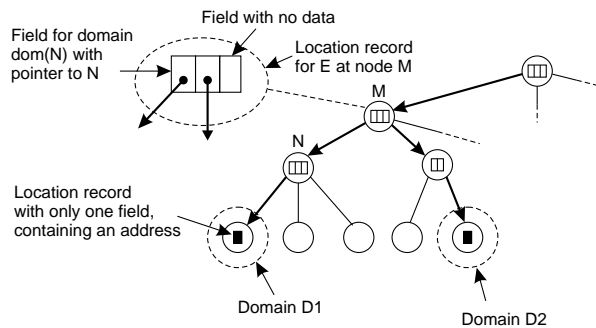
Hierarchical Location Services (HLS)

Basic idea: Build a large-scale search tree for which the underlying network is divided into hierarchical domains. Each domain is represented by a separate directory node.



HLS: Tree Organization

- The address of an entity is stored in a leaf node, or in an intermediate node
- Intermediate nodes contain a pointer to a child if and only if the subtree rooted at the child stores an address of the entity
- The root knows about all entities



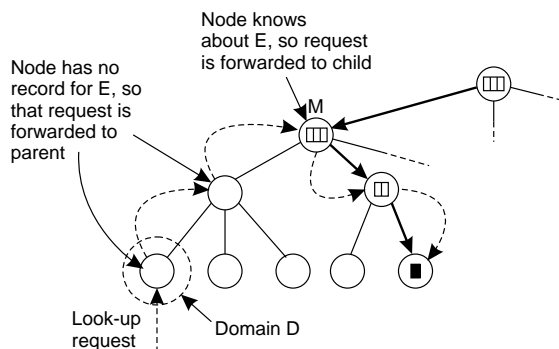
04 – 26

Naming/4.2 Locating Mobile Entities

HLS: Lookup Operation

Basic principles:

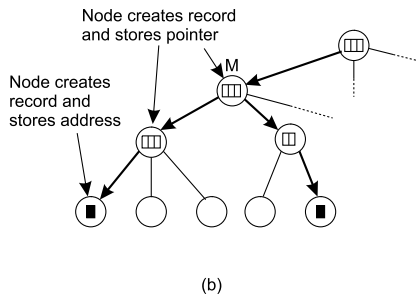
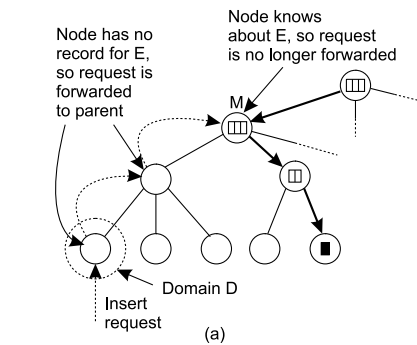
- Start lookup at local leaf node
- If node knows about the entity, follow downward pointer, otherwise go one level up
- Upward lookup always stops at root



04 – 27

Naming/4.2 Locating Mobile Entities

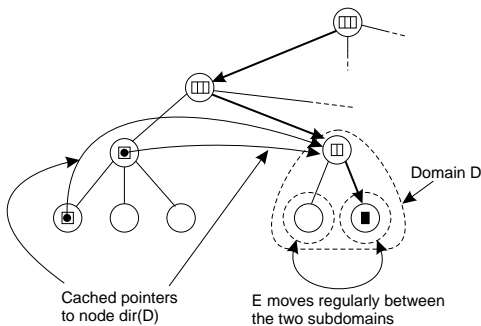
HLS: Insert Operation



HLS: Record Placement

Observation: If an entity E moves regularly between leaf domains D_1 and D_2 , it may be more efficient to store E 's contact record at the least common ancestor LCA of $\text{dir}(D_1)$ and $\text{dir}(D_2)$.

- Lookup operations from either D_1 or D_2 are on average cheaper
- Update operations (i.e., changing the current address) can be done directly at LCA
- Note: assuming that E generally stays in $\text{dom}(\text{LCA})$, it does make sense to cache a **pointer** to LCA



HLS: Scalability Issues

Size scalability: Again, we have a problem of overloading higher-level nodes:

- Only solution is to partition a node into a number of subnodes and evenly assign entities to subnodes
- Naive partitioning may introduce a node management problem, as a subnode may have to know how its parent and children are partitioned.

Geographical scalability: We have to ensure that lookup operations generally proceed monotonically in the direction of where we'll find an address:

- If entity E generally resides in California, we should not let a root subnode located in France store E 's contact record.
- Unfortunately, subnode placement is not that easy, and only a few tentative solutions are known

Reclaiming References

- Reference counting
- Reference listing
- Scalability issues

Unreferenced Objects: Problem

Assumption: Objects may exist only if it is known that they can be contacted:

- Each object should be named
- Each object can be located
- A reference can be resolved to client–object communication

Problem: Removing unreferenced objects:

- How do we know when an object is no longer referenced (think of cyclic references)?
- Who is responsible for (deciding on) removing an object?

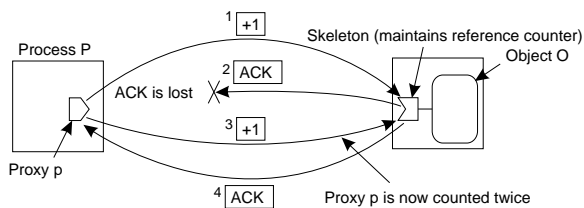
Reference Counting (1/2)

Principle: Each time a client creates (removes) a reference to an object O , a reference counter local to O is incremented (decremented)

Problem 1: Dealing with lost (and duplicated) messages:

- An increment is lost so that the object may be prematurely removed
- A decrement is lost so that the object is never removed
- An ACK is lost, so that the increment/decrement is resent.

Solution: Keep track of duplicate requests.



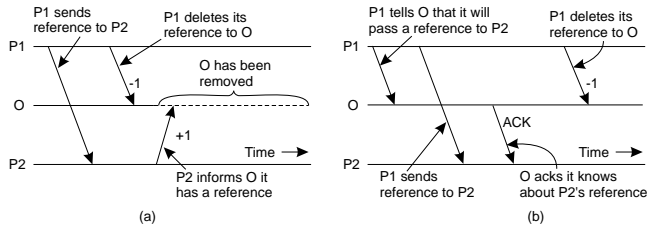
Reference Counting (2/2)

Problem 2: Dealing with duplicated references – client P_1 tells client P_2 about object O :

- Client P_2 creates a reference to O , but dereferencing (communicating with O) may take a long time
- If the last reference known to O is removed before P_2 talks to O , the object is removed prematurely

Solution 1: Ensure that P_2 talks to O on time:

- Let P_1 tell O it will pass a reference to P_2
- Let O contact P_2 immediately
- A reference may never be removed before O has acked that reference to the holder



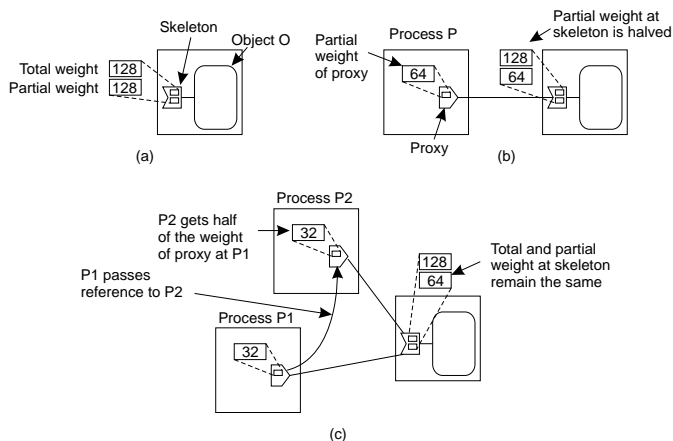
04 – 34

Naming/4.3 Reclaiming References

Weighted Reference Counting

Solution 2: Avoid increment and decrement messages:

- Let O allow a maximum M of references
- Client P_1 creates reference \Rightarrow grant it $M/2$ credit
- Client P_1 tells P_2 about O , it passes half of its credit grant to P_2
- Pass *current* credit grant back to O upon reference deletion



04 – 35

Naming/4.3 Reclaiming References

This image shows a single sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

[illegible]

This image shows a single sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.