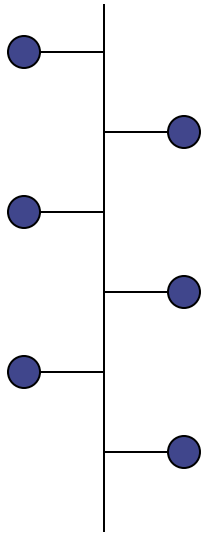# TechGI IV - ComDiS

**ComDis**

Introduction to
**Com**munication Networks
and **Dis**tributed Systems

*DHTs, Naming Summary*

Prof. Dr.-Ing. Adam Wolisz

TKN **Telecommunication Networks Group**

# Reminder:

- **We have discussed the mapping**

  *Name* ➔ *Address*
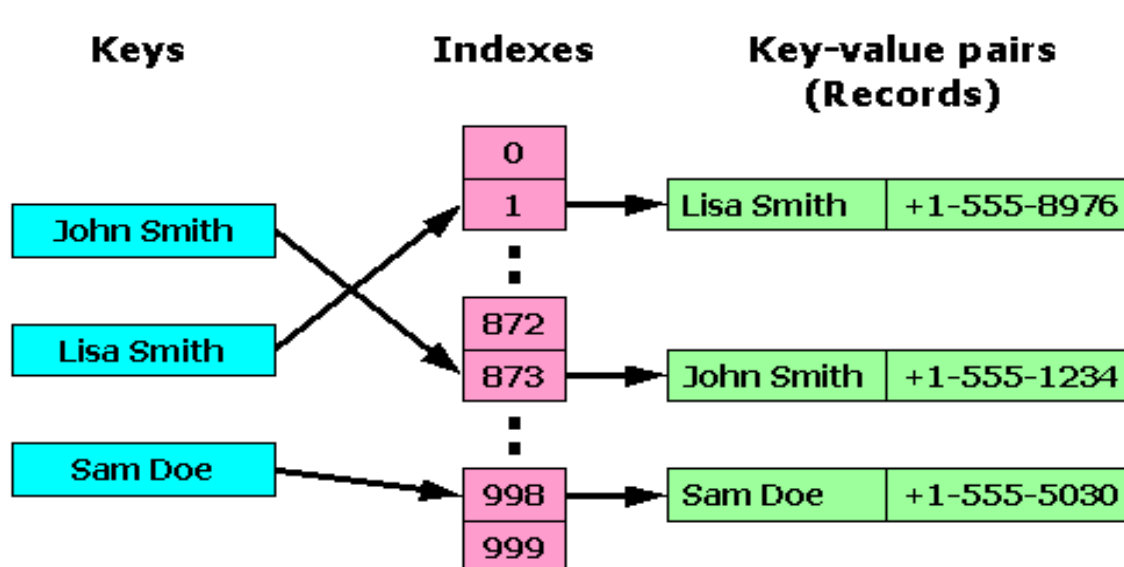
- **The solution has been DNS**
  - **Hierarchicall organization**
  - **Distributed**
  - **Using redundancy**

- **Challenge:**

  **How to provide storage of pairs  - like the above one – in a distributed way, even if there is no strong hierarchy?**

# Hash Tables (a classic)

- Items: [*Key*, *Value*] are stored
- The key is *hashed,* i.e. transformed (using a *hash function*) so that the result – **the hash** – can be used to locate a **bucket** in which the pair is stored. The bucket is identified by an **index.**
- .

| Keys | Indexes | Key-value pairs (Records) |
|------|---------|---------------------------|
| John Smith | 0 | |
| | 1 | Lisa Smith +1-555-8976 |
| Lisa Smith | 872 | |
| | 873 | John Smith +1-555-1234 |
| Sam Doe | 998 | Sam Doe +1-555-5030 |
| | 999 | |

In this example the index is simply the number of the record.

- The bucket might contain multiple such items (pairs)!

# A Distributed Hash Table (DHT)

- Remember the mapping of NAMES to IP Addresses? Could we use Hash tables? Remember the scaling issue…

- Distributed Hash Tables spread the pairs across a number of computers (buckets) located arbitrarily across the world.
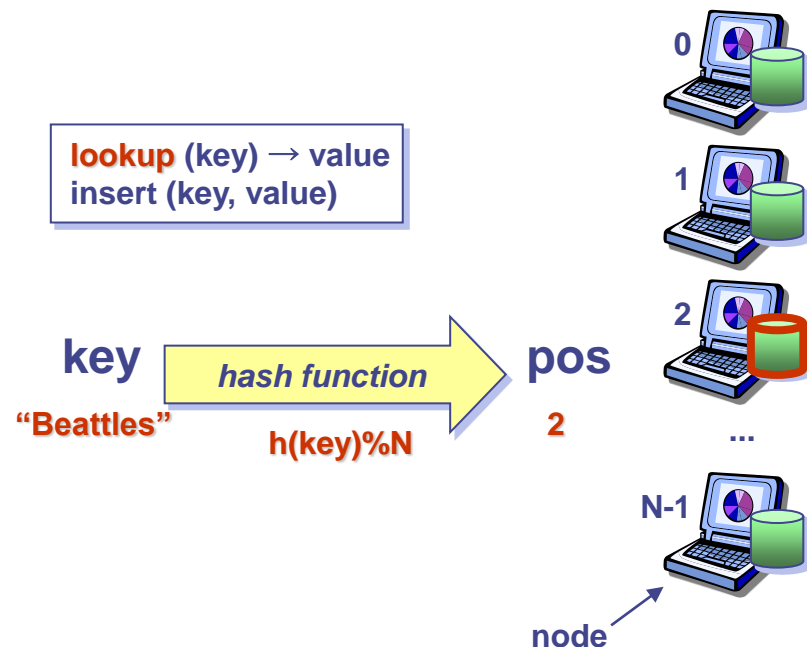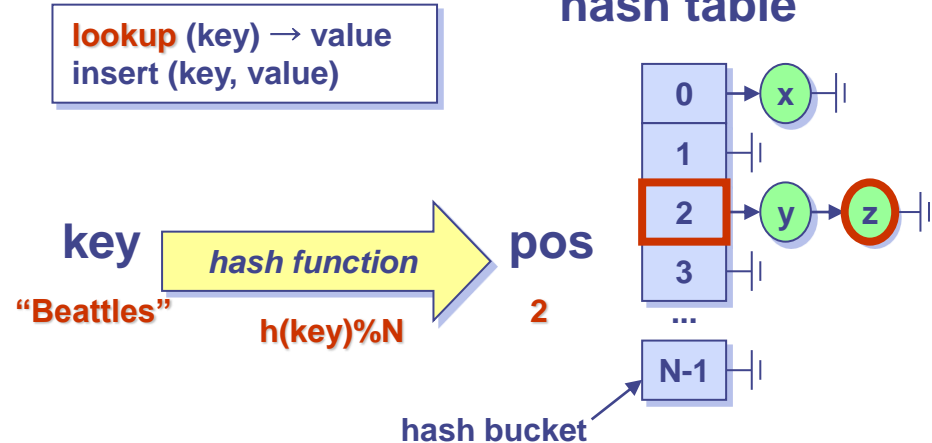
  *Note: Copies of a single pair can be stored in one or in multiple locations!*

- When a user queries the system, i.e. provides the key, the system uses the hash to find the pair from one of the computers where it's stored and returns the result.

- All the nodes are assumed to be reachable by some kind of unicast communication.

- DHT posses the features of : scaling, robustness, self-organization.

# The hash table vs. DHT

[Ala Khalifeh, UCI]

- The key is hashed to find the proper bucket in a hash table

**hash table**

lookup (key) → value
insert (key, value)

| 0 | → x |
| 1 | |
| 2 | → y → z |
| 3 | |
| ... | |
| N-1 | |

**key**  *hash function*  **pos**

"Beattles"  h(key)%N  2

hash bucket

- In a Distributed Hash Table (DHT), nodes are the hash buckets

  – Key is hashed to find the responsible Node

  – Pairs are distributed among the nodes with respect to load balancing

lookup (key) → value
insert (key, value)

**key**  *hash function*  **pos**

"Beattles"  h(key)%N  2

0
1
2
...
N-1

node

# DHT Interface

- Minimal interface (data-centric)

  Lookup(key) $\rightarrow$ value

  Insert(key, value)

  Delete (key)

- Supports a wide range of applications, because few restrictions

  - Value is application dependent

  - Keys have no semantic meaning

***Note**: DHTs do not have to store data useful to end users, e.g. data files…*
*Data storage can be build on top of DHTs*

# DHTs: Problems

- **Problem 1 (dynamicity):** adding or removing nodes
  - With hash mod N, virtually every key will change its location!

$$h(k) \bmod m \neq h(k) \bmod (m+1) \neq h(k) \bmod (m\text{-}1)$$

- **Solution:** use consistent hashing
  - Define a fixed hash space
  - All hash values fall within that space and do not depend on the number of peers (hash bucket)
  - Each key goes to peer closest to its ID in hash space (according to some proximity metric)

---

- **Problem 2 (size):** all nodes must be known ( in order to insert or lookup items!)
  - Works with *small* and *static* server populations

- **Solution:** each peer knows of only a few "neighbors"
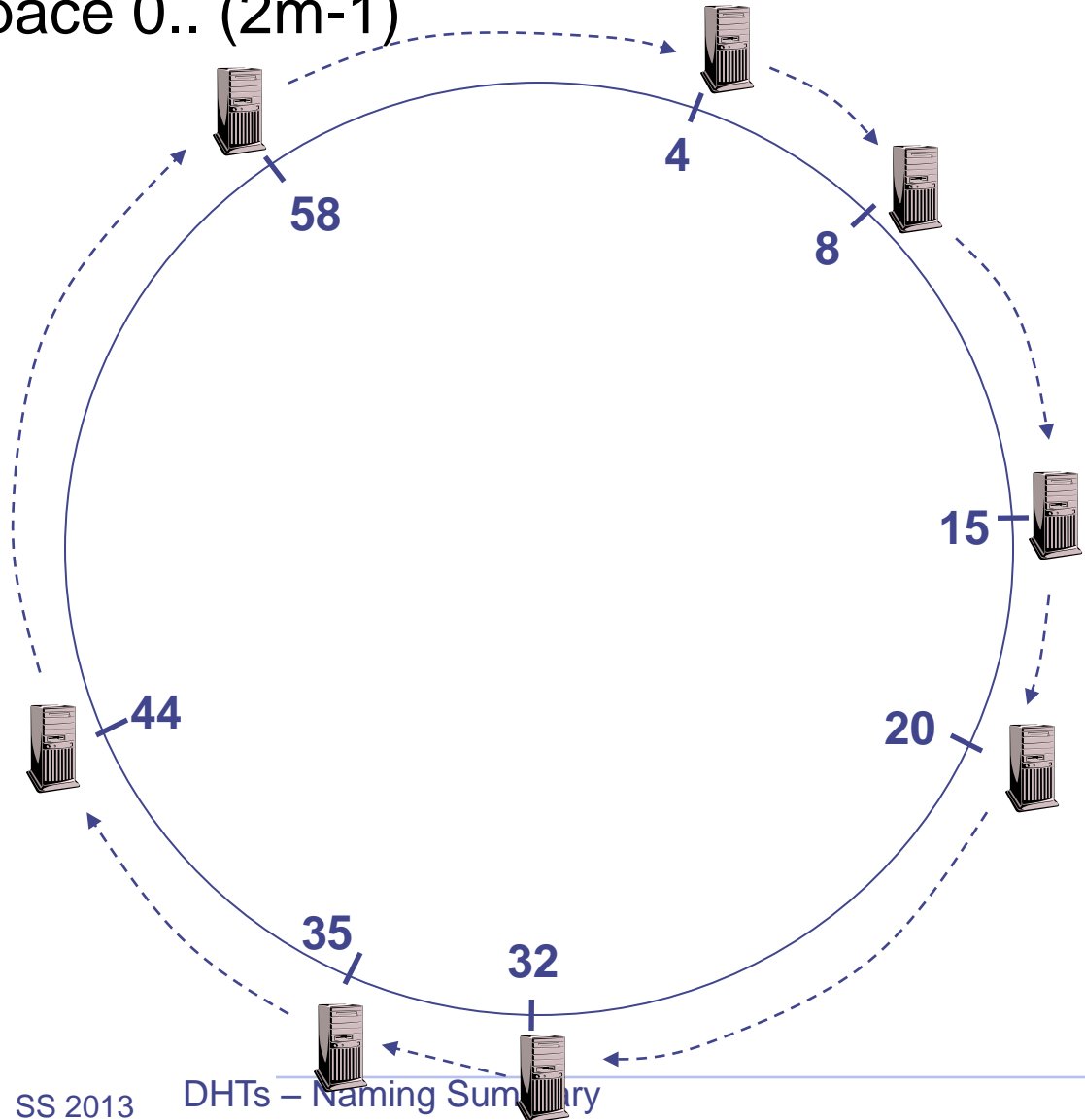  - Messages are routed through neighbors via multiple hops

# What Makes a Good DHT Design

- The number of neighbors for each node should remain "reasonable" (small degree)

- DHT routing mechanisms should be decentralized (no single point of failure or bottleneck)

- Should gracefully handle nodes joining and leaving
  - Repartition the affected keys over existing nodes
  - Reorganize the neighbor sets
  - Bootstrap mechanisms to connect new nodes into the DHT

- DHT must provide low stretch
  - Minimize ratio of DHT routing vs. unicast latency between two nodes

# Identifier to Node Mapping Example [S.Shenker and I.Stoica, UCB]

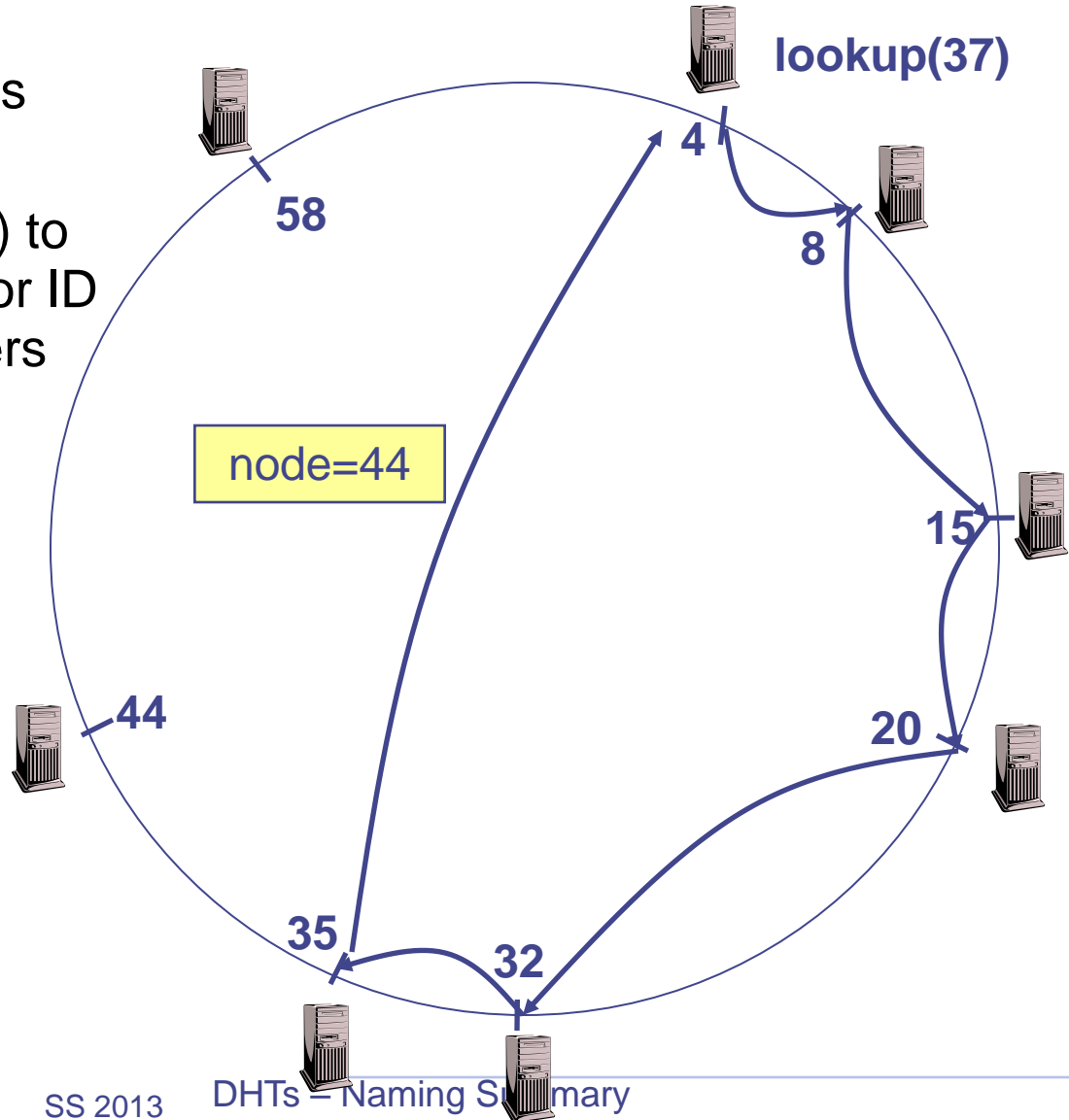• Associate to each node and item a unique *id* in an *uni*-dimensional space 0.. (2m-1)

- Node  8 maps [5,8]

- Node 15 maps [9,15]

- Node 20 maps [16, 20]

- …

- Node 4 maps [59, 4]

- Each node maintains a pointer to its successor

# Chord Lookup

- Each node maintains its successor

- Route packet (ID, data) to the node responsible for ID using successor pointers

**lookup(37)**

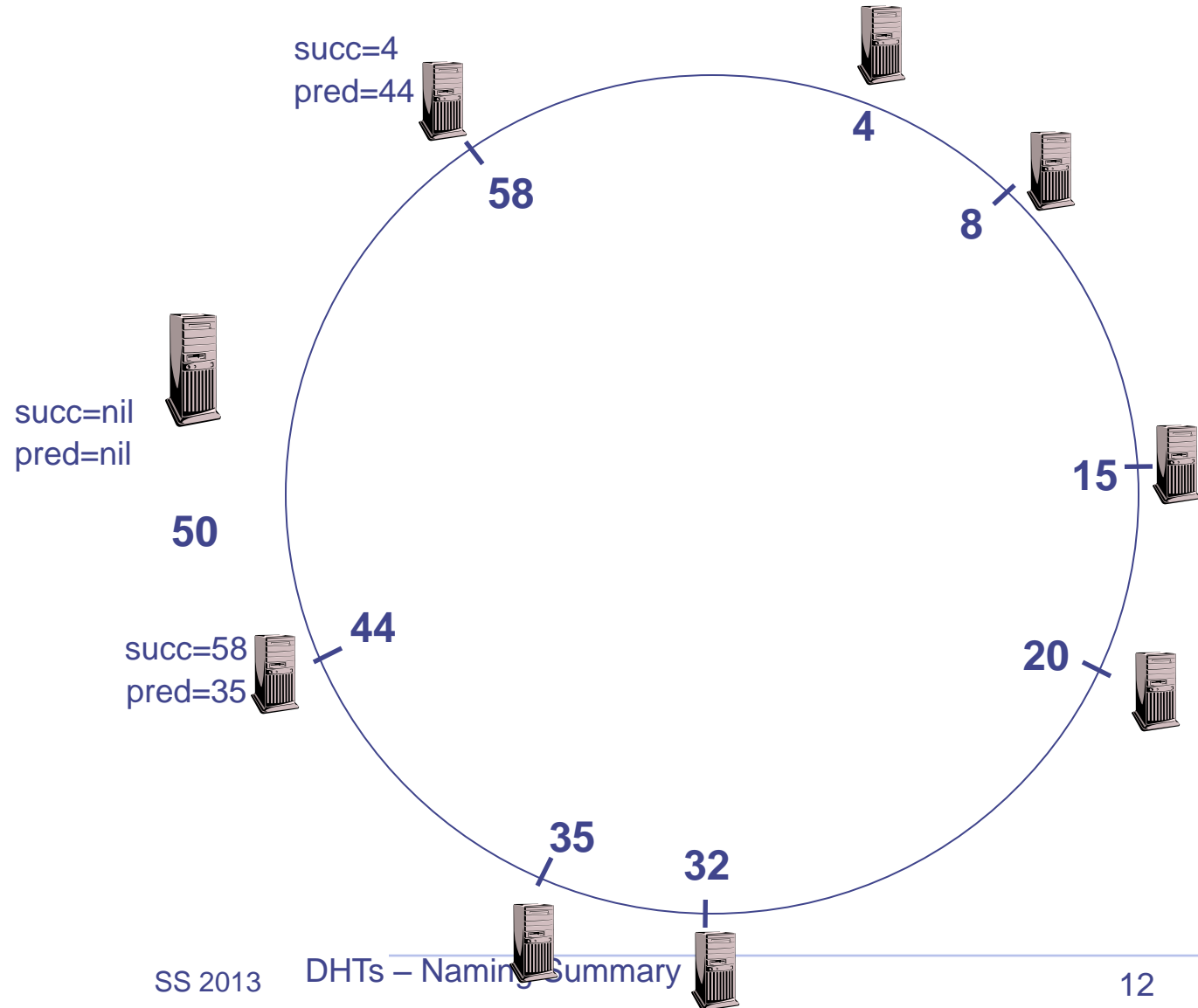node=44

4

58

8

15

44

20

35

32

# Chord Joining Operation

- Each node A periodically sends a stabilize() message to its successor B

- Upon receiving a stabilize() message node B
  - returns its predecessor B'=pred(B) to A by sending a notify(B') message

- Upon receiving notify(B') from B,
  - if B' is between A and B, A updates its successor to B'
  - A doesn't do anything, otherwise

- Node with id=50 joins the ring

- Node 50 needs to know at least one node already in the system

  – Assume known node is 15

**succ=4**
**pred=44**
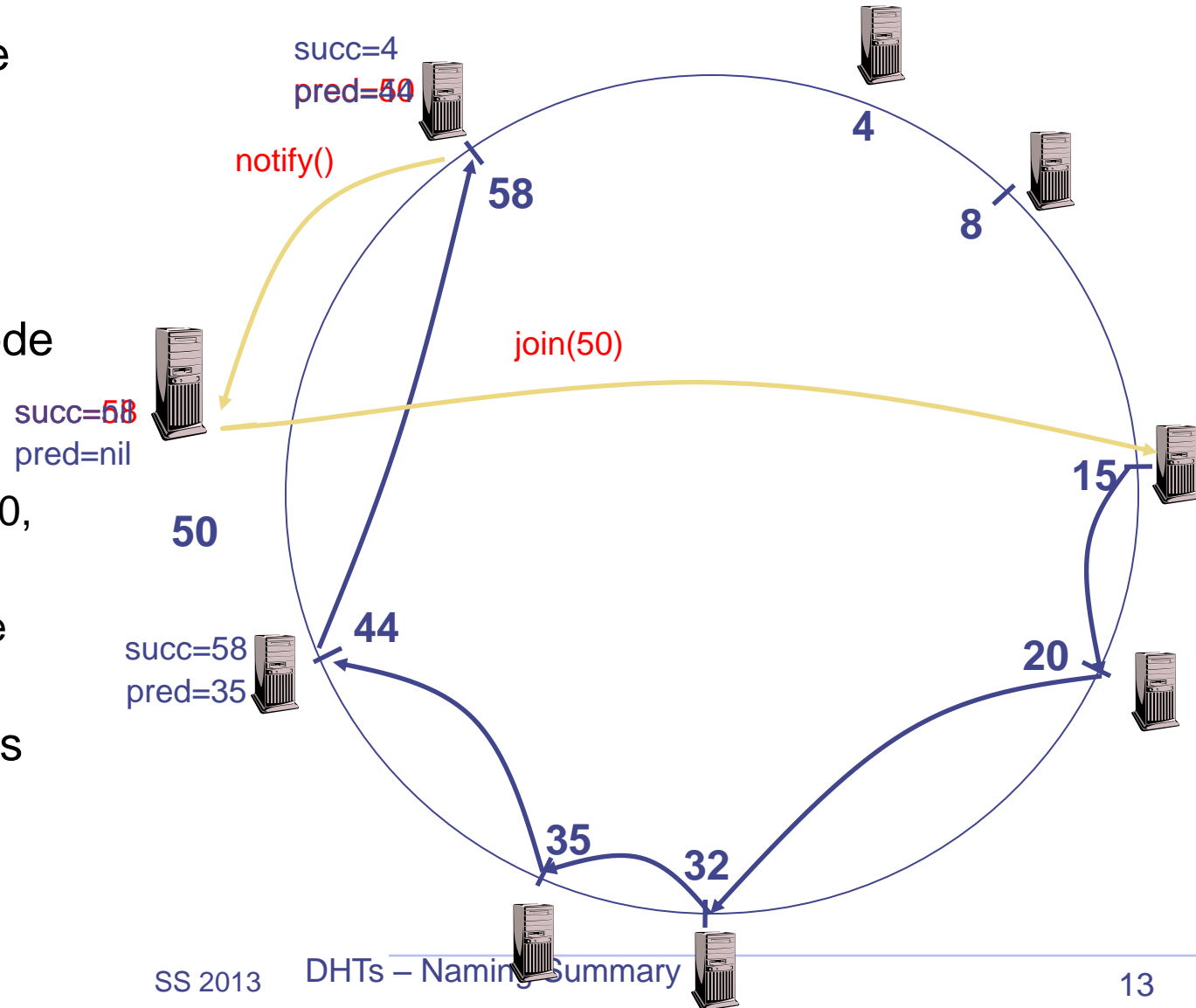
**4**

**58**

**8**

**succ=nil**
**pred=nil**

**15**

**50**

**succ=58**
**pred=35**

**44**

**20**

**35**

**32**

# Joining Operation

- Node 50 asks node 15 to forward join message

- When join(50) reaches the destination (i.e., node 58), node 58

  1) updates its predecessor to 50,

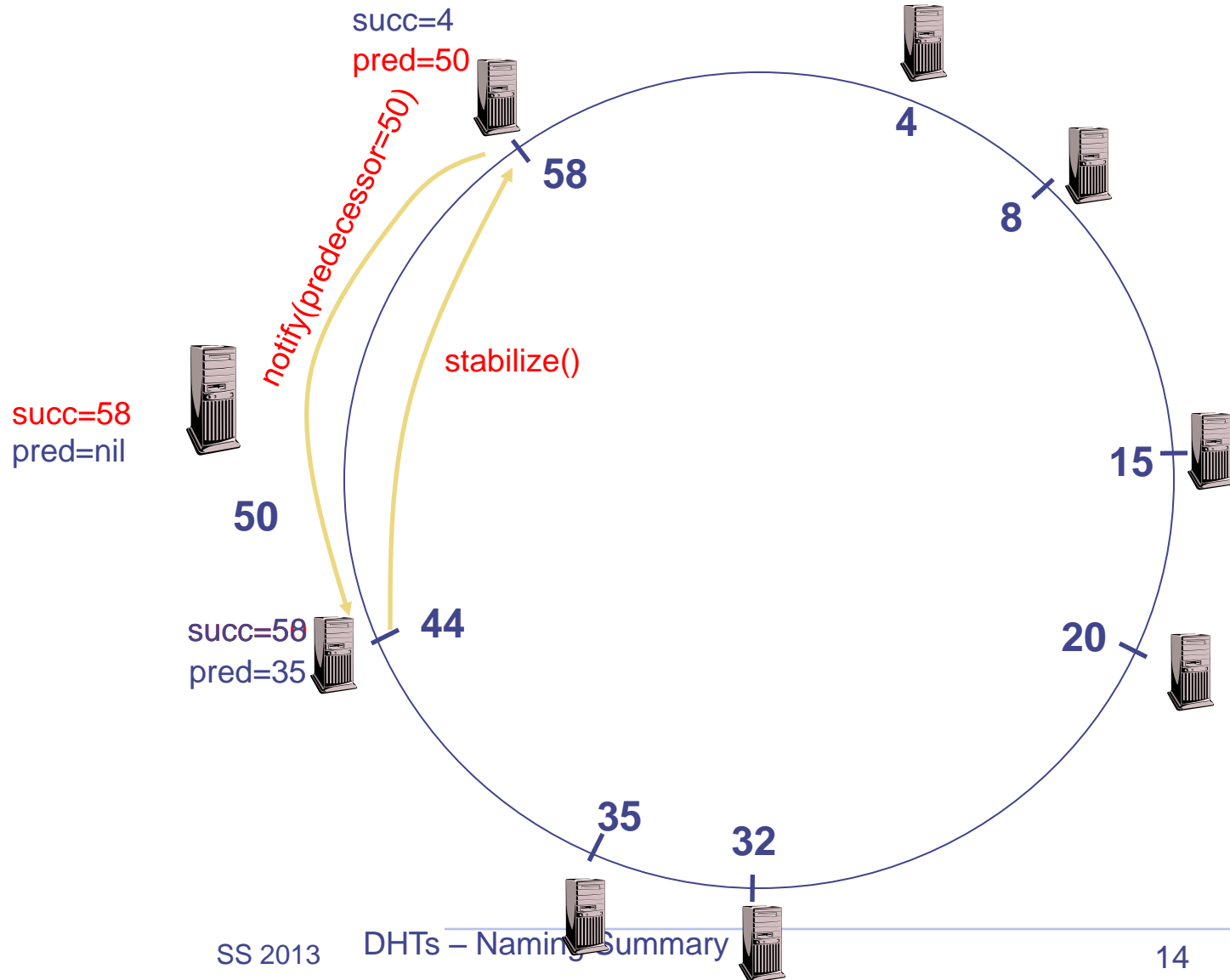  2) returns a notify message to node 50

- Node 50 updates its successor to 58

succ=4
pred=50

**4**

notify()

**58**

**8**

join(50)

succ=58
pred=nil

**15**

**50**

**44**

**20**

succ=58
pred=35

**35**

**32**

Node 44 sends a stabilize message to its successor, node 58

Node 58 reply with a notify message

Node 44 updates its successor to 50

succ=4
pred=50

**58**

**4**

**8**

notify(predecessor=50)

stabilize()

**15**

succ=58
pred=nil

**50**

succ=58
pred=35

**44**

**20**

**35**

**32**

# Joining Operation (cont'd)

- Node 44 sends a stabilize message to its new successor, node 50

- Node 50 sets its predecessor to node 44

succ=4
pred=50

**4**

**58**

**8**

succ=58
~~pred=44~~ pred=44

**50**

Stabilize()

**15**

**44**

**20**

succ=50
pred=35

**35**

**32**

# Joining Operation (cont'd)

- This completes the joining operation!



pred=50

**4**

**58**

**8**

succ=58

pred=44

**50**

**15**

**44**

**20**

succ=50

**35**

**32**

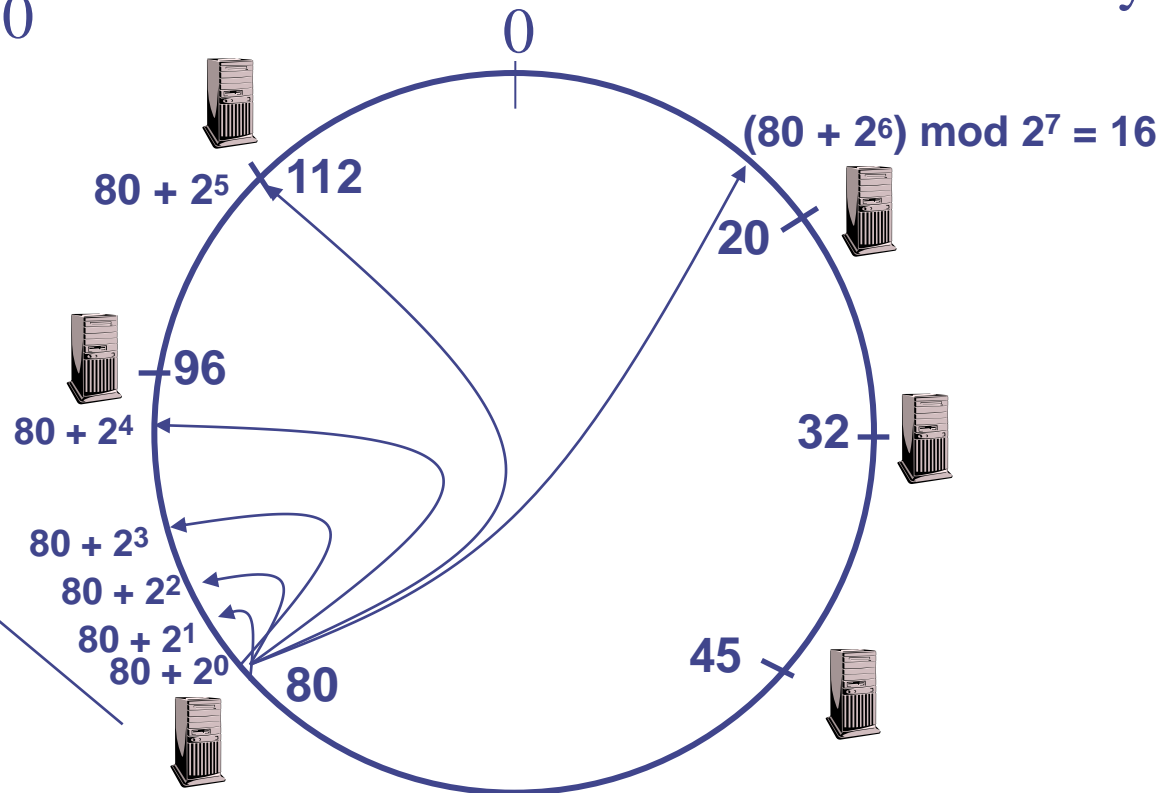# Achieving Efficiency: *finger tables*

Say *m=7*

Finger Table at 80

| $i$ | $ft[i]$ |
|---|---|
| 0 | 96 |
| 1 | 96 |
| 2 | 96 |
| 3 | 96 |
| 4 | 96 |
| 5 | 112 |
| 6 | 20 |

$(80 + 2^6) \bmod 2^7 = 16$

0

80 + 2^5  112

20

96

80 + 2^4

32

80 + 2^3

80 + 2^2

80 + 2^1

80 + 2^0   80

45

$i$th entry at peer with id $n$ is first peer with id $>= \quad n + 2^{i} \,(\bmod \quad 2^{m})$

# Chord  Performance (improvements)

- ## Chord Properties
  - Routing table size $O(\log(N))$ , where $N$ is the total number of nodes
  - Guarantees that a file is found in $O(\log(N))$ steps

- ## Reducing latency
  - Chose finger that reduces expected time to reach destination
  - Chose the closest node from range $[N+2^{i-1}, N+2^i)$ as successor

- ## Stretch is another parameter....

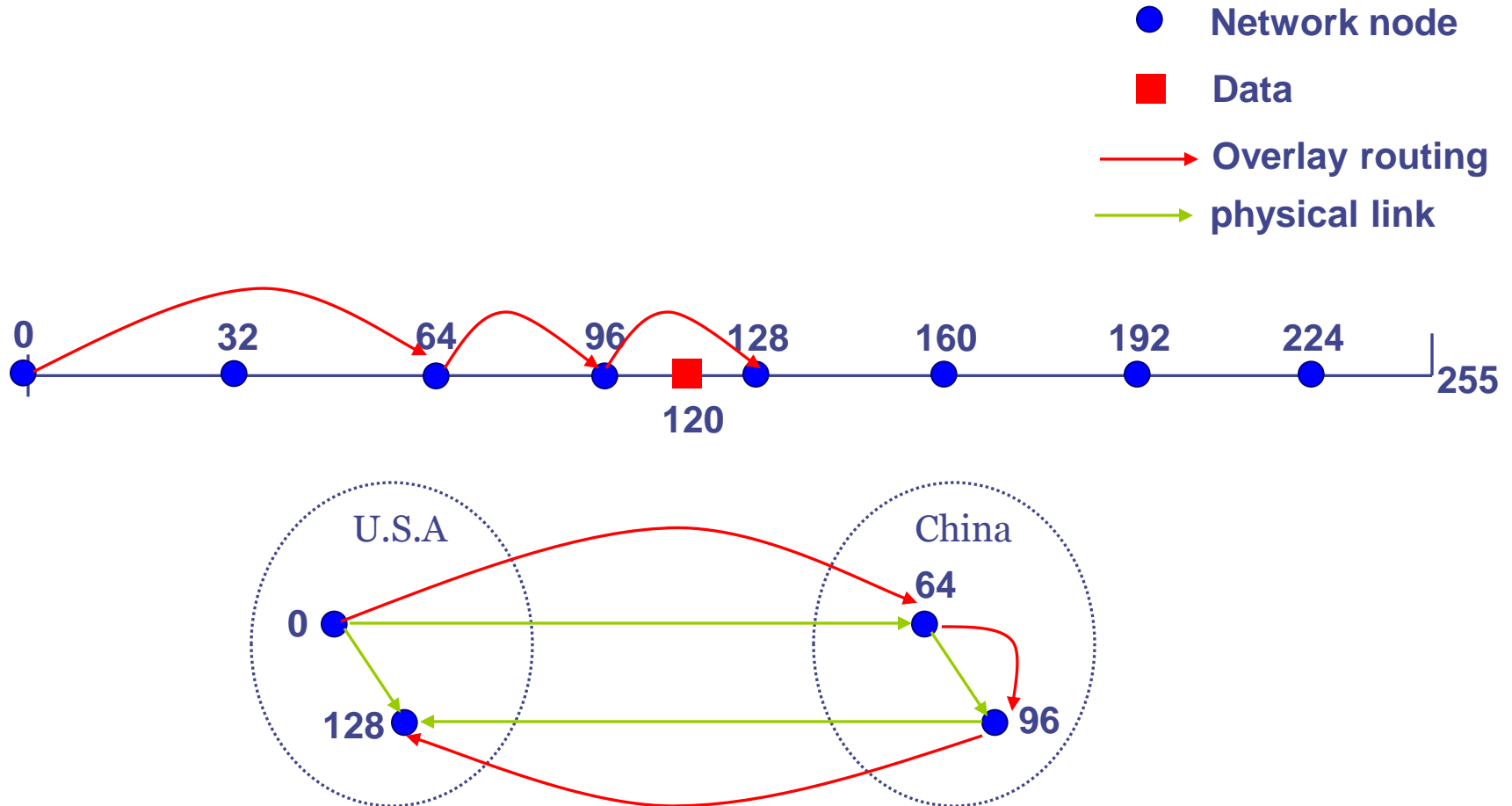- $$= \frac{\text{latency for each lookup on the overlay topology}}{\text{average latency on the underlying topology}}$$

  - Nodes <u>close</u> on ring, but <u>far away</u> in Internet
  - Goal: put nodes in routing table that result in few hops <u>and</u> low latency

# Latency stretch in Chord

[Ratnasamy et al. 2001]

**Network node**

**Data**

**Overlay routing**

**physical link**



A Chord network with N(=8) nodes and m(=8)-bit key space

# Achieving Robustness

- To improve robustness each node maintains the k (> 1) immediate successors instead of only one successor

- In the notify() message, node A can send its k-1 successors to its predecessor B

- Upon receiving notify() message, B can update its successor list by concatenating the successor list received from A with A itself

# The Problem of Membership Churn [Sean C. Rhea, Intel Res]

- In a system with 1,000s of machines, some machines failing / recovering at all times

- This process is called *churn*

- Without repair, quality of overlay network degrades over time

- A significant problem deployed DHTs systems

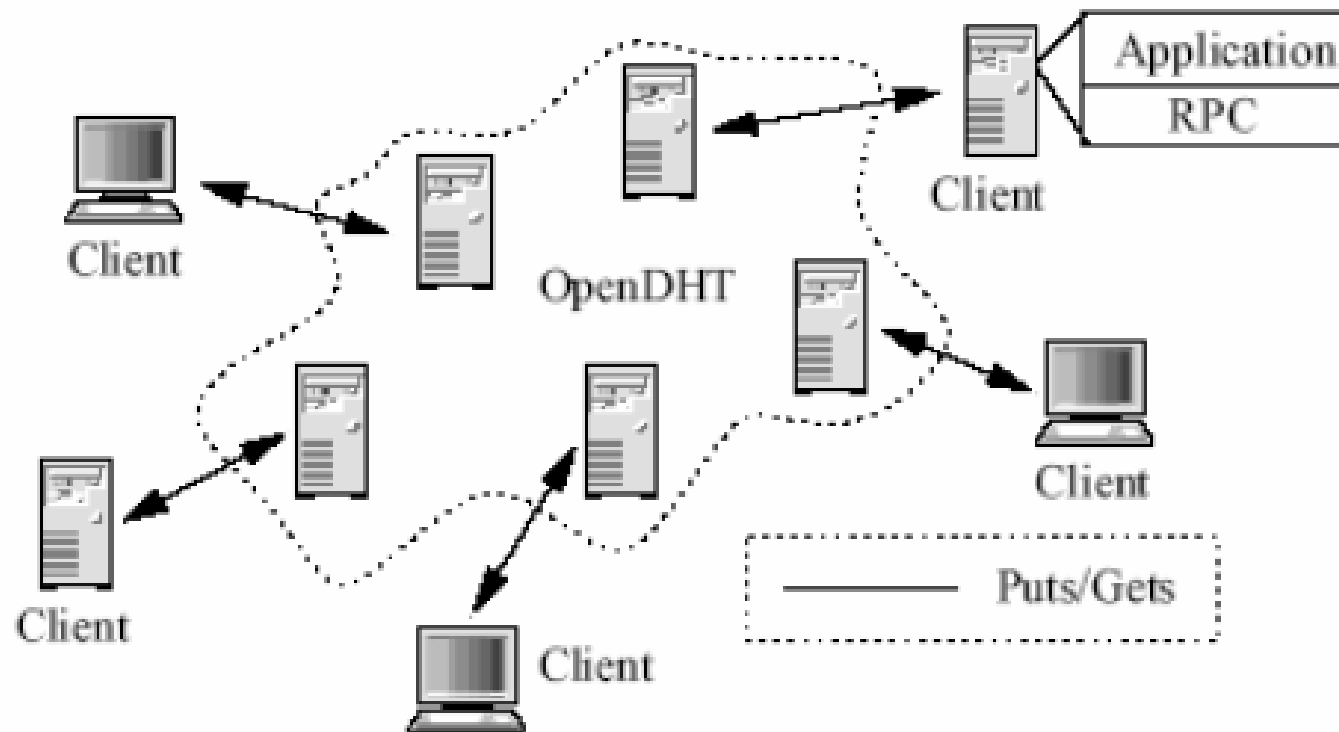*Observation: in >50 % cases, MTBF in order of minutes..*

# Multiple solutions…

- Chord

- Tapestry

- Pastry

- CAN

- ….

# DHTs Support Many Applications [Felber, Eurecom]

- File sharing [CFS, OceanStore, PAST, …]

- Web cache [Squirrel, …]

- Censor-resistant stores [Eternity, FreeNet, …]

- Application-layer multicast [Narada, …]

- Event notification [Scribe]

- Naming systems [ChordDNS, INS, …]

- Query and indexing [Kademlia, …]

- Communication primitives  [I3, …]

- Backup store [HiveNet]

- Web archive [Herodotus]

# DHT is a good shared infrastructure

- A single DHT (namely, Open DHT) is shared across multiple applications, thus amortizing the cost of deployment.

- Applications inherit <u>some</u> security and robustness from DHT
  - DHT replicates data
  - Resistant to malicious participants

- Low-cost deployment
  - Self-organizing across administrative domains
  - Allows to be shared among applications

# DHT as an Infrastructure



**Open DHT Architecture**

# Open DHT Deployment Model

- A single DHT (namely, Open DHT) is shared across multiple applications, thus amortizing the cost of deployment.

- Each DHT node serves as a *gateway* into the DHT for clients.

- Any Internet-connected computer can act as client:
  - Clients of Open DHT do not need to run a DHT node
  - using DHT services, i.e. can store or *put* key-value pairs in Open DHT, and can retrieve or *get* the value stored under a particular key

- Each DHT node serves as a *gateway* into the DHT for clients.

- An Open DHT client communicates with the DHT through the gateway of its choice using an RPC over TCP. **The gateway** processes the operations on client´s behalf.

- Because of this, the service is easy to access from virtually every programming language

# Naming – summary…

- We have discussed the Duality:
  - Name (object id) vs: address
  - Naming: finding the address out of the name…
  - **Solutions?**

- We have already discussed DNS
  - Structuring of the name SPACE necessary!

- We have just discussed DHTs:
  - A solution for a FLAT NAME SPACE!

- Are there any other options possible?
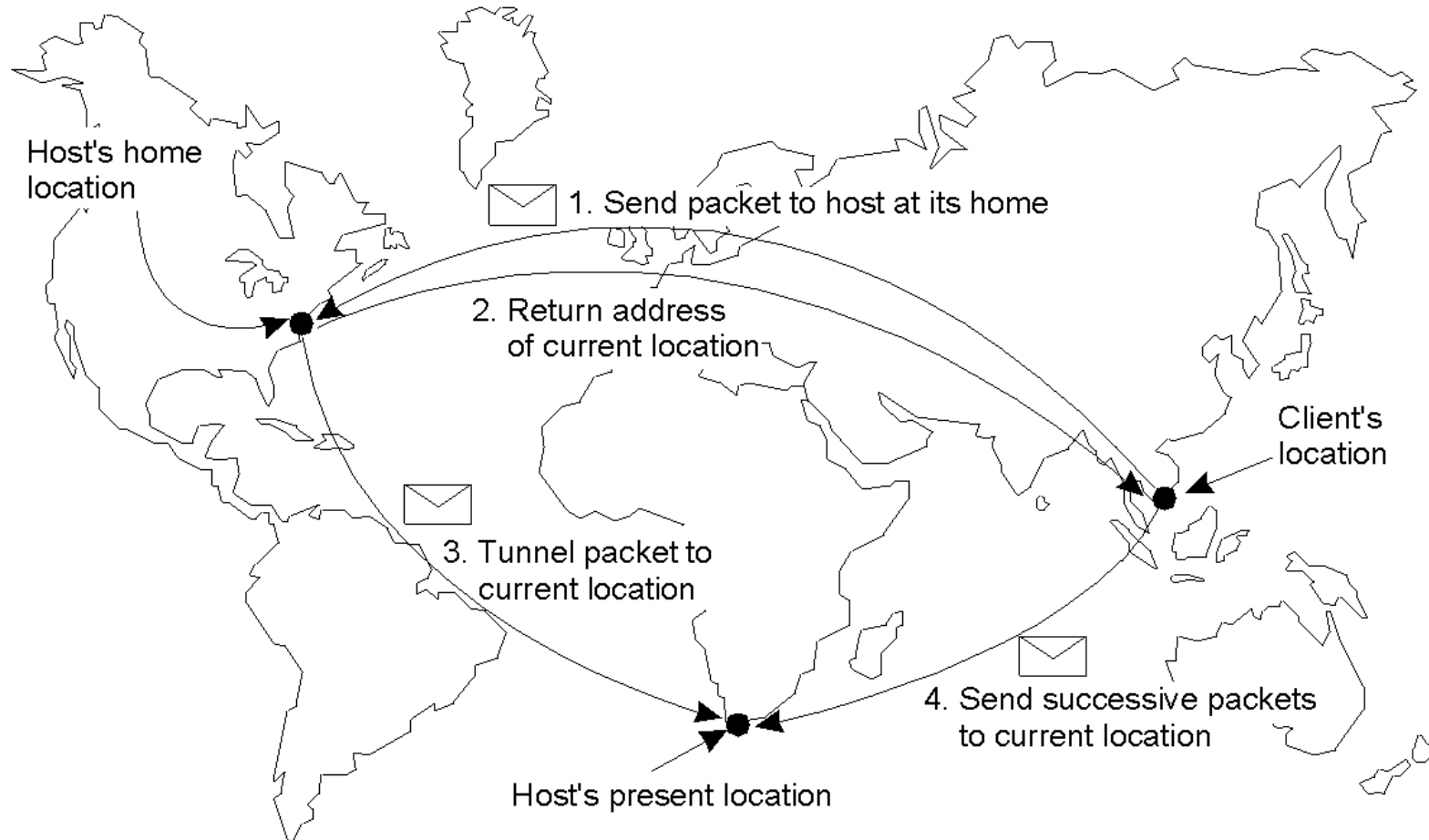
# Using a broadcast…

- To "owner of the name  XXX"
        "please post me your address"


- Works only for RELATIVELY SMALL CONFIGURATIONS!


- Has "side effects"
  - What if the owner/administrator  od the name XXX disappears ?
  - A third party can pretend to have the Name XXX
    (delegation of the responsibility!)


- *Look ahead: compare the ARP Service later….*

# Using the „Home Address"
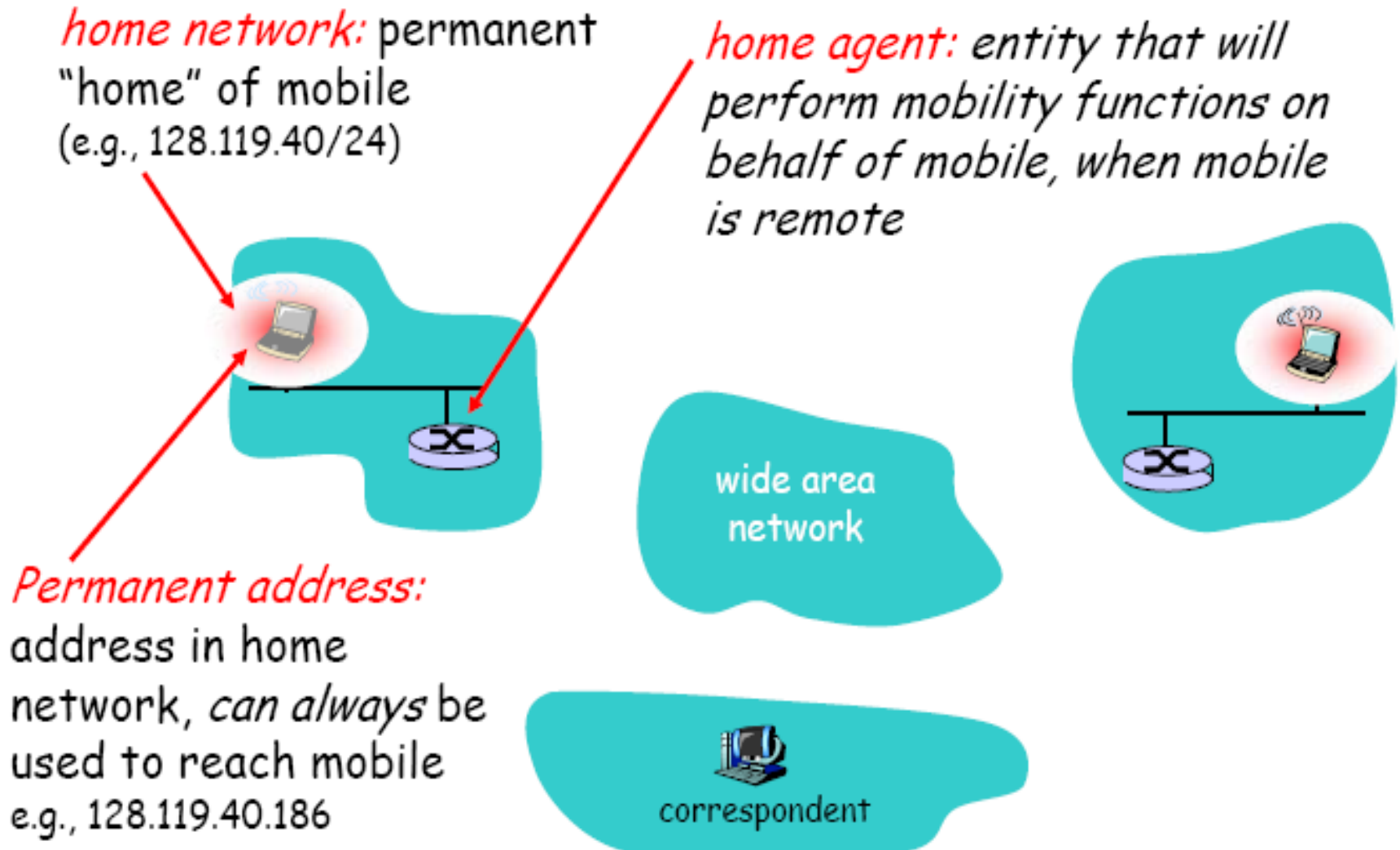
- The Name XXX is binded to a „home Address"

- Under the Home address a HOME REGISTER is maintained
  - a Data structure representing the the ACTUAL ADDRESS

- Resolving the actual address..
  - The "corresponding entity" uses the Home Address
  - Forwarding vs. answering with the "actual address"

- In both cases multiple Intermediate steps are possible

- Look ahead: Mobility support

# Home-Based Approaches



Host's home location

1. Send packet to host at its home

2. Return address of current location

3. Tunnel packet to current location

Client's location

4. Send successive packets to current location

Host's present location

## Disclosing the current address….

# Forwarding (see: cellular telephony)



*home network:* permanent "home" of mobile (e.g., 128.119.40/24)

*home agent:* entity that will perform mobility functions on behalf of mobile, when mobile is remote

*Permanent address:* address in home network, *can always* be used to reach mobile e.g., 128.119.40.186

wide area network

correspondent

# Mobility support

**Permanent address:** remains constant (e.g., 128.119.40.186)

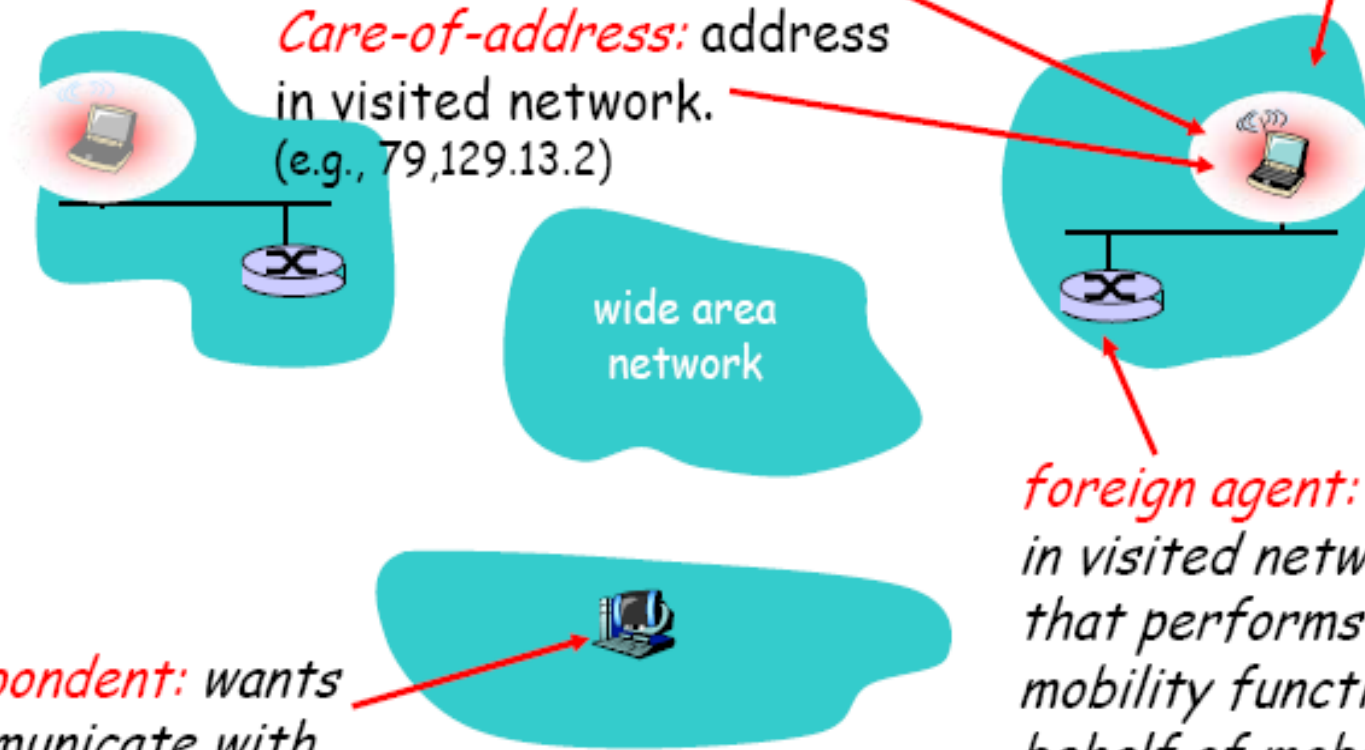**visited network:** network in which mobile currently resides (e.g., 79.129.13/24)

**Care-of-address:** address in visited network. (e.g., 79,129.13.2)

**wide area network**
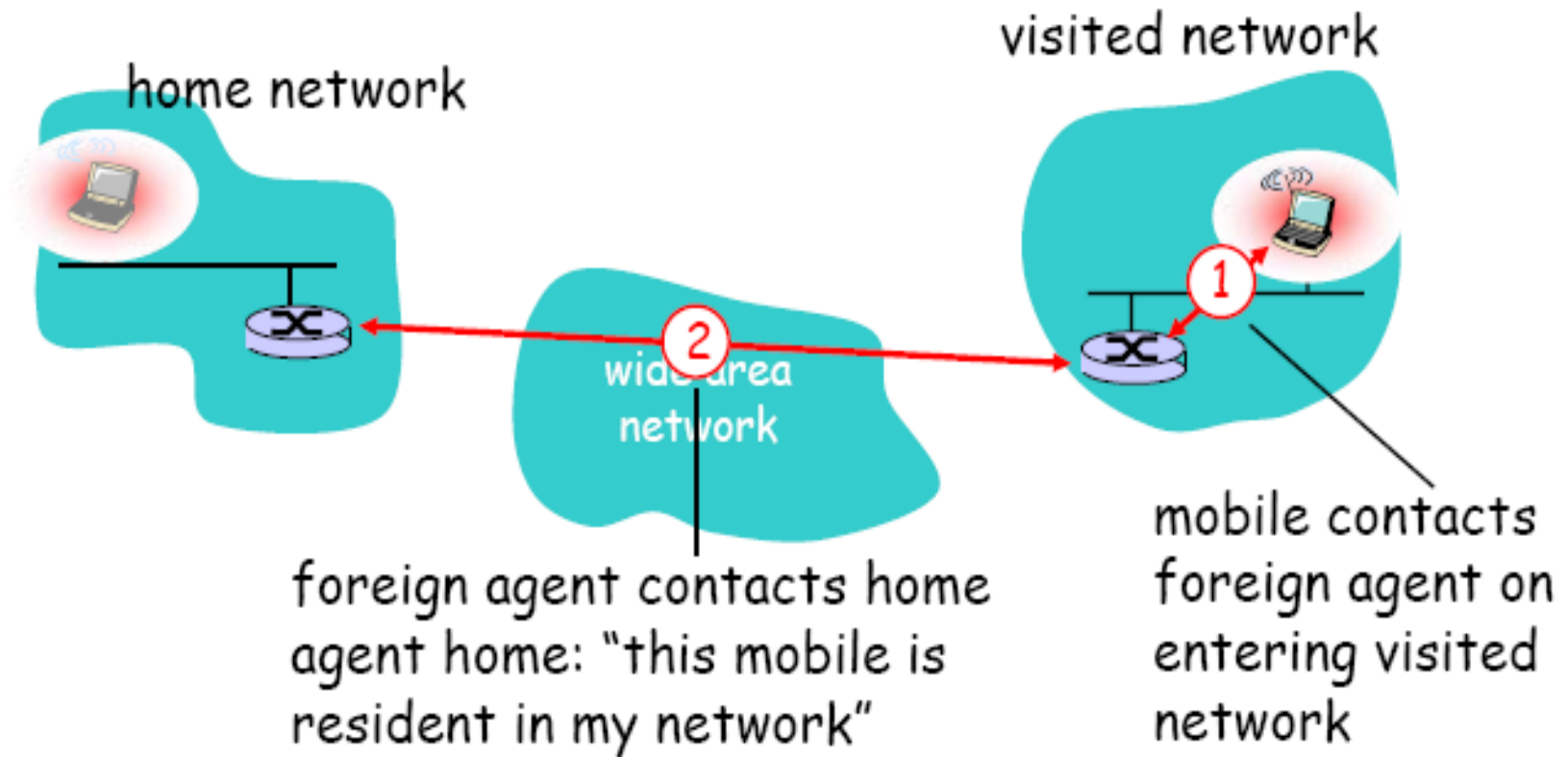
**correspondent:** wants to communicate with mobile

**foreign agent:** entity in visited network that performs mobility functions on behalf of mobile.

4c: Network Layer    13

# Updates



home network

visited network

wide area network

foreign agent contacts home agent home: "this mobile is resident in my network"

mobile contacts foreign agent on entering visited network

End result:
- Foreign agent knows about mobile
- Home agent knows location of mobile

# Indirect Routing



home agent intercepts packets, forwards to foreign agent

foreign agent receives packets, forwards to mobile

visited network

home network

wide area network

correspondent addresses packets using home address of mobile

mobile replies directly to correspondent