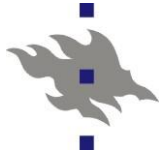**HELSINGIN YLIOPISTO**
**HELSINGFORS UNIVERSITET**
**UNIVERSITY OF HELSINKI**

# Peer-to-Peer and Grid Computing

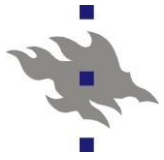Chapter 3: Networks, Searching and
Distributed Hash Tables

# Chapter Outline

n Networks and graphs

    n Graph theory meets networking

    n Different types of graphs and their properties

n Searching and addressing

    n Structured and unstructured networks

n Distributed Hash Tables

    n What they are?

    n How they work?

    n What are they good for?

    n Examples: Chord, CAN, Plaxton/Pastry/Tapestry

# Networks and Graphs

n Refresher of graph theory

n Graph families and models

    n Random graphs

    n Small world graphs

    n Scale-free graphs

n Graph theory and P2P

    n How are the graph properties reflected in real systems?
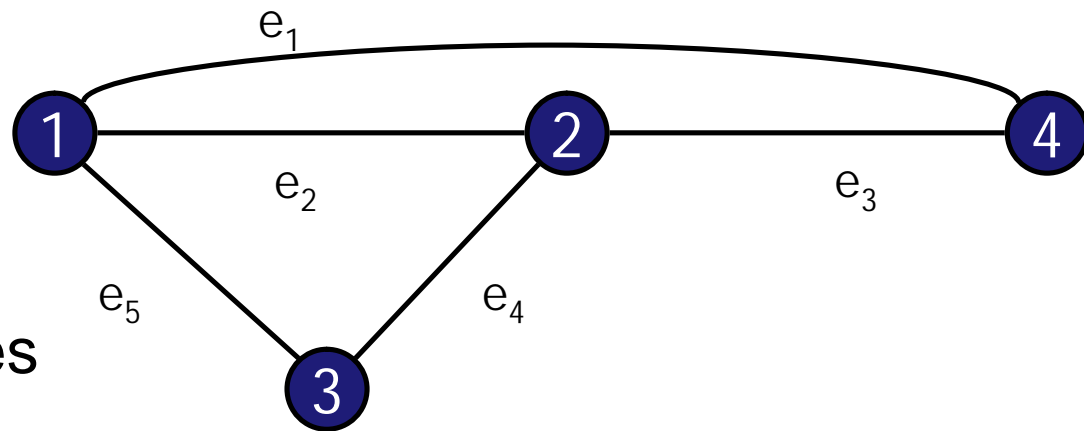
# What Is a Graph?

n Definition of a graph:

Graph $G = (V, E)$ consists of two finite sets, set $V$ of vertices (nodes) and set $E$ of edges (arcs) for which the following applies:

1. If $e \in E$, then exists $(v, u) \in V \times V$, such that $v \in e$ and $u \in e$

2. If $e \in E$ and above $(v, u)$ exists, and further for $(x, y) \in V \times V$ applies $x \in e$ and $y \in e$, then $\{v, u\} = \{x, y\}$
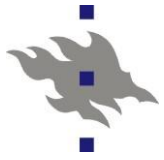
Example graph with 4 vertices and 5 edges

# Properties of Graphs

- n An edge $e \in E$ is directed if the start and end vertices in condition 2 above are identical: $v = x$ and $y = u$
- n An edge $e \in E$ is undirected if $v = x$ and $y = u$ as well as $v = y$ and $u = x$ are possible
- n A graph $G$ is directed (undirected) if the above property holds for all edges
- n A *loop* is an edge with identical endpoints
- n Graph $G_1 = (V_1, E_1)$ is a subgraph of $G = (V, E)$, if $V_1 \subseteq V$ and $E_1 \subseteq E$ (such that conditions 1 and 2 are met)
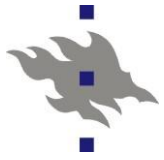
# Important Types of Graphs

- Vertices $v, u \in V$ are connected if there is a path from $v$ to $u$: $(v, v_2), (v_2, v_3), \ldots, (v_{k-1}, u) \in E$
- Graph $G$ is connected if all $v, u \in V$ are connected
- Undirected, connected, acyclic graph is called a tree
  - Sidenote: Undirected, acyclic graph which is not connected is called a forest
- Directed, connected, acyclic graph is also called DAG
  - DAG = directed, acyclic graph (connected is "assumed")
- An induced graph $G(V_C) = (V_C, E_C)$ is a graph $V_C \subseteq V$ and with edges $E_C = \{e = (i, j) \mid i, j \in V_C\}$
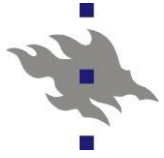- An induced graph is a component if it is connected

# Vertex Degree

n In graph *G = (V, E)*, the degree of vertex *v* ∈ *V* is the total number of edges *(v, u)* ∈ *E* and *(u, v)* ∈ *E*

    n Degree is the number of edges which touch a vertex

n For directed graph, we distinguish between in-degree and out-degree

    n In-degree is number of edges coming to a vertex

    n Out-degree is number of edges going away from a vertex

n The degree of a vertex can be obtained as:

    n Sum of the elements in its row in the incidence matrix

    n Length of its vertex incidence list

# Important Graph Metrics

∩ Distance: *d(v, u)* between vertices *v* and *u* is the length of the shortest path between *v* and *u*

∩ Average path length: Sum of the distances over all pairs of nodes divided by the number of pairs

∩ Diameter: *d(G)* of graph *G* is the maximum of *d(v, u)* for all *v, u* $\in$ *V*

# Six Degrees of Separation

- Famous experiment from 1960's (S. Milgram)
- Send a letter to random people in Kansas and Nebraska and ask people to forward letter to a person in Boston
    - Person identified by name, profession, and city
- Rule: Give letter only to people you know by first name and ask them to pass it on according to same rule
- Some letters reached their goal
- Letter needed six steps on average to reach the person
- Graph theoretically: Social networks have dense local structure, but (apparently) small diameter
- How to model such networks?

# Random Graphs

n Random graphs are first widely studied graph family

    n Many P2P networks choose neighbors more or less randomly

n Two different notations generally used:

    n Erdös and Renyi

    n Gilbert (we will use this)

n Gilbert's definition: Graph $G_{n,p}$ (with $n$ nodes) is a graph where the probability of an edge $e = (v, w)$ is $p$

Construction algorithm:

n For each possible edge, draw a random number

n If the number is smaller than $p$, then the edge exists

n $p$ can be function of $n$ or constant

# Basic Results for Random Graphs

Giant Connected Component:

*Let c > 0 be a constant and p = c/n. If c < 1 every component of $G_{n,p}$ has order O(log N) with high probability. If c > 1 then there will be one component of size n\*(f(c) + O(1)) where f(c) > 0, with high probability. All other components have size O(log N)*

n In plain English: Giant connected component emerges with high probability when average degree is about 1

Node degree distribution

n If we take random node, how high is probability *P(k)* that node has degree *k*?

n Node degree is Poisson distributed

$$P(k) = \frac{c^k e^{-c}}{k!}$$

# More Basic Results

Clustering coefficient

n Clustering coefficient measures number of edges between neighbors divided by maximum number of edges between them (clique-like)

n Clustering coefficient C(i) is defined as

$$C(i) = \frac{E(N(i))}{d(i)(d(i)-1)}$$

    n $E(N(i))$ = number of edges between neighbors of $i$

    n $d(i)$ = degree of $i$

n Clustering coefficient of a random graph is asymptotically equal to $p$ with high probability

# Random Graphs: Summary

n   Before random graphs, regular graphs were popular

   n   Regular: Every node has same degree

n   Random graphs have two advantages over regular graphs

1.   Many interesting properties analytically solvable

2.   Much better for applications, e.g., social networks

n   Note: Does not mean social networks are random graphs; just that the properties of social networks are well-described by random graphs

n   Question: How to model networks with local clusters and small diameter?

n   Answer: Small-world networks

# Small-World Networks

n Developed/discovered by Watts and Strogatz (1998)

  n Over 30 years after Milgram's experiment!

n Watts and Strogatz looked at three networks

  n Film collaboration between actors

  n US power grid

  n Neural network of worm *C. elegans*

n Results:

  n Compared to a random graph with same number of nodes

  n Diameters similar, slightly higher for real graph
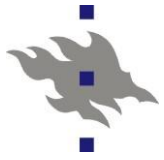
  n Clustering coefficient orders of magnitude higher

Definition of small-worlds network:

n Dense local clustering structure and small diameter
comparable to that of a same-sized random graph

# Constructing Small-World Graphs

n Put all *n* nodes on a ring, number them consecutively from 1 to *n*

n Connect each node with its *k* clockwise neighbors

n Traverse around ring in clockwise order

n For every edge:

  n Draw random number *r*

  n If *r < p*, then re-wire edge by selecting a random target node from the set of all nodes (no duplicates)

  n Otherwise keep old edge

n Different values of *p* give different graphs

  n If *p* is close to 0, then original structure mostly preserved

  n If *p* is close to 1, then new graph is random

  n Interesting things happen when *p* is somewhere in-between

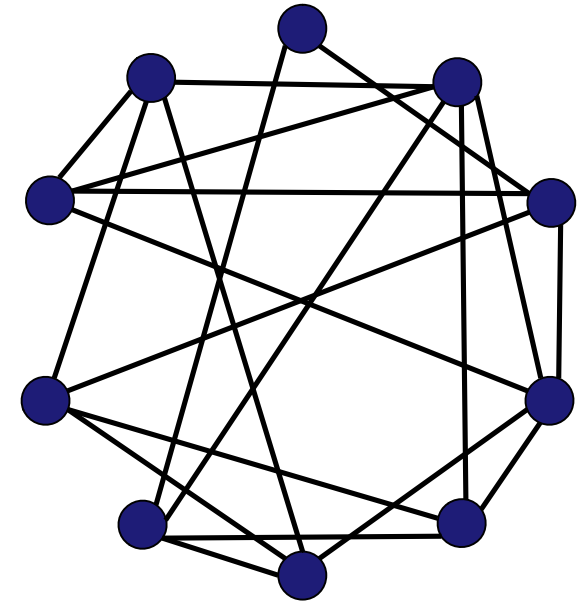# Regular, Small-World, Random

Regular           Small-World           Random

p = 0           p = 1

# Problems with Small-World Graphs

Small-world graphs explain why:

n Highly clustered graphs can have short average path lengths

Small-world graphs do *NOT* explain why:

n This property emerges in real networks
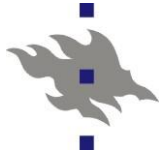
    n Real networks are practically never ring-like

Further problem with small-world graphs:

n Nearly all nodes have same degree

n Not true for random graphs ($k$ edges $\sim c^k/k!$)

n Is same true for real networks too?

n Let's look at the Internet…

# Internet

n Famous study by Faloutsos et al. (3 brothers! ;-) in 1999

n They examined Internet topology during 1998

n AS-level topology, during 1998 Internet grew 45%

Motivation for work:

n What does the Internet look like?

n Are there any topological properties that don't change over time?

n How can I generate Internet-like graphs for simulations?

# Faloutsos Results

n   4 key properties, each follows a power-law

n   Sort nodes according to their (out)degree

1. *Outdegree of a node is proportional to its rank to the power of a constant*

2. *Number of nodes with same outdegree is proportional to the outdegree to the power of a constant*

3. *Eigenvalues of a graph are proportional to the order to the power of a constant*

4. *Total number of pairs of nodes within a distance d is proportional to d to the power of a constant*

•   Why would Internet obey such laws?

# Answer: Power-Law Networks

n   Also known as scale-free networks

n   Barabasi-Albert-Model

1.  Network grows in time

2.  New node has preferences to whom it wants to connect

n   Preferential connectivity modeled as

    n   Each new node wants to connect to $m$ other nodes

    n   Probability that an existing node $j$ gets one of the $m$ connections is proportional to its degree $d(j)$

n   New nodes tend to connect to well-connected nodes

n   Another way to express this is "rich get richer"

# Applications to Peer-to-Peer

n Small-world model explains why short paths exist

n Why can we find these paths?

    n Each node has only local information

    n Milgram's results showed first steps were the largest

n How to model this?

n Kleinberg's Small-World Model

    n Set of points in an $n$ x $n$ grid

    n Distance is the number of "steps" separating points

      - $d(i, j) = |x_i - x_j| + |y_i - y_j|$

n Construct graph as follows:

    n Every node $i$ is connected to node $j$ within distance $q$

    n For every node $i$, additional $q$ edges are added. Probability that node $j$ is selected is proportional to $d(i, j)^{-r}$, for some constant $r$

# Navigation in Kleinberg's Model

n We want to send a message to another node

n Algorithm is decentralized if sending node only knows:

    n Its local neighbors

    n Position of the target node on the grid

    n Locations and long-range contacts of all nodes who come in contact of the message (not needed below, actually)

n Can be shown: Number of messages needed is proportional to $O(log\ n)$ (only one correct $r$ per case)

n Practical algorithm: Forward message to contact who is closest to target

n Note: Kleinberg's model assumes some way of associating nodes with points in grid

    n Compare with CAN DHT

# Power Law Networks and P2P

n Robustness comparison of random and power-law graphs

n Take network of 10000 nodes (random and power-law) and remove nodes randomly

n Random graph:

  n Take out 5% of nodes: Biggest component 9000 nodes

  n Take out 18% of nodes: No biggest component, all components between 1 and 100 nodes

  n Take out 45% of nodes: Only groups of 1 or 2 survive

n Power-law graph:

  n Take out 5% of nodes: Only isolated nodes break off

  n Take out 18% of nodes: Biggest component 8000 nodes

  n Take out 45% of nodes: Large cluster persists, fragments small

n Recall Gnutella: *Applies ONLY for random failures*

# Summary of Graphs

n Three kinds of graph models:

  n Random graph

  n Small-World

  n Power-Law (Scale-Free)

n Small-world graphs explain why we can have high clustering and short average paths

n Power-law graphs explain how graphs are built in many real networks

# Searching and Addressing

n  Two basic ways to find objects:

1. Search for them

2. Address them using their unique name

n  Both have pros and cons (see below)

n  Most existing P2P networks built on searching, but some
   networks are based on addressing objects

n  Difference between searching and addressing is a very
   fundamental difference

   n  Determines how network is constructed

   n  Determines how objects are placed

   n  "Determines" efficiency of object location

n  Let's compare searching and addressing

# Addressing vs. Searching

• "Addressing" networks find objects by addressing them with their unique name (cf. URLs in Web)
• "Searching" networks find objects by searching with keywords that match objects's description (cf. Google)

## Addressing

n  Pros:

    n  Each object uniquely identifiable

    n  Object location can be made efficient

n  Cons:

    n  Need to know unique name

    n  Need to maintain structure required

       by addresses

## Searching

n  Pros:

    n  No need to know unique names

      -  More user friendly

n  Cons:

    n  Hard to make efficient

      -  Can solve with money, see Google

    n  Need to compare actual objects to know

       if they are same

# Addressing vs. Searching: Examples

|  | Searching | Addressing |
|---|---|---|
| Physical name of object | Searching in P2P networks, Searching in filesystem (Desktop searches) (Search components of URL with Google?) | URLs in Web |
| Logical name of object | ? (Search components of URNs) | Object names in DHT, URNs |
| Content or metadata of object | Searching in P2P networks, Standard Google search Desktop searches | N/A |

# Searching, Addressing, and P2P

n We can distinguish two main P2P network types

Unstructured networks/systems

n Based on searching

n Unstructured does NOT mean complete lack of structure

  n Network has graph structure, e.g., scale-free

n Network has structure, but peers are free to join anywhere and objects can be stored anywhere

n So far we have seen unstructured networks

Structured networks/systems

n Based on addressing

n Network structure determines where peers belong in the network and where objects are stored

n How to build structured networks?

# Another Classification of P2P Systems

n Sometimes P2P systems classified in generations

n No 100% consensus on what is in which generation

n 1st generation

    n Typically: Napster

n 2nd generation

    n Typically: Gnutella

n 3rd generation

    n Typically: Superpeer networks

n 4th generation

    n Typically: Distributed hash tables

    n Note: For DHTs, no division into generations yet

# Distributed Hash Tables

n What are they?

n How they work?

n What are they good for?

n Examples:

    n Chord

    n CAN

    n Plaxton/Pastry/Tapestry

# DHT: Motivation

n Why we need DHTs?

n Searching in P2P networks is not efficient

  n Either centralized system with all its problems

  n Or distributed system with all its problems

  n Hybrid systems cannot guarantee discovery either

n Actual file transfer process in P2P network is scalable

  n File transfers directly between peers

n Searching does not scale in same way

n Original motivation for DHTs: More efficient searching and object location in P2P networks

n Put another way: Use addressing instead of searching

# Recall: Hash Tables

n Hash tables are a well-known data structure

n Hash tables allow insertions, deletions, and finds in
   constant (average) time

n Hash table is a fixed-size array

   n Elements of array also called *hash buckets*

n *Hash function* maps keys to elements in the array

n Properties of good hash functions:

   n Fast to compute

   n Good distribution of keys into hash table

   n Example: SHA-1 algorithm

# Hash Tables: Example

| | |
|---|---|
| 0 | → 0 |
| 1 | → 1 |
| 2 | |
| 3 | |
| 4 | → 4 |
| 5 | → 25 |
| 6 | → 16 |
| 7 | |
| 8 | |
| 9 | → 9 |

n Hash function:

$hash(x) = x \bmod 10$

n Insert numbers 0, 1, 4, 9, 16, and 25

n Easy to find if a given key is present in the table

# Distributed Hash Table: Idea

n Hash tables are fast for lookups

n Idea: Distribute hash buckets to peers

n Result is Distributed Hash Table (DHT)

n Need efficient mechanism for finding which peer is responsible for which bucket and routing between them

# DHT: Principle

n In a DHT, each node is responsible for one or more hash buckets

  n As nodes join and leave, the responsibilities change

n Nodes communicate among themselves to find the responsible node

  n Scalable communications make DHTs efficient

n DHTs support all the normal hash table operations

# Summary of DHT Principles

- Hash buckets distributed over nodes
- Nodes form an overlay network
    - Route messages in overlay to find responsible node
- Routing scheme in the overlay network is the difference between different DHTs
- DHT behavior and usage:
    - Node knows "object" name and wants to find it
        - Unique and known object names assumed
    - Node routes a message in overlay to the responsible node
    - Responsible node replies with "object"
        - Semantics of "object" are application defined

# DHT Examples

n In the following look at some example DHTs

    n Chord

    n CAN

    n Tapestry

n Several others exist too

    n Pastry, Plaxton, Kademlia, Koorde, Symphony, P-Grid, CARP, …

n All DHTs provide the same abstraction:

    n DHT stores key-value pairs

    n When given a key, DHT can retrieve/store the value

    n No semantics associated with key or value

n Difference is in overlay routing scheme

# Chord

n Chord was developed at MIT

n Originally published in 2001 at Sigcomm conference

n Chord's overlay routing principle quite easy to understand

    n Paper has mathematical proofs of correctness and performance

n Many projects at MIT around Chord

    n CFS storage system

    n Ivy storage system

    n Plus many others…

# Chord: Basics

n Chord uses SHA-1 hash function

 n Results in a 160-bit object/node identifier

 n Same hash function for objects and nodes

n Node ID hashed from IP address

n Object ID hashed from object name

 n Object names somehow assumed to be known by everyone

n SHA-1 gives a 160-bit identifier space

n Organized in a ring which wraps around

 n Nodes keep track of predecessor and successor

 n Node responsible for objects between its predecessor and itself

 n Overlay is often called "Chord ring" or "Chord circle"

# Chord: Examples

n Below examples for:

    n How to join the Chord ring

    n How to store and retrieve values

# Joining: Step-By-Step Example

n Setup: Existing network with nodes on 0, 1 and 4

n Note: Protocol messages simply examples

n Many different ways to implement Chord

    n Here only conceptual example

    n Covers all important aspects

# Joining: Step-By-Step Example: Start

n New node wants to join

n Hash of the new node: 6

n Known node in network: Node1

n Contact Node1

   n Include own hash

# Joining: Step-By-Step Example: Situation Before Join



Data for ]4;0]

pred0

pred1

0

succ0

Data for ]0;1]

7

1

6

2

succ4

succ1

5

3

No data

4

pred4

Data for ]1;4]

# Joining: Step-By-Step Example: Contact known node

- n Arrows indicate open connections
- n Example assumes connections are kept open, i.e., messages processed recursively
- n Iterative processing is also possible

0

7

1

6

2

5

3

4

JOIN 6

**Joining: Step-By-Step Example:**
**Join gets routed along the network**



0

7

1

6

2

JOIN 6

5

4

# Joining: Step-By-Step Example:
# Successor of New Node Found



JOIN 6

0
1
2
3
4
5
6
7

# Joining: Step-By-Step Example: Joining Successful + Transfer

Joining is successful

Old responsible node transfers data that should be in new node

New node informs Node4 about new successor (not shown)

TRANSFER
Data in range ]4;6]

0

7

1

6

2

5

3

4

Note: Transferring can happen also later

# Joining: Step-By-Step Example:
# All Is Done



Data for ]6;0]

pred0    pred1

0    succ0    Data for ]0;1]

pred0
7    1

Data for ]4;6]    succ6

succ4    6    2    succ1

pred6    5    4    pred4

Data for ]1;4]

# Storing a Value

n Node 6 wants to store
object with name "Foo"
and value 5

n hash(Foo) = 2

# Storing a Value

STORE 2 5

0

7

1

6

2

5

3

4

# Storing a Value

STORE 2 5

0

7

1

6

2

5

3

4

# Storing a Value



0

7

1

6

2

3

5

4

STORE 2 5

Value is now stored
in node 4.

# Retrieving a Value

n Node 1 wants to get object with name "Foo"
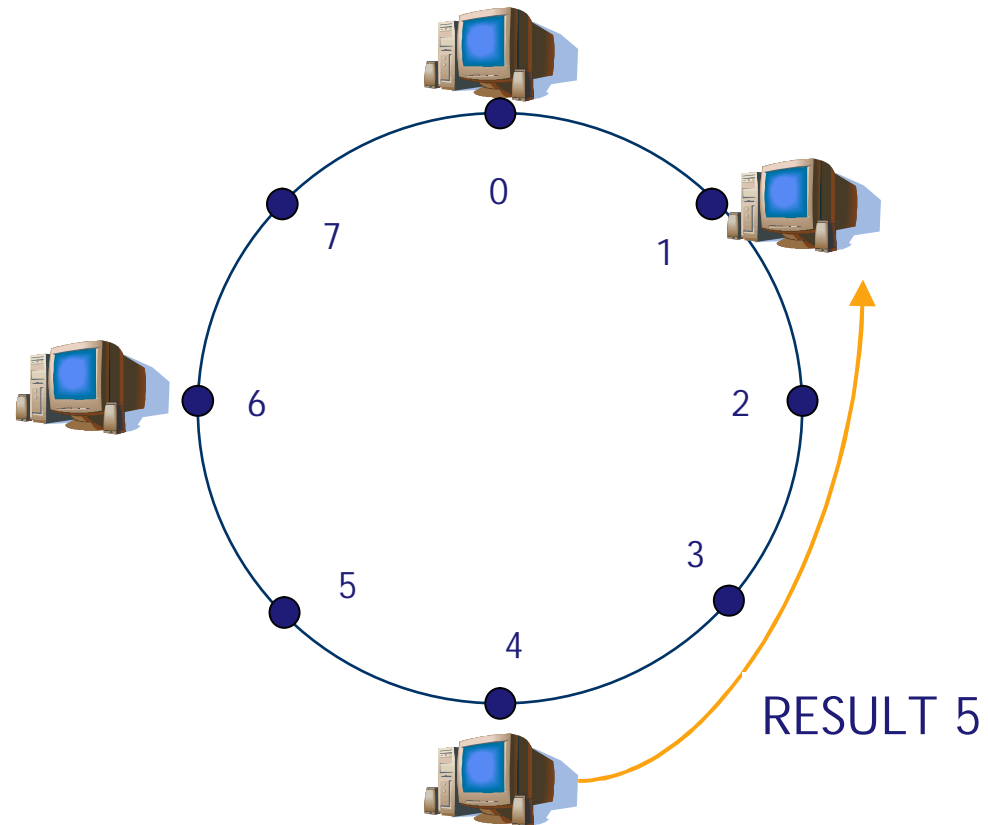
n hash(Foo) = 2

à Foo is stored on node 4

0

7

1
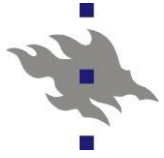
6

2

5

3

4

# Retrieving a Value



RETRIEVE 2

# Retrieving a Value



0

7

1

6

2

5

3

4

RESULT 5

# Chord: Scalable Routing

n Routing happens by passing message to successor

n What happens when there are 1 million nodes?

  n On average, need to route 1/2-way across the ring

  n In other words, 0.5 million hops! Complexity *O(n)*

n How to make routing scalable?

n Answer: Finger tables

n Basic Chord keeps track of predecessor and successor

n Finger tables keep track of more nodes

  n Allow for faster routing by jumping long way across the ring

  n Routing scales well, but need more state information

n Finger tables not needed for correctness, only performance improvement

# Chord: Finger Tables

n In *m*-bit identifier space, node has up to *m* fingers

n Fingers are stored in the finger table

n Row *i* in finger table at node *n* contains first node *s* that succeeds *n* by at least $2^{i-1}$ on the ring

n In other words:

$$finger[i] = successor(n + 2^{i-1})$$

n First finger is the successor

n Distance to *finger[i]* is at least $2^{i-1}$

# Chord: Scalable Routing

n Finger intervals increase with distance from node n

    n If close, short hops and if far, long hops

Two key properties:

n Each node only stores information about a small number of nodes

n Cannot in general determine the successor of an arbitrary ID

n Example has three nodes at 0, 1, and 4

n 3-bit ID space --> 3 rows of fingers

| Start | Int. | Succ. |
|---|---|---|
| 1 | [1,2) | 1 |
| 2 | [2,4) | 4 |
| 4 | [4,0) | 4 |

| Start | Int. | Succ. |
|---|---|---|
| 2 | [2,3) | 4 |
| 3 | [3,5) | 4 |
| 5 | [5,1) | 0 |

| Start | Int. | Succ. |
|---|---|---|
| 5 | [5,6) | 0 |
| 6 | [6,0) | 0 |
| 0 | [0,4) | 0 |

# Chord: Performance

n Search performance of "pure" Chord $O(n)$

  n Number of nodes is $n$

n With finger tables, need $O(\log n)$ hops to find the correct node

  n Fingers separated by at least $2^{i-1}$

  n With high probability, distance to target halves at each step

  n In beginning, distance is at most $2^m$

  n Hence, we need at most $m$ hops

n For state information, "pure" Chord has only successor and predecessor, $O(1)$ state

n For finger tables, need $m$ entries

  n Actually, only $O(\log n)$ are distinct

  n Proof is in the paper

# CAN: Content Addressable Network

n CAN developed at UC Berkeley

n Originally published in 2001 at Sigcomm conference(!)

n CANs overlay routing easy to understand

  n Paper concentrates more on performance evaluation

  n Also discussion on how to improve performance by tweaking

n CAN project did not have much of a follow-up
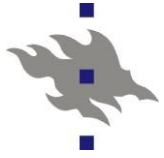
  n Only overlay was developed, no bigger follow-ups

# CAN: Basics

- CAN based on N-dimensional Cartesian coordinate space
    - Our examples: N = 2
    - One hash function for each dimension
- Entire space is partitioned amongst all the nodes
    - Each node owns a zone in the overall space

- Abstractions provided by CAN:
    - Can store data at points in the space
    - Can route from one point to another

- Point = Node that owns the zone in which the point (coordinates) is located
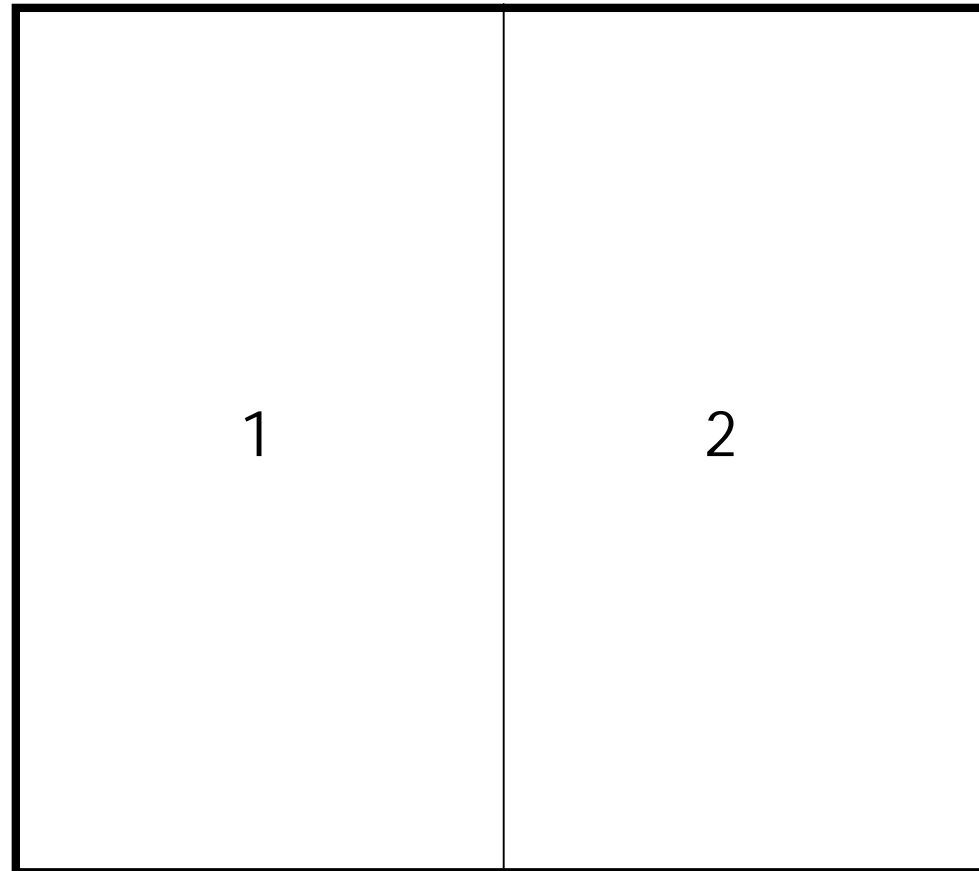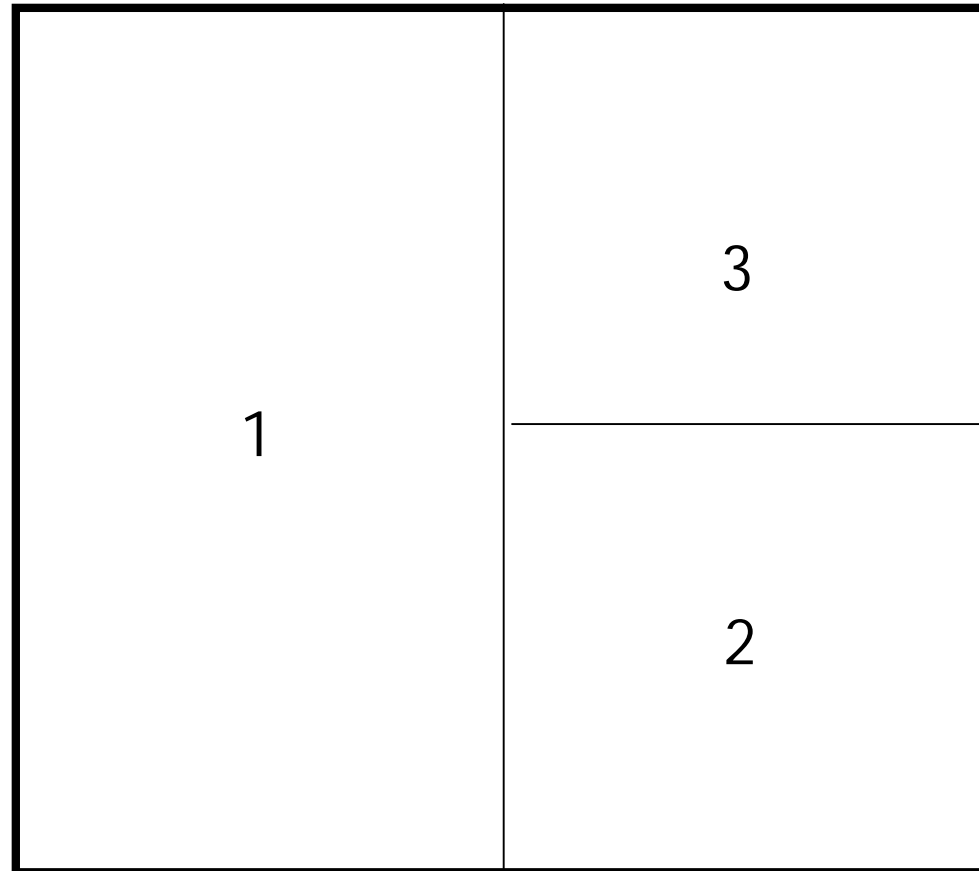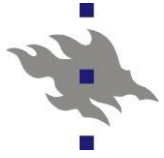- Order in which nodes join is important

# CAN: Partitioning

```
┌─────────────────────────────────┐
│                                 │
│                                 │
│                                 │
│                                 │
│                1                │
│                                 │
│                                 │
│                                 │
│                                 │
└─────────────────────────────────┘
```
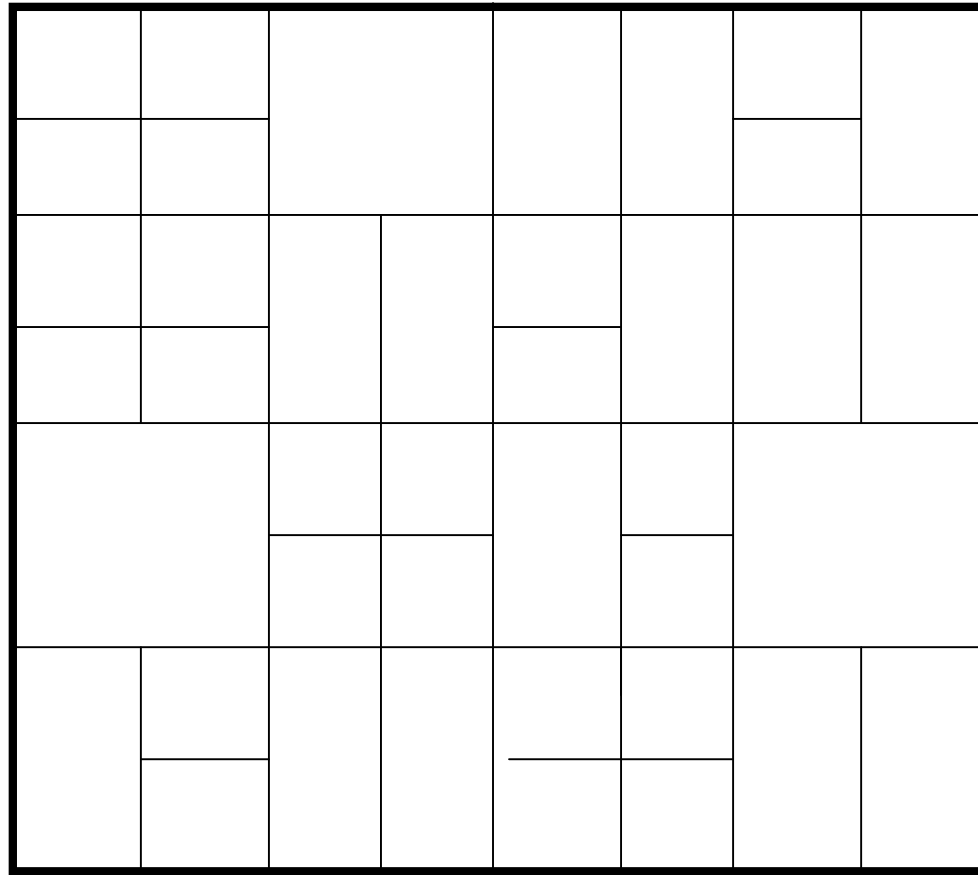
# CAN: Partitioning

# CAN: Partitioning

# CAN: Partitioning

# CAN: Partitioning

n CAN forms a d-
dimensional
torus

# CAN: Examples

n Below examples for:

    n How to join the network

    n How routing tables are managed
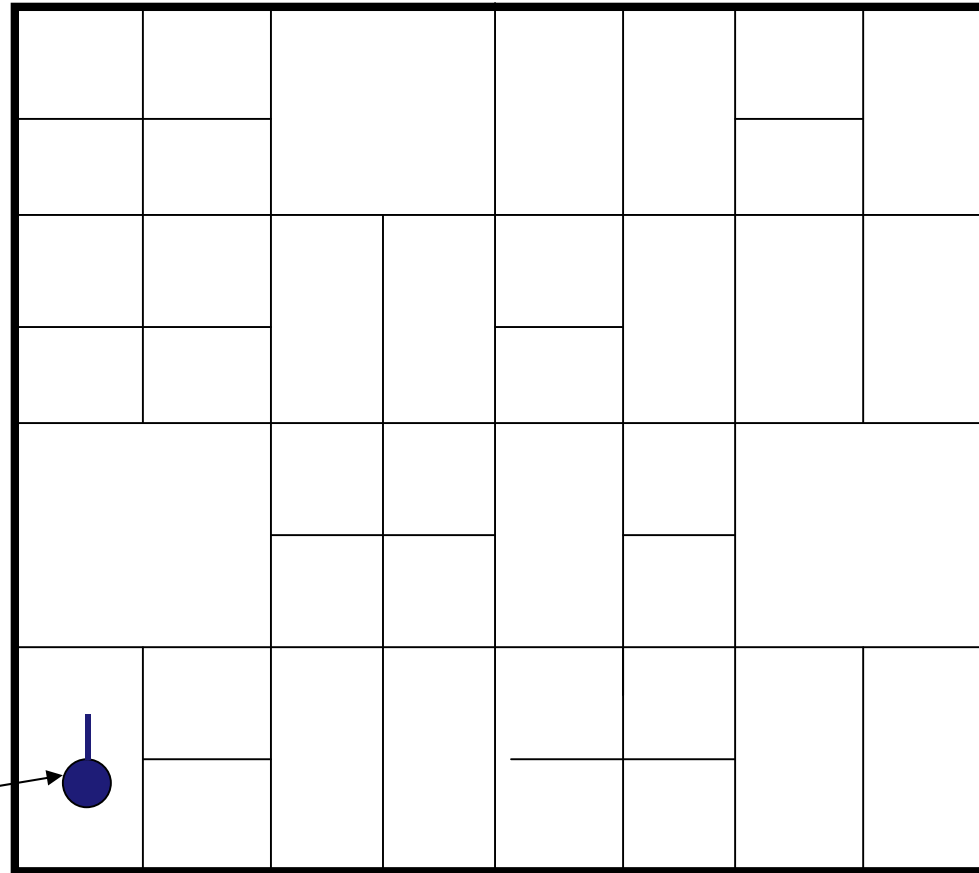
    n How to store and retrieve values

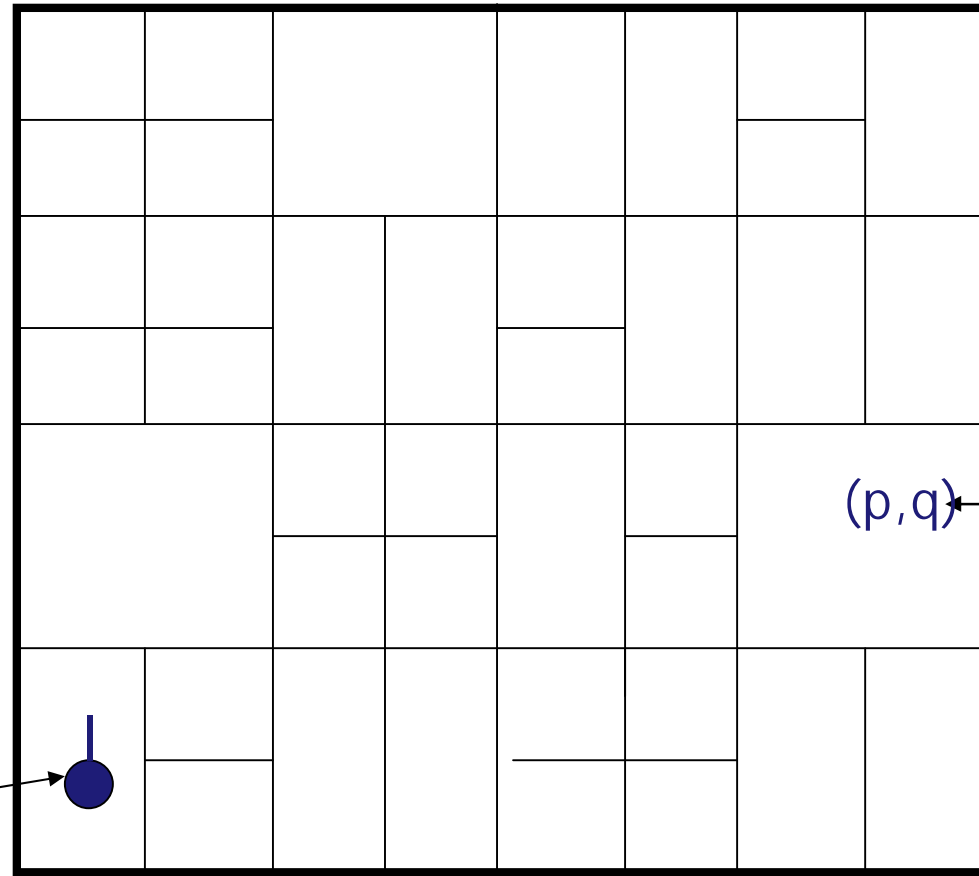# CAN: Node Insertion

Discover some
node "I" already
in CAN



New node

# CAN: Node Insertion
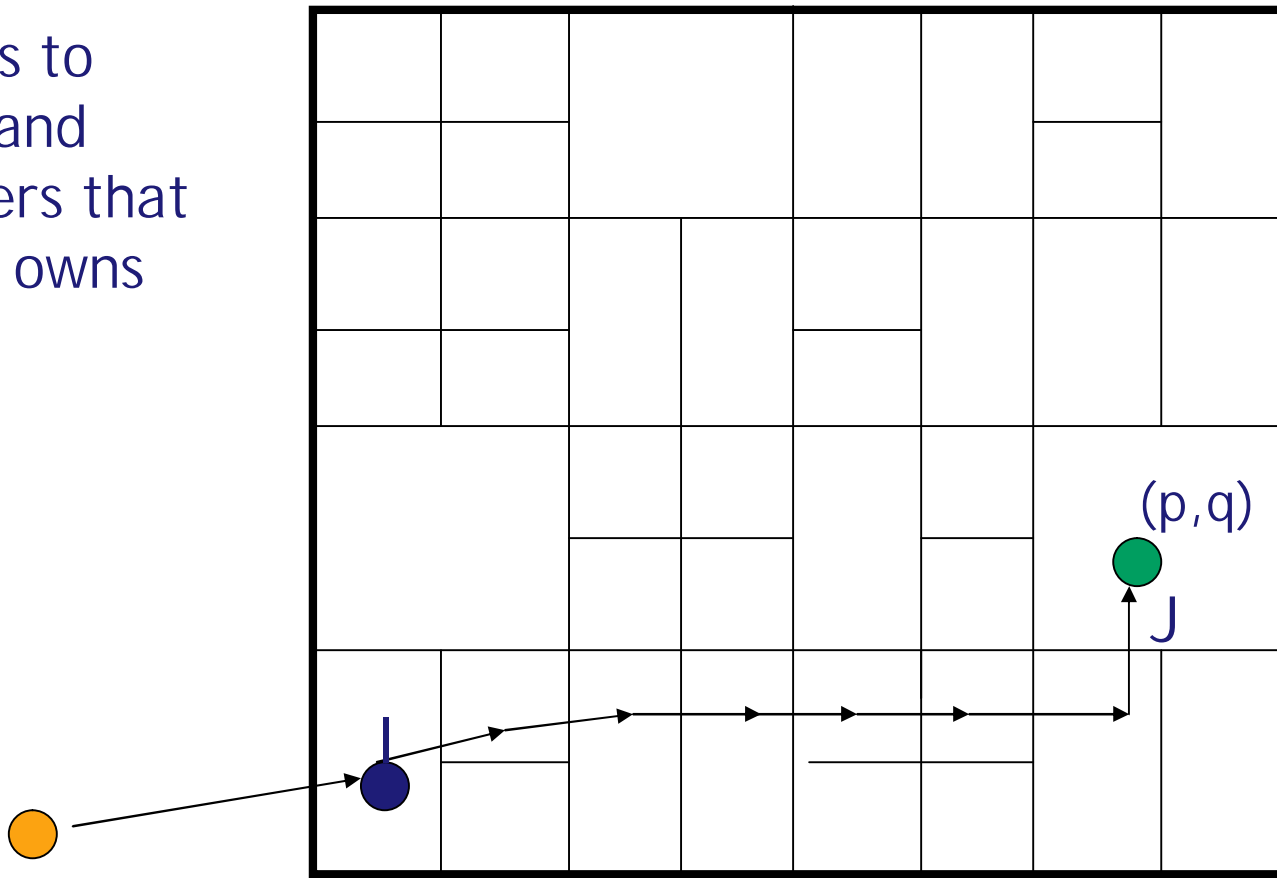
New node picks its coordinates in space

(p,q) ← pick random point in space

New node

# CAN: Node Insertion
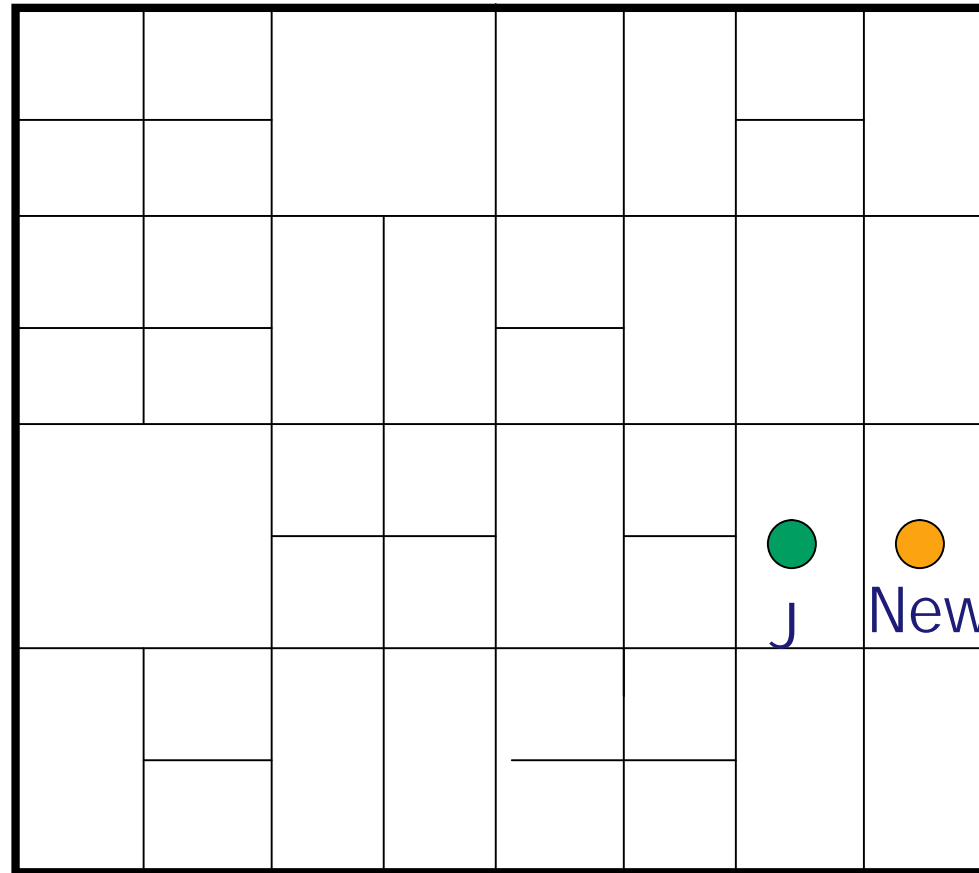
I routes to (p,q), and discovers that node J owns (p,q)


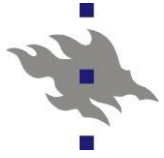
(p,q)
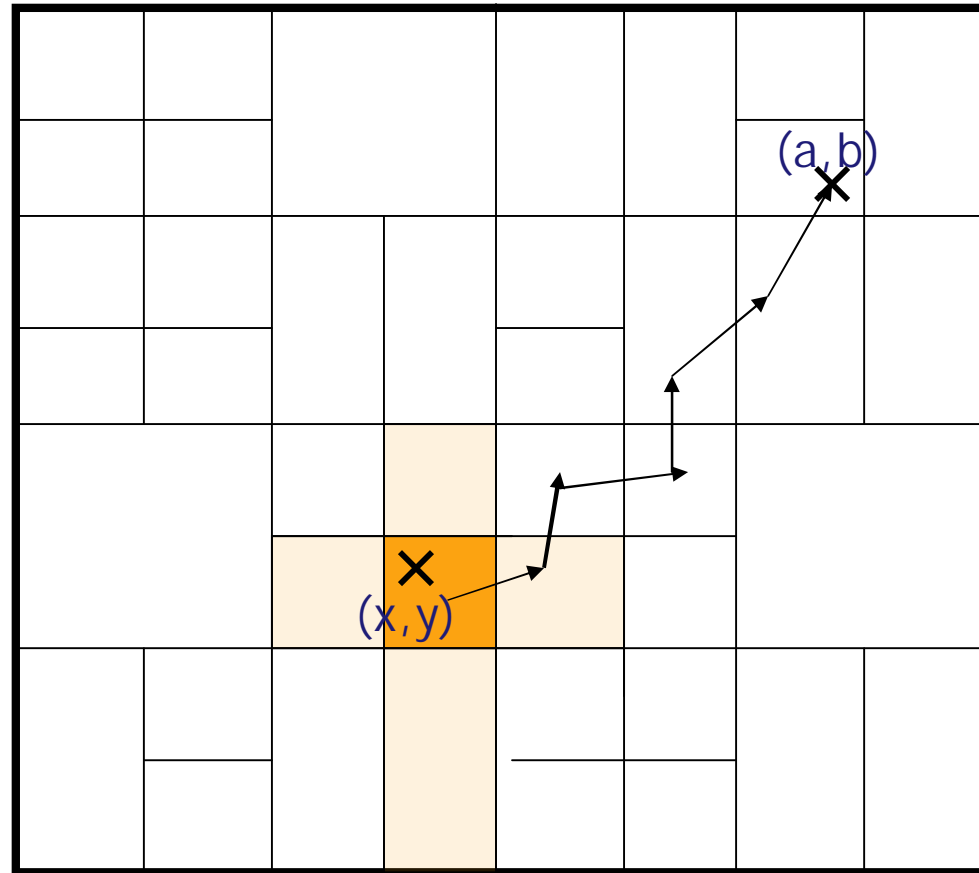
J

I

new node

Split J's zone in half. New owns one half



J    New

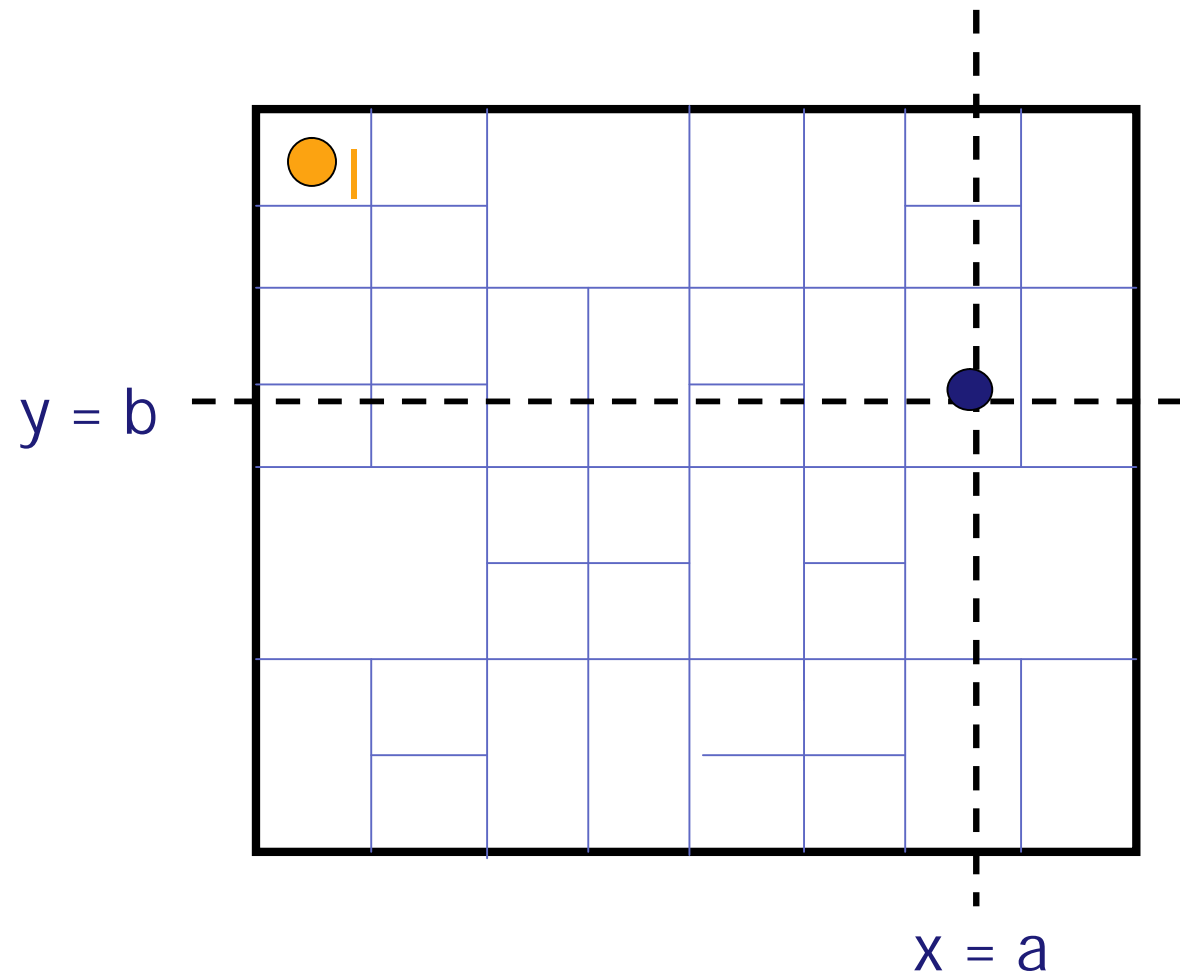# CAN: Routing Table

# CAN: Routing

# CAN: Storing Values

node $I::insert(K,V)$
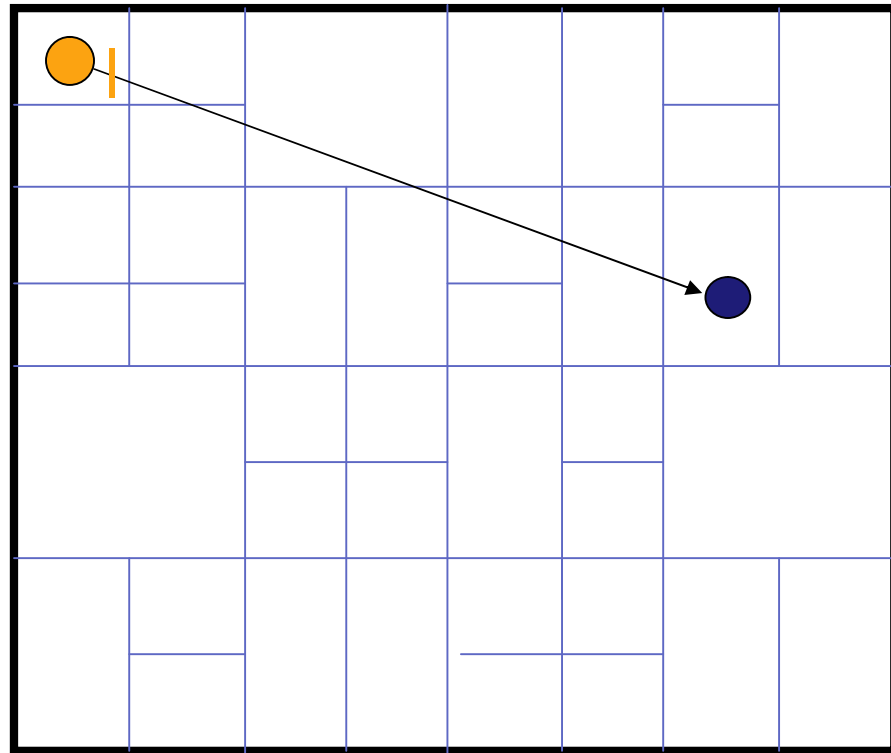
$a = h_x(K)$

$b = h_y(K)$

$y = b$

$x = a$

# CAN: Storing Values
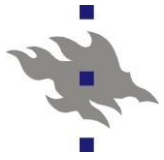
node I::insert(K,V)

(1)  $a = h_x(K)$

 $b = h_y(K)$
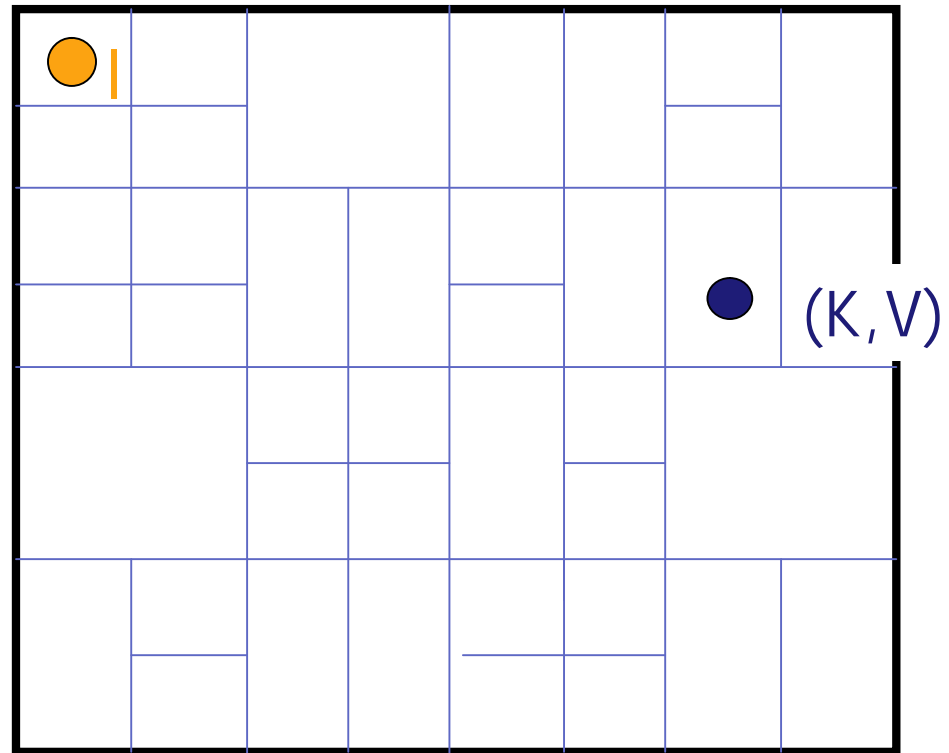
(2)  route(K,V) ->  (a,b)

# CAN: Storing Values

node I::insert(K,V)

(1)  $a = h_x(K)$
     $b = h_y(K)$

(2)  route(K,V) ->  (a,b)

(3)  (a,b) stores (K,V)

# CAN: Retrieving Values

node J::retrieve(K)

(1)  $a = h_x(K)$
     $b = h_y(K)$

(2)  route "retrieve(K)" to (a,b)

(K,V)

J

# CAN: Improvements

n Possible to increase number of dimensions *d*

  n Small increase in routing table size

  - Shorter routing path, more neighbors for fault tolerance

n Multiple realities (= coordinate spaces)

  n Use more hash functions

  n Same properties as increased dimensions

n Routing weighted by round-trip times

  n Take into account network topology

  n Forward to the "best" neighbor

# CAN: More Improvements

n Use well-known landmark servers (e.g., DNS roots)
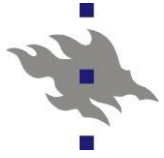
    n Nodes join CAN in different areas, depending on distance to landmarks

       - Pick points "near" landmark

    n Idea: Geographically close nodes see same landmarks

n Uniform partitioning

    n New node splits the largest zone in the neighborhood instead of the zone of the responsible node

# CAN: Performance

n State information at node $O(d)$

  n Number of dimensions is $d$

  n Need two neighbors in all coordinate axis

  n Independent of the number of nodes!

n Routing takes $O(dn^{1/d})$ hops

  n Network has $n$ nodes

  n Multiple dimensions and realities improve this

  n For routing: multiple dimensions are better

  n But: multiple realities improve availability and fault tolerance

# Tapestry

n Tapestry developed at UC Berkeley(!)

  n Different group from CAN developers

n Tapestry developed in 2000, but published in 2004

  n Originally only as technical report, 2004 as journal article

n Many follow-up projects on Tapestry

  n Example: OceanStore

n Tapestry based on work by Plaxton et al.

n Plaxton network has also been used by Pastry

n Pastry was developed at Microsoft Research and Rice University

  n Difference between Pastry and Tapestry minimal

  n Tapestry and Pastry add dynamics and fault tolerance to Plaxton network

# Tapestry: Plaxton Network

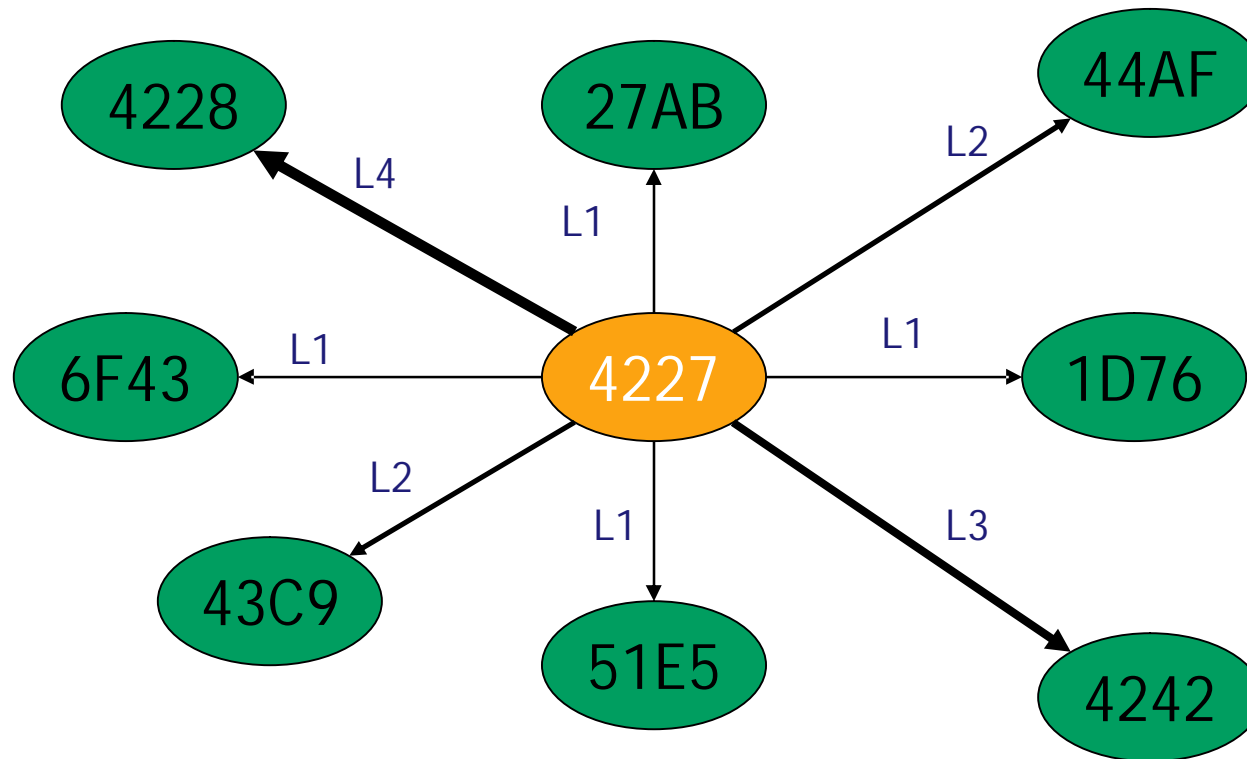n Plaxton network (or Plaxton mesh) based on prefix routing (similar to IP address allocation)

  n Prefix and postfix are functionally identical

  n Tapestry originally postfix, now prefix?!?

n Node ID and object ID hashed with SHA-1

  n Expressed as hexadecimal (base 16) numbers (40 digits)

  n Base is very important, here we use base 16

n Each node has a neighbor map with multiple levels

  n Each level represents a matching prefix up to digit position in ID

  n A given level has number of entries equal to the base of ID

  n $i^{th}$ entry in $j^{th}$ level is closest node which starts *prefix(N,j-1)+"i"*

  n Example: 9th entry of 4th level for node 325AE is the closest node with ID beginning with 3259

# Tapestry: Routing Mesh

n (Partial) routing mesh for a single node 4227

n Neighbors on higher levels match more digits

# Tapestry: Neighbor Map for 4227

| Level | 1 | 2 | 3 | 4 | 5 | 6 | 8 | A |
|-------|------|------|------|------|------|------|------|------|
| 1 | 1D76 | 27AB | | | 51E5 | 6F43 | | |
| 2 | | | 43C9 | 44AF | | | | |
| 3 | | | | | | | | 42A2 |
| 4 | | | | | | | 4228 | |

- There are actually 16 columns in the map (base 16)
- Normally more (most?) entries would be filled
- Tapestry has neighbor maps of size 40 x 16

# Tapestry: Routing Example



- n  Route message from 5230 to 42AD
- n  Always route to node closer to target
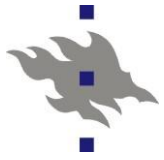    - n  At $n^{th}$ hop, look at $n+1^{th}$ level in neighbor map --> "always" one digit more
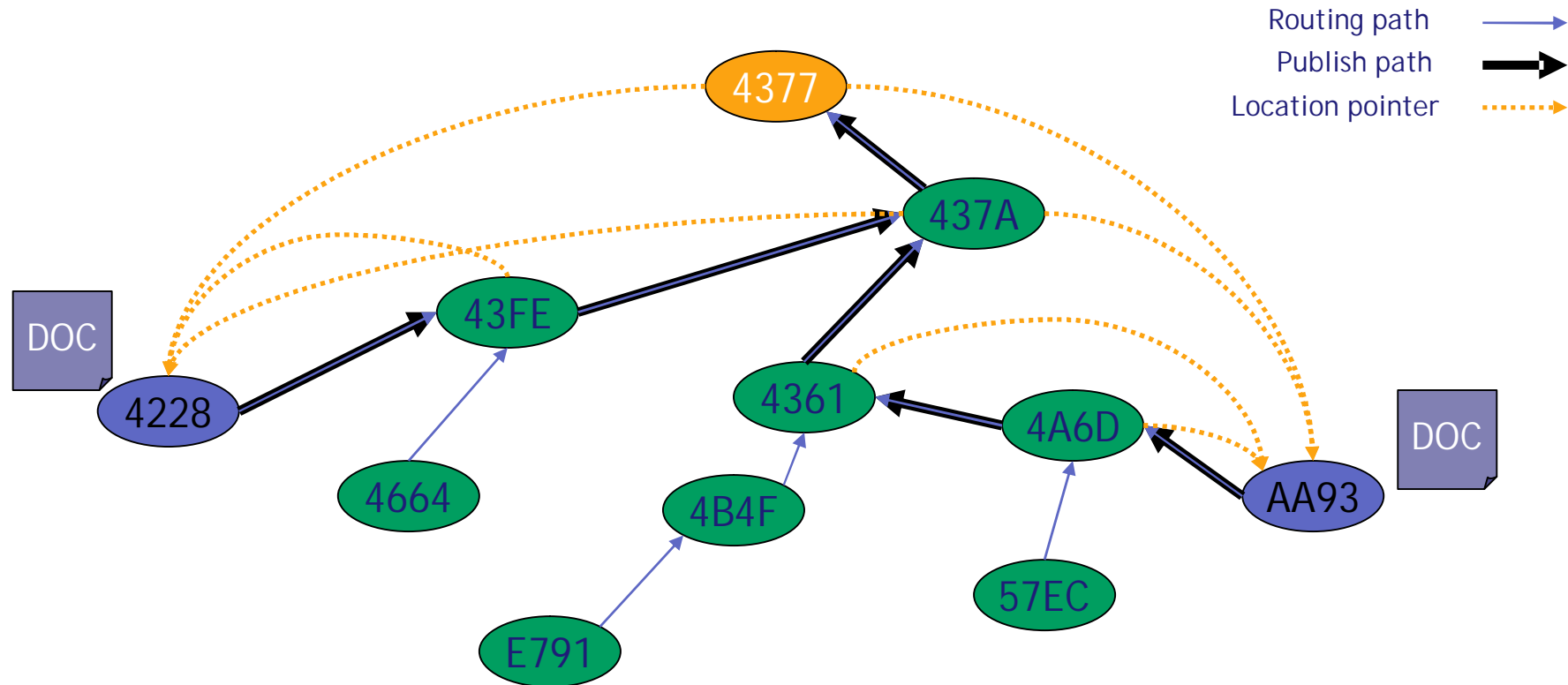- n  Not all nodes and links are shown

# Tapestry: Properties

- Node responsible for objects which have same ID
  - Unlikely to find such node for every object
  - Node responsible also for "nearby" objects (surrogate routing, see below)
- Object publishing:
  - Responsible nodes store only pointers
    - Multiple copies of object possible
    - Each copy must publish itself
  - Pointers cached along the publish path
  - Queries routed towards responsible node
  - Queries "often" hit cached pointers
    - Queries for same object go (soon) to same nodes
- Note: Tapestry focuses on storing objects
  - Chord and CAN focus on values, but in practice no difference

# Tapestry: Publishing Example



Routing path
Publish path
Location pointer

n Two copies of object "DOC" with ID 4377 created at AA93 and 4228
n AA93 and 4228 publish object DOC, messages routed to 4377
    n Publish messages create location pointers on the way
n Any subsequent query can use location pointers

# Tapestry: Querying Example

- **n** Requests initially route towards 4377

- **n** When they encounter the publish path, use location pointers to find object

- **n** Often, no need to go to responsible node

- **n** Downside: Must keep location pointers up-to-date

# Tapestry: Making It Work

n Previous examples show a Plaxton network

  n Requires global knowledge at creation time

  n No fault tolerance, no dynamics

n Tapestry adds fault tolerance and dynamics

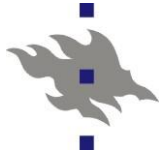  n Nodes join and leave the network

  n Nodes may crash

  n Global knowledge is impossible to achieve

n Tapestry picks closest nodes for neighbor table

  n Closest in IP network sense (= shortest RTT)

  n Network distance (usually) transitive

    - If A is close to B, then B is also close to A

  n Idea: Gives best performance

# Tapestry: Fault-Tolerant Routing

n Tapestry keeps mesh connected with keep-alives

  n Both TCP timeouts and UDP "hello" messages

  n Requires extra state information at each node

n Neighbor table has backup neighbors

  n For each entry, Tapestry keeps 2 backup neighbors

  n If primary fails, use secondary

  - Works well for uncorrelated failures

n When node notices a failed node, it marks it as invalid

  n Most link/connection failures short-lived

  n Second chance period (e.g., day) during which failed node can come back and old route is valid again

  n If node does not come back, one backup neighbor is promoted and a new backup is chosen

# Tapestry: Fault-Tolerant Location

n Responsible node is a single point of failure

n Solution: Assign multiple roots per object

n Add *"salt"* to object name and hash as usual

n Salt = globally constant sequence of values (e.g., 1, 2, 3, …)

n Same idea as CAN's multiple realities

n This process makes data more available, even if the network is partitioned

n With $s$ roots, availability is $P \approx 1 - (1/2)^s$

n Depends on partition

n These two mechanisms "guarantee" fault-tolerance

n In most cases :-)

n Problem: If the only out-going link fails…

# Tapestry: Surrogate Routing

n Responsible node is node with same ID as object

    n Such a node is unlikely to exist

n Solution: surrogate routing

n What happens when there is no matching entry in neighbor map for forwarding a message?

n Node picks (deterministically) one entry in neighbor map

    n Details are not explained in the paper :(

n Idea: If "missing links" are deterministically picked, any message for that ID will end up at same node

    n This node is the surrogate

n If new nodes join, surrogate may change

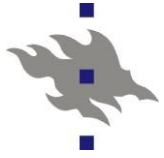    n New node is neighbor of surrogate

# Tapestry: Performance

n **Messages routed in** $O(log_b N)$ **hops**

   n At each step, we resolve one more digit in ID

   n $N$ is the size of the namespace (e.g, SHA-1 = 40 digits)

   n Surrogate routing adds a bit to this, but not significantly

n **State required at a node is** $O(b\ log_b N)$

   n Tapestry has $c$ backup links per neighbor, $O(cb\ log_b N)$

   n Additionally, same number of backpointers

# DHT: Comparison

| | **Chord** | **CAN** | **Tapestry** |
|---|---|---|---|
| Type of network | Ring | N-dimensional | Prefix routing |
| Routing | $O(\log n)$ | $O(d \cdot n^{1/d})$ | $O(\log_b N)$ |
| State | $O(\log n)$ | $O(d)$ | $O(b \cdot \log_b N)$ |
| Caching efficient | + | ++ | ++ |
| Robustness | -/+ | +++ | ++ |
| IP Topology-Aware | N | N/Y | Y |
| Used for other projects | +++ | -- | ++ |

Note: *n* is number of nodes, *N* is size of Tapestry's namespace

# Other DHTs

n Many other DHTs exist too

    n Pastry, similar to Tapestry

    n Kademlia, uses XOR metric

    n Kelips, group nodes into $k$ groups, similar to KaZaA

    n Plus some others…

n Overnet P2P network (also eDonkey) uses Kademlia

    n Wide-spread deployed DHT

n All DHTs provide same API

    n In principle, DHT-layer is interchangeable

# Chapter Summary

**n** Different networks and graphs

  **n** Random, small world, scale-free networks

**n** Searching and addressing

  **n** Fundamental difference

  **n** Unstructured vs. structured networks

**n** Distributed Hash Tables

  **n** DHT provides a key to value mapping

  **n** Three examples: Chord, CAN, Tapestry