

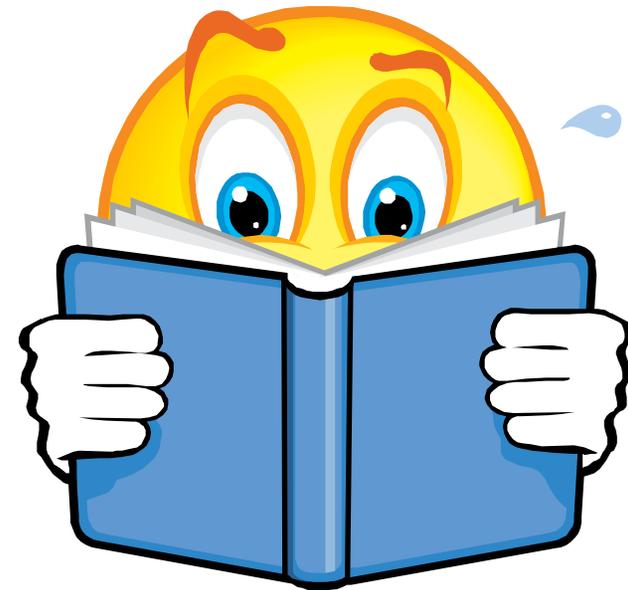


7680: Distributed Systems

Lookup services. Chord. CAN. Pastry. Kademlia.

Required Reading

- ▶ I. Stoica, R. Morris, D. Karger, M. F. Kaashoek, H. Balakrishnan, Chord: A Scalable Peer-to-peer Lookup Service for Internet Applications, SIGCOMM 2001.
- ▶ A Scalable Content-Addressable Network S.a Ratnasamy, P. Francis, M. Handley, R. Karp, S. Shenker, SIGCOMM 2001
- ▶ A. Rowstron and P. Druschel. "Pastry: Scalable, decentralized object location and routing for large-scale peer-to-peer systems". IFIP/ACM International Conference on Distributed Systems Platforms (Middleware), 2001
- ▶ Kademlia: A Peer-to-peer Information System Based on the XOR Metric. P. Maymounkov and D. Mazieres, IPTPS '02





1: Lookup services

Peer-to-Peer (P2P) Systems

- ▶ Applications that take advantage of resources (storage, cycles, content, human presence) available at the edges of the Internet.
- ▶ Characteristics:
 - ▶ System consists of clients connected through Internet and acting as peers
 - ▶ System is designed to work in the presence of variable connectivity
 - ▶ Nodes at the edges of the network have significant autonomy; no centralized control
 - ▶ Nodes are symmetric in function

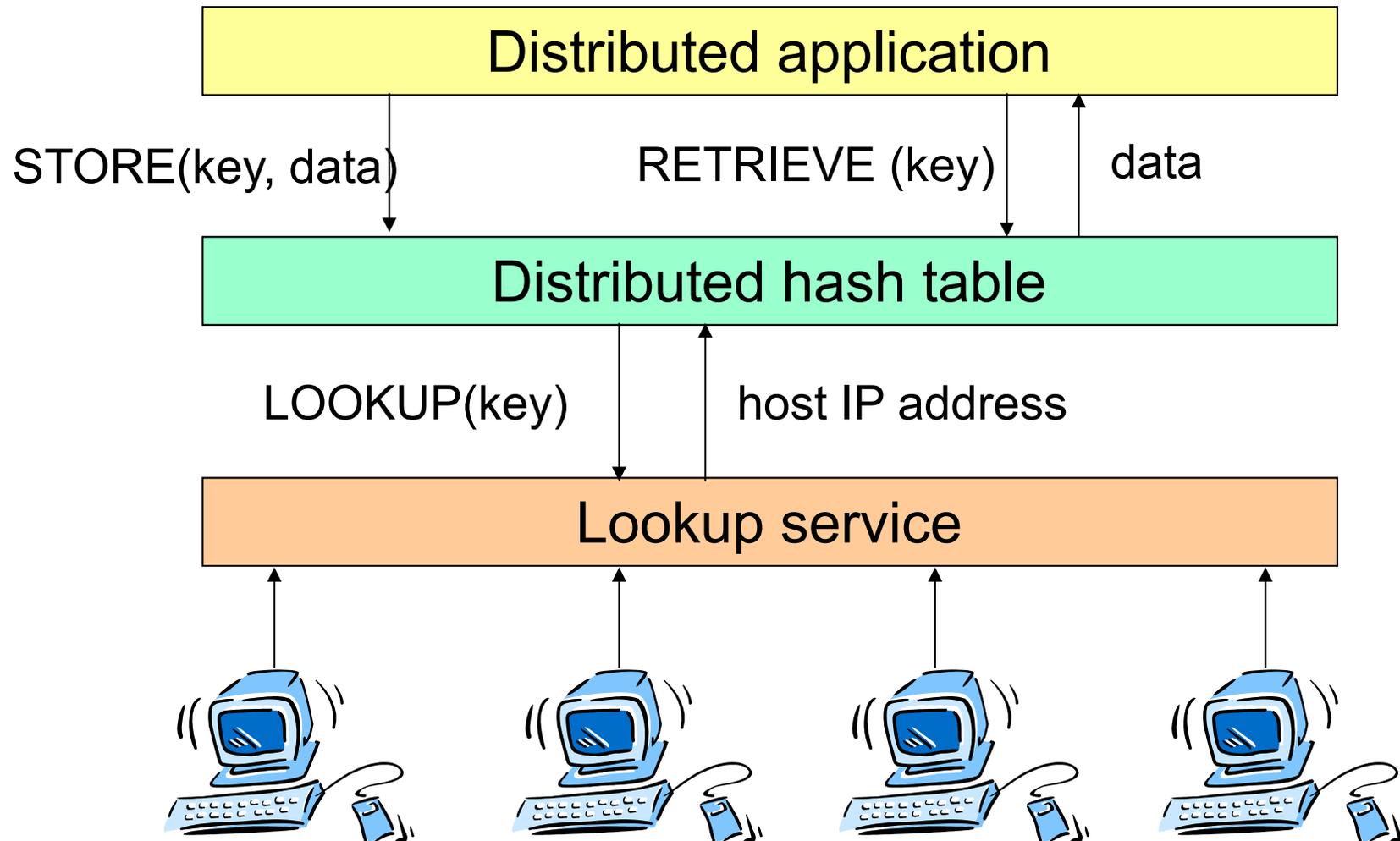
Benefits of P2P and Applications

- ▶ **High capacity:** all clients provide resources (bandwidth, storage space, and computing power). The capacity of the system increases as more nodes become part of the system.
- ▶ **Increased reliability:** achieved by replicating data over multiple peers, and by enabling peers to find the data without relying on a centralized index server.
- ▶ **Applications:**
 - ▶ File sharing: Napster, Gnutella, Freenet, BitTorrent
 - ▶ Distributed file systems: Ivy
 - ▶ Multicast overlays: ESM, NICE, AIML

Issues in P2P Systems Design

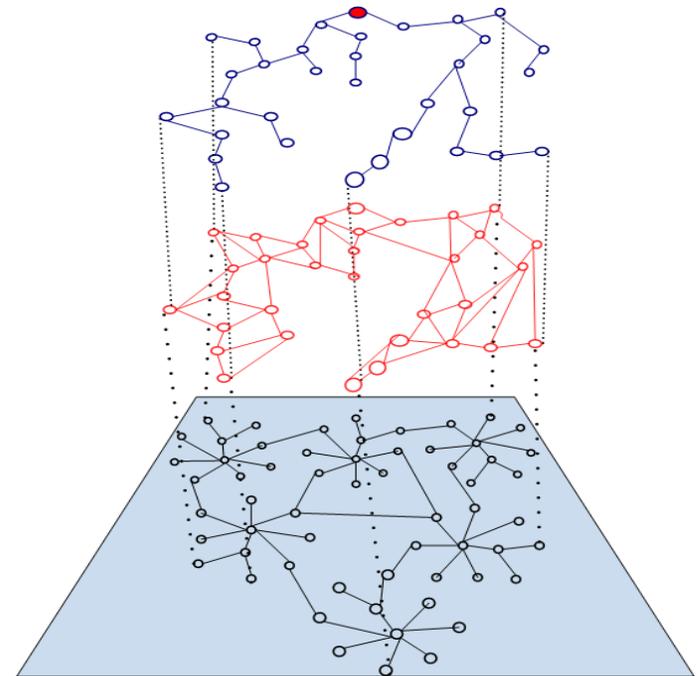
- ▶ How do nodes self-organize, what are appropriate structures?
- ▶ How to search efficiently or perform more complex queries?
- ▶ How to route efficiently on such structures?
- ▶ How to maintain performance in spite of crashes, transient failures?
- ▶ How to maintain availability in spite of failures and partitions?

Structure of P2P File Sharing Systems



Structure of P2P Multicast Systems

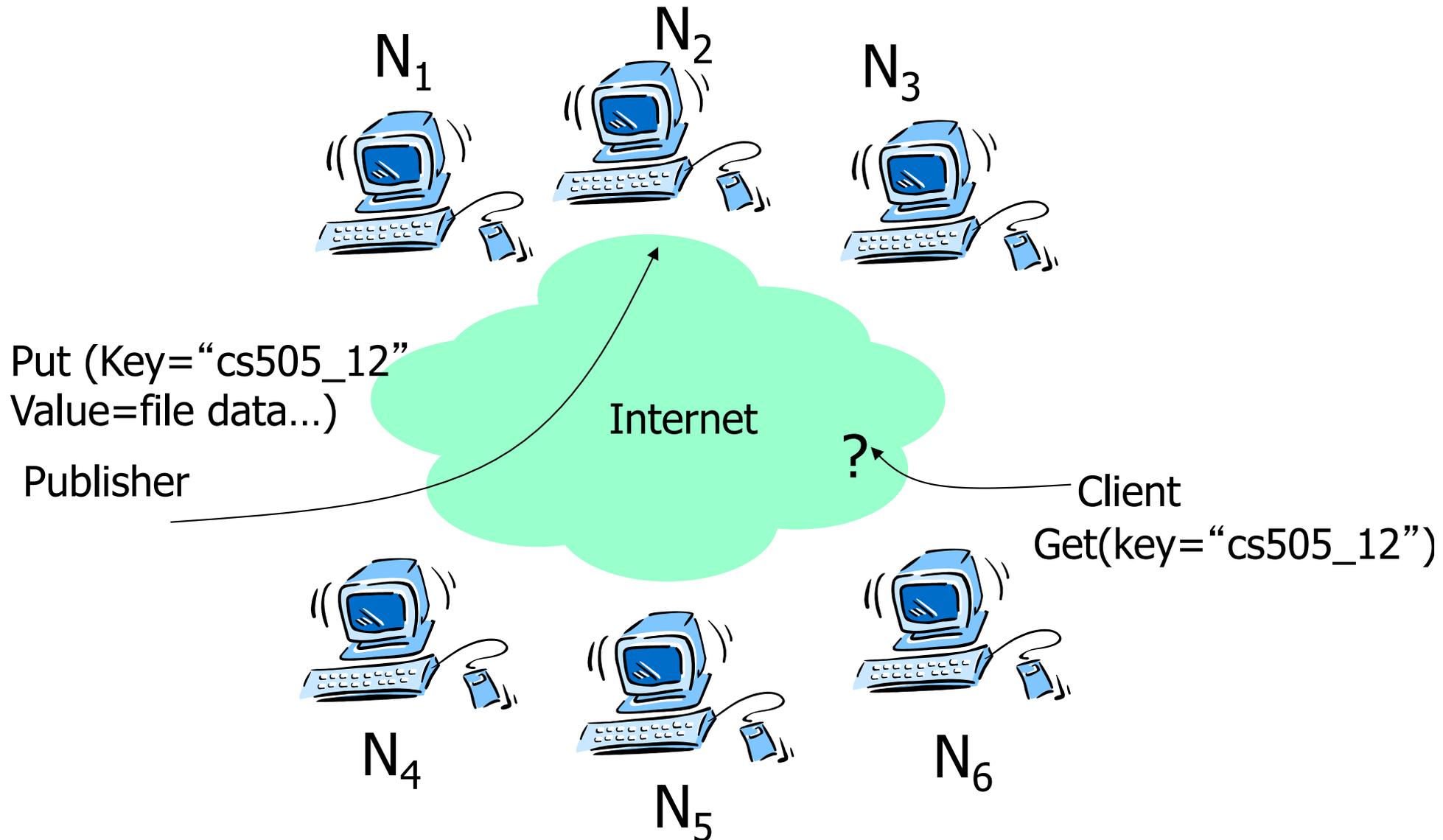
- ▶ Extend network functionality by providing multicast services
- ▶ Usually build a multicast tree that dynamically adapts to improve suboptimal overlay meshes.
- ▶ Overlay is unstructured and optimizations are done by using measurement-based heuristics
- ▶ ESM, Nice, Overcast, ALMI



Structured vs. Unstructured P2P

- ▶ **Many file sharing P2P systems are structured:**
 - ▶ A small subset of nodes meeting presubscribed conditions are eligible to become neighbors
 - ▶ The goal here is to bound the cost of locating objects and the number of network hops
- ▶ **Many multicast/broadcast P2P systems are not structured:**
 - ▶ The goal here is maximizing performance in terms of throughput and latency

Why Lookup Services



Challenges for Lookup Services

- ▶ **Availability**
- ▶ **Scalability**
- ▶ **Complexity**
- ▶ **Exact-match searching vs. approximate matching**
- ▶ **General lookup vs specialized lookup**

Centralized

- ▶ Central index server maintaining the list of files available in the system
- ▶ Upon joining, a node sends the list of files it stores locally, to the central index server
- ▶ When performing a search, a node contacts the central index server to find out the location of the file
- ▶ Vulnerable to single point of failures
- ▶ Maintains $O(N)$ state, costly to maintain the state
- ▶ Example: Napster

Flooded Query

- ▶ There is no centralized index server
- ▶ Each node stores the list of the files it stores locally, no cost on join
- ▶ When performing a search, a node floods the query to every other machine in the network
- ▶ More robust than the centralized approach, avoids the single point of failure
- ▶ Inefficient, worst case $O(N)$ messages per lookup
- ▶ Example: Gnutella

Rooted Query

- ▶ Completely distributed
- ▶ Use a more efficient key-based routing in order to bound the cost of lookup
- ▶ Less robust than flooded query approach, but more efficient
- ▶ Example: Chord, Pastry, Tapestry, Kademlia

Distributed Hash Tables

- ▶ Decentralized distributed systems that partition a set of keys among participating nodes
- ▶ Goal is to efficiently route messages to the unique owner of any given key
- ▶ Typically designed to scale to large numbers of nodes and to handle continual node arrivals and failures
- ▶ Examples: Chord, CAN, Pastry, Tapestry

DHT Design Goals

- ▶ **Decentralized system:**
 - ▶ One node needs to coordinate with a limited set of participants to find the location of a file; should work well in the presence of dynamic membership
- ▶ **Scalability:**
 - ▶ The system should function efficiently even with thousands or millions of nodes
- ▶ **Fault tolerance:**
 - ▶ The system should be reliable even with nodes continuously joining, leaving, and failing

DHT: Keys and Overlays

- ▶ **Key space:**
 - ▶ Ownership of keys is split among the nodes according to some partitioning scheme that maps nodes to keys
- ▶ **Overlay network:**
 - ▶ Nodes self organize in an overlay network; each node maintains a set of links to other nodes (its neighbors or routing table).
 - ▶ Overlay and routing information is used to locate an object based on the associated key

DHT: Storing an Object

- ▶ Compute key according to the object-key mapping method
- ▶ Send message `store(k,data)` to any node participating in the DHT
- ▶ Message is forwarded from node to node through the overlay network until it reaches the node S responsible for key k as specified by the keyspace partitioning method
- ▶ Store the pair $(k,data)$ at node S (sometimes the object is stored at several nodes to deal with node failures)

DHT: Retrieving an Object

- ▶ Compute key according to the object-key mapping method
- ▶ Send a message to any DHT node to find the data associated with k with a message `retrieve(k)`
- ▶ Message is routed through the overlay to the node S responsible for k
- ▶ Retrieve object from node S

Key Partitioning

- ▶ Key partitioning: defines what node “owns what keys” \Leftrightarrow “stores what objects”
- ▶ Removal or addition of nodes should not result in entire remapping of key space since this will result in a high cost in moving the objects around
- ▶ Use consistent hashing to map keys to nodes. A function $d(k_1, k_2)$ defines the distance between keys k_1 to key k_2 . Each node is assigned an identifier (ID). A node with ID i owns all the keys for which i is the closest ID, measured according to distance function d .
- ▶ Consistent hashing has the property that removal or addition of one node changes only the set of keys owned by the nodes with adjacent IDs, and leaves all other nodes unaffected.

Overlay Networks and Routing

- ▶ Nodes self-organize in a logical network defined by the set of links to other nodes each node must maintain
- ▶ Routing:
 - ▶ Greedy algorithm, at each step, forward the message to the neighbor whose ID is closest to k .
 - ▶ When there is no such neighbor, then this is the closest node, which must be the owner of key k



2: Chord

CHORD

- ▶ Efficient lookup of a node which stores data items for a particular search key.
- ▶ Provides only one operation: given a key, it maps the key onto a node.
- ▶ Example applications:
 - ▶ Co-operative mirroring
 - ▶ Time-shared storage
 - ▶ Distributed indexes
 - ▶ Large-scale combinatorial search

Design Goals

- ▶ **Load balance:** distributed hash function, spreading keys evenly over nodes
- ▶ **Decentralization:** CHORD is fully distributed, nodes have symmetric functionality, improves robustness
- ▶ **Scalability:** logarithmic growth of lookup costs with number of nodes in network
- ▶ **Availability:** CHORD guarantees correctness, it automatically adjusts its internal tables to ensure that the node responsible for a key can always be found

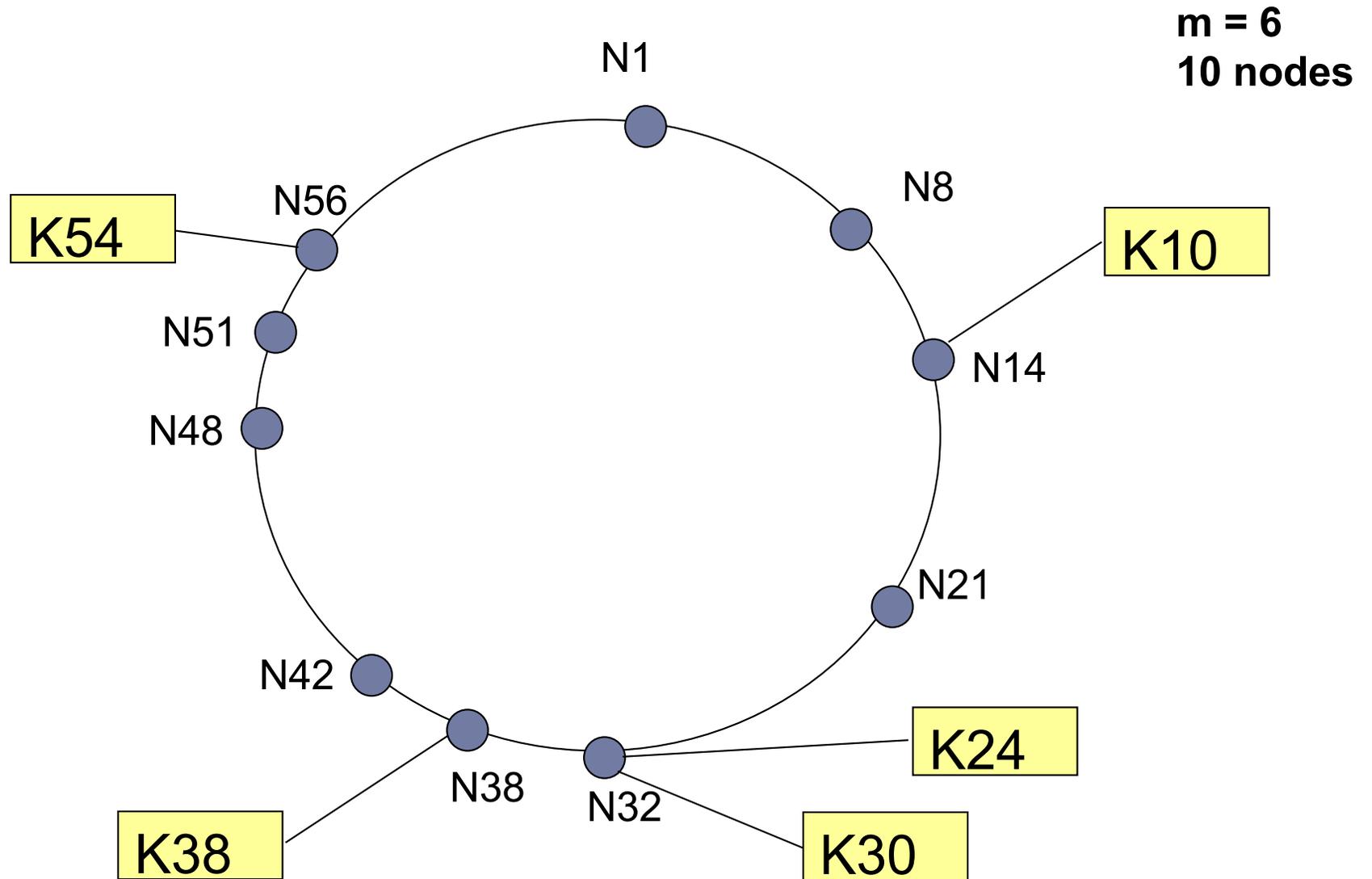
Assumptions

- ▶ Communication in underlying network is both symmetric and transitive
- ▶ Assigns keys to nodes with consistent hashing
- ▶ Hash function balances the load
- ▶ Participants are correct, nodes can join and leave at any time
- ▶ Nodes can fail

Chord Rings

- ▶ Key identifier = SHA-1(key)
- ▶ Node identifier = SHA-1(IP address)
- ▶ Consistent hashing function assigns each node and key an m-bit identifier using SHA-1
- ▶ Mapping key identifiers to node identifiers:
 - ▶ Identifiers are ordered on a circle modulo 2^m called a chord ring.
 - ▶ **The circle is split into contiguous segments whose endpoints are the node identifiers. If i_1 and i_2 are two adjacent IDs, then the node with ID greater identifier i_2 owns all the keys that fall between i_1 and i_2 .**

Example of Key Partitioning in Chord



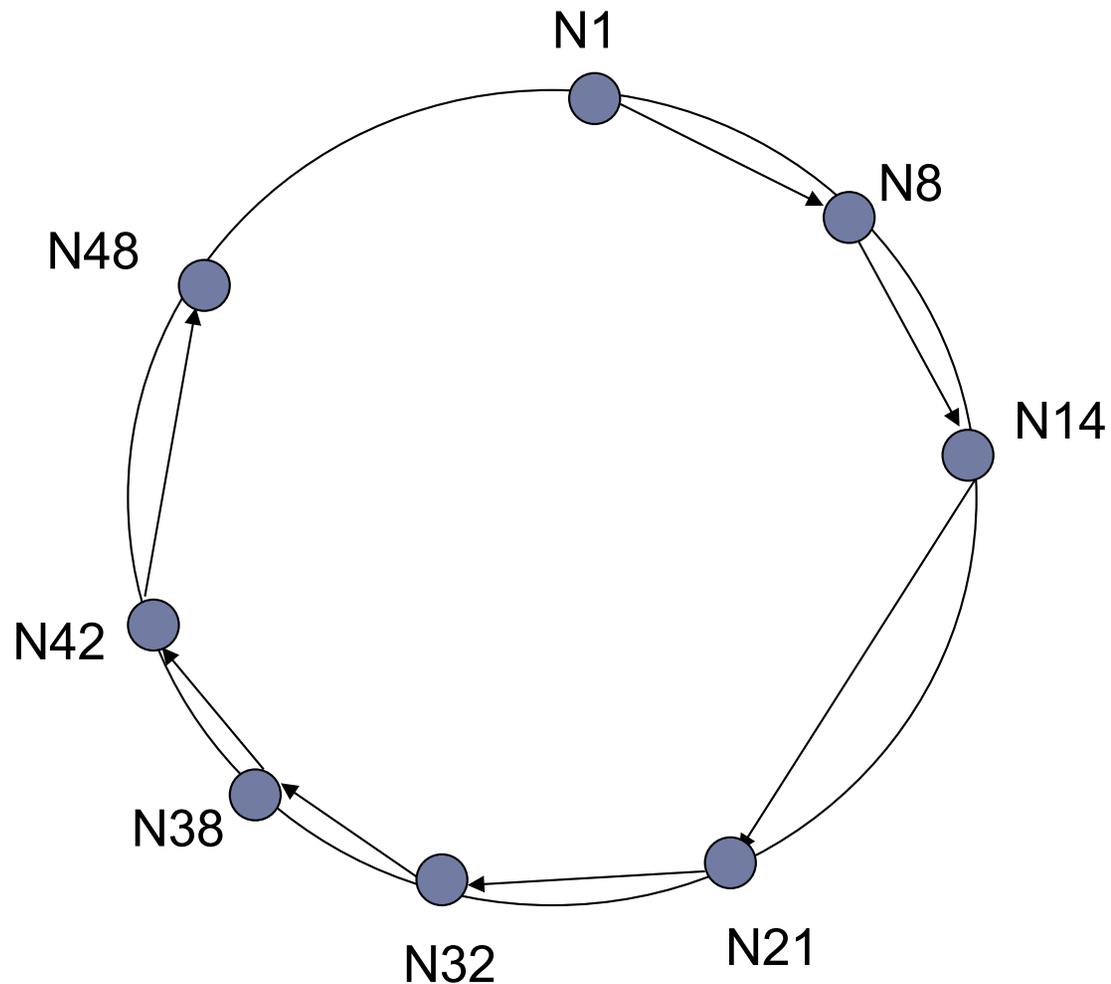
How to Perform Key Lookup

- ▶ Assume that each node knows only how to contact its current successor node on the identifier circle, then all nodes can be visited in linear order.
- ▶ When performing a search, the query for a given identifier could be passed around the circle via these successor pointers until they encounter the node that contains the key corresponding to the search.

Example of Key Lookup Scheme

$successor(k)$ = first node whose ID is \geq ID of k in identifier space

K45

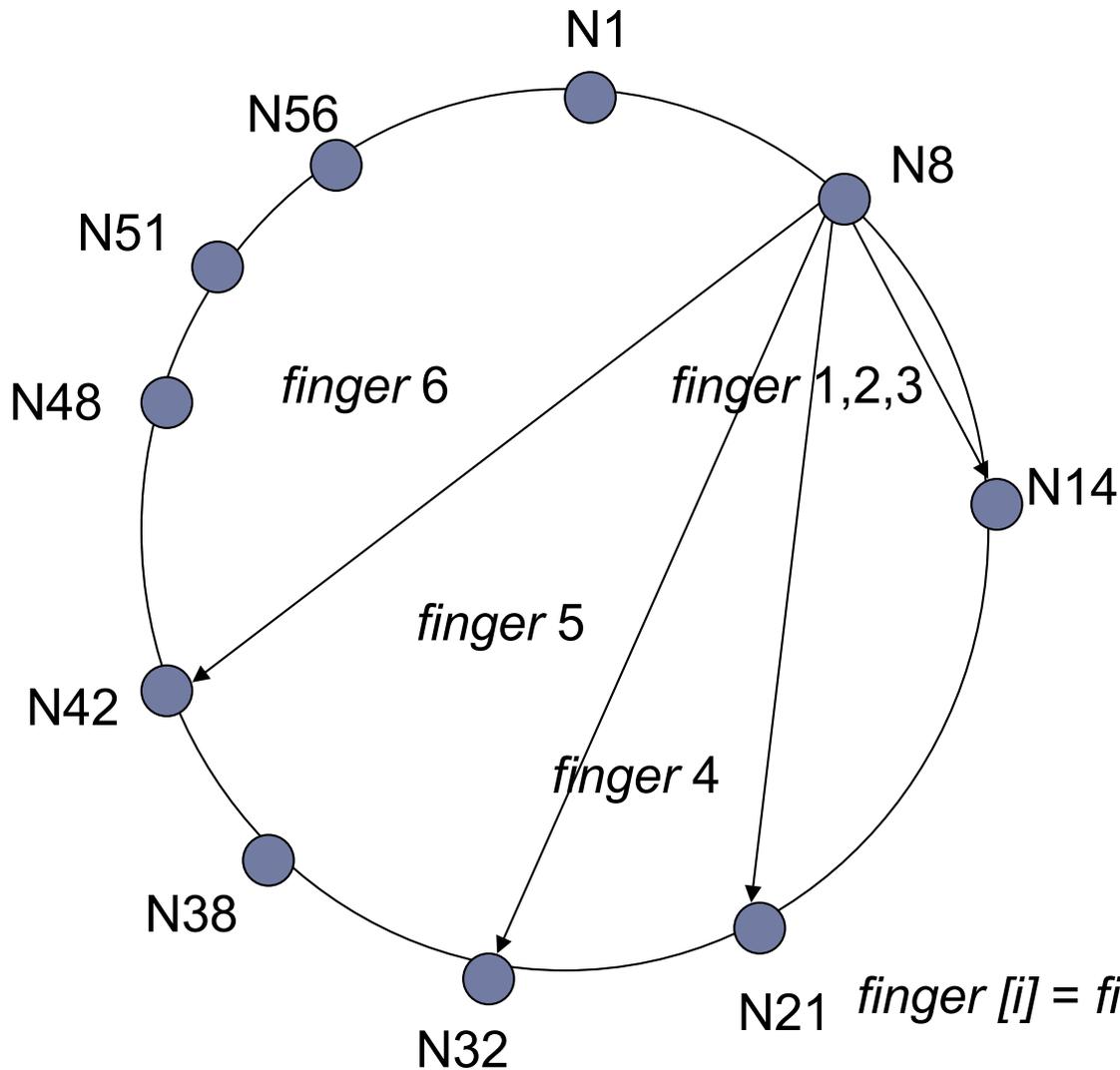


Scalable Key Location

- ▶ To accelerate lookups, Chord maintains additional routing information (m entries): finger table
- ▶ The i^{th} entry in the table at node n contains the identity of the first node s that succeeds n by at least 2^{i-1} on the identifier circle.
- ▶ $s = \text{successor}(n + 2^{i-1})$.
- ▶ s is called the i^{th} finger of node n

Scalable Lookup Scheme

$m = 6$



Finger Table for N8

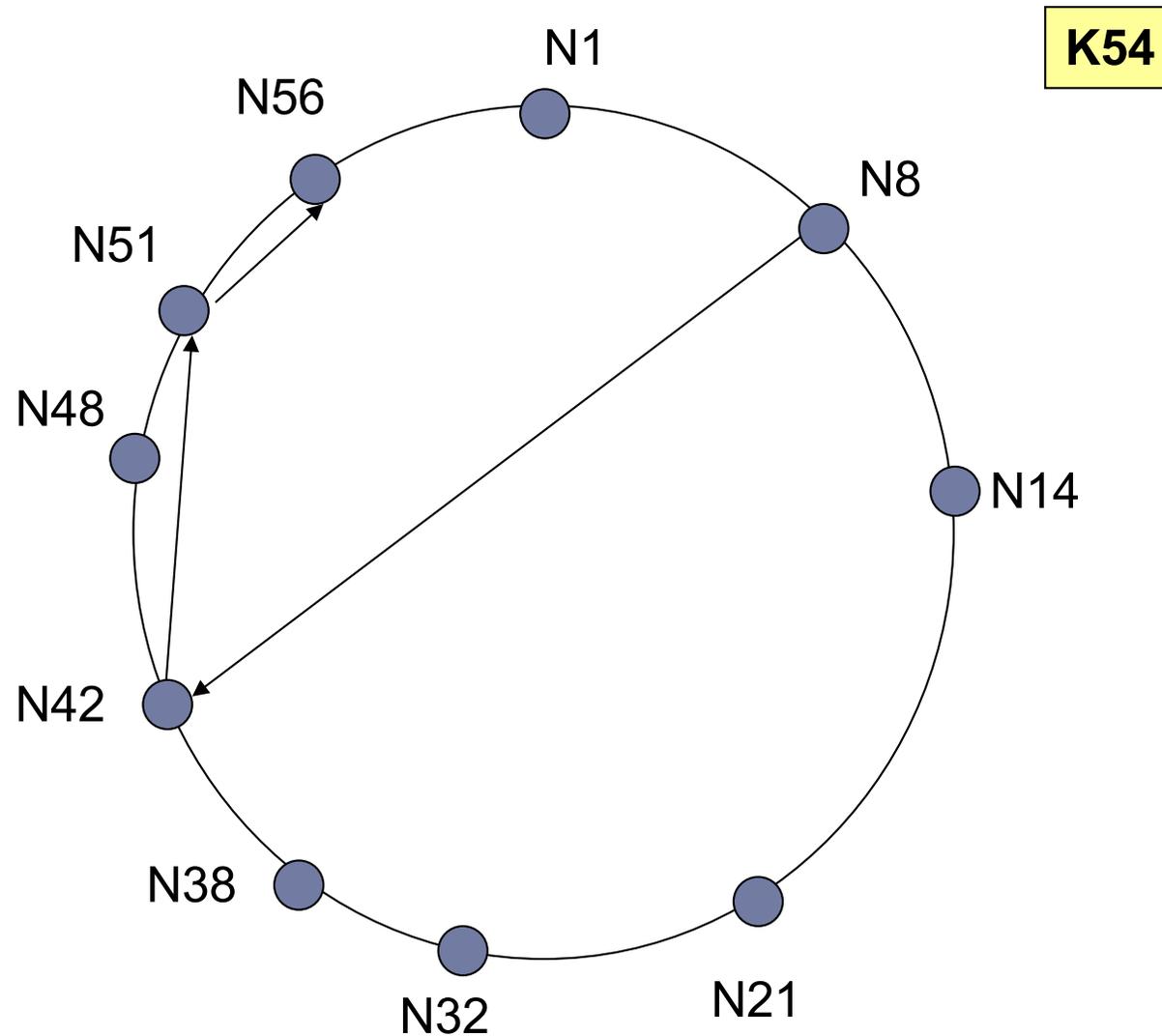
N8+1	N14
N8+2	N14
N8+4	N14
N8+8	N21
N8+16	N32
N8+32	N42

$\text{finger}[i] = \text{first node that succeeds } (n+2^{i-1}) \bmod 2^m$

Scalable Lookup

- ▶ Each node has finger entries at power of two intervals around the identifier circle
- ▶ Each node can forward a query at least halfway along the remaining distance between the node and the target identifier.

Lookup Using Finger Table



Node Joins and Failures/Leaves

- ▶ When a node N joins the network, some of the keys previously assigned to N' 's successor should become assigned to N .
- ▶ When node N leaves the network, all of its assigned keys should be reassigned to N' 's successor.
- ▶ How to deal with these cases?

Node Joins and Stabilizations

- ▶ Everything relies on successor pointer.
- ▶ Up to date successor pointer is sufficient to guarantee correctness of lookups
- ▶ Idea: run a “stabilization” protocol periodically in the background to update successor pointer and finger table.

Stabilization Protocol

- ▶ Guarantees to add nodes in a fashion to preserve reachability
- ▶ Does not address the cases when a Chord system has split into multiple disjoint cycles, or a single cycle that loops multiple times around the identifier space

Stabilization Protocol (cont.)

- ▶ Each time node N runs stabilize protocol, it asks its successor for its predecessor p , and decides whether p should be N 's successor instead.
- ▶ Stabilize protocol notifies node N 's successor of N 's existence, giving the successor the chance to change its predecessor to N .
- ▶ The successor does this only if it knows of no closer predecessor than N .

Impact of Node Joins on Lookups

- ▶ If finger table entries are current then lookup finds the correct successor in $O(\log N)$ steps
- ▶ If successor pointers are correct but finger tables are incorrect, correct lookup but slower
- ▶ If incorrect successor pointers, then lookup may fail

Voluntary Node Departures

- ▶ Leaving node may transfer all its keys to its successor
- ▶ Leaving node may notify its predecessor and successor about each other so that they can update their links

Node Failures

- ▶ **Stabilize successor lists:**

- ▶ Node N reconciles its list with its successor S by copying S 's successor list, removing its last entry, and prepending S to it.
- ▶ If node N notices that its successor has failed, it replaces it with the first live entry in its successor list and reconciles its successor list with its new successor.

CHORD Summary

- ▶ Efficient location of the node that stores a desired data item is a fundamental problem in P2P networks
- ▶ Separates correctness (successor) from performance (finger table)
- ▶ Chord protocol solves it in an efficient decentralized manner
 - ▶ Routing information: $O(\log N)$ nodes
 - ▶ Lookup: $O(\log N)$ nodes
 - ▶ Update: $O(\log^2 N)$ messages
- ▶ It also adapts dynamically to the topology changes introduced during the run



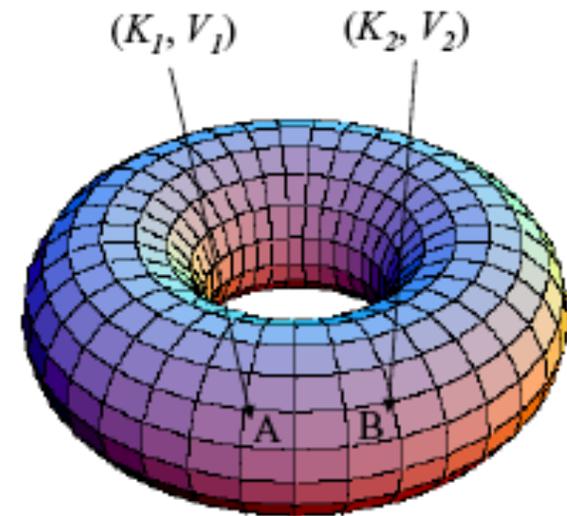
3: CAN

Content Addressable Network

- ▶ Each node stores a chunk (zone) of the hash table
- ▶ Each node stores state information about neighbor zones
- ▶ Requests (insert, lookup, or delete) for a key are routed by intermediate nodes using a greedy routing algorithm
- ▶ Completely distributed
- ▶ Nodes can crash
- ▶ Keys mapped to a d dimensional space

Design of CAN

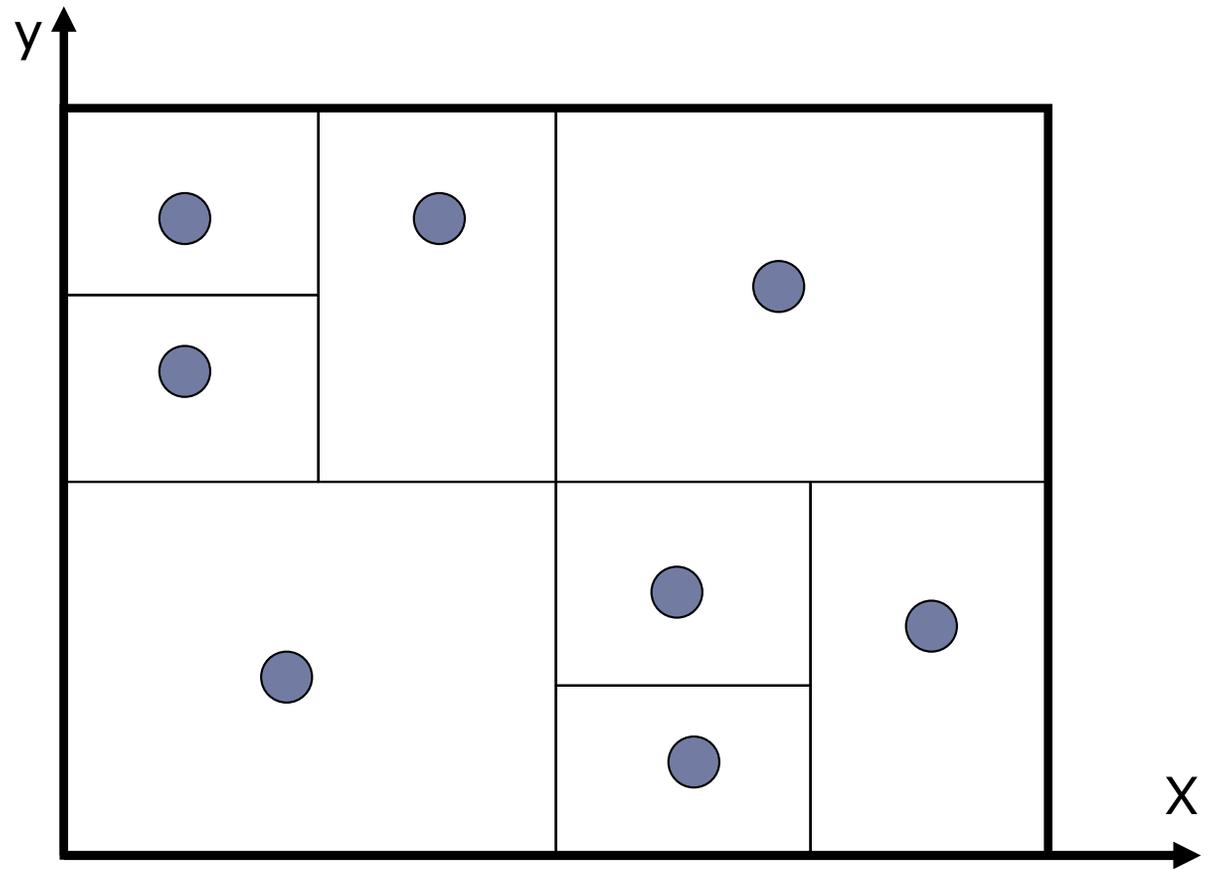
- ▶ d-dimensional Cartesian coordinate space (d-torus)
- ▶ Each node owns a zone on the torus
- ▶ To store key value pair (K_I, V_I) ,
 - ▶ K_I mapped to point P_I using uniform
 - ▶ (K_I, V_I) stored at the node N that owns P_I



2-torus

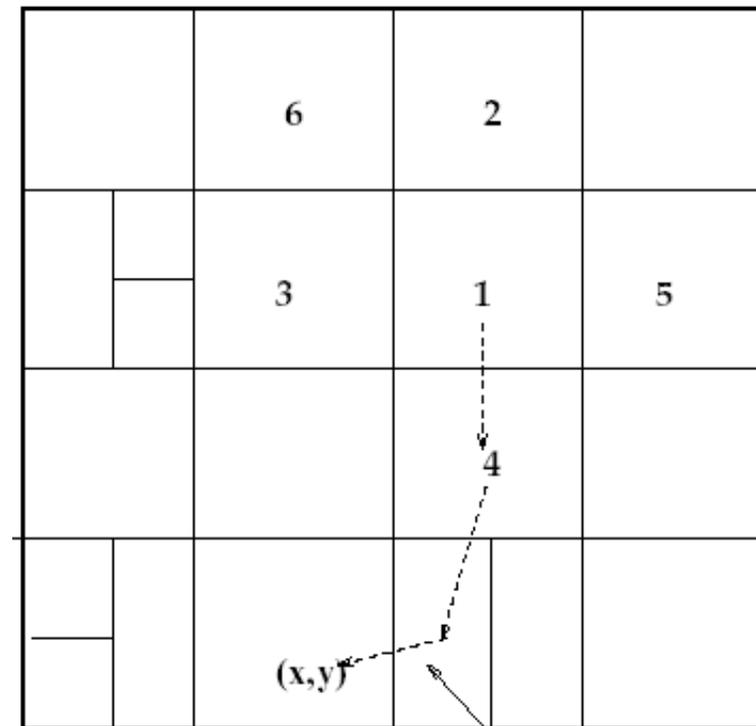
Key Partitioning

- ▶ Virtual d-dimensional Cartesian coordinate dynamically partitioned among all nodes
- ▶ Key partitioning using a uniform hash function



Routing in CAN

- ▶ Each node maintains a table of the IP address and virtual coordinate zone of each local neighbor
- ▶ Use greedy routing to neighbor closest to destination

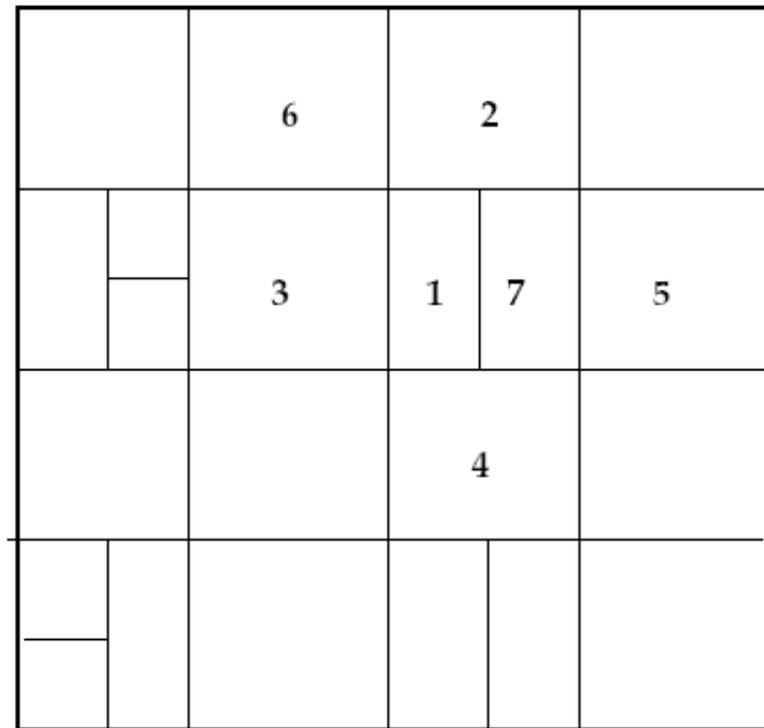


sample routing path from node 1 to point (x,y)

1's coordinate neighbor set = {2,3,4,5}
7's coordinate neighbor set = { }

Node Join

- ▶ Joining node locates a bootstrap node B (random member) using the CAN DNS entry
- ▶ Node in zone containing B splits the zone and allocates “half” to joining node
- ▶ Keys and objects are transferred to new node
- ▶ Joining node and previous node zone update their neighbor set



1's coordinate neighbor set = {2,3,4,7}

7's coordinate neighbor set = {1,2,4,5}

Node Voluntarily Leave and Failure

- ▶ Leave: node transfers its zone and the pairs of keys and data to a neighbor
- ▶ Failure: unreachable node(s) trigger an immediate takeover algorithm that allocate failed node's zone to a neighbor
- ▶ Multiple concentrated failures: neighbors may need to perform an expanding ring search to build sufficient neighbor state prior to initiating the takeover algorithm
- ▶ Background zone-reassignment algorithm to prevent space fragmentation
 - ▶ Over time nodes may takeover zones that cannot be merged with its own zone



4: Pastry

Pastry

- ▶ Each node in Pastry has a unique, uniform random identifier (nodeID) in a circular 128-bit identifier space (as in Chord)
- ▶ Each object has a unique, uniform random identifier (objectID)
- ▶ Node with numerically closest nodeID maintains the object
- ▶ Routing is done to the numerically closest to the key that is searched
- ▶ Network of N nodes
 - ▶ Expected number of forwarding steps in the Pastry overlay network is $O(\log N)$
 - ▶ Size of the routing table maintained at each node is $O(\log N)$

Routing

- ▶ Address blocks: 128-bit key is divided up into digits, each digit being b bits long
 - ▶ level 0 represents a zero-digit common prefix between two addresses, level 1 a one-digit common prefix, and so on.
- ▶ Leaf set: IP addresses of the nodes with the $L/2$ numerically closest peers by nodeID in each direction around the circle
- ▶ Neighbor list: M closest peers with respect to the routing metric.
- ▶ Routing table: address of the closest known peer for each possible digit at each address level, except for the digit that belongs to the peer itself at that particular level.
- ▶ Storage of $2^b - 1$ contacts per level, number of levels scaling as $(\log 2N) / b$, $b=4$, $L = 2^b$ and $M=2^b$ are typical settings

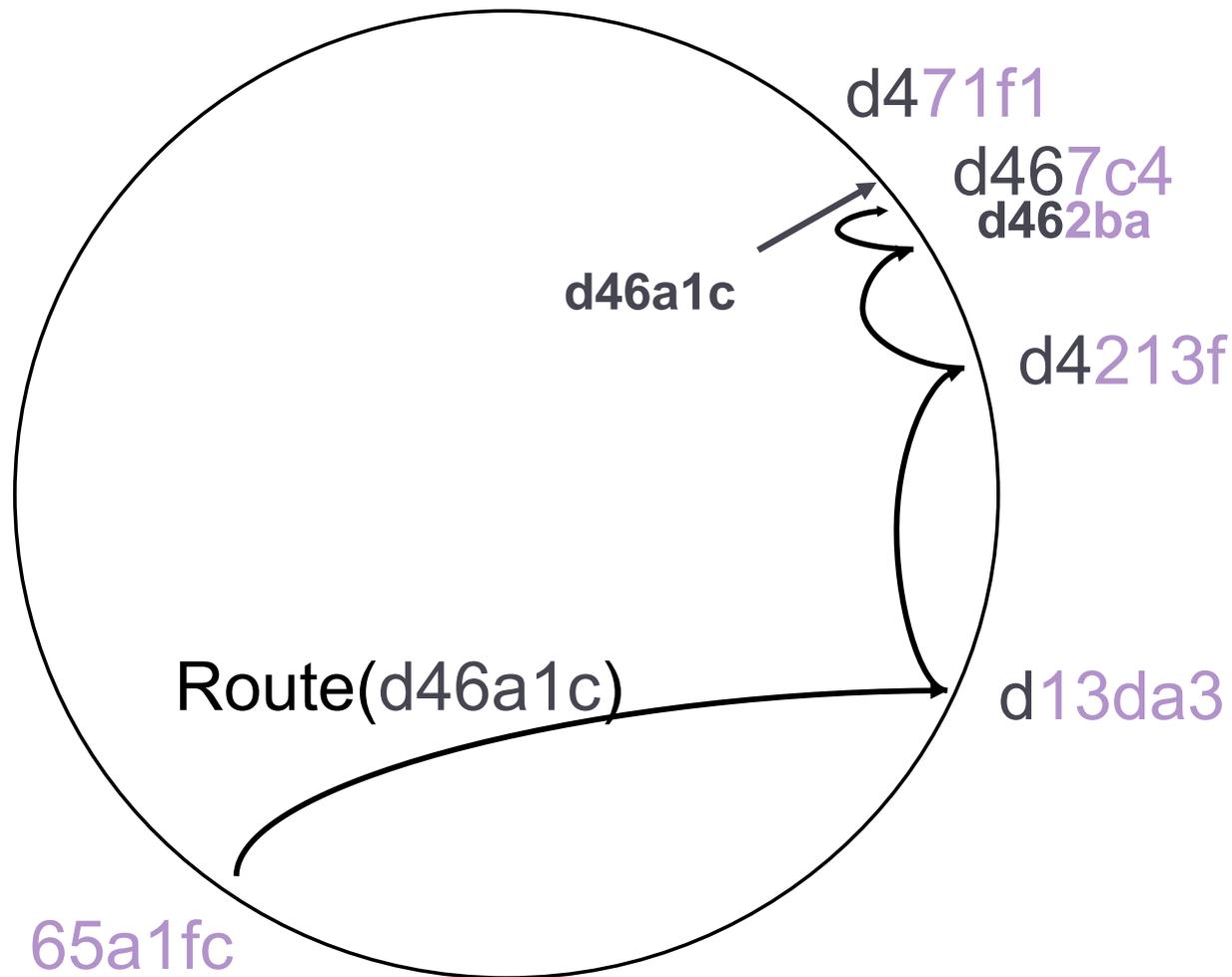
Routing Table of Node 65a1fcx

Row 0	<i>0</i>	<i>1</i>	<i>2</i>	<i>3</i>	<i>4</i>	<i>5</i>		<i>7</i>	<i>8</i>	<i>9</i>	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
	<i>x</i>	<i>x</i>	<i>x</i>	<i>x</i>	<i>x</i>	<i>x</i>		<i>x</i>								
Row 1	<i>6</i>	<i>6</i>	<i>6</i>	<i>6</i>	<i>6</i>		<i>6</i>									
	<i>0</i>	<i>1</i>	<i>2</i>	<i>3</i>	<i>4</i>		<i>6</i>	<i>7</i>	<i>8</i>	<i>9</i>	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
	<i>x</i>	<i>x</i>	<i>x</i>	<i>x</i>	<i>x</i>		<i>x</i>									
Row 2	<i>6</i>		<i>6</i>	<i>6</i>	<i>6</i>	<i>6</i>	<i>6</i>									
	<i>5</i>		<i>5</i>	<i>5</i>	<i>5</i>	<i>5</i>	<i>5</i>									
	<i>0</i>	<i>1</i>	<i>2</i>	<i>3</i>	<i>4</i>	<i>5</i>	<i>6</i>	<i>7</i>	<i>8</i>	<i>9</i>		<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
	<i>x</i>		<i>x</i>	<i>x</i>	<i>x</i>	<i>x</i>	<i>x</i>									
Row 3	<i>6</i>		<i>6</i>													
	<i>5</i>		<i>5</i>													
	<i>a</i>		<i>a</i>													
	<i>0</i>		<i>2</i>	<i>3</i>	<i>4</i>	<i>5</i>	<i>6</i>	<i>7</i>	<i>8</i>	<i>9</i>	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
	<i>x</i>		<i>x</i>													

$\log_{16} N$
rows

DHTs

Routing intuition



Properties

- ▶ $\log_{16} N$ steps
- ▶ $O(\log N)$ state

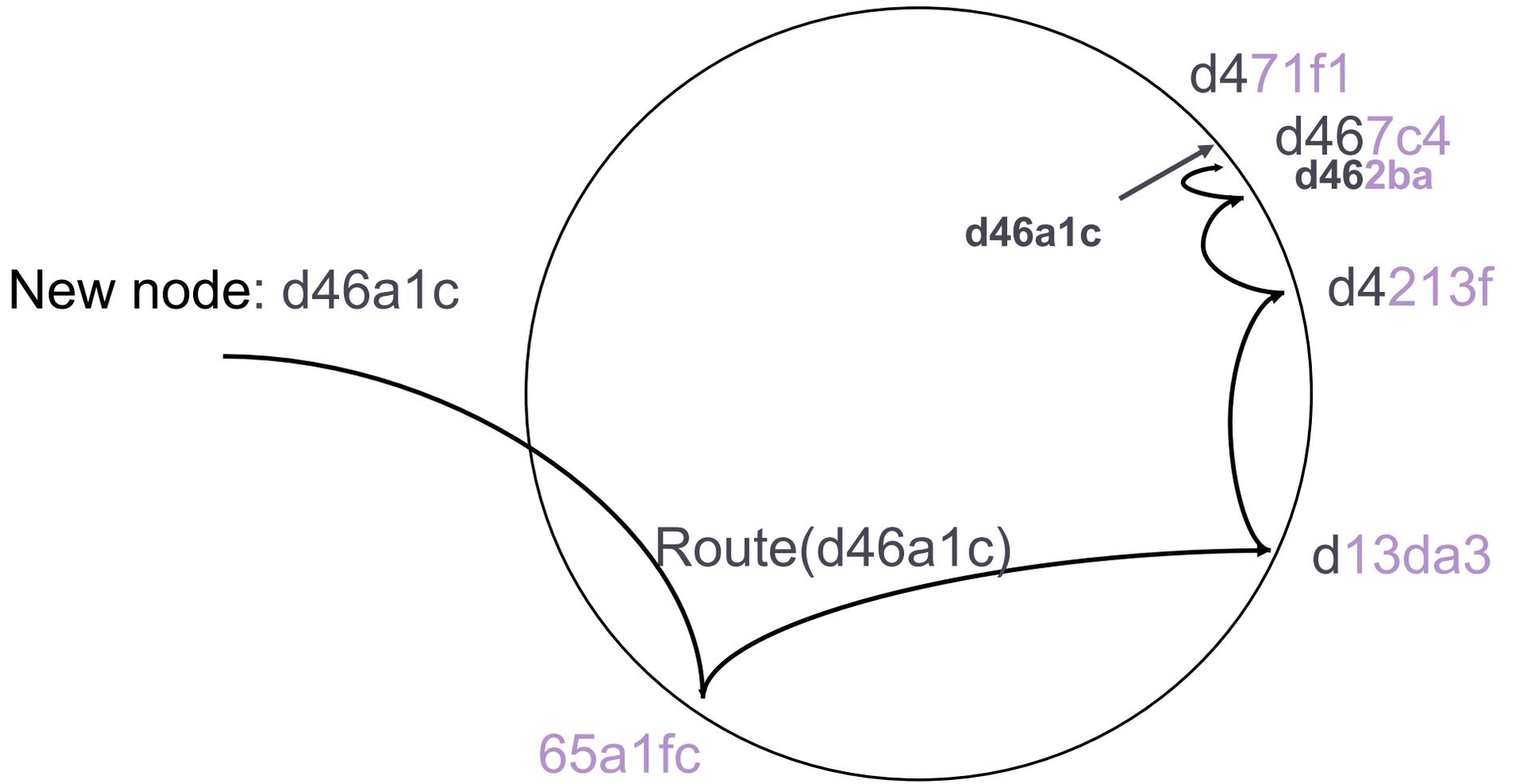
Mapping Objects to Nodes

- ▶ Objects are assigned unique, uniform random identifiers (objIDs) and mapped to the k nodes with nodeIDs numerically closest to the objID.
- ▶ Inserting objects : the message reaches a node with one of the k closest nodeIDs to the objID, that node replicates the object among the other $k-1$ nodes with closest nodeIDs (which are, by definition, in the same leaf set for $k \leq L/2$).
- ▶ Searching for objects: Application-specific objects can be looked up, contacted, or retrieved by routing a Pastry message, using the objID as the key. By definition, the message is guaranteed to reach a node that maintains a replica of the requested object unless all k nodes with nodeIDs closest to the objID have failed.

Routing

- ▶ A peer first examines its leaf set and routes directly to the correct node if one is found.
- ▶ If this fails, the peer checks its routing table to find the address of a node which shares a longer prefix with the destination address than the peer itself.
- ▶ If the peer does not have any contacts with a longer prefix or the contact has died it will pick a peer from its contact list with the same length prefix whose node ID is numerically closer to the destination and send the packet to that peer.

Node Addition



Node departure (failure)

- ▶ Leaf set members exchange keep-alive messages
- ▶ Leaf set repair (eager): request set from farthest live node in set
- ▶ Routing table repair (lazy): get table from peers in the same row, then higher rows

From www.cs.rice.edu/~druschel/comp413/lectures/Pastry.ppt



5: Kademia

Kademlia in a Nutshell

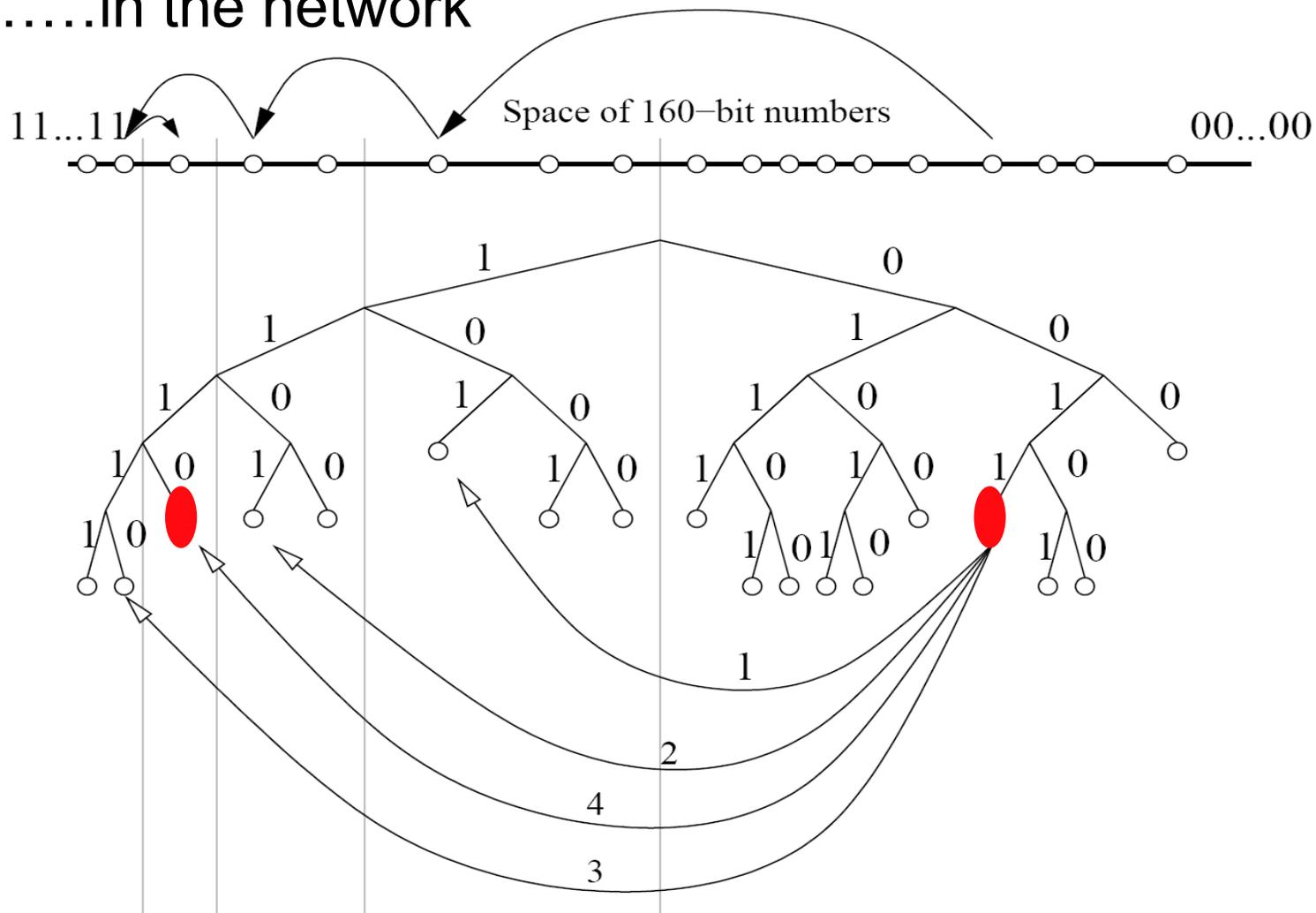
- ▶ Similar with other services, IDs based on SHA-1 hash into a 160 bits space.
- ▶ Closeness between two objects measured as their bitwise XOR interpreted as an integer.
- ▶ $\text{distance}(a, b) = a \text{ XOR } b$
- ▶ Distance is symmetric, $\text{dist}(a, b) = \text{dist}(b, a)$
- ▶ Uses parallel asynchronous queries to avoid timeout delays of the failed nodes. Routes are selected based on latency
- ▶ Kademlia uses tree-based routing

Kademlia Binary Tree

- ▶ Start from root, for any given node, dividing the binary tree into a series of successively lower subtrees that do not contain the node and each correspond to a k-bucket
- ▶ Every node keeps track of at least one node from each of its subtrees.
- ▶ Every node keeps a list of (IP, Port, NodeID) triples, and (key, value) tuples for further exchanging information with others.

Kademlia Search

An example of lookup: node 0011 is searching for 1110.....in the network



Kademlia Lookup

- ▶ Locate the k closest nodes to a given `nodeID`.
- ▶ Uses a recursive algorithm for node lookups.
 - ▶ The lookup initiator starts by picking a node from its closest non-empty k -bucket.
 - ▶ The initiator then sends parallel, asynchronous `FIND_NODE` to the α nodes it has chosen.
 - ▶ The initiator resends the `FIND_NODE` to nodes it has learned about from previous requests.
 - ▶ If a round of `FIND_NODES` fails to return a node any closer than the closest already seen, the initiator resends the `FIND_NODE` to all of the k closest nodes it has not already queried.

Kademlia Keys Store

- ▶ To store a (key,value) pair, a participant locates the k closest nodes to the key.
- ▶ Additionally, each node re-publishes (key,value) pairs as necessary to keep them alive
- ▶ Kademlia's current application (file sharing), requires the original publisher of a (key,value) pair to republish it every 24 hours. Otherwise, (key,value) pairs expire 24 hours after publication.

Kademlia Cost

- ▶ **Operation cost**
 - ▶ As low as other popular protocols
 - ▶ Look up, $O(\log N)$
 - ▶ Join or leave, $O(\log^2 N)$
- ▶ **Fault tolerance and concurrent change**
 - ▶ Handles well, for the use of k-buckets
- ▶ **Proximity routing**
 - ▶ Kademlia can choose from α nodes that has lower latency