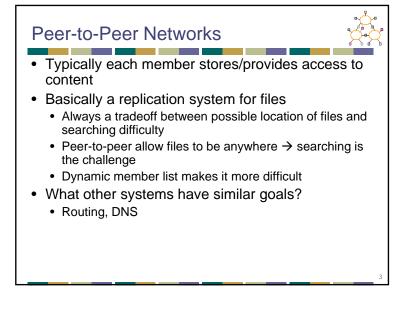
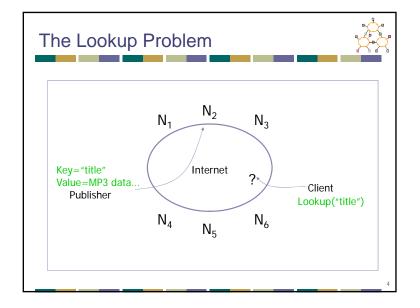
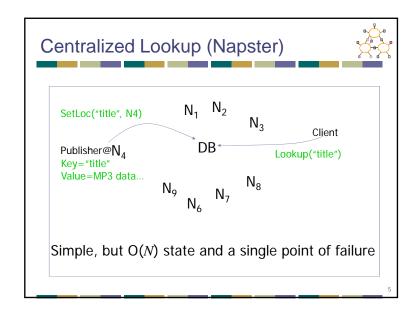


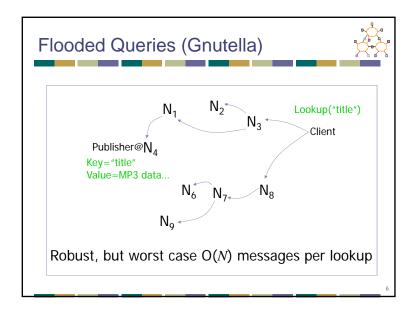
### P2P LookupOverviewNapsterGnutella

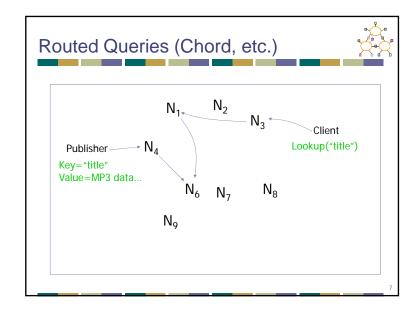
- Routed Lookups Chord
- Comparison of DHTs
- I3











### Centralized: Napster



- Simple centralized scheme → motivated by ability to sell/control
- How to find a file:
  - On startup, client contacts central server and reports list of files
  - Query the index system → return a machine that stores the required file
    - Ideally this is the closest/least-loaded machine
  - Fetch the file directly from peer

### Centralized: Napster



- Advantages:
  - Simple
  - Easy to implement sophisticated search engines on top of the index system
- Disadvantages:
  - Robustness, scalability
  - Easy to sue!

### Flooding: Old Gnutella



- On startup, client contacts any servent (server + client) in network
  - Servent interconnection used to forward control (queries, hits, etc)
- Idea: broadcast the request
- How to find a file:
  - Send request to all neighbors
  - Neighbors recursively forward the request
  - Eventually a machine that has the file receives the request, and it sends back the answer
  - Transfers are done with HTTP between peers

10

### Flooding: Old Gnutella



- Advantages:
  - Totally decentralized, highly robust
- Disadvantages:
  - Not scalable; the entire network can be swamped with request (to alleviate this problem, each request has a TTL)
  - Especially hard on slow clients
    - At some point broadcast traffic on Gnutella exceeded 56kbps – what happened?
    - Modem users were effectively cut off!

### Flooding: Old Gnutella Details



- Basic message header
  - Unique ID, TTL, Hops
- Message types
  - Ping probes network for other servents
  - Pong response to ping, contains IP addr, # of files, # of Kbytes shared
  - Query search criteria + speed requirement of servent
  - QueryHit successful response to Query, contains addr + port to transfer from, speed of servent, number of hits, hit results, servent ID
  - Push request to servent ID to initiate connection, used to traverse firewalls
- · Ping, Queries are flooded
- QueryHit, Pong, Push reverse path of previous message

# Flooding: Old Gnutella Example Assume: m1's neighbors are m2 and m3; m3's neighbors are m4 and m5;...

### Flooding: Gnutella, Kazaa



- Modifies the Gnutella protocol into two-level hierarchy
  - · Hybrid of Gnutella and Napster
- Supernodes
  - · Nodes that have better connection to Internet
  - · Act as temporary indexing servers for other nodes
  - · Help improve the stability of the network
- Standard nodes
  - · Connect to supernodes and report list of files
  - · Allows slower nodes to participate
- Search
  - · Broadcast (Gnutella-style) search across supernodes
- Disadvantages
  - Kept a centralized registration → allowed for law suits ☺

14

### Overview



- P2P Lookup
  - Overview
  - Napster
  - Gnutella
- Routed Lookups Chord
- Comparison of DHTs
- I3

### **Routing: Structured Approaches**



- Goal: make sure that an item (file) identified is always found in a reasonable # of steps
- Abstraction: a distributed hash-table (DHT) data structure
  - insert(id, item);
  - item = query(id);
  - Note: item can be anything: a data object, document, file, pointer to a file...
- Proposals
  - CAN (ICIR/Berkeley)
  - Chord (MIT/Berkeley)
  - Pastry (Rice)
  - Tapestry (Berkeley)
  - ...

### Routing: Chord



- Associate to each node and item a unique id in an uni-dimensional space
- 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

### Aside: Hashing



- Advantages
  - Let nodes be numbered 1..m
  - Client uses a good hash function to map a URL to 1..m
  - Say hash (url) = x, so, client fetches content from node
  - No duplication not being fault tolerant.
  - · One hop access
  - · Any problems?
    - · What happens if a node goes down?
    - What happens if a node comes back up?
    - · What if different nodes have different views?

### Robust hashing



- Let 90 documents, node 1..9, node 10 which was dead is alive again
- % of documents in the wrong node?
  - 10, 19-20, 28-30, 37-40, 46-50, 55-60, 64-70, 73-80, 82-90
  - Disruption coefficient = 1/2
  - Unacceptable, use consistent hashing idea behind Akamai!

### **Consistent Hash**



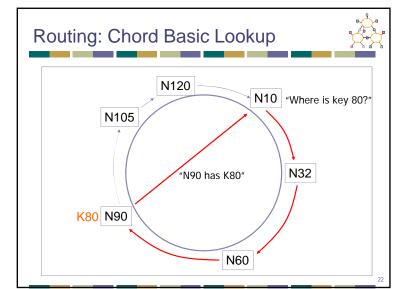
- "view" = subset of all hash buckets that are visible
- Desired features
  - Balanced in any one view, load is equal across buckets
  - Smoothness little impact on hash bucket contents when buckets are added/removed
  - Spread small set of hash buckets that may hold an object regardless of views
  - Load across all views # of objects assigned to hash bucket is small

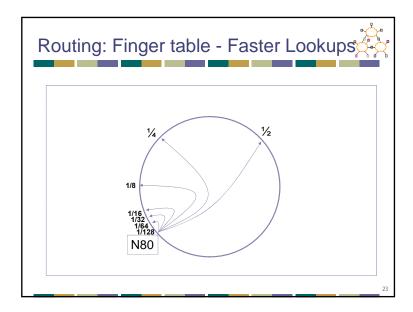
0!

### Consistent Hash - Example



- Construction
  - Assign each of C hash buckets to random points on mod 2<sup>n</sup> circle, where, hash key size = n.
  - Map object to random position on circle
  - Hash of object = closest clockwise bucket
- Smoothness → addition of bucket does not cause much movement between existing buckets
- Spread & Load → small set of buckets that lie near object
- Balance → no bucket is responsible for large number of objects

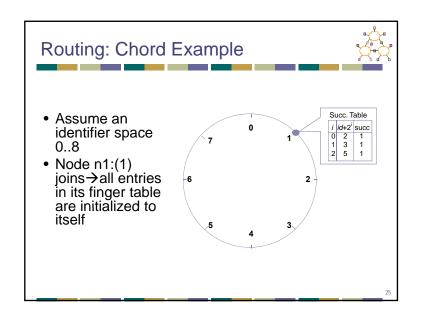


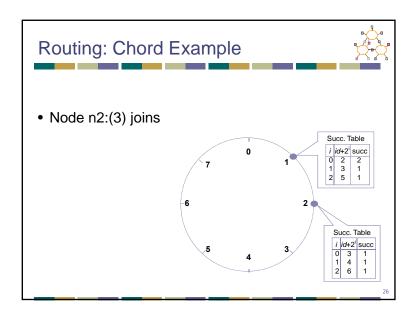


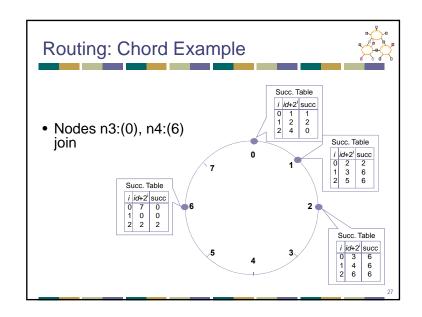
### Routing: Chord Summary

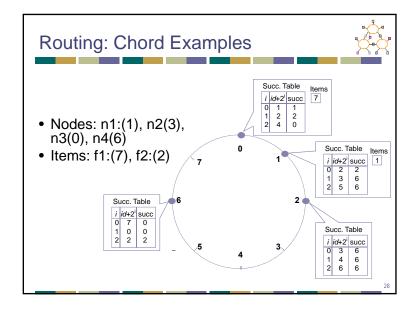


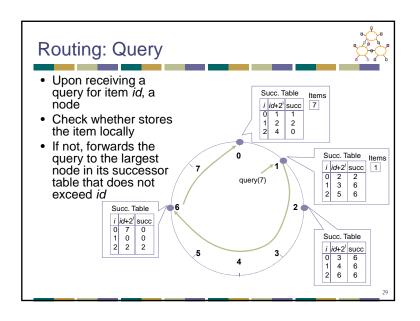
- Assume identifier space is 0...2<sup>m</sup>
- Each node maintains
  - Finger table
    - Entry *i* in the finger table of *n* is the first node that succeeds or equals  $n + 2^{i}$
  - Predecessor node
- An item identified by id is stored on the successor node of id















- Distributed object lookup
  - · Based on object ID
- De-centralized file systems
  - CFS, PAST, Ivy
- Application Layer Multicast
  - Scribe, Bayeux, Splitstream
- Databases
  - PIER

Overview



- P2P Lookup Overview
- Centralized/Flooded Lookups
- Routed Lookups Chord
- Comparison of DHTs
  - Slides Michael Stroucken
- I3

Comparison



- Many proposals for DHTs
  - Tapestry (UCB)
- ...
  - -- Symphony (Stanford) -- 1hop (MIT)
  - · Pastry (MSR, Rice)
- -- Tangle (UCB)
- -- conChord (MIT)

- Chord (MIT, UCB)
- -- SkipNet (MSR,UW)
- -- Apocrypha (Stanford)

- CAN (UCB, ICSI)Viceroy (Technion)
- -- Bamboo (UCB)
  -- Hieras (U.Cinn)
- -- LAND (Hebrew Univ.)
  -- ODRI (TexasA&M)

- Kademlia (NYU)
- -- Sprout (Stanford)
- Kelips (Cornell)
- -- Calot (Rochester)
- Koorde (MIT)
- -- JXTA's (Sun)
- What are the right design choices? Effect on performance?

### **Deconstructing DHTs**



### Two observations:

- 1. Common approach
  - N nodes; each labeled with a virtual identifier (128 bits)
  - · define "distance" function on the identifiers
  - routing works to reduce the distance to the destination
- 2. DHTs differ primarily in their definition of "distance"
  - typically derived from (loose) notion of a routing geometry

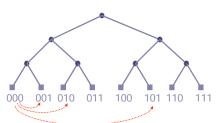
### **DHT** Routing Geometries



- · Geometries:
  - Tree (Plaxton, Tapestry)
  - Ring (Chord)
  - Hypercube (CAN)
  - XOR (Kademlia)
  - Hybrid (Pastry)
- What is the impact of geometry on routing?

Tree (Plaxton, Tapestry)





### Geometry

- nodes are leaves in a binary tree
- distance = height of the smallest common subtree
- logN neighbors in subtrees at distance 1,2,...,logN

0.5

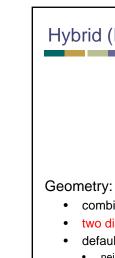
### Hypercube (CAN)



### Geometry

- nodes are the corners of a hypercube
- distance = #matching bits in the IDs of two nodes
- logN neighbors per node; each at distance=1 away

## Ring (Chord) Geometry • nodes are points on a ring • distance = numeric distance between two node IDs • logN neighbors exponentially spaced over 0...N



### Geometry: • combination of a tree and ring

- two distance metricsdefault routing uses tree; fallback to ring under failures
  - · neighbors picked as on the tree

XOR (Kademlia)



00 ←→01 ←→ 10 ←→ 11

01 ←→00 ←→11 ←→ 10

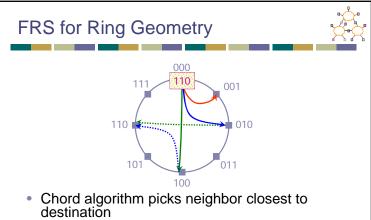
### Geometry:

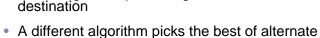
- distance(A,B) = A XOR B
- logN neighbors per node spaced exponentially
- <u>not</u> a ring because there is no single consistent ordering of all the nodes

Geometry's Impact on Routing

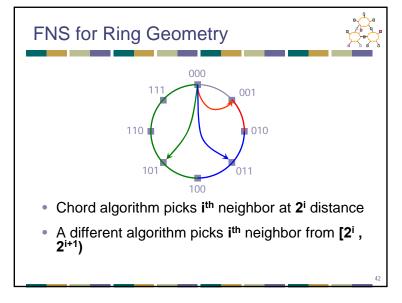


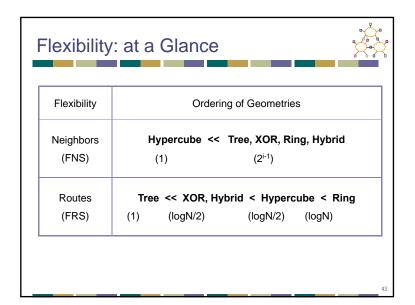
- Routing
  - Neighbor selection: how a node picks its routing entries
  - Route selection: how a node picks the next hop
- Proposed metric: flexibility
  - amount of freedom to choose neighbors and next-hop paths
    - FNS: flexibility in neighbor selection
    - FRS: flexibility in route selection
  - intuition: captures ability to "tune" DHT performance
  - single predictor metric dependent only on routing issues





paths





### Geometry → Flexibility → Performance?



### Validate over three performance metrics:

- 1. resilience
- 2. path latency
- 3. path convergence

### Metrics address two typical concerns:

- ability to handle node failure
- ability to incorporate proximity into overlay routing

### Analysis of Static Resilience



### Two aspects of robust routing

- Dynamic Recovery: how quickly routing state is recovered after failures
- Static Resilience: how well the network routes before recovery finishes
  - · captures how quickly recovery algorithms need to work
  - · depends on FRS

### **Evaluation:**

- Fail a fraction of nodes, without recovering any state
- Metric: % Paths Failed

