



## 15-440 Distributed Systems

### Lecture 21 – CDN & Peer-to-Peer

## Last Lecture: DNS (Summary)



- Motivations → large distributed database
  - Scalability
  - Independent update
  - Robustness
- Hierarchical database structure
  - Zones
  - How is a lookup done
- Caching/prefetching and **TTLs**
- Reverse name lookup
- What are the steps to creating your own domain?

2

## Outline



- **Content Distribution Networks**
- P2P Lookup Overview
- Centralized/Flooded Lookups
- Routed Lookups – Chord

3

## Typical Workload (Web Pages)



- Multiple (typically small) objects per page
- File sizes are heavy-tailed
- Embedded references
- This plays havoc with performance. Why?
- Solutions?

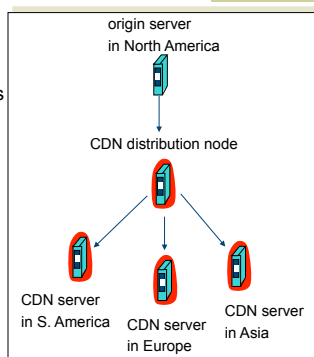
- Lots of small objects & TCP
  - 3-way handshake
  - Lots of slow starts
  - Extra connection state

4

## Content Distribution Networks (CDNs)



- The content providers are the CDN customers.
- Content replication
- CDN company installs hundreds of CDN servers throughout Internet
  - Close to users
- CDN replicates its customers' content in CDN servers. When provider updates content, CDN updates servers



5

## How Akamai Works



- Clients fetch html document from primary server
  - E.g. fetch index.html from cnn.com
- URLs for replicated content are replaced in html
  - E.g. `` replaced with ``
- Client is forced to resolve `aXYZ.g.akamaitech.net` hostname

Note: Nice presentation on Akamai at [www.cs.odu.edu/~mukka/cs775s07/Presentations/mklein.pdf](http://www.cs.odu.edu/~mukka/cs775s07/Presentations/mklein.pdf)

6

## How Akamai Works



- How is content replicated?
- Akamai only replicates static content (\*)
- Modified name contains original file name
- Akamai server is asked for content
  - First checks local cache
  - If not in cache, requests file from primary server and caches file

\* (At least, the version we're talking about today. Akamai actually lets sites write code that can run on Akamai's servers, but that's a pretty different beast)

7

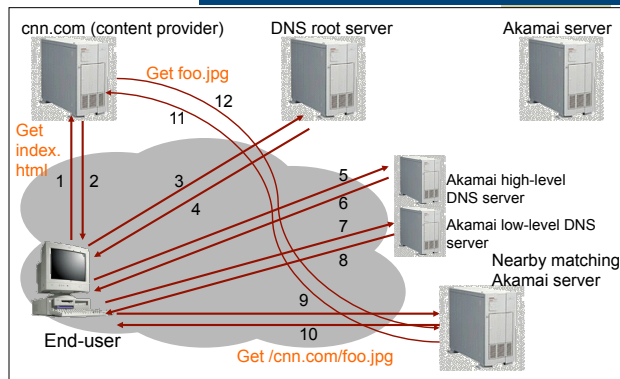
## How Akamai Works



- Root server gives NS record for akamai.net
- Akamai.net name server returns NS record for g.akamaitech.net
  - Name server chosen to be in region of client's name server
  - TTL is large
- G.akamaitech.net nameserver chooses server in region
  - Should try to choose server that has file in cache - How to choose?
  - Uses aXYZ name and hash
  - TTL is small → why?

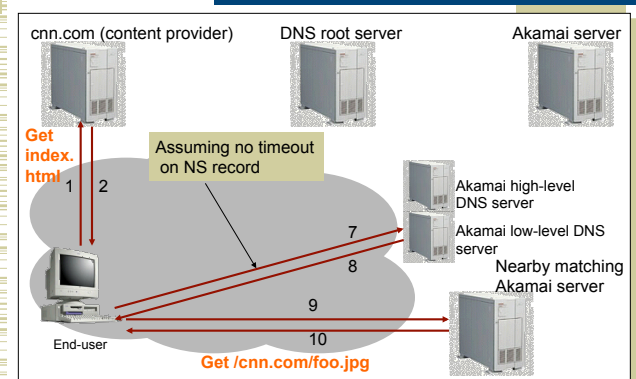
8

## How Akamai Works



9

## Akamai – Subsequent Requests



10

## Simple Hashing



- Given document XYZ, we need to choose a server to use
- Suppose we use modulo
- Number servers from 1...n
  - Place document XYZ on server (XYZ mod n)
  - What happens when a server fails?  $n \rightarrow n-1$ 
    - Same if different people have different measures of n
  - Why might this be bad?

11

## Consistent Hash



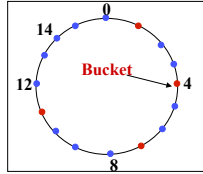
- “view” = subset of all hash buckets that are visible
- Desired features
  - Smoothness – little impact on hash bucket contents when buckets are added/removed
  - Spread – small set of hash buckets that may hold an object regardless of views
  - Load – across all views # of objects assigned to hash bucket is small

12

## Consistent Hash – Example



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



13

## Consistent Hashing not just for CDN



- Finding a nearby server for an object in a CDN uses centralized knowledge.
- Consistent hashing can also be used in a distributed setting
- P2P systems like BitTorrent need a way of finding files.
- Consistent Hashing to the rescue.

14

## Summary



- Content Delivery Networks move data closer to user, maintain consistency, balance load
- Consistent hashing maps keys AND buckets into the same space
- Consistent hashing can be fully distributed, useful in P2P systems using structured overlays

15

## Outline



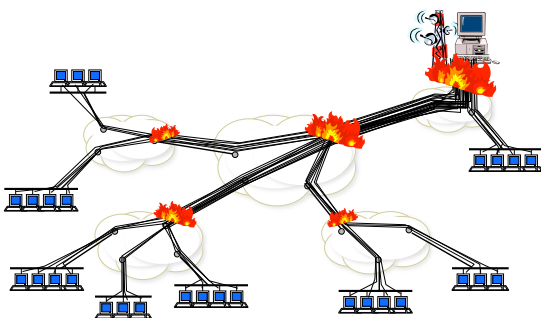
- Content Distribution Networks
- **P2P Lookup Overview**
- Centralized/Flooded Lookups
- Routed Lookups – Chord

16

## Scaling Problem

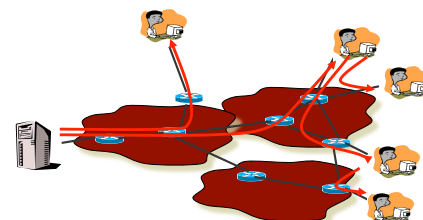


- Millions of clients ⇒ server and network meltdown



17

## P2P System



- Leverage the resources of client machines (peers)
  - Computation, storage, bandwidth

18

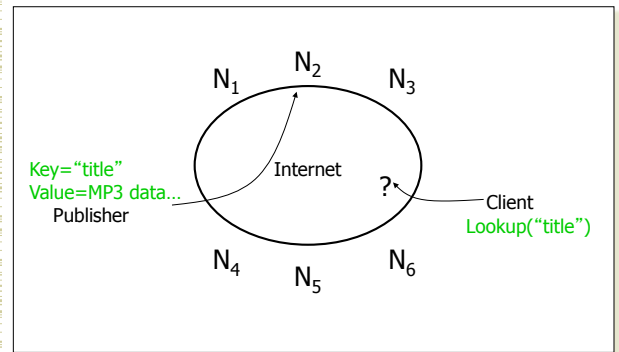
## Peer-to-Peer Networks



- Typically each member stores/provides access to content
- Basically a replication system for files
  - Always a tradeoff between possible location of files and searching difficulty
  - Peer-to-peer allow files to be anywhere → searching is the challenge
  - Dynamic member list makes it more difficult
- What other systems have similar goals?
  - Routing, DNS

19

## The Lookup Problem



20

## Searching



- Needles vs. Haystacks
  - Searching for top 40, or an obscure punk track from 1981 that nobody's heard of?
- Search expressiveness
  - Whole word? Regular expressions? File names? Attributes? Whole-text search?
    - (e.g., p2p gnutella or p2p google?)

21

## Framework



- Common Primitives:
  - **Join**: how do I begin participating?
  - **Publish**: how do I advertise my file?
  - **Search**: how do I find a file?
  - **Fetch**: how do I retrieve a file?

22

## Outline



- Content Distribution Networks
- P2P Lookup Overview
- **Centralized/Flooded Lookups**
- Routed Lookups – Chord

23

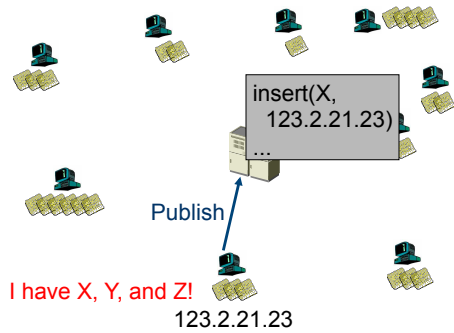
## Napster: Overview



- Centralized Database:
  - **Join**: on startup, client contacts central server
  - **Publish**: reports list of files to central server
  - **Search**: query the server => return someone that stores the requested file
  - **Fetch**: get the file directly from peer

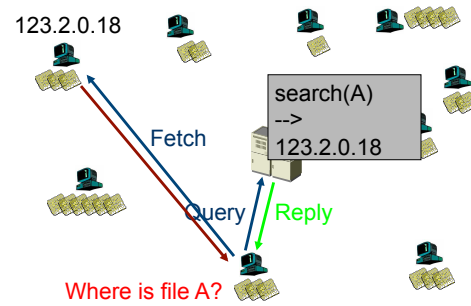
24

## Napster: Publish



25

## Napster: Search



26

## Napster: Discussion

- Pros:
  - Simple
  - Search scope is  $O(1)$
  - Controllable (pro or con?)
- Cons:
  - Server maintains  $O(N)$  State
  - Server does all processing
  - Single point of failure

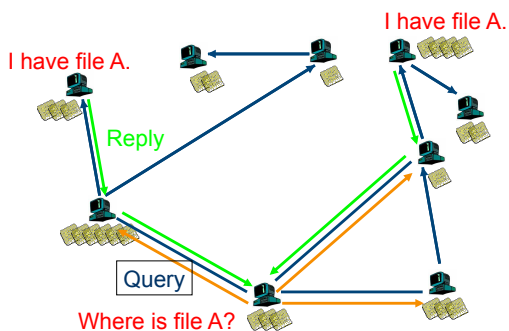
27

## "Old" Gnutella: Overview

- Query Flooding:
  - **Join**: on startup, client contacts a few other nodes; these become its "neighbors"
  - **Publish**: no need
  - **Search**: ask neighbors, who ask their neighbors, and so on... when/if found, reply to sender.
    - TTL limits propagation
  - **Fetch**: get the file directly from peer

28

## Gnutella: Search



29

## Gnutella: Discussion

- Pros:
  - Fully de-centralized
  - Search cost distributed
  - Processing @ each node permits powerful search semantics
- Cons:
  - Search scope is  $O(N)$
  - Search time is  $O(???)$
  - Nodes leave often, network unstable
- TTL-limited search works well for haystacks.
  - For scalability, does NOT search every node. May have to re-issue query later

30

## 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 ☹

31

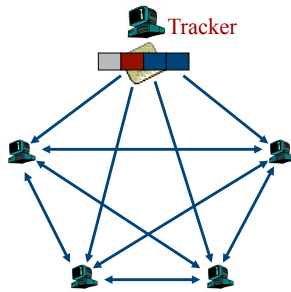
## BitTorrent: Overview



- Swarming:
  - **Join:** contact centralized “tracker” server, get a list of peers.
  - **Publish:** Run a tracker server.
  - **Search:** Out-of-band. E.g., use Google to find a tracker for the file you want.
  - **Fetch:** Download chunks of the file from your peers. Upload chunks you have to them.
- Big differences from Napster:
  - Chunk based downloading
  - “few large files” focus
  - Anti-freeloading mechanisms

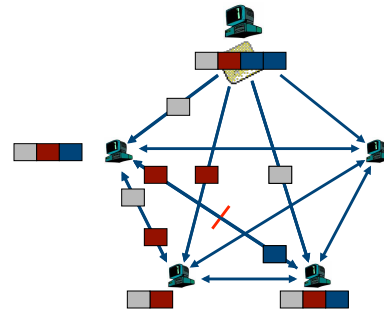
32

## BitTorrent: Publish/Join



33

## BitTorrent: Fetch



34

## BitTorrent: Sharing Strategy



- Employ “Tit-for-tat” sharing strategy
  - A is downloading from some other people
    - A will let the fastest N of those download from him
  - Be optimistic: occasionally let freeloaders download
    - Otherwise no one would ever start!
    - Also allows you to discover better peers to download from when they reciprocate
- Goal: Pareto Efficiency
  - Game Theory: “No change can make anyone better off without making others worse off”
  - Does it work? (not perfectly, but perhaps good enough?)

35

## BitTorrent: Summary



- Pros:
  - Works reasonably well in practice
  - Gives peers incentive to share resources; avoids freeloaders
- Cons:
  - Pareto Efficiency relative weak condition
  - Central tracker server needed to bootstrap swarm
    - Alternate tracker designs exist (e.g. DHT based)

36

## Outline



- Content Distribution Networks
- P2P Lookup Overview
- Centralized/Flooded Lookups
- Routed Lookups – Chord

37

## DHT: Overview (1)



- 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...
- Implementation: nodes in system form a distributed data structure
  - Can be Ring, Tree, Hypercube, Skip List, Butterfly Network, ...

38

## DHT: Overview (2)



- Structured Overlay Routing:
  - **Join:** On startup, contact a “bootstrap” node and integrate yourself into the distributed data structure; get a *node id*
  - **Publish:** Route publication for *file id* toward a close *node id* along the data structure
  - **Search:** Route a query for file id toward a close node id. Data structure guarantees that query will meet the publication.
  - **Fetch:** Two options:
    - Publication contains actual file => fetch from where query stops
    - Publication says “I have file X” => query tells you 128.2.1.3 has X, use IP routing to get X from 128.2.1.3

39

## DHT: Example - Chord



- Associate to each node and file a unique *id* in an *uni*-dimensional space (a Ring)
  - E.g., pick from the range  $[0 \dots 2^m]$
  - Usually the hash of the file or IP address
- Properties:
  - Routing table size is  $O(\log N)$ , where  $N$  is the total number of nodes
  - Guarantees that a file is found in  $O(\log N)$  hops

from MIT in 2001

40

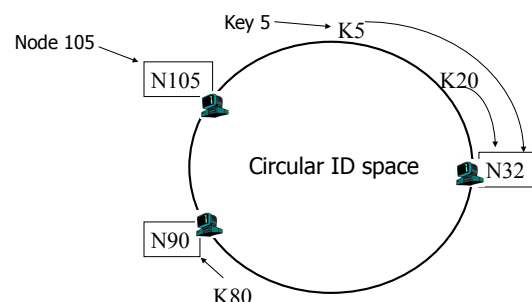
## Routing: Chord



- 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

41

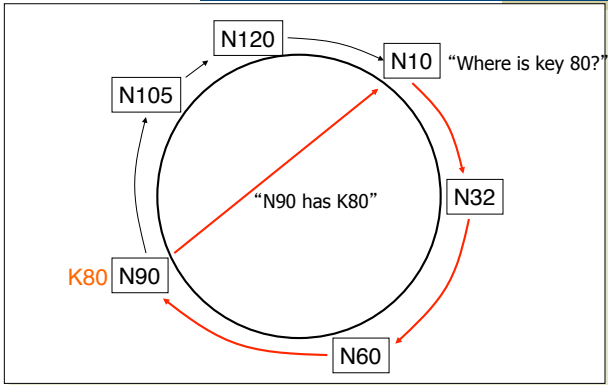
## DHT: Consistent Hashing



A key is stored at its successor: node with next higher ID

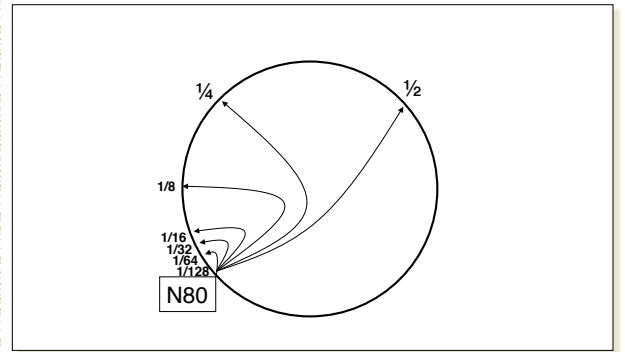
42

## Routing: Chord Basic Lookup



43

## Routing: Finger table - Faster Lookups



44

## Routing: Chord Summary



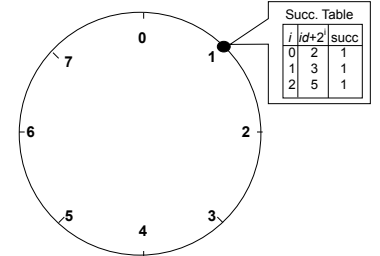
- Assume identifier space is  $0 \dots 2^m$
- Each node maintains
  - Finger table
    - Entry  $i$  in the finger table of  $n$  is the first node that succeeds or equals  $n + 2^i$
  - Predecessor node
- An item identified by  $id$  is stored on the successor node of  $id$

45

## Routing: Chord Example



- Assume an identifier space  $0 \dots 7$
- Node  $n1:(1)$  joins  $\rightarrow$  all entries in its finger table are initialized to itself

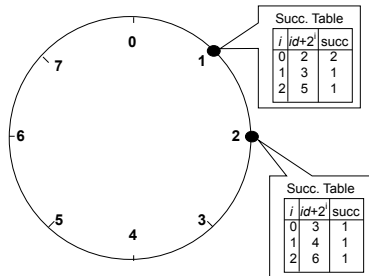


46

## Routing: Chord Example



- Node  $n2:(3)$  joins

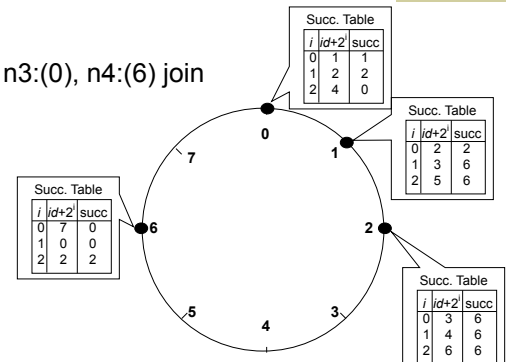


47

## Routing: Chord Example



- Nodes  $n3:(0)$ ,  $n4:(6)$  join



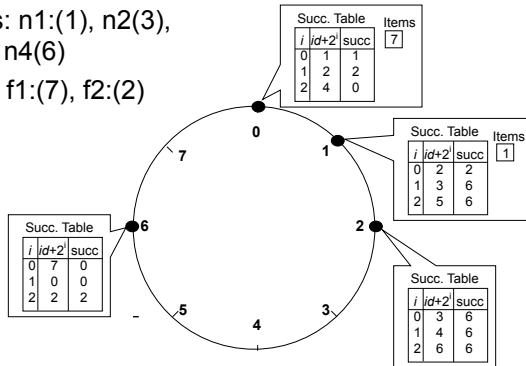
48



## Routing: Chord Examples



- Nodes:  $n_1:(1)$ ,  $n_2(3)$ ,  $n_3(0)$ ,  $n_4(6)$
- Items:  $f_1:(7)$ ,  $f_2:(2)$

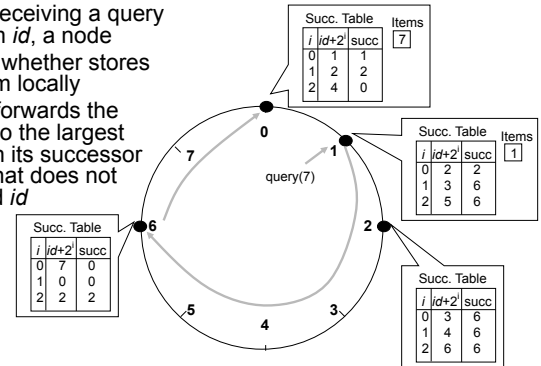


49

## Routing: Query



- Upon receiving a query for item  $id$ , a node
- Check whether stores the item locally
- If not, forwards the query to the largest node in its successor table that does not exceed  $id$



50

## DHT: Chord Summary



- Routing table size?
  - Log  $N$  fingers
- Routing time?
  - Each hop expects to  $1/2$  the distance to the desired  $id \Rightarrow$  expect  $O(\log N)$  hops.

51

## DHT: Discussion



- Pros:
  - Guaranteed Lookup
  - $O(\log N)$  per node state and search scope
- Cons:
  - No one uses them? (only one file sharing app)
  - Supporting non-exact match search is hard

52

## What can DHTs do for us?



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

53

## When are p2p / DHTs useful?



- Caching and “soft-state” data
  - Works well! BitTorrent, KaZaA, etc., all use peers as caches for hot data
- Finding read-only data
  - Limited flooding finds hay
  - DHTs find needles
- BUT

54

## A Peer-to-peer Google?



- Complex intersection queries (“the” + “who”)
  - Billions of hits for each term alone
- Sophisticated ranking
  - Must compare many results before returning a subset to user
- Very, very hard for a DHT / p2p system
  - Need high inter-node bandwidth
  - (This is exactly what Google does - massive clusters)

55

## Writable, persistent p2p



- Do you trust your data to 100,000 monkeys?
- Node availability hurts
  - Ex: Store 5 copies of data on different nodes
  - When someone goes away, you must replicate the data they held
  - Hard drives are \*huge\*, but cable modem upload bandwidth is tiny - perhaps 10 Gbytes/day
  - Takes many days to upload contents of 200GB hard drive. Very expensive leave/replication situation!

56

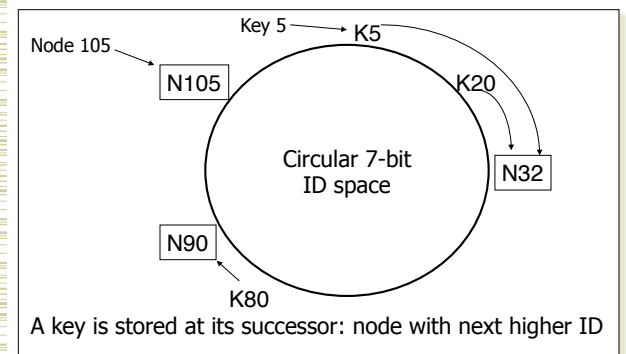
## P2P: Summary



- Many different styles; remember pros and cons of each
  - centralized, flooding, swarming, unstructured and structured routing
- Lessons learned:
  - Single points of failure are very bad
  - Flooding messages to everyone is bad
  - Underlying network topology is important
  - Not all nodes are equal
  - Need incentives to discourage freeloading
  - Privacy and security are important
  - Structure can provide theoretical bounds and guarantees

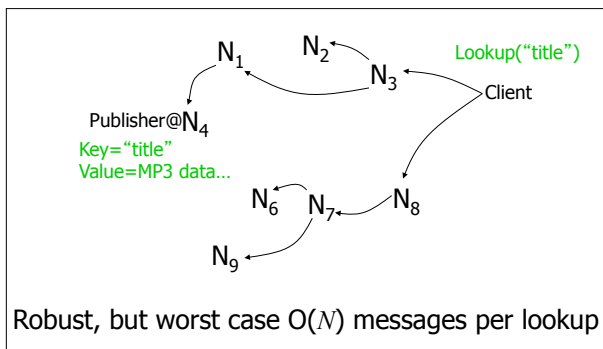
57

## Aside: Consistent Hashing [Karger 97]



58

## Flooded Queries (Gnutella)



59

## Flooding: Old Gnutella



- On startup, client contacts any server (**server + client**) in network
  - Server 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

60

## 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!

61

## Flooding: Old Gnutella Details



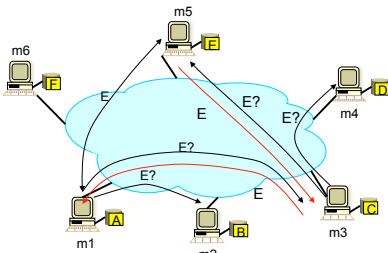
- Basic message header
  - Unique ID, TTL, Hops
- Message types
  - Ping – probes network for other servents
  - Pong – response to ping, contains IP addr, # of files, # of Kbytes shared
  - Query – search criteria + speed requirement of servent
  - QueryHit – successful response to Query, contains addr + 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

62

## Flooding: Old Gnutella Example

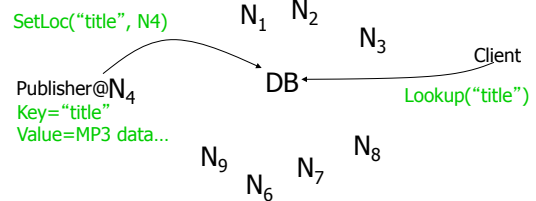


Assume: m1's neighbors are m2 and m3;  
m3's neighbors are m4 and m5;...



63

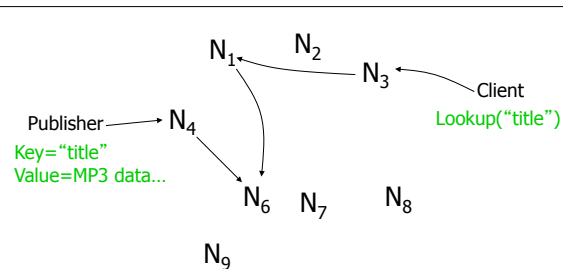
## Centralized Lookup (Napster)



Simple, but  $O(N)$  state and a single point of failure

64

## Routed Queries (Chord, etc.)



65

## Slashdot

NEWS FOR NERDS. STUFF THAT MATTERS.

Log In | Create Account | Help | Subscribe | Firehose

Sections  
Main  
Apple  
AskSlashdot  
Books  
Developers  
Games  
Hardware  
IT  
Idle  
Index  
Interviews

### Political Sites Scale Up For Election Traffic

Posted by [timothy](#) on Tuesday November 04, @ 12:15PM  
from the [hearken-are-those-trumpets](#) dept.

[miller@](#) writes

"News sites and political blogs are expecting extraordinary traffic tonight as Americans track results of the Presidential election, and are scaling their infrastructure to meet the challenge. Yahoo anticipates its Election Night traffic may be three times the volume seen in 2004, when it had 80 million page views on Election Day and 142 million more visits the following day. Hosting companies say customers have been ordering extra servers and load balancing services, while content delivery networks are also expecting a busy night. Will traffic approach record levels? Akamai's [Net Usage Index](#), which tracks traffic to its customer news sites, is one metric to watch."



<http://www.akamai.com/html/technology/nui/news/index.html>

66

## Content Distribution Networks & Server Selection



- Replicate content on many servers
- Challenges
  - How to replicate content
  - Where to replicate content
  - How to find replicated content
  - How to choose among known replicas
  - How to direct clients towards replica

67

## Server Selection



- Which server?
  - Lowest load → to balance load on servers
  - Best performance → to improve client performance
    - Based on Geography? RTT? Throughput? Load?
  - Any alive node → to provide fault tolerance
- How to direct clients to a particular server?
  - As part of routing → anycast, cluster load balancing
    - Not covered ☹
  - As part of application → HTTP redirect
  - As part of naming → DNS

68

## Application Based



- HTTP supports simple way to indicate that Web page has moved (30X responses)
- Server receives Get request from client
  - Decides which server is best suited for particular client and object
  - Returns HTTP redirect to that server
- Can make informed application specific decision
- May introduce additional overhead → multiple connection setup, name lookups, etc.
- While good solution in general, but...
  - HTTP Redirect has some design flaws – especially with current browsers

69

## Naming Based



- Client does name lookup for service
- Name server chooses appropriate server address
  - A-record returned is “best” one for the client
- What information can name server base decision on?
  - Server load/location → must be collected
  - Information in the name lookup request
    - Name service client → typically the local name server for client

70