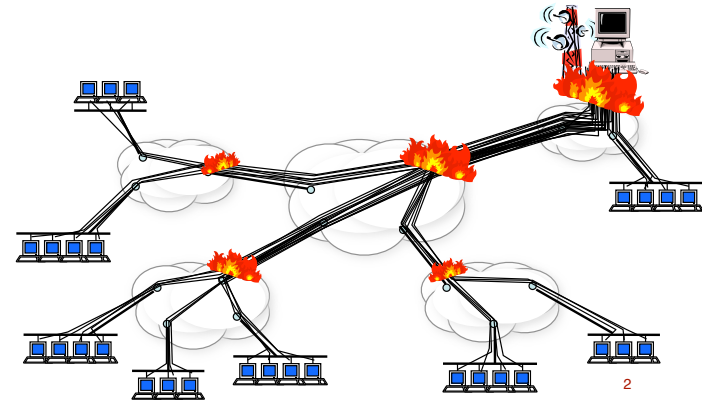


## Peer-to-Peer

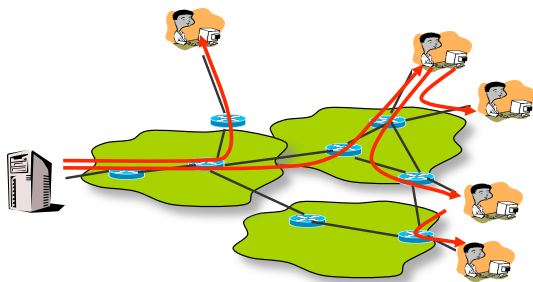
15-441

## Scaling Problem

- Millions of clients  $\Rightarrow$  server and network meltdown



## P2P System



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

3

## Why p2p?

- Scaling: Create system whose capacity grows with # of clients - automatically!
- Self-managing
  - This aspect attractive for corporate/datacenter needs
  - e.g., Amazon's 100,000-ish machines, google's 300k+
- Harness lots of "spare" capacity at end-hosts
- Eliminate centralization
  - Robust to failures, etc.
  - Robust to censorship, politics & legislation??

4

## Today's Goal

- p2p is hot.
- There are tons and tons of instances
- But that's not the point
- Identify fundamental techniques useful in p2p settings
- Understand the challenges
- Look at the (current!) boundaries of where 2p is particularly useful

5

## Outline

- p2p file sharing techniques
  - Downloading: Whole-file vs. chunks
  - Searching
    - Centralized index (Napster, etc.)
    - Flooding (Gnutella, etc.)
    - Smarter flooding (KaZaA, ...)
    - Routing (Freenet, etc.)
- Uses of p2p - what works well, what doesn't?
  - servers vs. arbitrary nodes
  - Hard state (backups!) vs soft-state (caches)
- Challenges

6

## Searching & Fetching

Human:

“I want to watch that great 80s cult classic ‘Better Off Dead’”

1. Search:

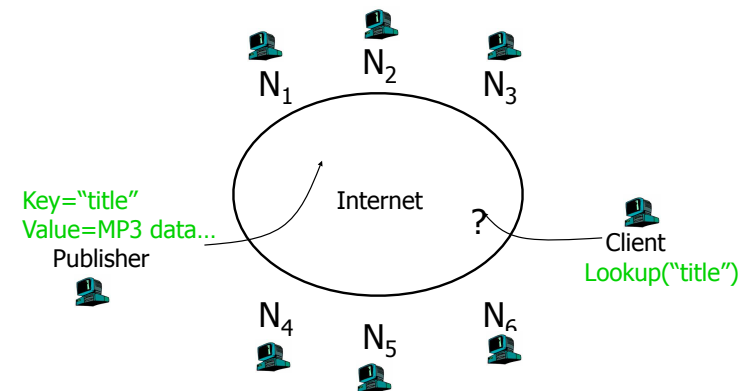
“better off dead” -> better\_off\_dead.mov  
or -> 0x539fba83ajdeadbeef

2. Locate sources of better\_off\_dead.mov

3. Download the file from them

7

## Searching



8

## Search Approaches

- Centralized
- Flooding
- A hybrid: Flooding between “Supernodes”
- Structured

9

## Different types of searches

- 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?)

10

## Framework

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

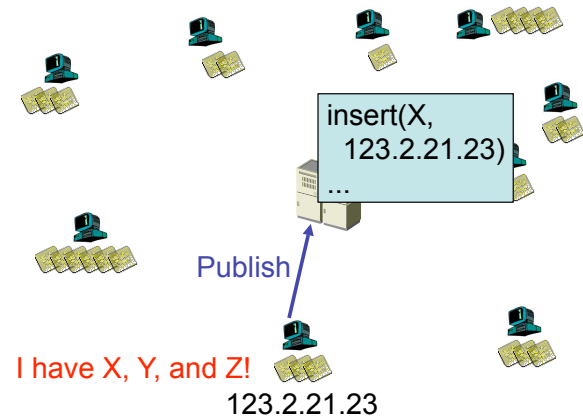
11

## Centralized

- Centralized Database:
  - **Join**: on startup, client contacts central server
  - **Publish**: reports list of files to central server
  - **Search**: query the server => return node(s) that store the requested file

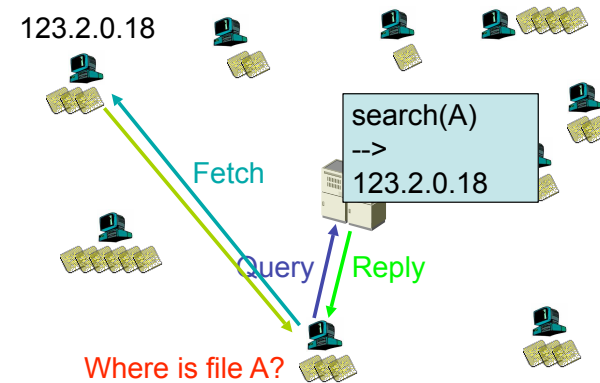
12

## Napster Example: Publish



13

## Napster: Search



14

## Napster: Discussion

- Pros:
  - Simple
  - Search scope is  $O(1)$  for even complex searches (one index, etc.)
  - Controllable (pro or con?)
- Cons:
  - Server maintains  $O(N)$  State
  - Server does all processing
  - Single point of failure
    - Technical failures + legal (napster shut down)

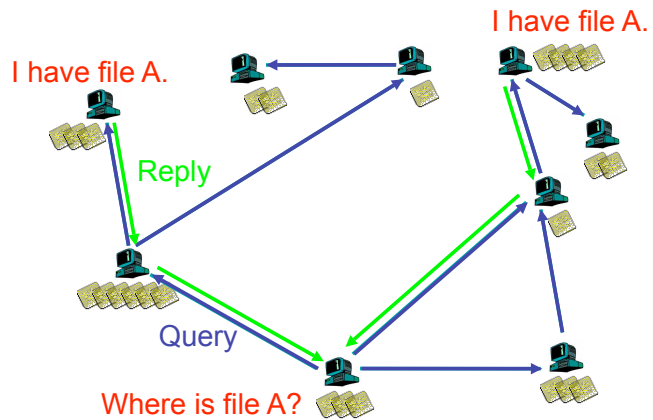
15

## Query Flooding

- **Join**: Must join a flooding network
  - Usually, establish peering with a few existing nodes
- **Publish**: no need, just reply
- **Search**: ask neighbors, who ask their neighbors, and so on... when/if found, reply to sender.
  - TTL limits propagation

16

## Example: Gnutella



17

## Flooding: Discussion

- Pros:
  - Fully de-centralized
  - Search cost distributed
  - Processing @ each node permits powerful search semantics
- Cons:
  - Search scope is  $O(M)$
  - 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

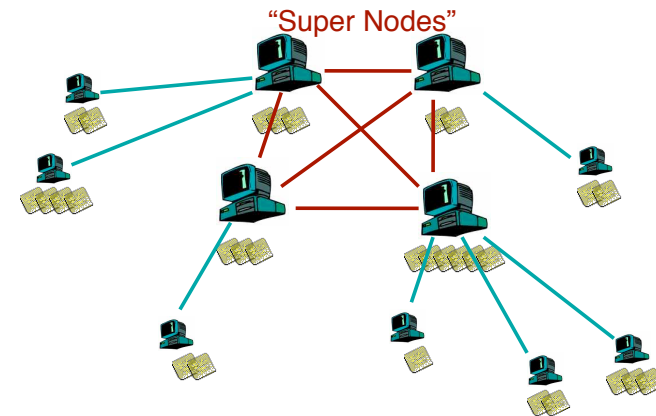
18

## Supernode Flooding

- **Join:** on startup, client contacts a "supernode" ... may at some point become one itself
- **Publish:** send list of files to supernode
- **Search:** send query to supernode, supernodes flood query amongst themselves.
  - Supernode network just like prior flooding net

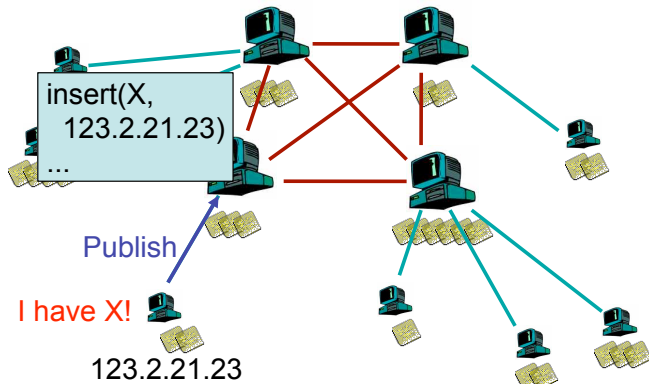
19

## Supernode Network Design



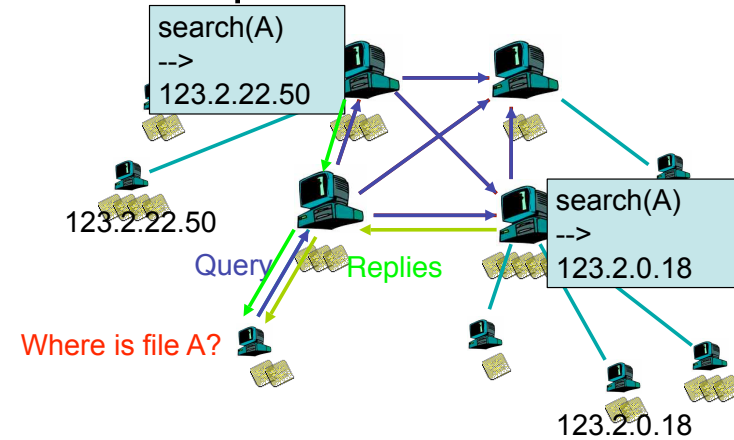
20

## Supernode: File Insert



21

## Supernode: File Search



22

## Supernode: Which nodes?

- Often, bias towards nodes with good:
  - Bandwidth
  - Computational Resources
  - Availability!

23

## Stability and Superpeers

- Why superpeers?
  - Query consolidation
    - Many connected nodes may have only a few files
    - Propagating a query to a sub-node would take more b/w than answering it yourself
  - Caching effect
    - Requires network stability
- Superpeer selection is time-based
  - How long you've been on is a good predictor of how long you'll be around.

24

## Superpeer results

- Basically, “just better” than flood to all
- Gets an order of magnitude or two better scaling
- But still fundamentally:  $o(\text{search}) * o(\text{per-node storage}) = O(N)$ 
  - central:  $O(1)$  search,  $O(N)$  storage
  - flood:  $O(N)$  search,  $O(1)$  storage
  - Superpeer: can trade between

25

## Structured Search: Distributed Hash Tables

- Academic answer to p2p
- Goals
  - Guaranteed lookup success
  - Provable bounds on search time
  - Provable scalability
- Makes some things harder
  - Fuzzy queries / full-text search / etc.
- Read-write, not read-only
- Hot Topic in networking since introduction in ~2000/2001

26

## Searching Wrap-Up

Type	$O(\text{search})$	storage	Fuzzy?
Central	$O(1)$	$O(N)$	Yes
Flood	$\sim O(N)$	$O(1)$	Yes
Super	$< O(N)$	$> O(1)$	Yes
Structured	$O(\log N)$	$O(\log N)$	not really

27

## DHT: Overview

- **Abstraction:** a distributed “hash-table” (DHT) data structure:
  - `put(id, item);`
  - `item = get(id);`
- **Implementation:** nodes in system form a distributed data structure
  - Can be Ring, Tree, Hypercube, Skip List, Butterfly Network, ...

28

## 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.
- Important difference: `get(key)` is for an *exact match* on key!
  - `search(“spars”)` will not find file(“briney spars”)
  - We can exploit this to be more efficient

29

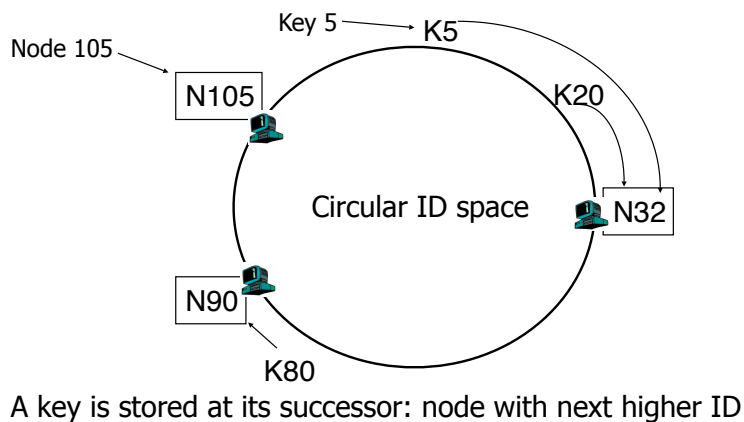
## DHT: Example - Chord

- Associate to each node and file a unique *id* in an *uni*-dimensional space (a Ring)
  - E.g., pick from the range  $[0...2^m]$
  - 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

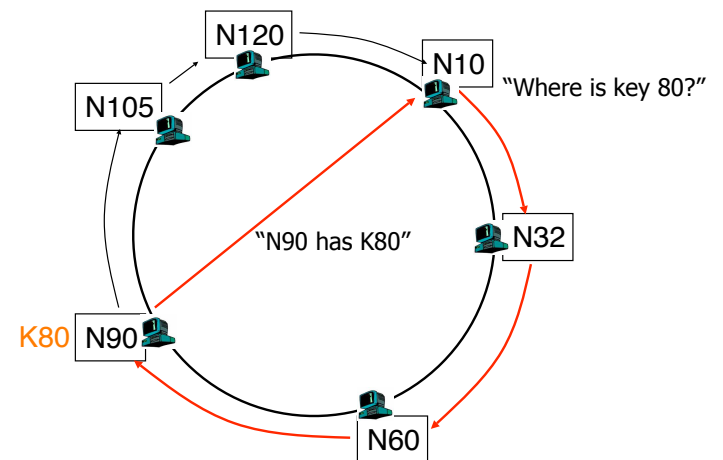
30

## DHT: Consistent Hashing



31

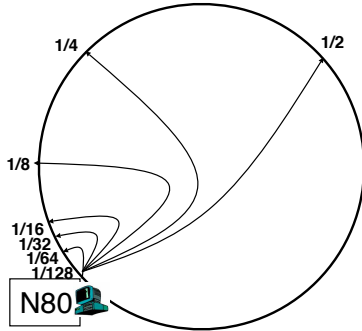
## DHT: Chord Basic Lookup



32



## DHT: Chord “Finger Table”



- Entry  $i$  in the finger table of node  $n$  is the first node that succeeds or equals  $n + 2^i$
- In other words, the  $i$ th finger points  $1/2^{n-i}$  way around the ring

33

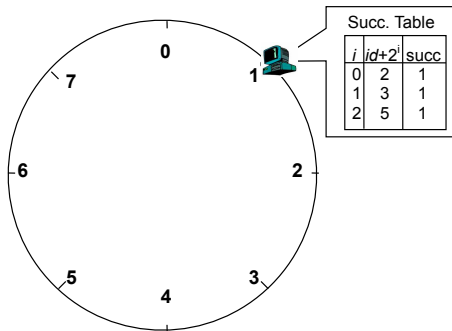
## Node Join

- Compute ID
- Use an existing node to route to that ID in the ring.
  - Finds  $s = \text{successor}(id)$
- ask  $s$  for its predecessor,  $p$
- Splice self into ring just like a linked list
  - $p \rightarrow \text{successor} = \text{me}$
  - $\text{me} \rightarrow \text{successor} = s$
  - $\text{me} \rightarrow \text{predecessor} = p$

34

## DHT: Chord Join

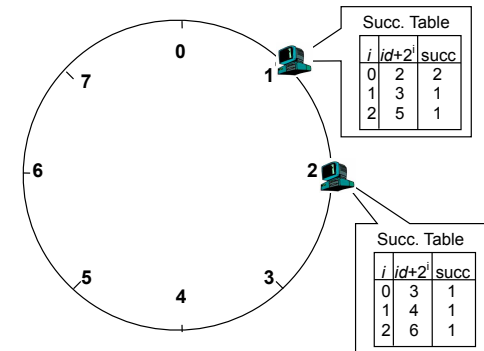
- Assume an identifier space  $[0..8]$
- Node  $n1$  joins



35

## DHT: Chord Join

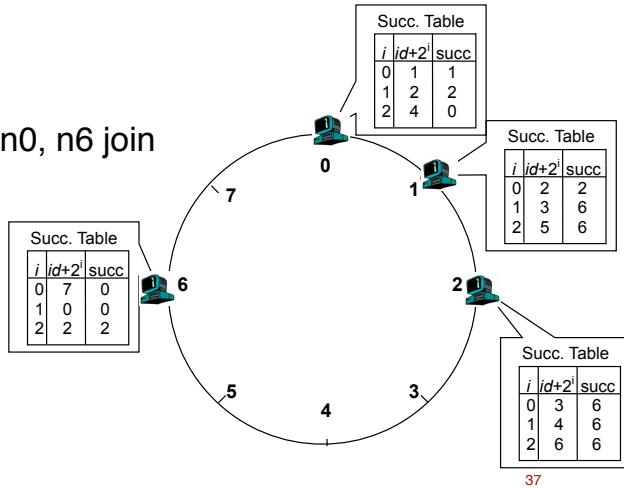
- Node  $n2$  joins



36

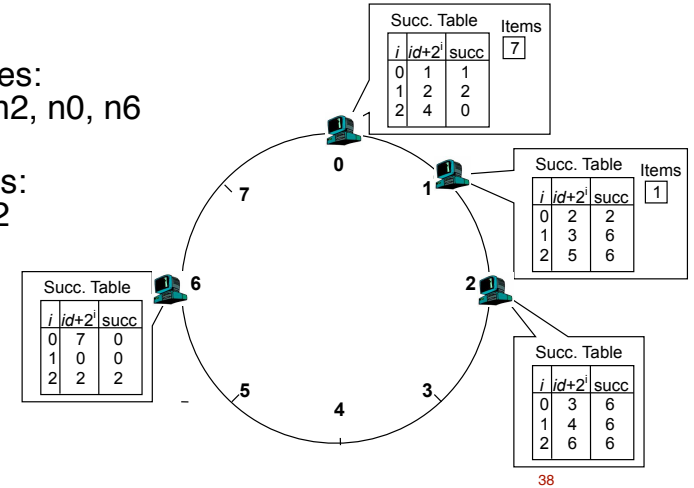
## DHT: Chord Join

- Nodes  $n_0, n_6$  join



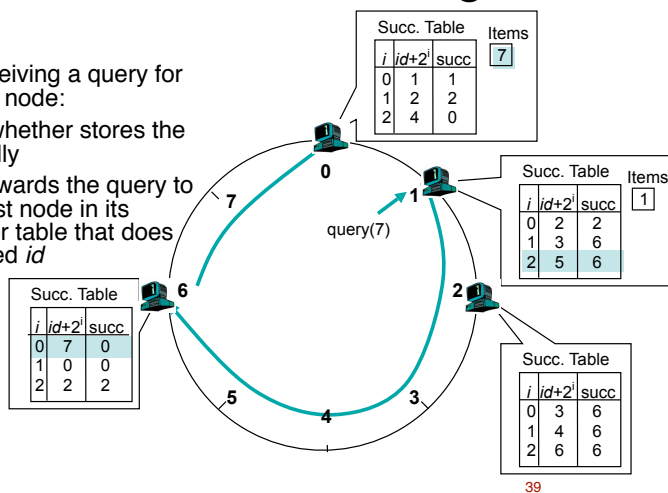
## DHT: Chord Join

- Nodes:  $n_1, n_2, n_0, n_6$
- Items:  $f_7, f_2$



## DHT: Chord Routing

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



## DHT: Chord Summary

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

## DHT: Discussion

- Pros:
  - Guaranteed Lookup
  - $O(\log M)$  per node state and search scope
- Cons:
  - This line used to say “not used.” But: Now being used in a few apps, including BitTorrent.
  - Supporting non-exact match search is (quite!) hard

41

## The limits of search: 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)
- But maybe many file sharing queries are okay..<sup>42</sup>

## Fetching Data

- Once we know which node(s) have the data we want...
- Option 1: Fetch from a single peer
  - Problem: Have to fetch from peer who has whole file.
    - Peers not useful sources until d/l whole file
    - At which point they probably log off. :)
  - How can we fix this?

43

## Chunk Fetching

- More than one node may have the file.
- How to tell?
  - Must be able to distinguish identical files
  - Not necessarily same filename
  - Same filename not necessarily same file...
- Use hash of file
  - Common: MD5, SHA-1, etc.
- How to fetch?
  - Get bytes [0..8000] from A, [8001...16000] from B
  - Alternative: Erasure Codes

44

## BitTorrent: Overview

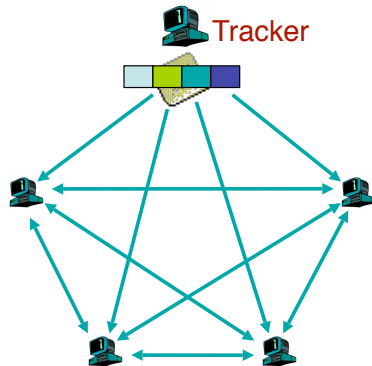
- **Swarming:**
  - **Join:** contact centralized “tracker” server, get a list of peers.
  - **Publish:** Run a tracker server.
  - **Search:** Out-of-band. E.g., use Google to find a tracker for the file you want.
  - **Fetch:** Download chunks of the file from your peers. Upload chunks you have to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

45

## BitTorrent

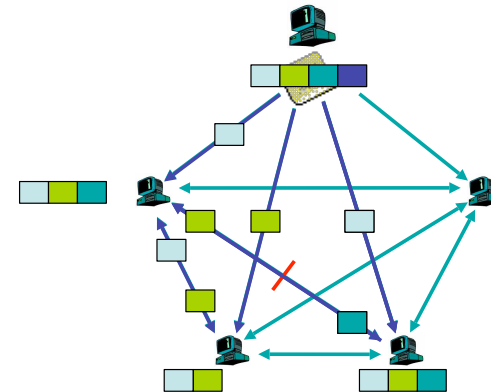
- Periodically get list of peers from tracker
- More often:
  - Ask each peer for what chunks it has
    - (Or have them update you)
- Request chunks from several peers at a time
- Peers will start downloading from you
- BT has some machinery to try to bias towards helping those who help you<sup>46</sup>

## BitTorrent: Publish/Join



47

## BitTorrent: Fetch



48

## BitTorrent: Summary

- Pros:
  - Works reasonably well in practice
  - Gives peers incentive to share resources; avoids freeloaders
- Cons:
  - Central tracker server needed to bootstrap swarm
  - (Tracker is a design choice, not a requirement, as you know from your projects. Modern BitTorrent can also use a DHT to locate peers. But approach still needs a “search” mechanism)

49

## 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!

50

## What's out there?

	Central	Flood	Super-node flood	Route
Whole File	Napster	Gnutella		Freenet
Chunk Based	BitTorrent		KaZaA (bytes, not chunks)	DHTs eDonkey 2000

51

## P2P: Summary

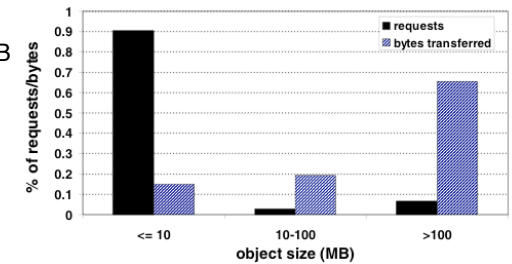
- Many different styles; remember pros and cons of each
  - centralized, flooding, swarming, unstructured and structured routing
- Lessons learned:
  - Single points of failure are bad
  - Flooding messages to everyone is bad
  - Underlying network topology is important
  - Not all nodes are equal
 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

52

## Extra Slides

## KaZaA: Usage Patterns

- KaZaA is more than one workload!
  - Many files < 10MB (e.g., Audio Files)
  - Many files > 100MB (e.g., Movies)

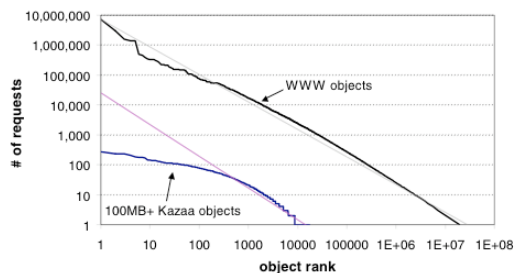


from Gummadi *et al.*, *SOSP 2003*

54

## KaZaA: Usage Patterns (2)

- KaZaA is not Zipf!
  - FileSharing: “Request-once”
  - Web: “Request-repeatedly”

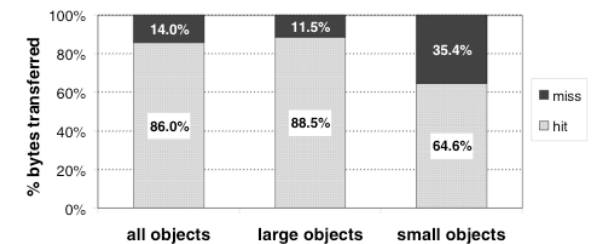


from Gummadi *et al.*, *SOSP 2003*

55

## KaZaA: Usage Patterns (3)

- What we saw:
  - A few big files consume most of the bandwidth
  - Many files are fetched once per client but still very popular
- Solution?
  - Caching!



from Gummadi *et al.*, *SOSP 2003*

56

## Freenet: History

- In 1999, I. Clarke started the Freenet project
- Basic Idea:
  - Employ Internet-like routing on the overlay network to publish and locate files
- Addition goals:
  - Provide anonymity and security
  - Make censorship difficult

57

## Freenet: Overview

- Routed Queries:
  - **Join**: on startup, client contacts a few other nodes it knows about; gets a unique *node id*
  - **Publish**: route file contents toward the *file id*. File is stored at node with *id* closest to *file id*
  - **Search**: route query for *file id* toward the closest *node id*
  - **Fetch**: when query reaches a node containing *file id*, it returns the file to the sender

58

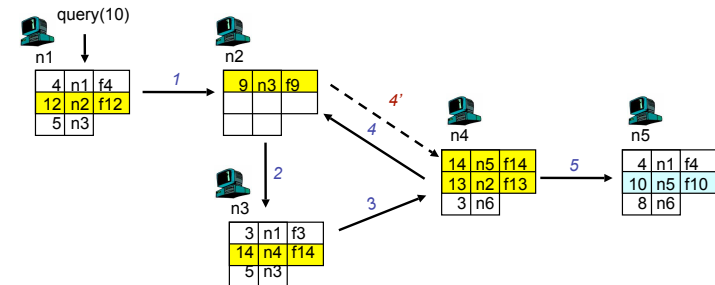
## Freenet: Routing Tables

- *id* – file identifier (e.g., hash of file)
- *next\_hop* – another node that stores the file id
- *file* – file identified by *id* being stored on the local node
- Forwarding of query for file *id*
  - If file *id* stored locally, then stop
    - Forward data back to upstream requestor
  - If not, search for the “closest” *id* in the table, and forward the message to the corresponding *next\_hop*
  - If data is not found, failure is reported back
    - Requestor then tries next closest match in routing table

<i>id</i>	<i>next_hop</i>	<i>file</i>
4	n1	f4
12	n2	f12
	⋮	
	⋮	
	⋮	
	⋮	

59

## Freenet: Routing



60

## Freenet: Routing Properties

- “Close” file ids tend to be stored on the same node
  - Why? Publications of similar file ids route toward the same place
- Network tend to be a “small world”
  - Small number of nodes have large number of neighbors (i.e., ~ “six-degrees of separation”)
- Consequence:
  - Most queries only traverse a small number of hops to find the file

61

## Freenet: Anonymity & Security

- Anonymity
  - Randomly modify source of packet as it traverses the network
  - Can use “mix-nets” or onion-routing
- Security & Censorship resistance
  - No constraints on how to choose *ids* for files => easy to have to files collide, creating “denial of service” (censorship)
  - Solution: have a *id* type that requires a private key signature that is verified when updating the file
  - Cache file on the reverse path of queries/publications => attempt to “replace” file with bogus data will just cause the file to be replicated more!

62

## Freenet: Discussion

- Pros:
  - Intelligent routing makes queries relatively short
  - Search scope small (only nodes along search path involved); no flooding
  - Anonymity properties may give you “plausible deniability”
- Cons:
  - Still no provable guarantees!
  - Anonymity features make it hard to measure, debug

63

## 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
  - Let N peop
- Goal: Pareto Efficiency
  - Game Theory: “No change can make anyone better off without making others worse off” <sup>64</sup>