

Lecture 21 (Apr 13, 2004)

Outline

- P2P: Napster, Gnutella, KaZaA
- Chord

CS 640

1

P2P: Napster, Gnutella, KaZaA
Chord

CS 640

1

What is P2P?

- Significant autonomy from central servers
- Exploits resources at the edges of Internet
 - Bandwidth
 - Storage
 - Processing
- Resources at edge have intermittent connectivity
 - Dynamic joins and leaves

CS 640

2

- CS 640

2

Applications

- P2P file sharing
 - Napster, Gnutella, KaZaA, etc.
- Storage and lookup
 - Chord, CAN, etc.
- P2P communication
 - Instant messaging
- P2P computation
 - seti@home

CS 640

3

- CS 640

3

P2P file sharing software

- Allows Alice to open up a directory in her file system
 - Anyone can retrieve a file from directory
 - Like a Web server
- Allows Alice to copy files from other users' open directories:
 - Like a Web client
- Allows users to search the peers for content based on keyword matches:
 - Like Google

CS 640

4

- CS 640

4

Napster

- the most (in)famous
- Instructive for what it gets right, and
- also wrong...
- also has a political message...and economic and legal...

CS 640 5

- CS 640

5

Napster

- popular for sharing files over the Internet
- a “disruptive” application/technology?
- history:
 - 5/99: Shawn Fanning (freshman, Northeastern U.) founds Napster Online music service
 - 12/99: first lawsuit
 - 3/00: 25% UWise traffic Napster
 - 2/01: US Circuit Court of Appeals: Napster knew users violating copyright laws
 - 7/01: # simultaneous online users:

Napster 160K, Gnutella: 40K,
Morpheus (KaZaA): 300K

CS 640

6

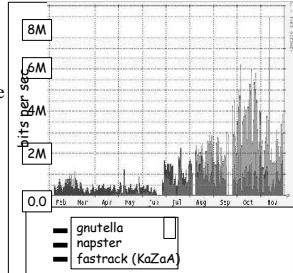
-
- Well Known Services MB/s**
- This chart displays network activity in Megabits per second (MB/s) for various well-known services over a 24-hour period. The Y-axis ranges from 0 to 120 MB/s. The X-axis shows time in hours, with labels at 18:00, 20:00, and 22:00. The chart includes a legend for the following services:
- Nagios* I/O
 - FTP-DATA I/O
 - SMTP SRC I/O
 - SMTP DST I/O
 - HTTP SRC I/O
 - HTTP DST I/O
 - Realserver I/O
 - IISW
 - TOTAL I/O
 - HTTP DST I/O
 - HTTP SRC I/O
 - SMTP SRC I/O
 - SMTP DST I/O
- The chart shows a significant peak in activity around 20:00, reaching approximately 120 MB/s. The total I/O is the sum of all individual service I/Os.
- Nagios* 22.12666628** **FTP-DATA 10.26610628**
HTTP 0.000000 **SMTP 1.93060128** **Real 0.76465216** **SMTP 0.40000000**
ICMP 0.18197171 **ether 34.25875872**

CS 640

6

Napster

- judge orders Napster to pull plug in July '01
- other file sharing apps take over!



CS 640

7

Napster: how does it work

- Application-level, client-server protocol over point-to-point TCP
- Centralized directory server

Steps:

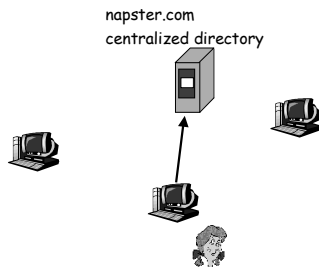
- connect to Napster server
- upload your list of files to server.
- give server keywords to search the full list with.
- select "best" of correct answers. (pings)

CS 640

8

Napster

1. File list and IP address is uploaded

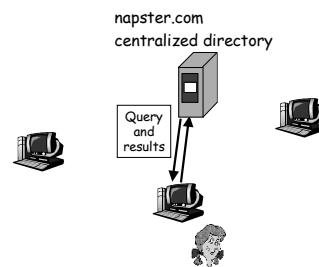


CS 640

9

Napster

2. User requests search at server.



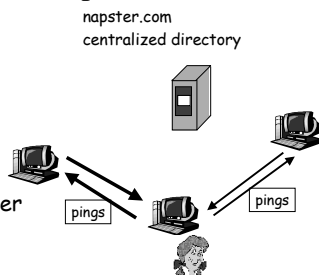
CS 640

10

Napster

3. User pings hosts that apparently have data.

Looks for **best** transfer rate.



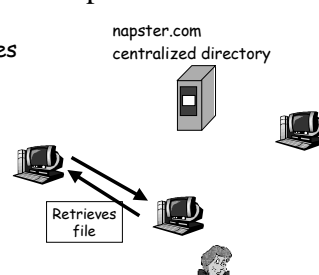
CS 640

11

Napster

4. User chooses server

Napster's centralized server farm had difficult time keeping up with traffic



CS 640

12

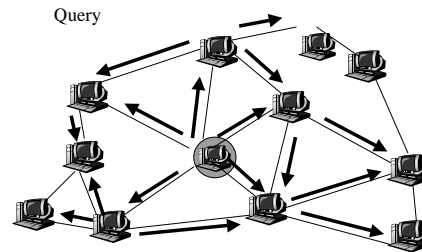
2. Unstructured P2P File Sharing

- Napster
- Gnutella
- KaZaA
- search theory
- dealing with flash crowds

CS 640

13

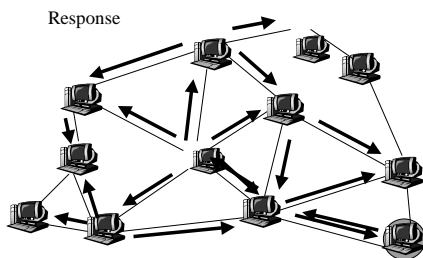
Distributed Search/Flooding



CS 640

14

Distributed Search/Flooding



CS 640

15

Gnutella

- focus: decentralized method of searching for files
 - central directory server no longer the bottleneck
 - more difficult to “pull plug”
- each application instance serves to:
 - store selected files
 - route queries from and to its neighboring peers
 - respond to queries if file stored locally
 - serve files

CS 640

16

Gnutella

- Gnutella history:
 - 3/14/00: release by AOL, almost immediately withdrawn
 - became open source
 - many iterations to fix poor initial design (poor design turned many people off)
- issues:
 - how much traffic does one query generate?
 - how many hosts can it support at once?
 - what is the latency associated with querying?
 - is there a bottleneck?

CS 640

17

Gnutella: limited scope query

Searching by flooding:

- if you don't have the file you want, query 7 of your neighbors.
- if they don't have it, they contact 7 of their neighbors, for a maximum hop count of 10.
- reverse path forwarding for responses (not files)

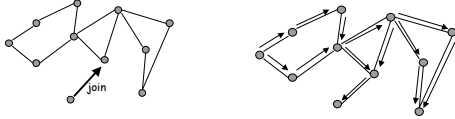
Note: Play gnutella animation at:
<http://www.limewire.com/index.jsp/p2p>

CS 640

18

Gnutella overlay management

- New node uses bootstrap node to get IP addresses of existing Gnutella nodes
- New node establishes neighboring relations by sending join messages



CS 640

19

Gnutella in practice

- Gnutella traffic \ll KaZaA traffic
- KaZaA:
 - hierarchy, queue management, parallel download,...

CS 640

20

Gnutella Discussion:

- researchers like it because it's open source
 - but is it truly representative?
- architectural lessons learned?
- good source for technical info/open questions:
 - http://www.limewire.com/index.jsp/tech_papers

CS 640

21

2. Unstructured P2P File Sharing

- Napster
- Gnutella
- KaZaA
- search theory
- dealing with flash crowds

CS 640

22

KaZaA: The service

- more than 3 million up peers sharing over 3,000 terabytes of content
- more popular than Napster ever was
- more than 50% of Internet traffic ?
- MP3s & entire albums, videos, games
- optional parallel downloading of files
- automatically switches to new download server when current server becomes unavailable
- provides estimated download times

CS 640

23

KaZaA: The service (2)

- User can configure max number of simultaneous uploads and max number of simultaneous downloads
- queue management at server and client
 - Frequent uploaders can get priority in server queue
- Keyword search
 - User can configure "up to x" responses to keywords
- Responses to keyword queries come in waves; stops when x responses are found
- From user's perspective, service resembles Google, but provides links to MP3s and videos rather than Web pages

CS 640

24

KaZaA: Technology

Software

- Proprietary
- files and control data encrypted
- Hints:
 - KaZaA Web site gives a few
 - Some reverse engineering attempts described in Web
- Everything in HTTP request and response messages

Architecture

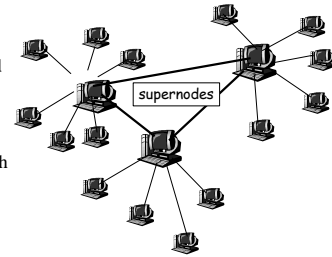
- hierarchical
- cross between Napster and Gnutella

CS 640

25

KaZaA: Architecture

- Each peer is either a supernode or is assigned to a supernode
- Each supernode knows about many other supernodes (almost mesh overlay)



CS 640

26

KaZaA: Architecture (2)

- Nodes that have more connection bandwidth and are more available are designated as supernodes
- Each supernode acts as a mini-Napster hub, tracking the content and IP addresses of its descendants
- Guess: supernode has (on average) 200-500 descendants; roughly 10,000 supernodes
- There is also dedicated user authentication server and supernode list server

CS 640

27

KaZaA: Overlay maintenance

- List of potential supernodes included within software download
- New peer goes through list until it finds operational supernode
 - Connects, obtains more up-to-date list
 - Node then pings 5 nodes on list and connects with the one with smallest RTT
- If supernode goes down, node obtains updated list and chooses new supernode

CS 640

28

KaZaA Queries

- Node first sends query to supernode
 - Supernode responds with matches
 - If x matches found, done.
- Otherwise, supernode forwards query to subset of supernodes
 - If total of x matches found, done.
- Otherwise, query further forwarded
 - Probably by original supernode rather than recursively

CS 640

29

Parallel Downloading; Recovery

- If file is found in multiple nodes, user can select parallel downloading
- Most likely HTTP byte-range header used to request different portions of the file from different nodes
- Automatic recovery when server peer stops sending file

CS 640

30

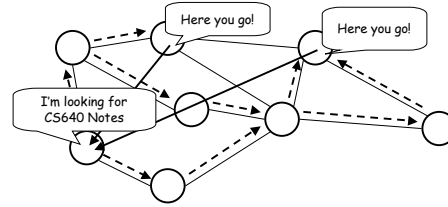
3. Structured P2P: DHT Approaches

- Want a storage and lookup service with better service guarantees and more efficient
- A Distributed Hash Table (DHT)
 - Chord
 - CAN
 - Pastry
 - Tapestry

CS 640

31

Challenge: Locating Content



- Simplest strategy: expanding ring search
- If K of N nodes have copy, expected search cost *at least* N/K , i.e., $O(N)$
- Need many cached copies to keep search overhead small

CS 640

32

Directed Searches

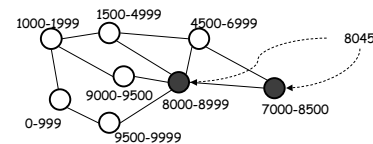
- Idea:
 - assign particular nodes to hold particular content (or pointers to it, like an information booth)
 - when a node wants that content, go to the node that is supposed to have or know about it
- Challenges:
 - Distributed: want to distribute responsibilities among existing nodes in the overlay
 - Adaptive: nodes join and leave the P2P overlay
 - distribute knowledge responsibility to joining nodes
 - redistribute responsibility knowledge from leaving nodes

CS 640

33

DHT Step 1: The Hash

- Introduce a hash function to map the object being searched for to a unique identifier:
 - e.g., $h(\text{"CS640 Class notes"}) \rightarrow 8045$
- Distribute the range of the hash function among all nodes in the network



- Each node must "know about" at least one copy of each object that hashes within its range (when one exists)

CS 640

34

DHT Step 2: Routing

- For each object, node(s) whose range(s) cover that object must be reachable via a "short" path
- by the querier node (assumed can be chosen arbitrarily)
- Different approaches (CAN, Chord, Pastry, Tapestry) differ fundamentally only in the routing approach
 - any "good" random hash function will suffice

CS 640

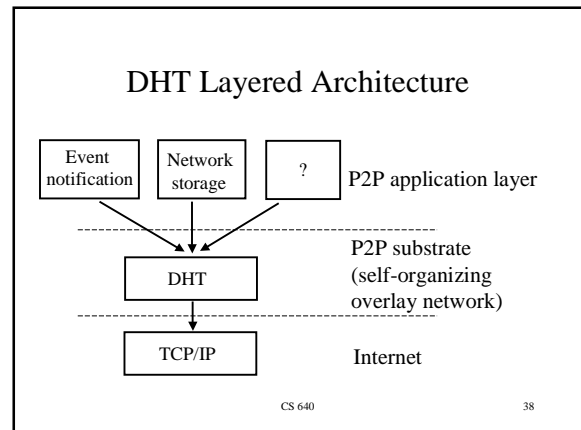
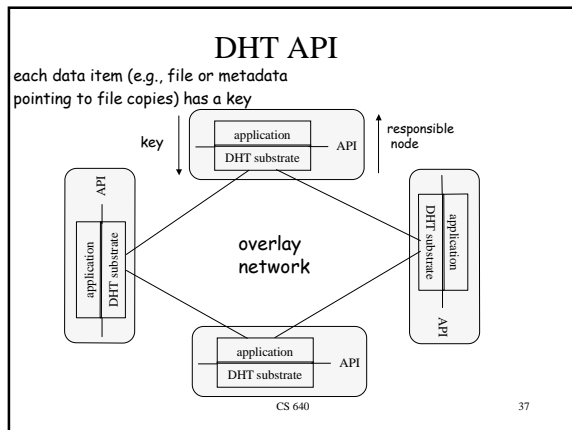
35

DHT API

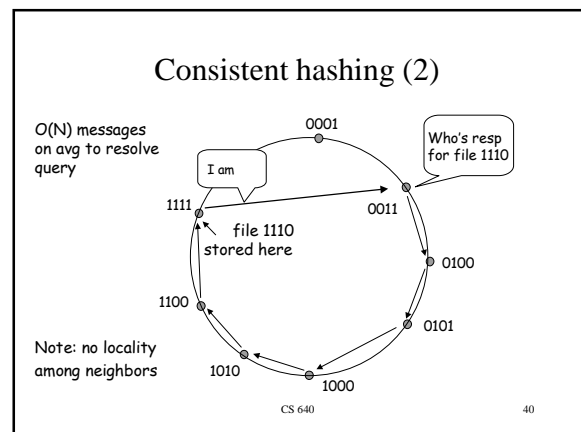
- each data item (e.g., file or metadata containing pointers) has a key in some ID space
- In each node, DHT software provides API:
 - Application gives API key k
 - API returns IP address of node that is responsible for k
- API is implemented with an underlying DHT overlay and distributed algorithms

CS 640

36



- ### Consistent hashing (1)
- Overlay network is a circle
 - Each node has randomly chosen id
 - Keys in same id space
 - Node's successor in circle is node with next largest id
 - Each node knows IP address of its successor
 - Key is stored in closest successor
- CS 640 39



- ### Consistent hashing (3)
- Node departures

 - Each node must track $s \geq 2$ successors
 - If your successor leaves, take next one
 - Ask your new successor for list of its successors; update your s successors

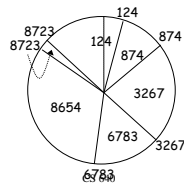
Node joins

 - You're new, node id k
 - ask any node n to find the node n' that is the successor for id k
 - Get successor list from n'
 - Tell your predecessors to update their successor lists
 - Thus, each node must track its predecessor
- CS 640 41

- ### Consistent hashing (4)
- Overlay is actually a circle with small chords for tracking predecessor and k successors
 - # of neighbors = $s+1$: $O(1)$
 - The ids of your neighbors along with their IP addresses is your "routing table"
 - average # of messages to find key is $O(N)$
- Can we do better?
- CS 640 42

Chord

- Nodes assigned 1-dimensional IDs in hash space at random (e.g., hash on IP address)
- Consistent hashing: Range covered by node is from previous ID up to its own ID (modulo the ID space)



43

Chord Routing

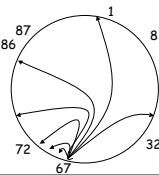
- A node s 's i^{th} neighbor has the ID that is equal to $s+2^i$ or is the next largest ID (mod ID space), $i \geq 0$
- To reach the node handling ID t , send the message to neighbor $\# \log_2(t-s)$
- Requirement: each node s must know about the next node that exists clockwise on the Chord (0^{th} neighbor)
- Set of known neighbors called a finger table

CS 640

44

Chord Routing (cont'd)

- A node s is node t 's neighbor if s is the closest node to $t+2^i \bmod H$ for some i . Thus,
 - each node has at most $\log_2 N$ neighbors
 - for any object, the node whose range contains the object is reachable from any node in no more than $\log_2 N$ overlay hops (each step can always traverse at least half the distance to the ID)
- Given K objects, with high probability each node has at most $(1 + \log_2 N) K / N$ in its range
- When a new node joins or leaves the overlay, objects move between nodes



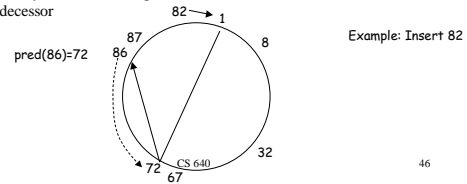
CS 640

i	Finger table for node 67	$(t-s)/N$
0	72	Closest node clockwise to
1	72	
2	72	
3	86	$67 \cdot 2^i \bmod 100$
4	86	
5	1	
6	32	

45

Chord Node Insertion

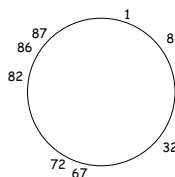
- One protocol addition: each node knows its closest counter-clockwise neighbor
- A node selects its unique (pseudo-random) ID and uses a bootstrapping process to find some node in the Chord
- Using Chord, the node identifies its successor in the clockwise direction
- An newly inserted node's predecessor is its successor's former predecessor



46

Chord Node Insertion (cont'd)

- First: set added node s 's fingers correctly
 - s 's predecessor t does the lookup for each distance of 2^i from s

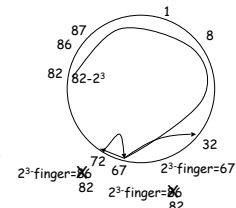


Lookups from node 72

i	Finger table for node 82
0	86
1	86
2	86
3	1
4	1
5	32
6	67

Chord Node Insertion (cont'd)

- Next, update other nodes' fingers about the entrance of s (when relevant). For each i :
 - Locate the closest node to s (counter-clockwise) whose 2^i -finger can point to s : largest possible is $s - 2^i$
 - Use Chord to go (clockwise) to largest node t before or at $s - 2^i$
 - route to $s - 2^i$; if arrived at a larger node, select its predecessor as t
 - If t 's 2^i -finger routes to a node larger than s
 - change t 's 2^i -finger to s
 - set $t = \text{predecessor of } t$ and repeat
 - Else $i++$, repeat from top
- $O(\log^2 N)$ time to find and update nodes



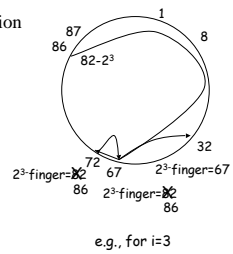
e.g., for $i=3$

CS 640

48

Chord Node Deletion

- Similar process can perform deletion



CS 640

49