

#### Typical Workload (Web Pages)



- Multiple (typically small) objects per page
- File sizes
- Heavy-tailed
  - · Pareto distribution for tail
  - Lognormal for body of distribution
- Embedded references
- •Lots of small objects means & TCP
- 3-way handshake
- Lots of slow starts
- Extra connection state
- Number of embedded objects also Pareto  $Pr(X>x) = (x/x_m)^{-k}$
- This plays havoc with performance. Why?
- · Solutions?

5

#### HTTP 0.9/1.0

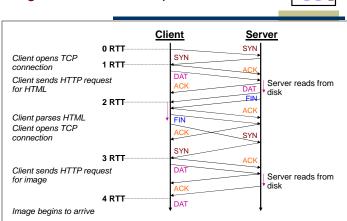


- One request/response per TCP connection
  - Simple to implement
- · Short transfers are very hard on TCP
  - Multiple connection setups → three-way handshake each time
    - Several extra round trips added to transfer
  - Many slow starts low throughput because of small window
    - Never leave slow start for short transfers
  - Loss recovery is poor when windows are small
  - · Lots of extra connections
    - Increases server state/processing

6

# Single Transfer Example

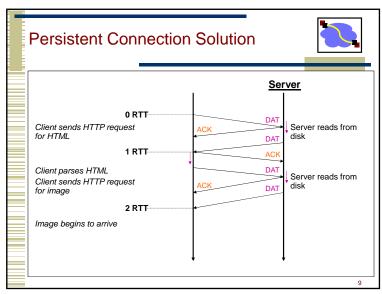


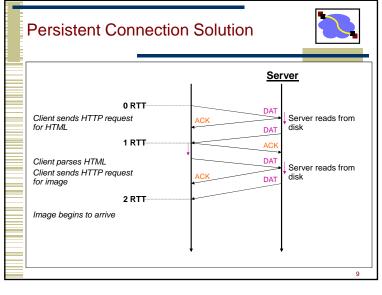


#### Improving Performance



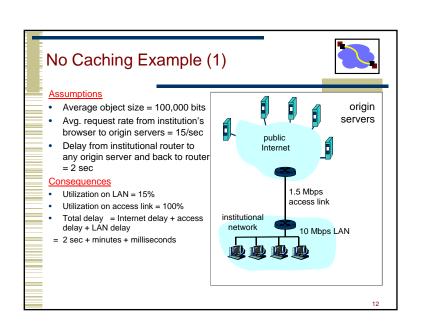
- Multiple concurrent connections (Netscape)
  - · Benefits are mixed: more state, timeouts, ...
- Multiplex multiple transfers onto one TCP connection (HTTP 1.1)
  - · Also allow pipelined transfers, i.e., multiple outstanding requests
- How to identify requests/responses
  - Delimiter → Server must examine response for delimiter string
  - Content-length and delimiter → Must know size of transfer in advance
  - Block-based transmission → send in multiple length delimited blocks
  - Store-and-forward → wait for entire response and then use content-length
  - Solution → use existing methods and close connection otherwise

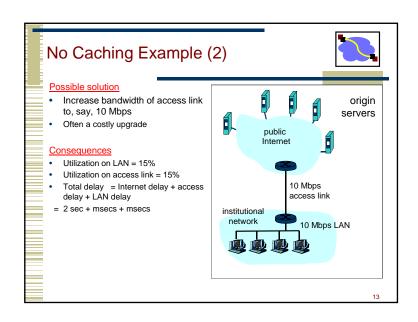


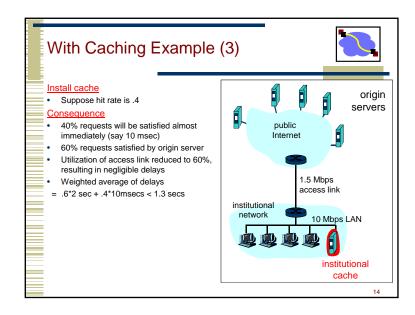


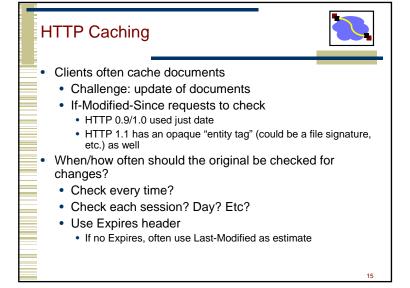
# Web Proxy Caches User configures browser: Web origin accesses via cache Browser sends all HTTP Proxy requests to cache · Object in cache: cache returns object Else cache requests object from origin server, then returns object to client origin server

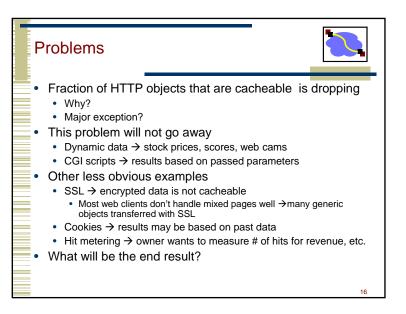
# Other Problems Serialized transmission but first bytes may be most useful May be better to get the 1st 1/4 of all images than one complete image (e.g., progressive JPEG) • Can "packetize" transfer over TCP, e.g., range requests Application specific solution to transport protocol problems. :( Could fix TCP so it works well with multiple simultaneous connections · More difficult to deploy Persistent connections can significantly increase state on the server Allow server to close HTTP session · Client can no longer send requests



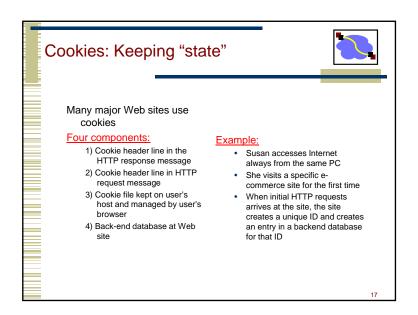


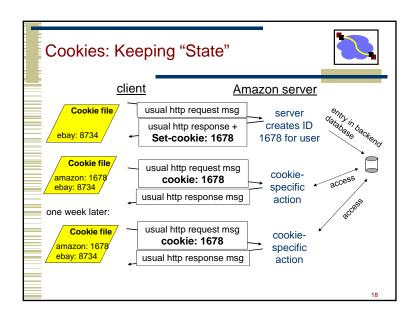


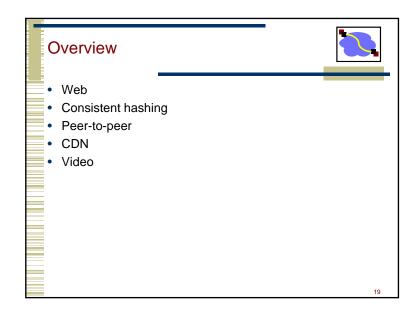


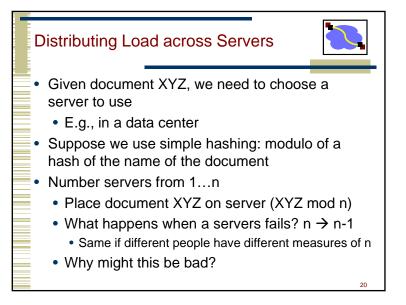


Δ









#### Consistent Hash: Goals



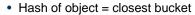
- "view" = subset of all hash buckets that are candidate locations
  - Correspond to a real server
- Desired features
  - Load all hash buckets have a similar number of objects assigned to them
  - Smoothness little impact on hash bucket contents when buckets are added/removed
  - Spread small set of hash buckets that may hold an object regardless of views

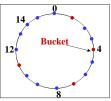
21

#### 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 unit interval





- 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

22

### Consistent Hashing: Ring



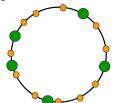
- Use consistent has to map both keys and nodes to an m-bit identifier in the same (metric) identifier space
  - · For example, use SHA-1 hashes
  - Node identifier: SHA-1 hash of IP address

IP="198.10.10.1" <u>SHA-1</u> ID=123

• Key identifier: SHA-1 hash of key

 $Key = \text{``LetItBe''} \quad \underline{SHA-1} \qquad ID = 60$ 

- Also need "rule" for assigning keys to nodes
  - For example: "closest", higher, lower, ..



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

RP="198.10.10.1"

O K5

Circular 7-bit
ID space

N90

Key="LetItBe"

#### **Consistent Hashing Properties**



- Load balance: all nodes receive roughly the same number of keys
  - For *N* nodes and *K* keys, with high probability
    - Each node holds at most (1+ε)K/N keys
    - · Provided that K is large compared to N
- When server is added, it receives its initial work load from "neighbors" on the ring
  - "Local" operation: no other servers are affected
  - · Similar property when a server is removed

25

#### Finer Grain Load Balancing



- · Redirector knows all server IDs
- It can also track approximate "load" for more precise load balancing
  - Need to define load and be able to track it
- To balance load:
  - W<sub>i</sub> = Hash(URL, ip of s<sub>i</sub>) for all i
  - Sort W<sub>i</sub> from high to low
  - Find first server with low enough load
- Benefits and drawbacks?

26

# Consistent Hashing Used in Many Contexts



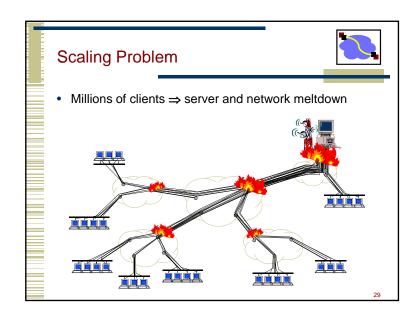
- Distribute load across servers in a data center
  - · The redirector sits in data center
- Finding storage cluster for an object in a CDN uses centralized knowledge
  - Why?
  - Can use consistent hashing in the cluster
- Consistent hashing can also be used in a distributed setting
  - P2P systems can use it find files

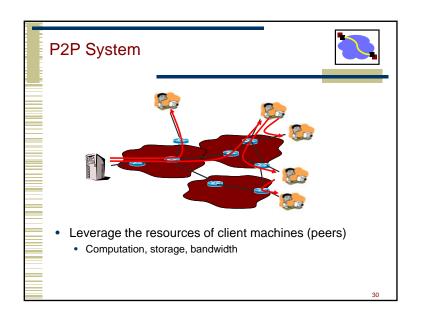
7

#### Overview



- Web
- · Consistent hashing
- Peer-to-peer
  - Motivation
  - Architectures
  - TOR
  - Skype
- CDN
- Video





## Why p2p?



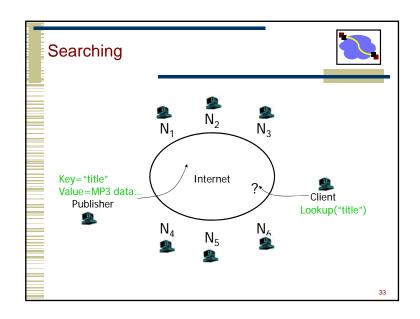
- Harness lots of spare capacity
  - 1 Big Fast Server: 1Gbit/s, \$10k/month++
  - 2,000 cable modems: 1Gbit/s, \$ ??
  - 1M end-hosts: Uh, wow.
  - Capacity grows with the number of users!
- Build very large-scale, self-managing systems
  - Same techniques useful for companies and p2p apps
    - E.g., Akamai's 14,000+ nodes, Google's 100,000+ nodes
  - · Many differences to consider
    - · Servers versus arbitrary nodes
    - Hard state (backups!) versus soft state (caches)
    - · Security, fairness, freeloading, ..

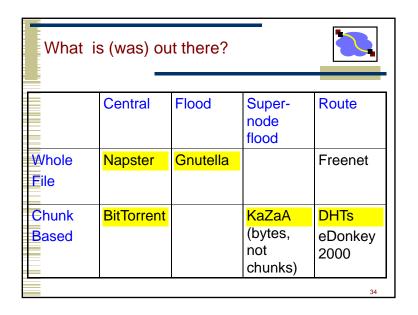
1

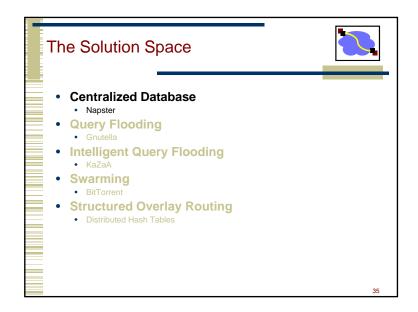
#### Common P2P Framework

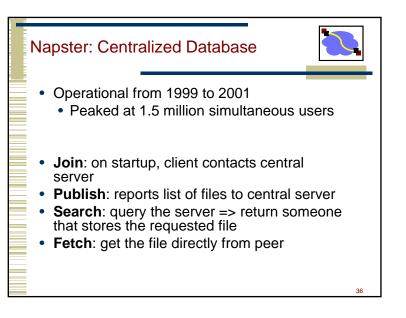


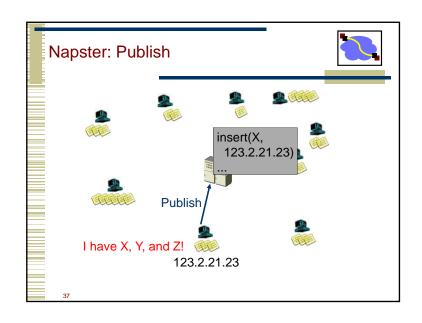
- Common Primitives:
  - Join: how to I begin participating?
  - Publish: how do I advertise my file?
  - Search: how to I find a file?
  - **Fetch**: how to I retrieve a file?
- · Search tends to be the most challenging:
  - Needles vs. Haystacks
    - Searching for top 40, or an obscure punk track from 1981 that nobody ever heard of?
  - Search expressiveness: Whole word? Regular expressions? File names? Attributes? Whole-text? ...
    - E.g., p2p gnutella or p2p google?

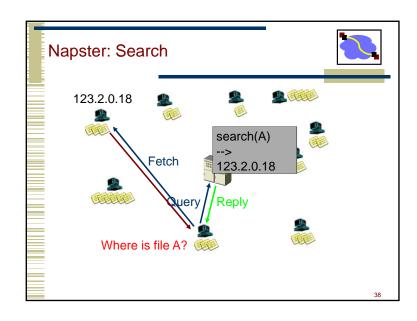


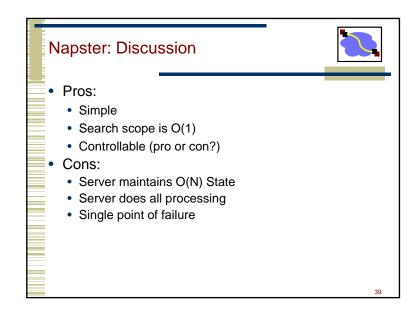














#### **Gnutella: Flooding**



- Released in 2000 and improved over time; still alive
  - · Many clients available, very popular
- 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

Gnutella: Search

I have file A.

Query

Where is file A?

#### **Gnutella: Discussion**



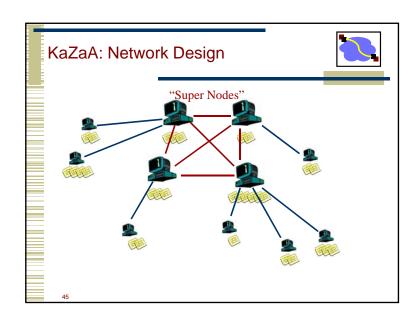
- Pros:
  - Fully de-centralized
  - Search cost distributed
  - Processing @ each node permits powerful search semantics
- Cons:
  - Search scope is O(N)
  - Search time is O(???)
  - · Nodes leave often, network unstable
- TTL-limited search works well for haystacks.
  - For scalability, does NOT search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

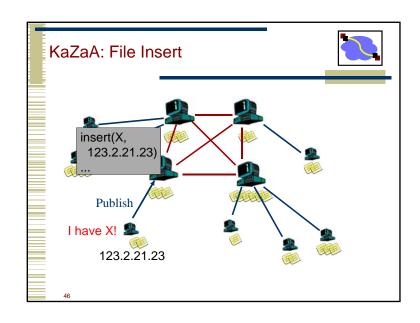
2

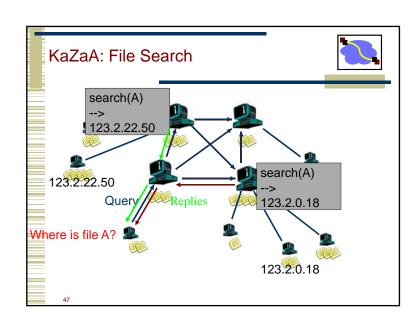
### KaZaA: Query Flooding

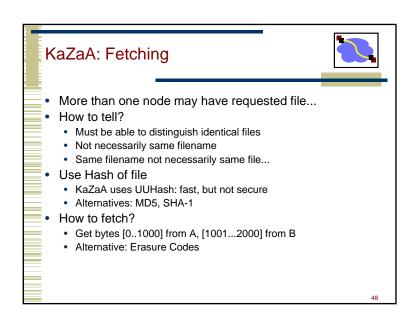


- First released in 2001 and still used today
  - Also very popular
- **Join**: on startup, client contacts a "supernode" ... may at some point become one itself
- Publish: send list of files to supernode
- Search: send query to supernode, supernodes flood query amongst themselves.
- Fetch: get the file directly from peer(s); can fetch simultaneously from multiple peers









#### KaZaA: Discussion



- Pros:
  - · Tries to take into account node heterogeneity:
    - Bandwidth
    - Host Computational Resources
    - Host Availability (?)
  - · Rumored to take into account network locality
- · Cons:
  - · Mechanisms easy to circumvent
  - · Still no real guarantees on search scope or search time
- Similar behavior to gnutella, but better.

49

# 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 have been on is a good predictor of how long you will be around

50

## The Solution Space



- Centralized Database
  - Napster
- Query Flooding
- Coutoll
- Intelligent Query Flooding
  - KaZaA
- Swarming
  - BitTorrent
- Structured Overlay Routing
  - Distributed Hash Tables
- More on Thursday ...

Napster: History



- 1999: Sean Fanning launches Napster
- Peaked at 1.5 million simultaneous users
- Jul 2001: Napster shuts down

# Gnutella: History



- In 2000, J. Frankel and T. Pepper from Nullsoft released Gnutella
- Soon many other clients: Bearshare, Morpheus, LimeWire, etc.
- In 2001, many protocol enhancements including "ultrapeers"

# KaZaA: History



- In 2001, KaZaA created by Dutch company Kazaa BV
- Single network called FastTrack used by other clients as well: Morpheus, giFT, etc.
- Eventually protocol changed so other clients could no longer talk to it
- Very popular file sharing network today with >10 million users (number varies)