

Peer-peer and Application-level Networking

Presented by Jon Crowcroft

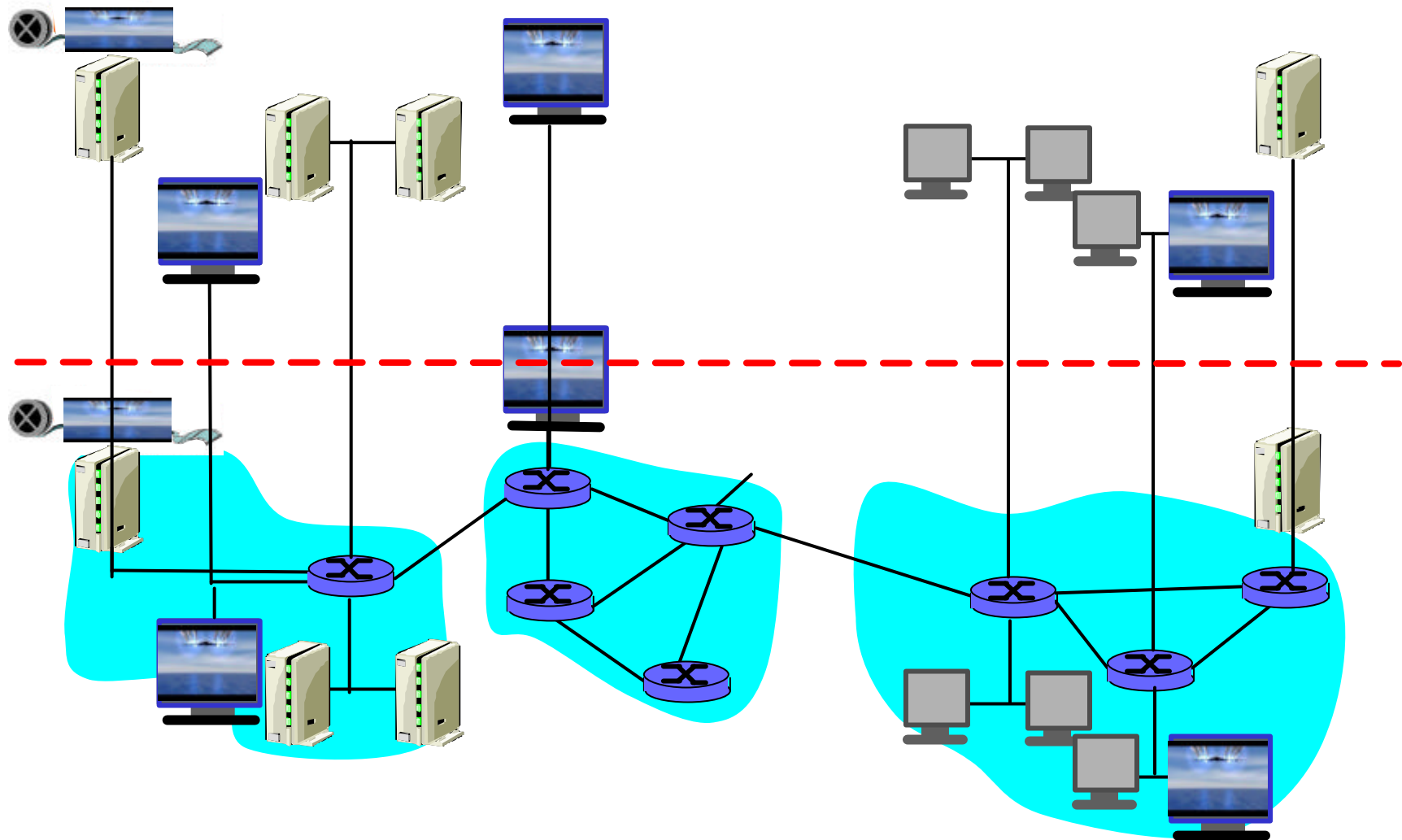
Based ***strongly*** on material by

Jim Kurose, Brian Levine, Don Towsley, and
the class of 2001 for the Umass Comp Sci
791N course..

0. Introduction

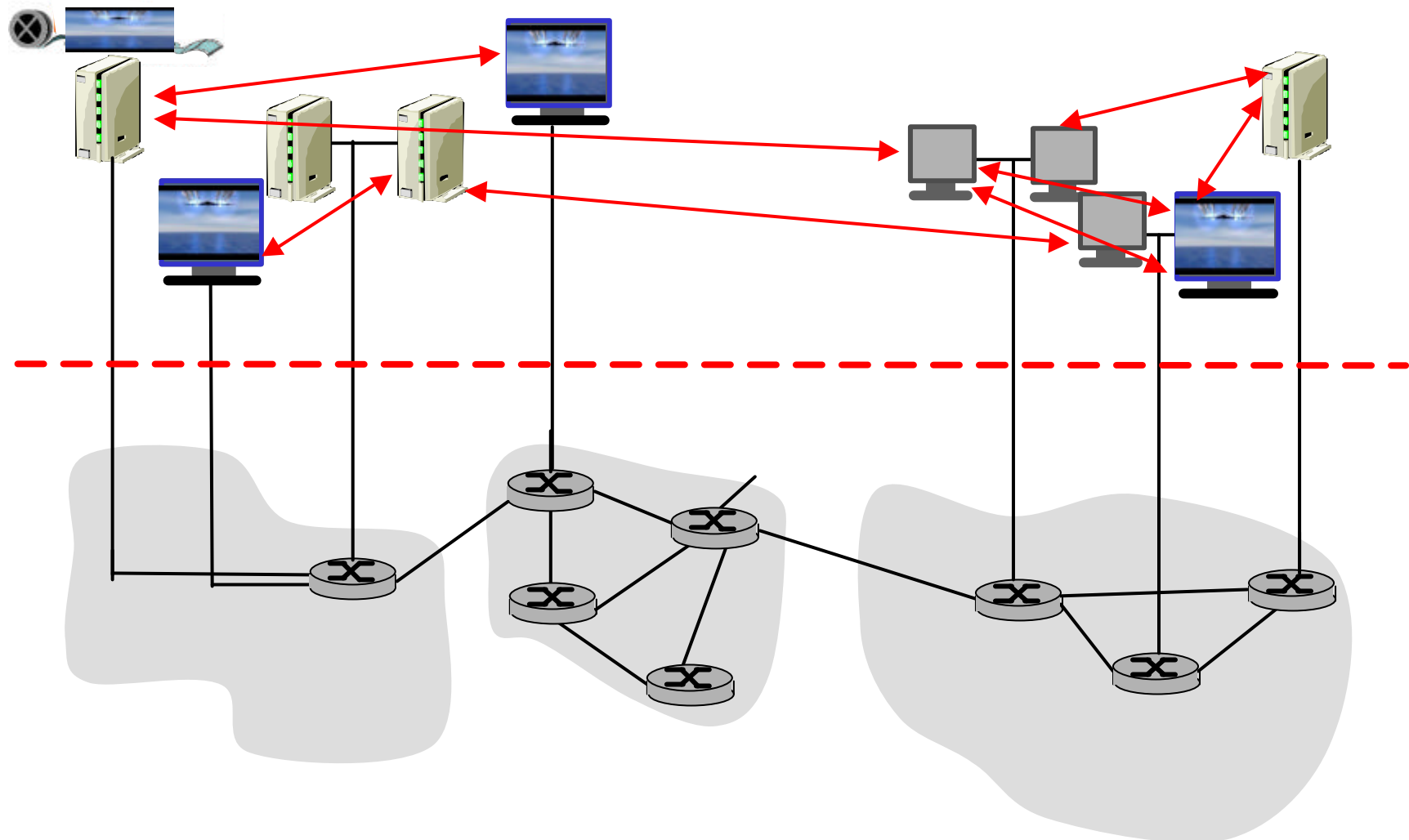
- r Background
- r Motivation
- r outline of the tutorial

Peer-peer networking



Peer-peer networking

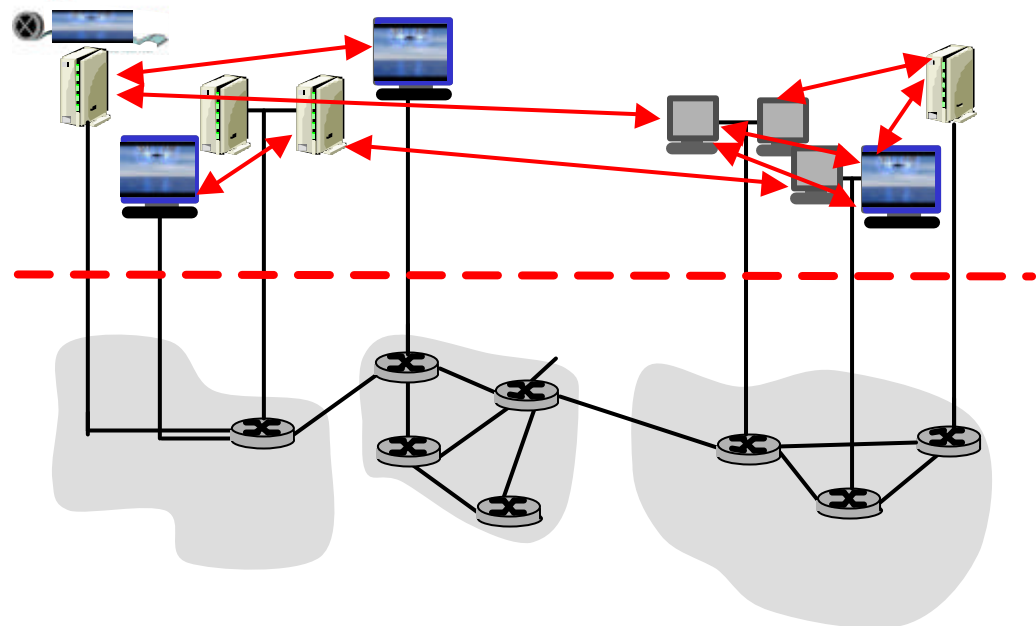
Focus at the application level



Peer-peer networking

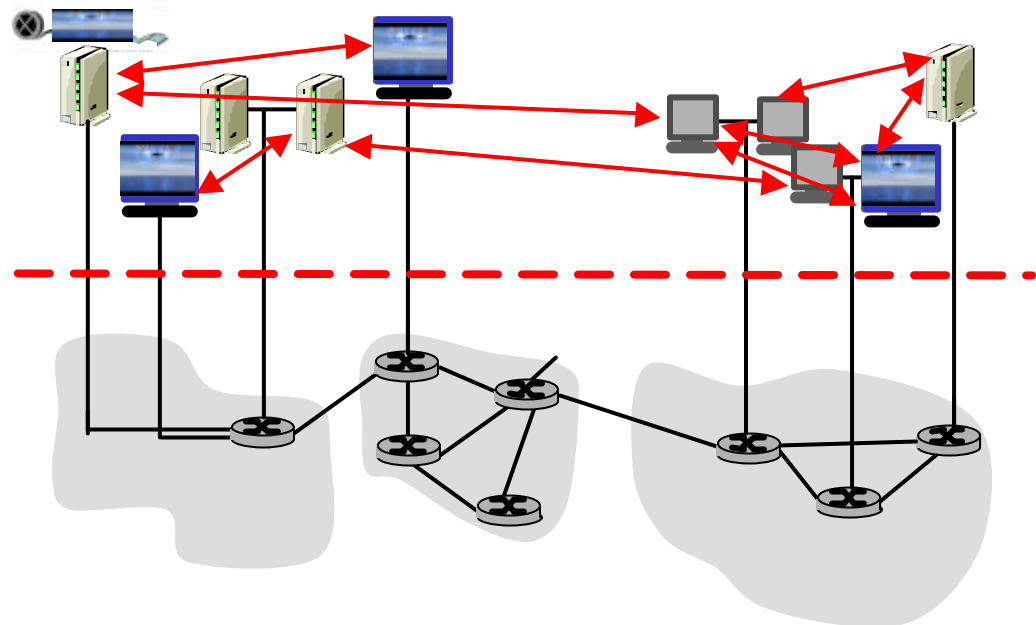
Peer-peer applications

- r Napster, Gnutella, Freenet: file sharing
- r ad hoc networks
- r multicast overlays (e.g., video distribution)



Peer-peer networking

- r Q: What are the new technical challenges?
- r Q: What new services/applications enabled?
- r Q: Is it just “networking at the application-level”?
 - m “There is nothing new under the Sun” (William Shakespeare)



Tutorial Contents

- r Introduction
- r Client-Server v. P2P
- r Case Studies
 - m Napster
 - m Gnutella
 - m RON
 - m Freenet
 - m Publius
- r Middleware
 - m Chord
 - m CAN
 - m Tapestry
 - m JXTA
 - m ESM
 - m Overcast
- r Applications
 - m Storage
 - m Conferencing

Client Server v. Peer to Peer(1)

- r RPC/RMI
- r Synchronous
- r Asymmetric
- r Emphasis on language integration and binding models (stub *IDL/XDR* compilers etc)
- r Kerberos style security – access control, crypto
- r Messages
- r Asynchronous
- r Symmetric
- r Emphasis on service location, content addressing, application layer routing.
- r Anonymity, high availability, integrity.
- r Harder to get right ✍

Client Server v. Peer to Peer(2)

RPC

```
Cli_call(args)
```

```
Srv_main_loop()  
{  
    while(true) {  
        deque(call)  
        switch(call.procid)  
        case 0:  
            call.ret=proc1(call.args)  
        case 1:  
            call.ret=proc2(call.args)  
        ... ..  
        default:  
            call.ret = exception  
        }  
    }  
}
```

Client Server v. Peer to Peer(3)

P2P

```
Peer_main_loop()  
{  
    while(true) {  
        await(event)  
        switch(event.type) {  
            case timer_expire: do some p2p work()  
                                randomize timers  
                                break;  
            case inbound message: handle it  
                                respond  
                                break;  
            default: do some book keeping  
                    break;  
        }  
    }  
}
```

Peer to peer systems actually old

- r IP routers are peer to peer.
- r Routers discover topology, and maintain it
- r Routers are neither client nor server
- r Routers continually chatter to each other
- r Routers are fault tolerant, inherently
- r Routers are autonomous

Peer to peer systems

- r Have no distinguished role
- r So no single point of bottleneck or failure.
- r However, this means they need distributed algorithms for
 - m Service discovery (name, address, route, metric, etc)
 - m Neighbour status tracking
 - m Application layer routing (based possibly on content, interest, etc)
 - m Resilience, handling link and node failures
 - m Etc etc etc

Ad hoc networks and peer2peer

- r Wireless ad hoc networks have many similarities to peer to peer systems
- r No *a priori knowledge*
- r No given infrastructure
- r Have to construct it from "thin air"!
- r Note for later – wireless ✍

Overlays and peer 2 peer systems

- r P2p technology is often used to create overlays which offer services that could be offered in the I P level
- r Useful **deployment** strategy
- r Often economically a way around other barriers to deployment
- r I P I tself was an overlay (on telephone core infrastructure)
- r Evolutionary path!!!

Next we look at some case studies

- r Piracy^{H^H^H^H^H}content sharing ✎
- r Napster
- r Gnutella
- r Freenet
- r Publius
- r etc

1. NAPSTER

- r The most (in)famous
- r Not the first (c.f. probably Eternity, from Ross Anderson in Cambridge)
- r But instructive for what it gets right, and
- r Also wrong...
- r Also has a political message...and economic and legal...etc etc etc

Napster

- r program for sharing files over the Internet
- r a “disruptive” application/technology?
- r history:

- m 5/99: Shawn Fanning (freshman, Northeastern U.) founds Napster Online music service

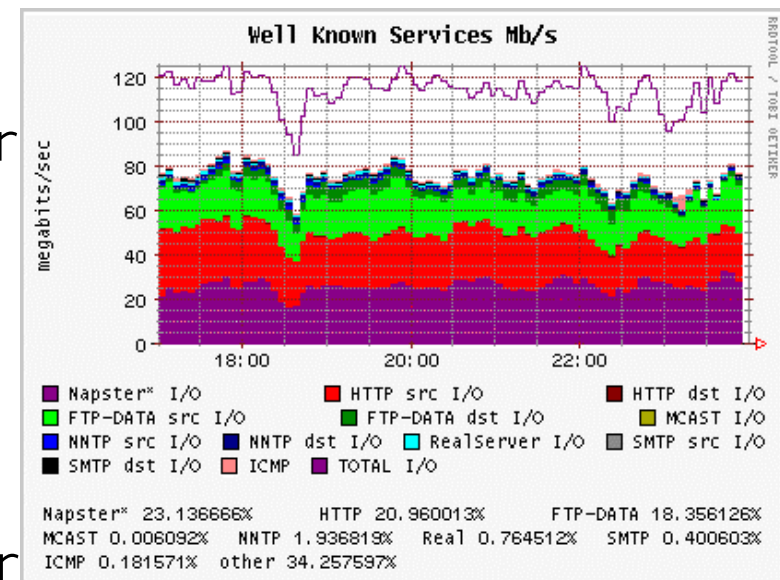
- m 12/99: first lawsuit

- m 3/00: 25% UWisc traffic Napster

- m 2000: est. 60M users

- m 2/01: US Circuit Court of Appeals: Napster knew users violating copyright laws

- m 7/01: # simultaneous online users: Napster 160K, Gnutella: 40K, Mor



Napster: how does it work

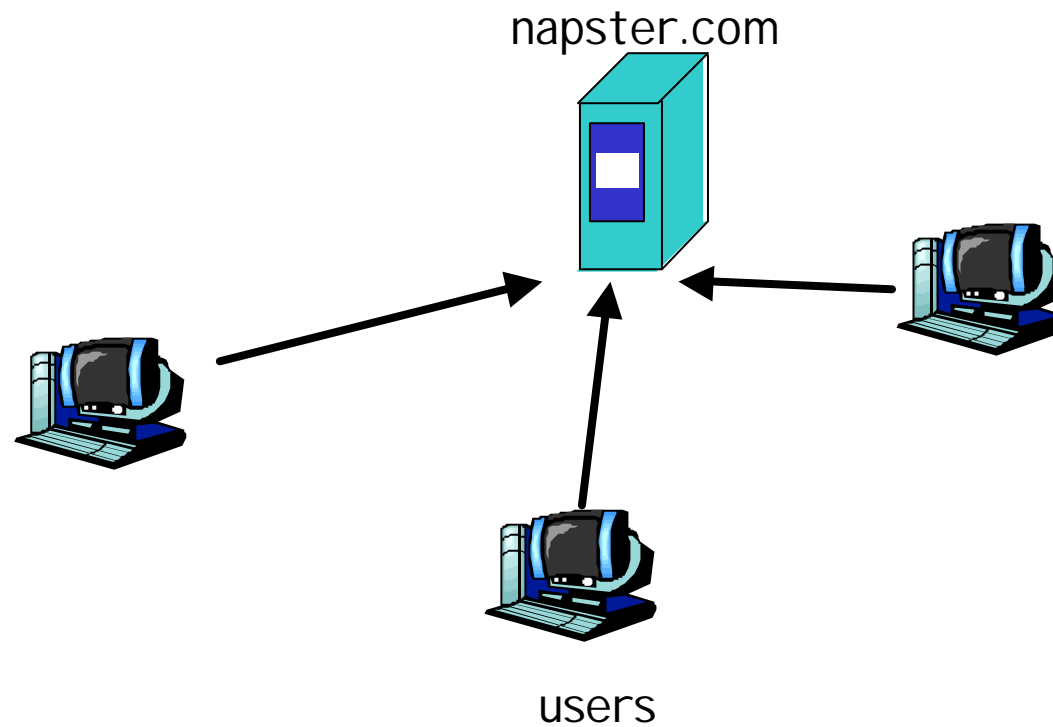
Application-level, client-server protocol over point-to-point TCP

Four steps:

- r Connect to Napster server
- r Upload your list of files (push) to server.
- r Give server keywords to search the full list with.
- r Select "best" of correct answers. (pings)

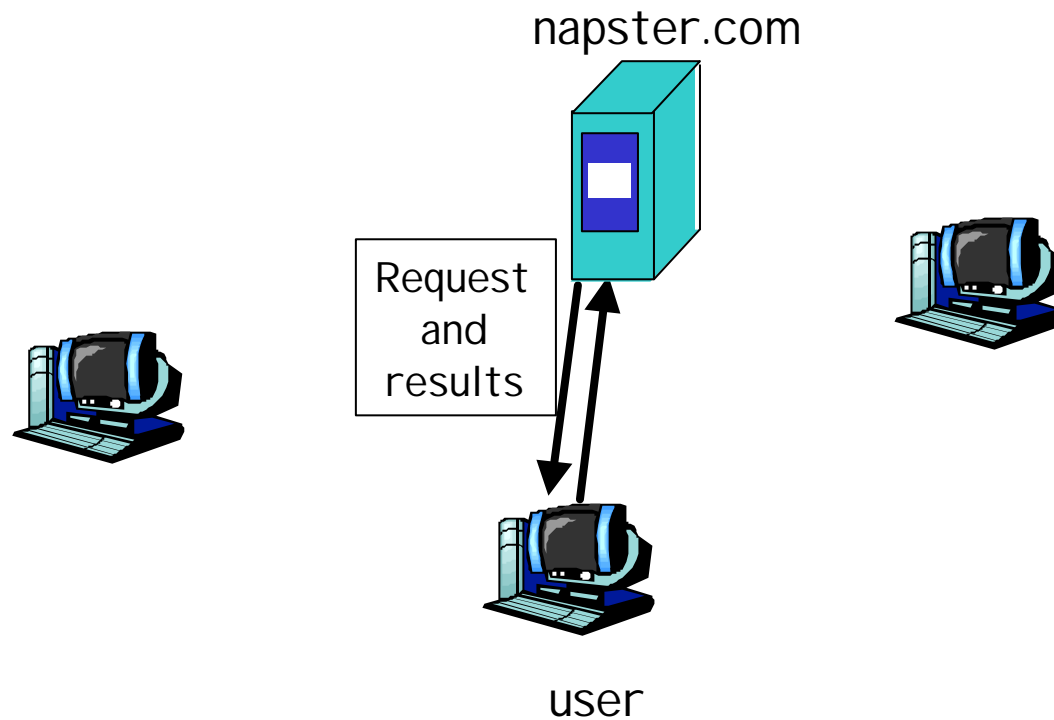
Napster

1. File list is uploaded



Napster

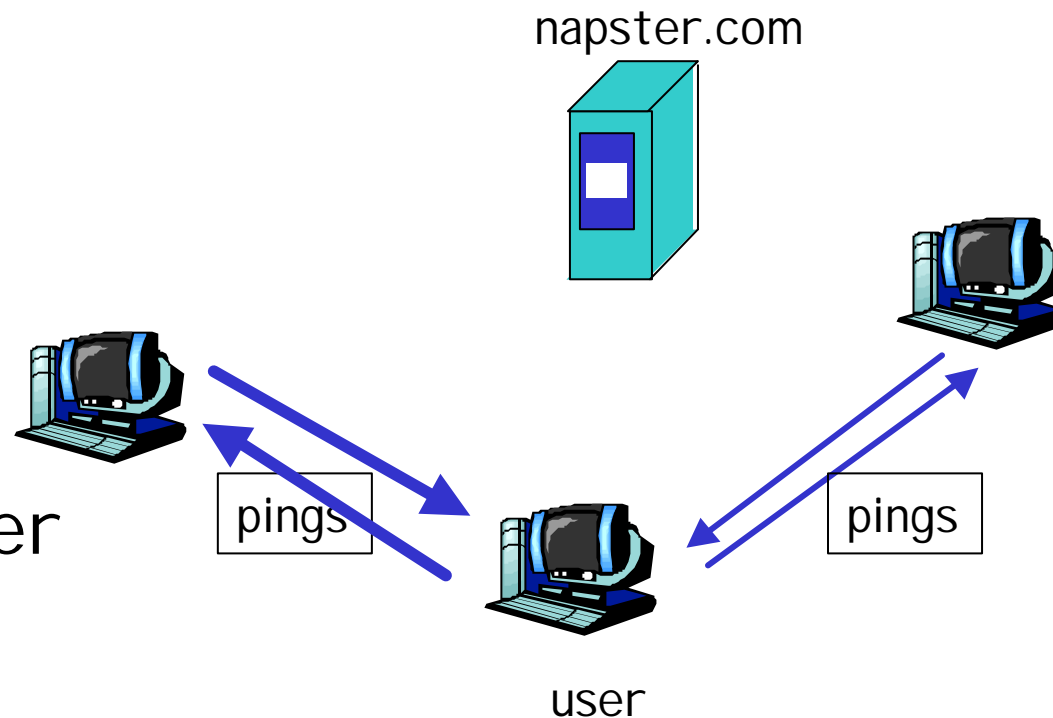
2. User requests search at server.



Napster

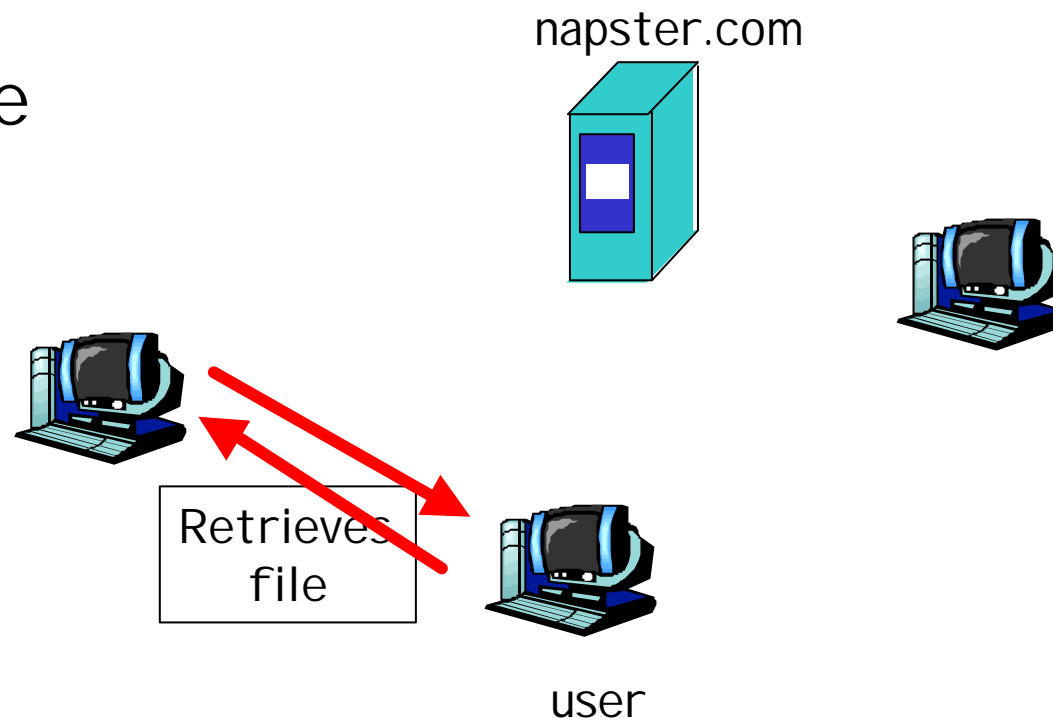
3. User pings hosts that apparently have data.

Looks for **best** transfer rate.



Napster

4. User retrieves file



Napster messages

General Packet Format

[chunksize] [chunkinfo] [data...]

CHUNKSIZE:

Intel-endian 16-bit integer
size of [data...] in bytes

CHUNKINFO: (hex)

Intel-endian 16-bit integer.

00 - login rejected	5B - whois query
02 - login requested	5C - whois result
03 - login accepted	5D - whois: user is offline!
0D - challenge? (nuprin1715)	69 - list all channels
2D - added to hotlist	6A - channel info
2E - browse error (user isn't online!)	90 - join channel
2F - user offline	91 - leave channel

.....

Napster: requesting a file

SENT to server (after logging in to server)

2A 00 CB 00 username

"C:\MP3\REM - Everybody Hurts.mp3"

RECEIVED

5D 00 CC 00 username

2965119704 (IP-address backward-form = A.B.C.D)

6699 (port)

"C:\MP3\REM - Everybody Hurts.mp3" (song)

(32-byte checksum)

(line speed)

[connect to A.B.C.D:6699]

RECEIVED from client

31 00 00 00 00 00

SENT to client

GET

RECEIVED from client

00 00 00 00 00 00

SENT to client

Myusername

"C:\MP3\REM - Everybody Hurts.mp3"

0 (port to connect to)

RECEIVED from client

(size in bytes)

SENT to server

00 00 DD 00 (give go-ahead thru server)

RECEIVED from client

[DATA]

Napster: architecture notes

- r centralized server:
 - m single logical point of failure
 - m can load balance among servers using DNS rotation
 - m potential for congestion
 - m Napster “in control” (freedom is an illusion)
- r no security:
 - m passwords in plain text
 - m no authentication
 - m no anonymity

2 Gnutella

- r Napster fixed
- r Open Source
- r Distributed
- r Still very political...

Gnutella

- r peer-to-peer networking: applications connect to peer applications
- r focus: decentralized method of searching for files
- r each application instance serves to:
 - m store selected files
 - m route queries (file searches) from and to its neighboring peers
 - m respond to queries (serve file) if file stored locally
- r Gnutella history:
 - m 3/14/00: release by AOL, almost immediately withdrawn
 - m too late: 23K users on Gnutella at 8 am this AM
 - m many iterations to fix poor initial design (poor design turned many people off)

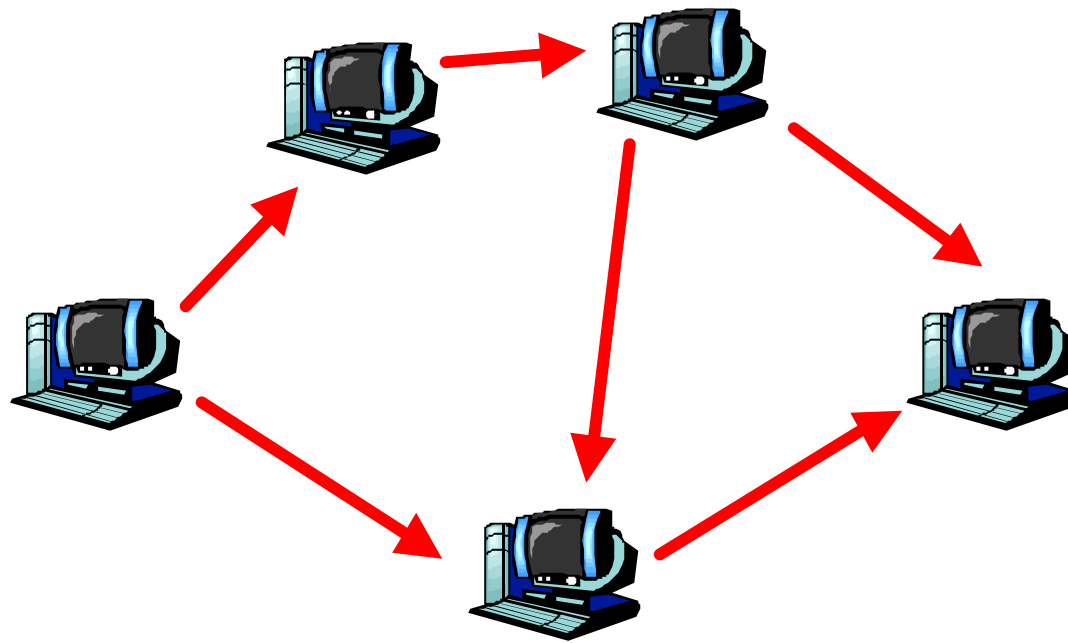
Gnutella: how it works

Searching by flooding:

- r If you don't have the file you want, query 7 of your partners.
- r If they don't have it, they contact 7 of their partners, for a maximum hop count of 10.
- r Requests are flooded, but there is no tree structure.
- r No looping but packets may be received twice.
- r Reverse path forwarding(?)

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

Flooding in Gnutella: loop prevention



Seen already list: "A"

Gnutella message format

- r **Message ID:** 16 bytes (yes bytes)
- r **FunctionID:** 1 byte indicating
 - m 00 ping: used to probe gnutella network for hosts
 - m 01 pong: used to reply to ping, return # files shared
 - m 80 query: search string, and desired minimum bandwidth
 - m 81: query hit: indicating matches to 80:query, my IP address/port, available bandwidth
- r **RemainingTTL:** decremented at each peer to prevent TTL-scoped flooding
- r **HopsTaken:** number of peer visited so far by this message
- r **DataLength:** length of data field

Gnutella: initial problems and fixes

- r Freeloading: WWW sites offering search/retrieval from Gnutella network without providing file sharing or query routing.
 - m Block file-serving to browser-based non-file-sharing users
- r Prematurely terminated downloads:
 - m long download times over modems
 - m modem users run gnutella peer only briefly (Napster problem also!) or any users becomes overloaded
 - m fix: peer can reply "I have it, but I am busy. Try again later"
 - m late 2000: only 10% of downloads succeed
 - m 2001: more than 25% downloads successful (is this success or failure?)

Gnutella: initial problems and fixes (more)

- r 2000: avg size of reachable network only 400-800 hosts. Why so small?
 - m **modem users**: not enough bandwidth to provide search routing capabilities: routing black holes
- r **Fix**: create peer hierarchy based on capabilities
 - m previously: all peers identical, most modem blackholes
 - m connection preferencing:
 - favors routing to well-connected peers
 - favors reply to clients that themselves serve large number of files: prevent freeloading
 - m Limewire gateway functions as Napster-like central server on behalf of other peers (for searching purposes)

Anonymous?

- r Not anymore than it's scalable.
- r The person you are getting the file from knows who you are. That's not anonymous.
- r Other protocols exist where the owner of the files doesn't know the requester.
- r Peer-to-peer anonymity exists.
- r See Eternity and, next, Freenet!

Gnutella Discussion:

- r Architectural lessons learned?
- r Do Gnutella's goals seem familiar? Does it work better than say squid or summary cache? Or multicast with carousel?
- r anonymity and security?
- r Other?
- r Good source for technical info/open questions:
http://www.limewire.com/index.jsp/tech_papers

3. Overlays

- r Next, we need to look at overlays in general, and more specifically, at
- r Routing...
- r RON is a good example...

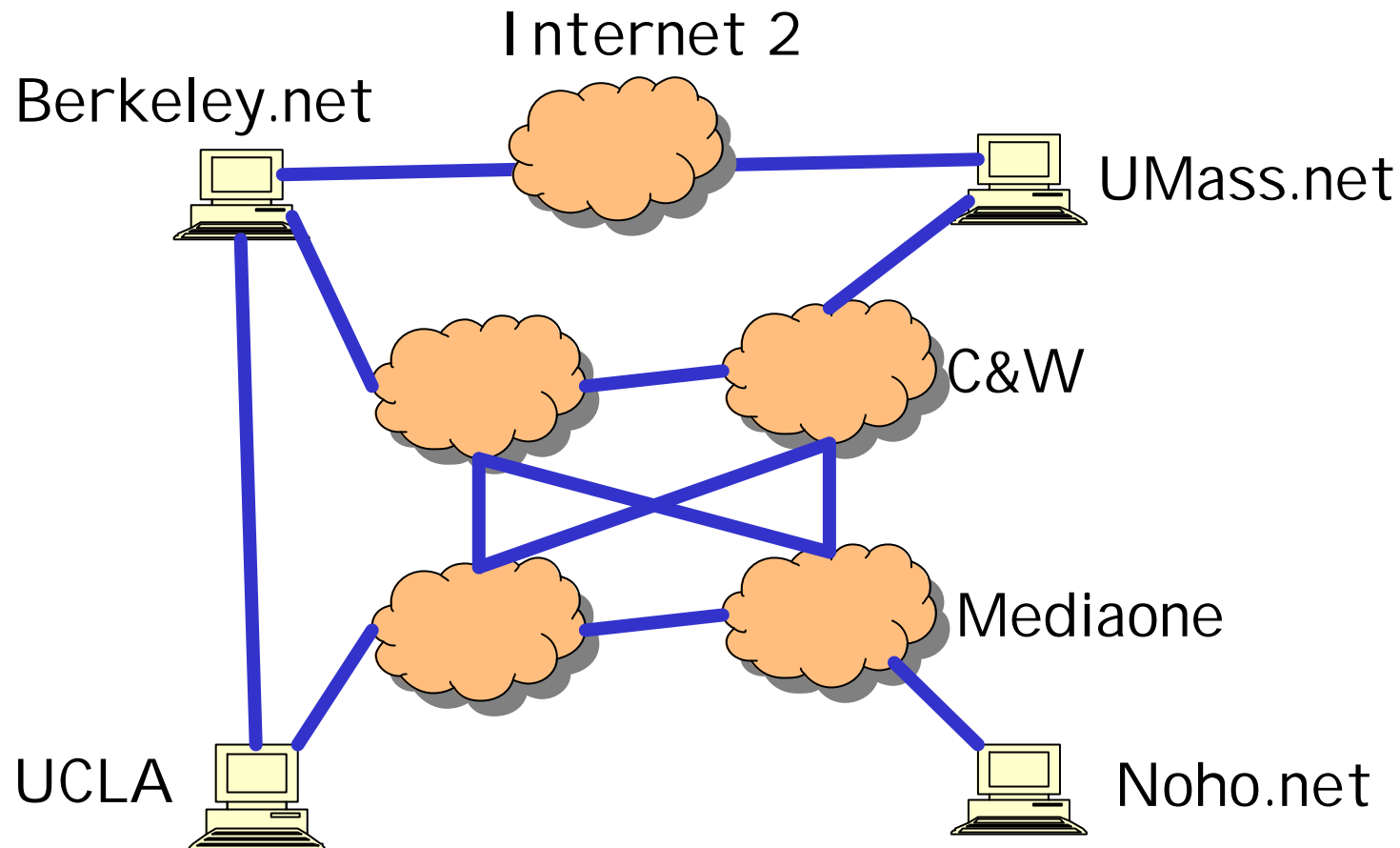
Resilient Overlay Networks

Overlay network:

- r applications, running at various sites
- r create “logical” links (e.g., TCP or UDP connections) pairwise between each other
- r each logical link: multiple physical links, routing defined by native Internet routing
- r let's look at an example, taken from:
 - r D. Anderson, H. Balakrishnan, F. Kaashoek, R. Morris, "The case for resilient overlay networks," Proc. HotOS VIII, May 2001, <http://nms.lcs.mit.edu/papers/ron-hotos2001.html>.

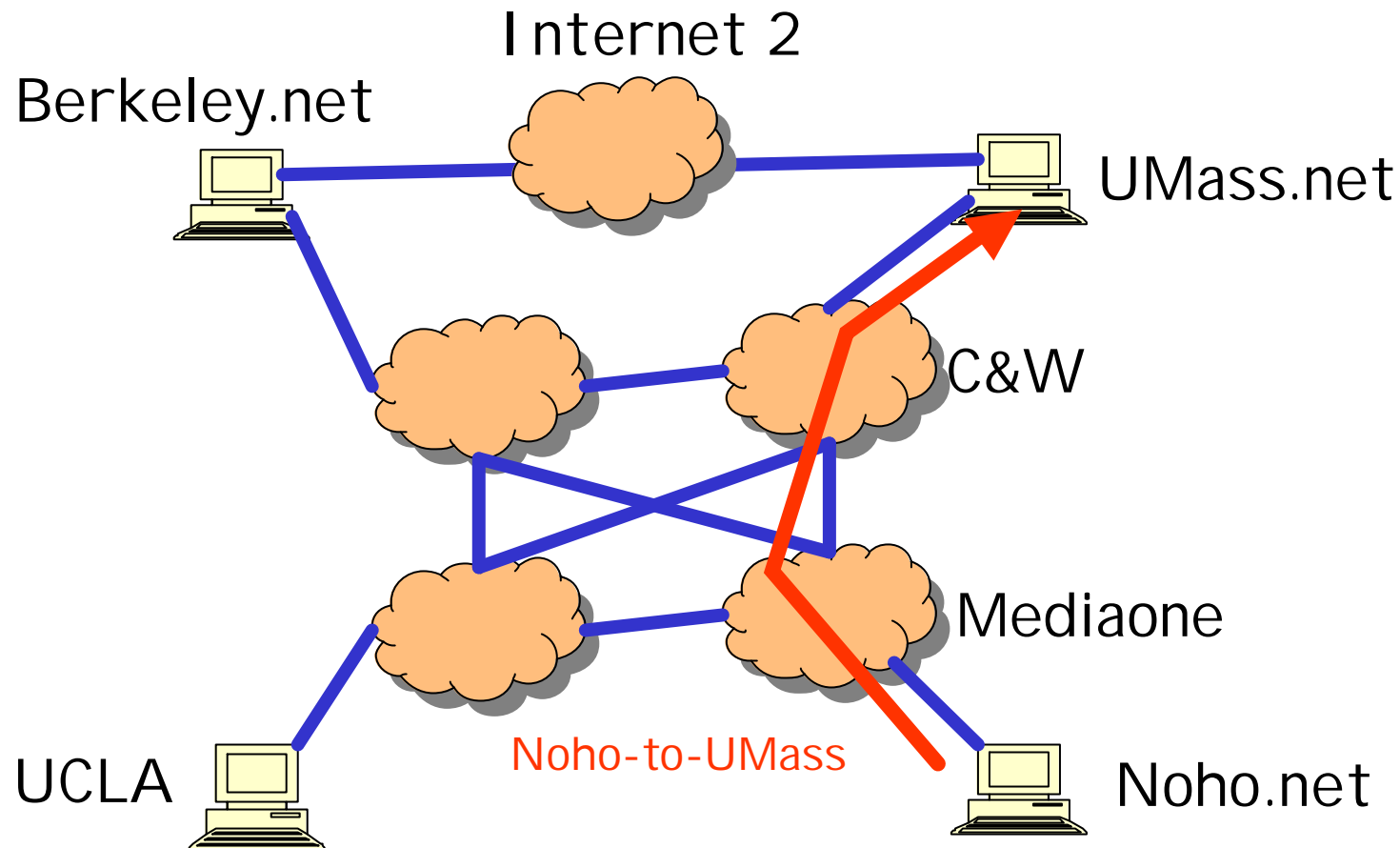
Internet Routing

- r BGP defines routes between stub networks



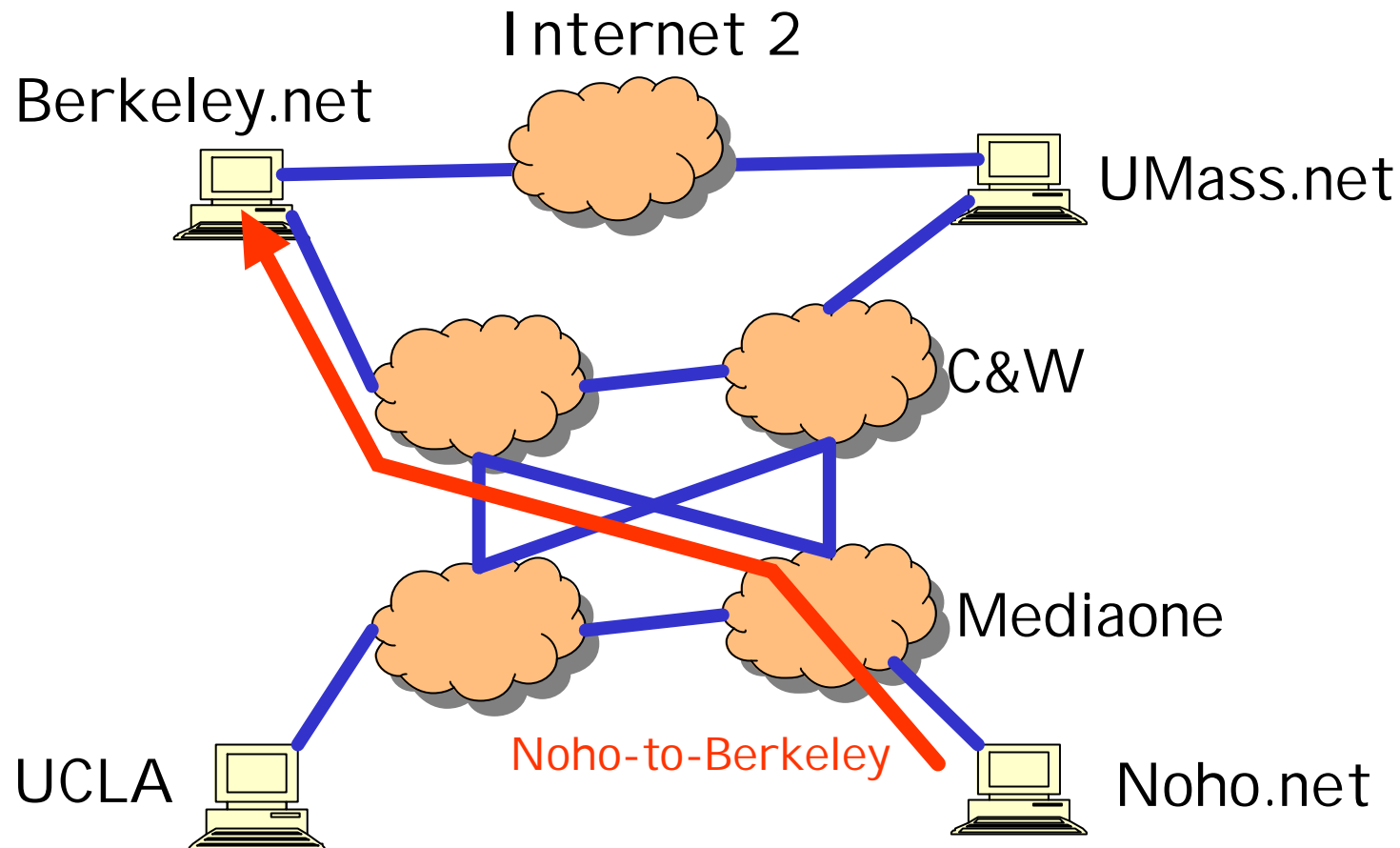
Internet Routing

- r BGP defines routes between stub networks

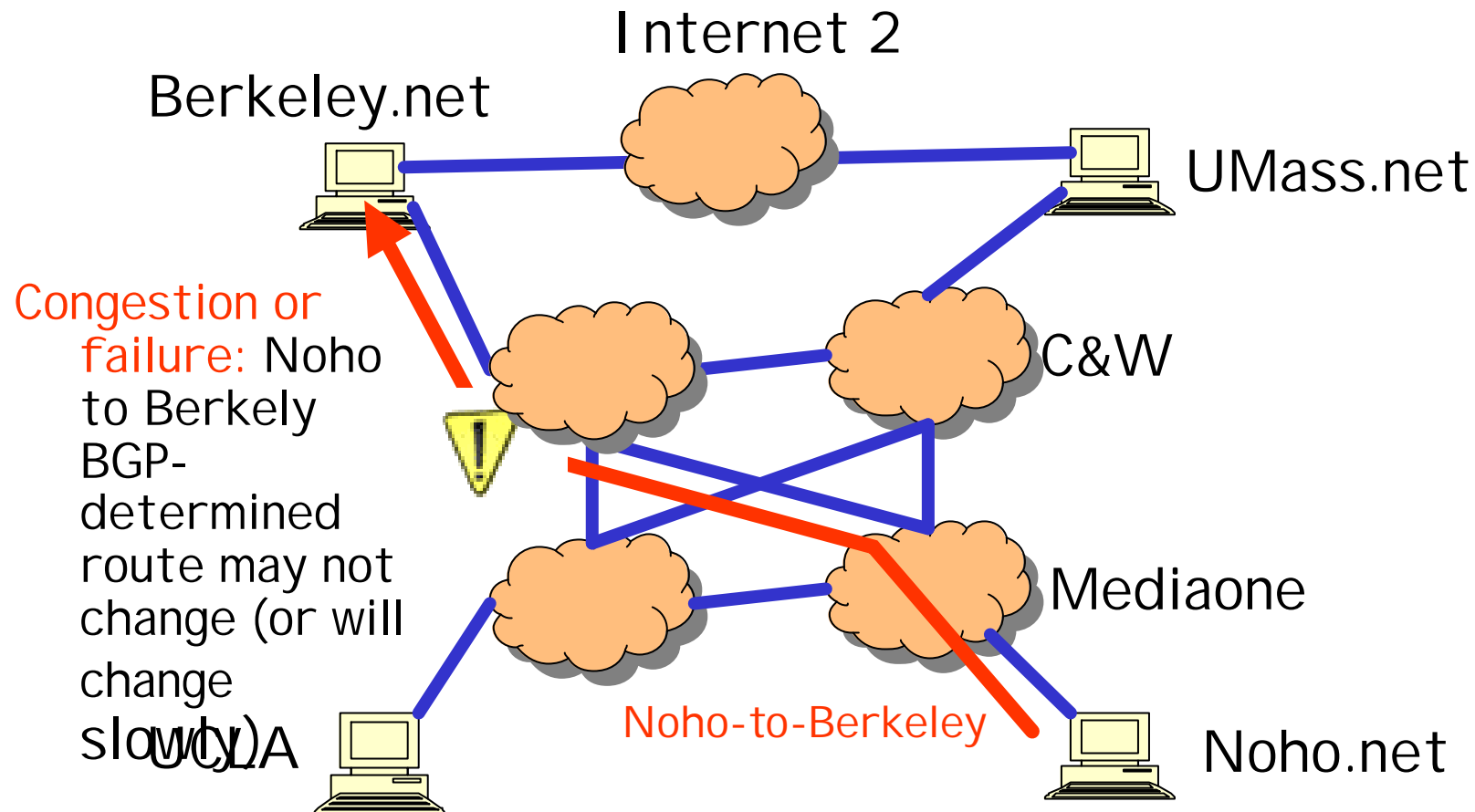


Internet Routing

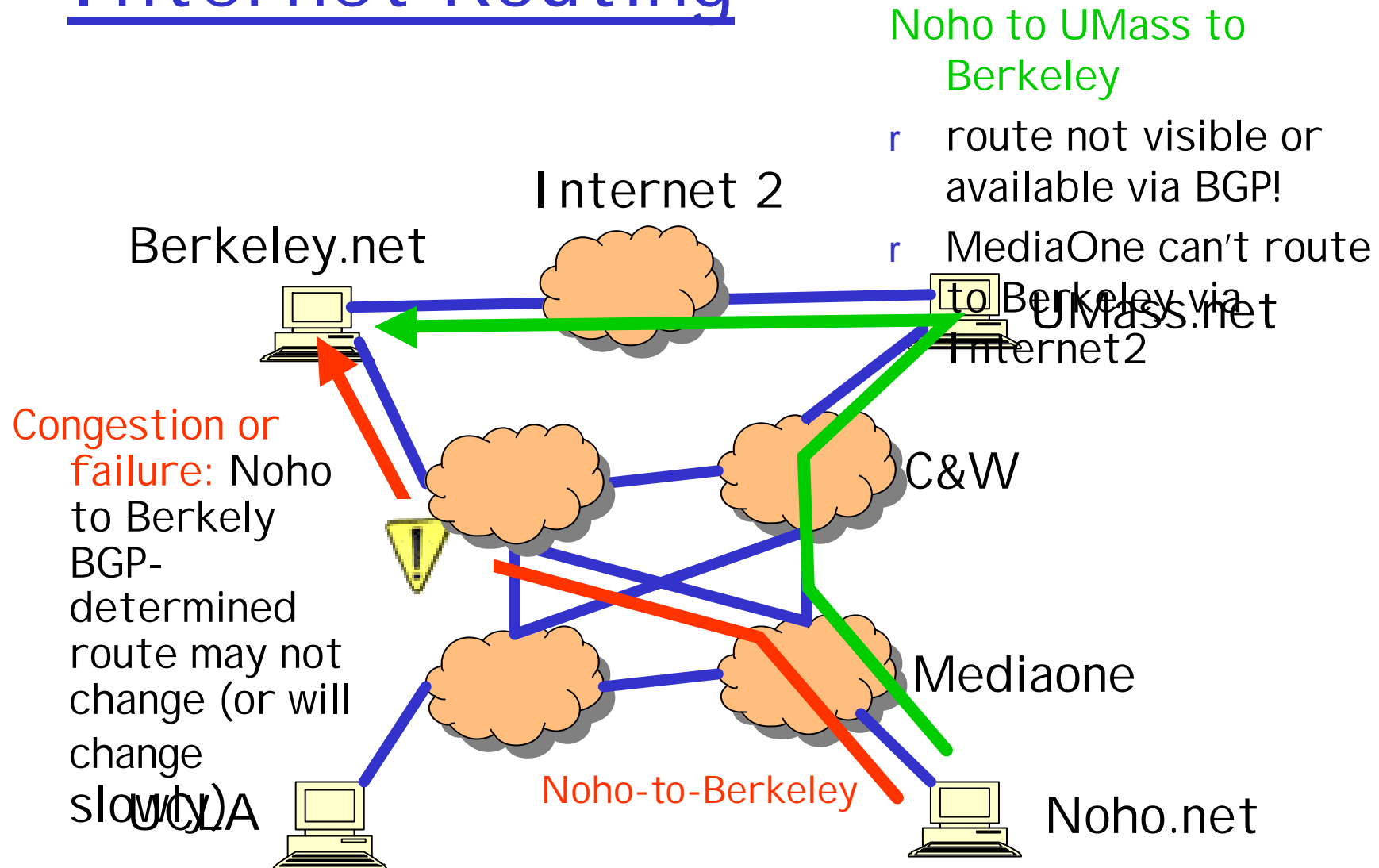
- r BGP defines routes between stub networks



Internet Routing

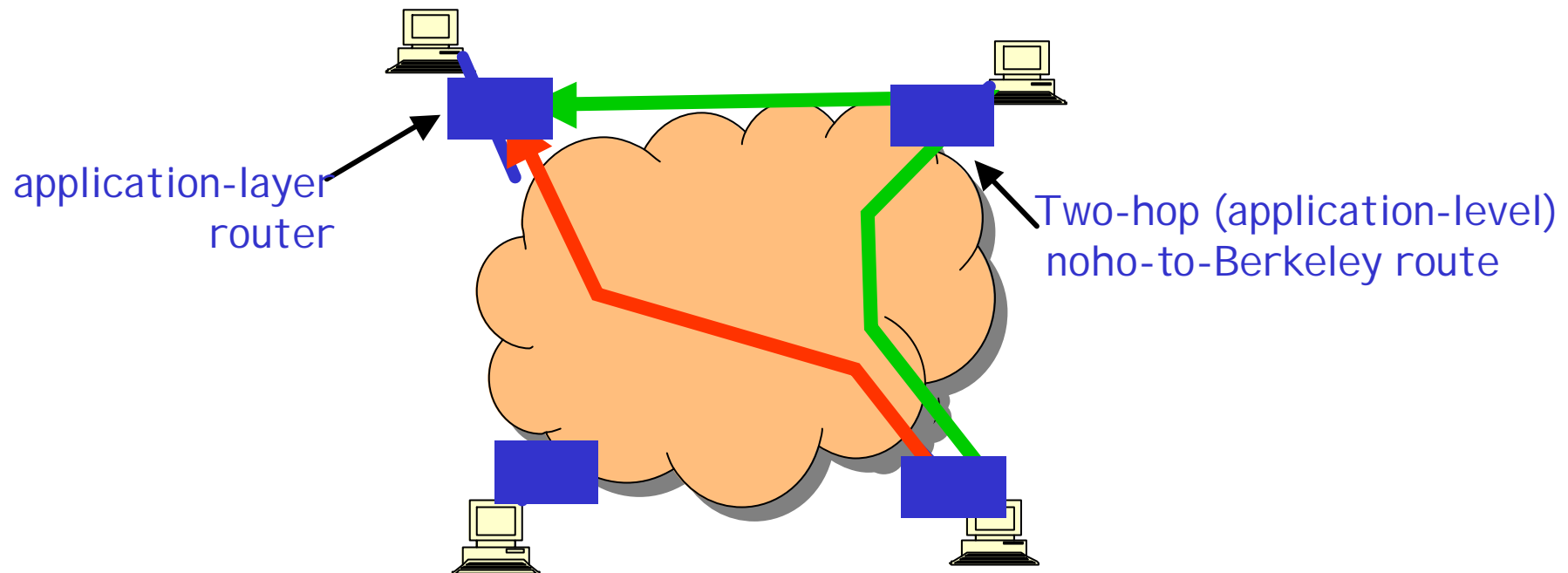


Internet Routing



RON: Resilient Overlay Networks

Premise: by building application overlay network, can increase performance, reliability of routing



RON Experiments

- r Measure loss, latency, and throughput with and without RON
- r 13 hosts in the US and Europe
- r 3 days of measurements from data collected in March 2001
- r 30-minute average loss rates
 - m A 30 minute outage is very serious!
- r Note: Experiments done with “No-Internet2-for-commercial-use” policy

An order-of-magnitude fewer failures

Loss Rate	<i>30-minute average loss rates</i>		
	RON Better	No Change	RON Worse
10%	479	57	47
20%	127	4	15
30%	32	0	0
50%	20	0	0
80%	14	0	0
100%	10	0	0

6,825 “path hours” represented here

12 “path hours” of essentially complete outage

76 “path hours” of TCP outage

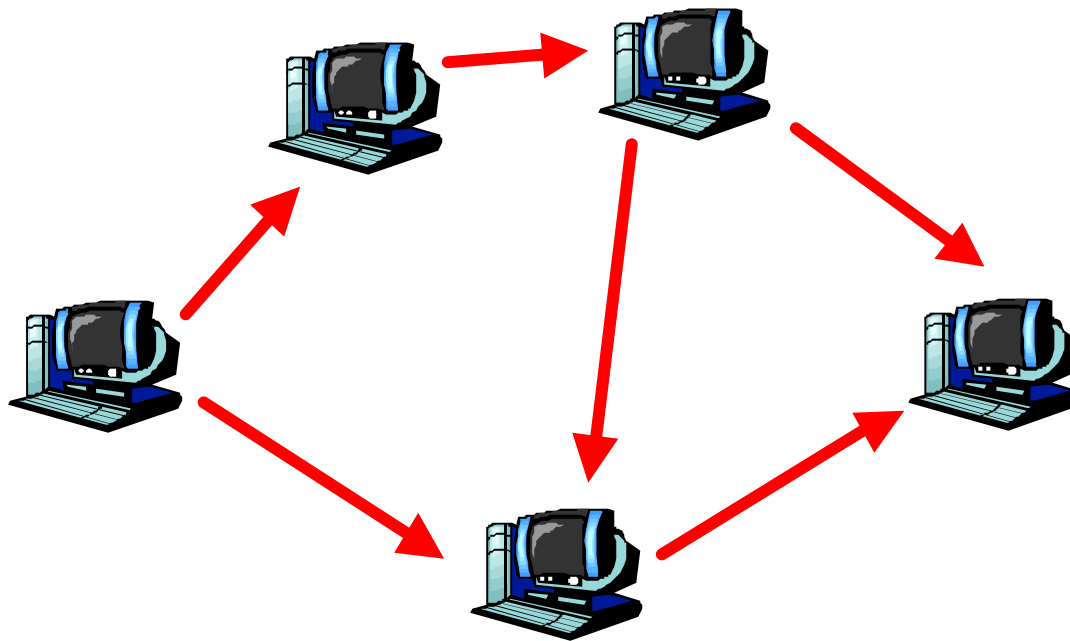
RON routed around all of these!

One indirection hop provides almost all the benefit!

RON Research Issues

- How to design overlay networks?
 - Measurement and self-configuration
 - Understanding performance of underlying net.
 - Fast fail-over.
 - Sophisticated metrics.
 - application-sensitive (e.g., delay versus throughput) path selection.
- Effect of RON on underlying network
 - If everyone does RON, are we better off?

4. Freenet



P2P Application

- r Centralized model

- m e.g. Napster

- m global index held by central authority

- m direct contact between requestors and providers

- r Decentralized model

- m e.g. Freenet, Gnutella, Chord

- m no global index – local knowledge only
(approximate answers)

- m contact mediated by chain of intermediaries

Freenet history

- r Final Year project [Ian Clarke](#) , [Edinburgh University](#), Scotland, June, 1999
- r Sourceforge Project, most active
- r V.0.1 (released March 2000)
- r Latest version(Sept, 2001): 0.4

What is Freenet and Why?

- r Distributed, Peer to Peer, file sharing system
- r Completely anonymous, for producers or consumers of information
- r Resistance to attempts by third parties to deny access to information

Freenet: How it works

- r Data structure
- r Key Management
- r Problems
 - m How can one node know about others
 - m How can it get data from remote nodes
 - m How to add new nodes to Freenet
 - m How does freenet manage its data
- r Protocol Details
 - m Header information

Data structure

r Routing Table

m Pair: node address: ip, tcp; corresponding key value

r Data Store

m Requirement:

- rapidly find the document given a certain key
- rapidly find the closest key to a given key
- keep track the popularity of documents and know which document to delete when under pressure

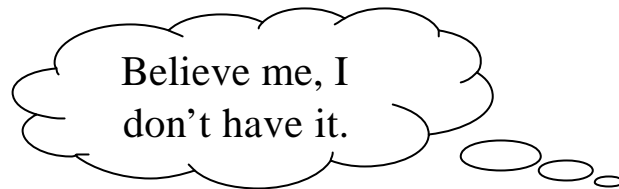
Key Management(1)

- r A way to locate a document anywhere
- r Keys are used to form a URI
- r Two similar keys don't mean the subjects of the file are similar!
- r Keyword-signed Key(KSK)
 - m *Based on a short descriptive string, usually a set of keywords that can describe the document*
 - m *Example: University/umass/cs/hzhang*
 - m *Uniquely identify a document*
 - m *Potential problem – global namespace*

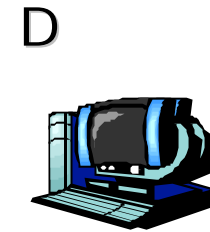
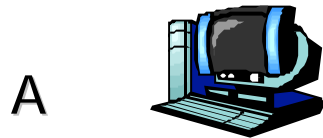
Key Management (2)

- r Signed-subspace Key(SSK)
 - m *Add sender information to avoid namespace conflict*
 - m *Private key to sign/ public key to varify*
- r Content-hash Key(CHK)
 - m *Message digest algorithm, Basically a hash of the document*

Freenet: Routing Algorithm: search or insert



Freenet: Routing Algorithm: search or insert

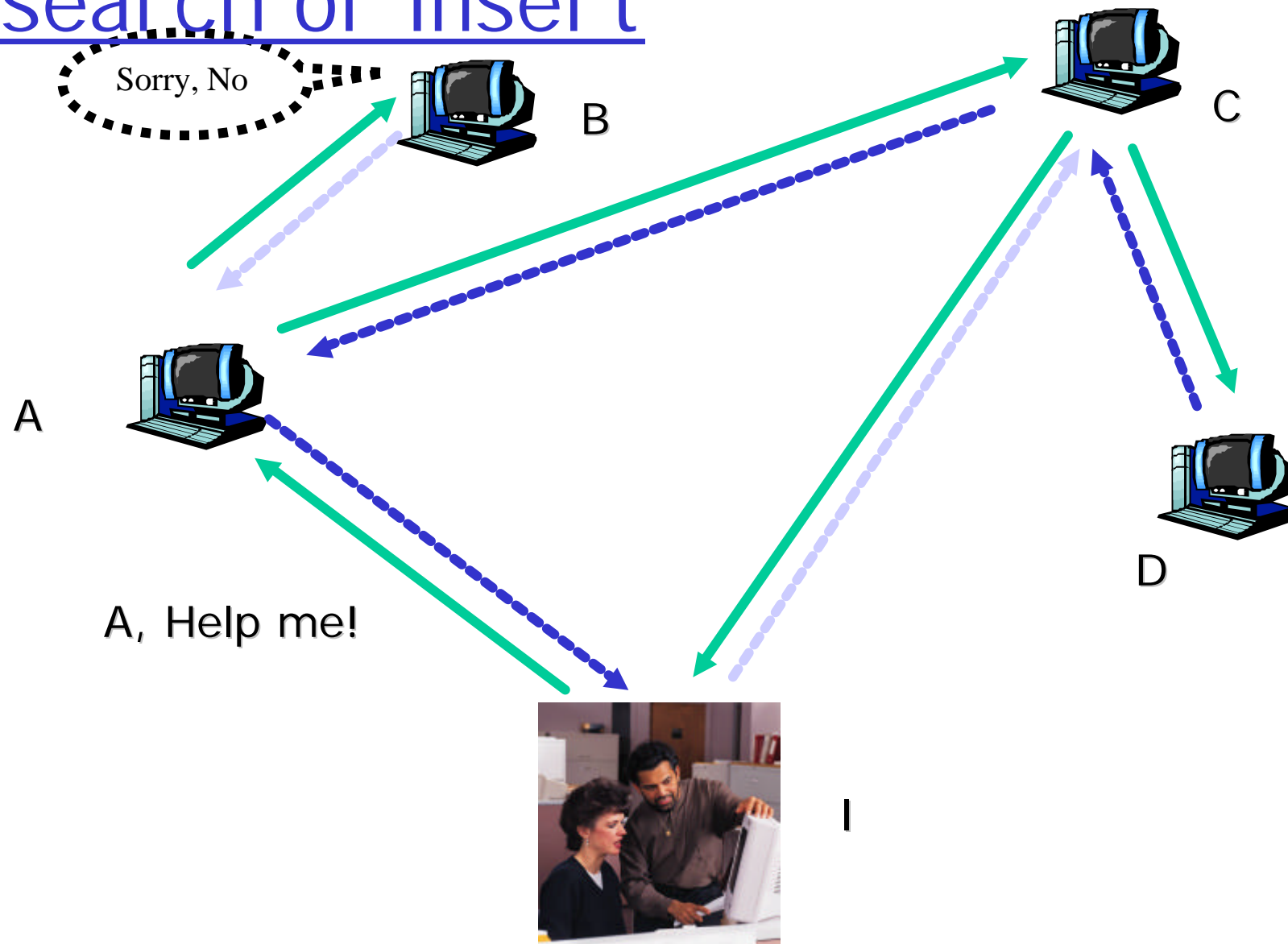


But I know Joe
may have it
since I
borrowed
similar stuff
him last time.



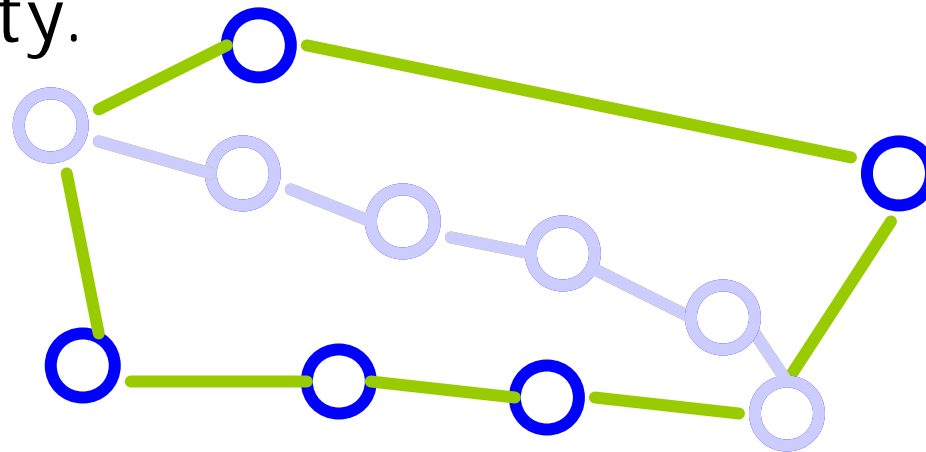
I

Freenet: Routing Algorithm: search or insert



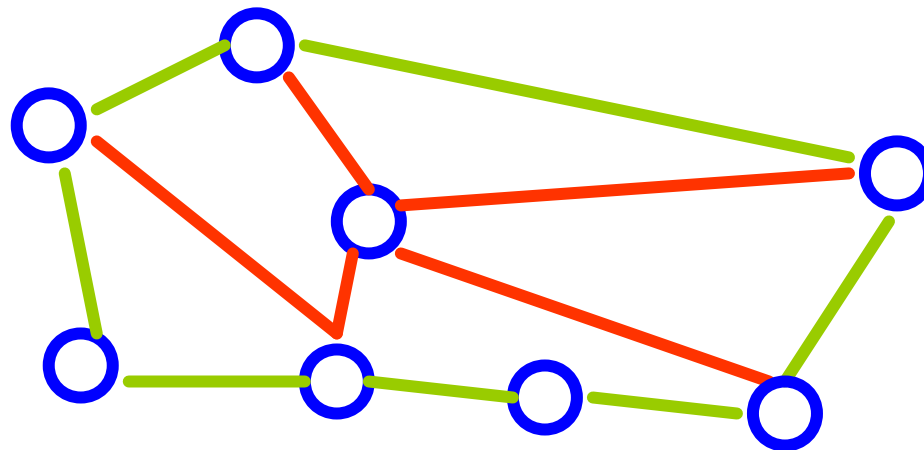
Strength of routing algorithm(1)

- r Replication of Data Clustering (1)
(Note: Not subject-clustering but key-clustering!)
- r Reasonable Redundancy: improve data availability.



Strength of routing algorithm(2)

- r New Entry in the Routing Table: the graph will be more and more connected. --- Node discovery



Protocol Details

r Header information

m DataReply

UniqueID=C24300FB7BEA06E3

Depth=a

* HopsToLive=2c

Source=tcp/127.0.0.1:2386

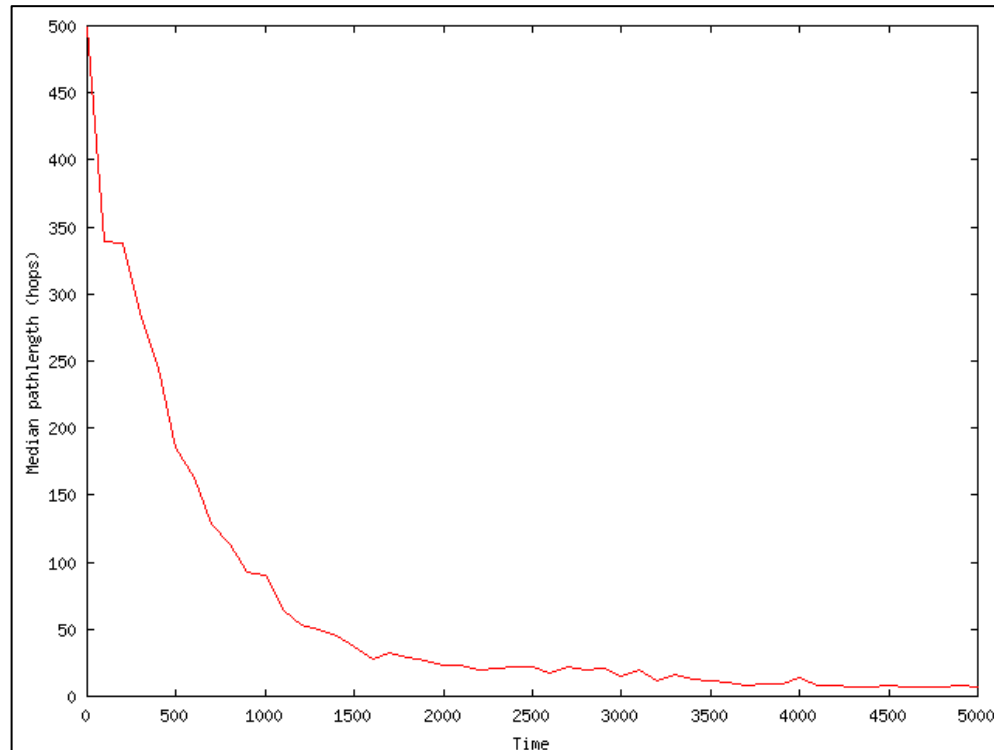
DataLength=131

Data 'Twas brillig, and the slithy toves Did
gyre and gimble in the wabe: All
mimsy were the borogoves And the
mome raths outgrabe

Some security and authentication issues

- r How to ensure anonymity:
 - m Nodes can lie randomly about the requests and claim to be the origin or the destination of a request
 - m Hop-To-Live values are fuzzy
 - m Then it's impossible to trace back a document to its original node
 - m Similarly, it's impossible to discover which node inserted a given document.

Network convergence



X-axis: time

Y-axis: # of pathlength

1000 Nodes, 50 items

datastore, 250 entries

routing table

the routing tables were

initialized to ring-lattice

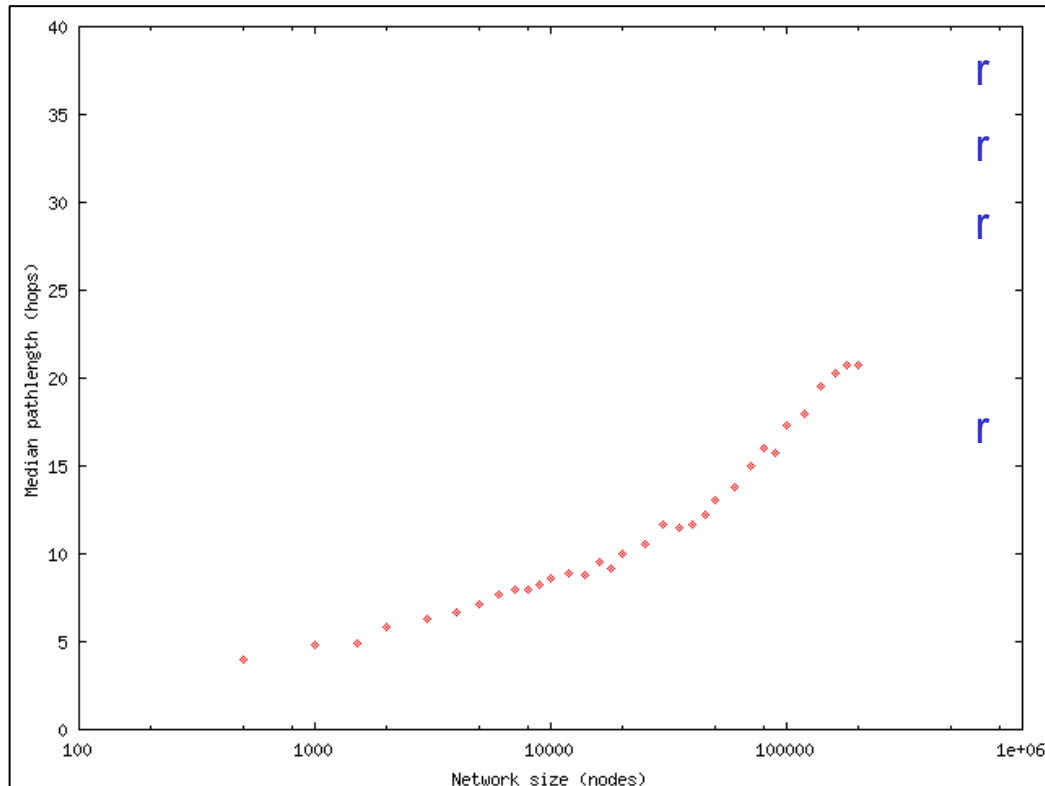
topology

Pathlength: the number

of hops actually taken

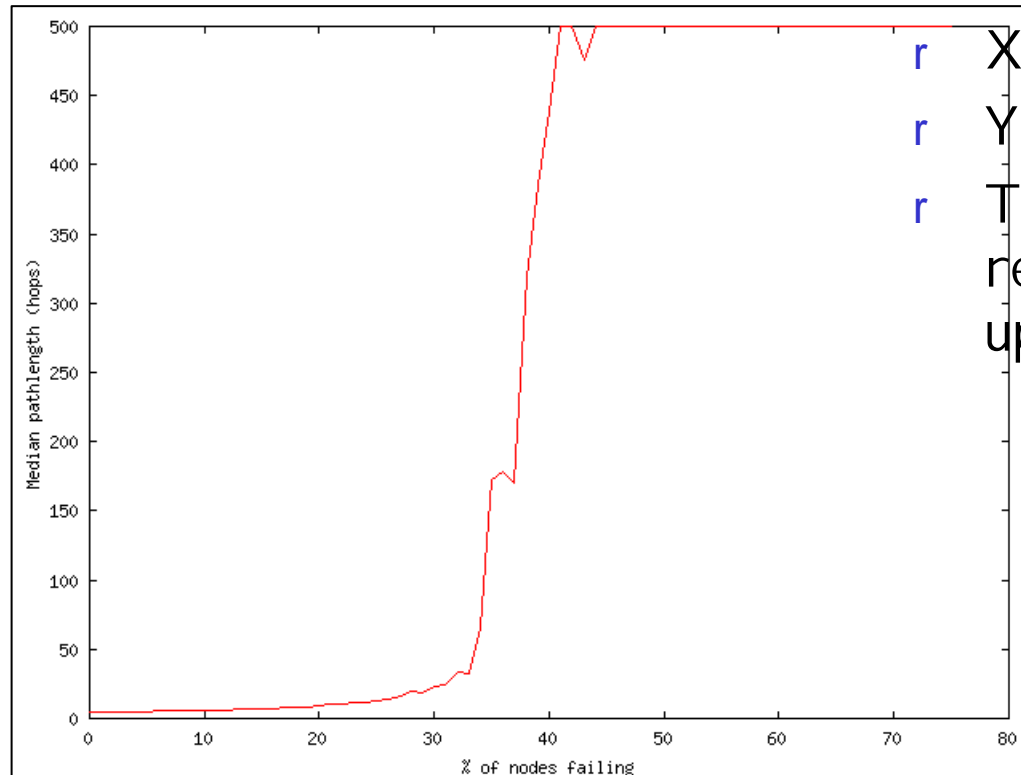
before finding the data.

Scalability



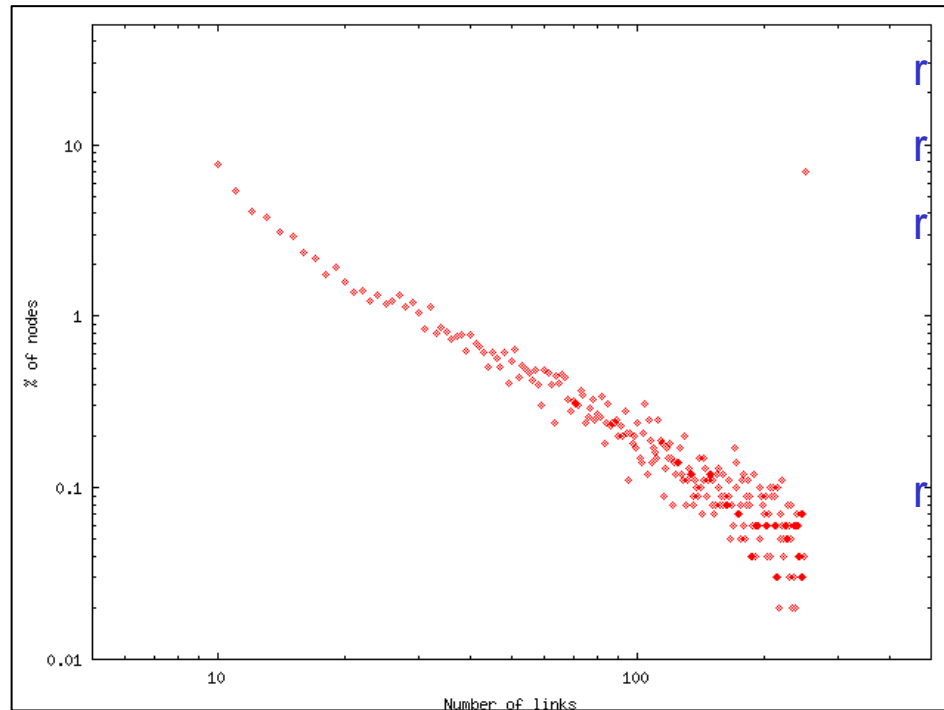
- r X-axis: # of nodes
- r Y-axis: # of pathlength
- r The relation between network size and average pathlength.
- r Initially, 20 nodes. Add nodes regularly.

Fault Tolerance



- r X-axis: # of nodes failing
- r Y-axis: # of pathlength
- r The median pathlength remains below 20 even when up to 30% nodes fails.

Small world Model



X-axis: # of nodes failing

Y-axis: # of pathlength

Most of nodes have only few connections while a small number of news have large set of connections.

The authors claim it follows power law.

So far, it's really a good model

- r Keep anonymity
- r Distributed model; data available
- r Converge fast
- r Adaptive

Is it Perfect?

- r How long will it take to search or insert?
 - m Trade off between anonymity and searching efforts:
Chord vs Freenet
 - m Can we come up a better algorithm? A good try: "Search in Power-Law Networks"
- r Have no idea about if search fails due to no such document or just didn't find it.
- r File lifetime. Freenet doesn't guarantee a document you submit today will exist tomorrow!!

Question??

- r Anonymity? Security?
- r Better search algorithm? Power law?
- r ...

5. Publius: A robust, tamper-evident, censorship-resistant web publishing system

Marc Waldman

Aviel Rubin

Lorrie Faith Cranor

Outline

- r Design Goals
- r Kinds of Anonymity
- r Publius Features
- r Publius Limitations and Threats
- r Questions

Design Goals

- r Censorship resistant
 - m Difficult for a third party to modify or delete content
- r Tamper evident
 - m Unauthorized changes should be detectable
- r Source anonymous
 - m No way to tell who published the content
- r Updateable
 - m Changes to or deletion of content should be possible for publishers

Design Goals

- r Deniable

- m Involved third parties should be able to deny knowledge of what is published

- r Fault Tolerant

- m System remains functional, even if some third parties are faulty or malicious

- r Persistent

- m No expiration date on published materials

Web Anonymity

r Connection Based

m Hides the identity of the individual requesting a page Examples:

- Anonymizing proxies, such as The Anonymizer or Proxymate
- Proxies utilizing Onion Routing, such as Freedom
- Crowds, where users in the Crowd probabilistically route or retrieve for other users in the Crowd

Web Anonymity

r Author Based

m Hides the location or author of a particular document Examples:

- Rewebber, which proxies requests for encrypted URLs
- The Eternity Service, which for a fee inserts a document into a random subset of servers, and guarantees its future existence
- Freenet

m Publius provides this sort of anonymity

Publius System Overview

- r Publishers

- m Post Publius content to the web

- r Servers

- m A static set which host random-looking content

- r Retrievers

- m Browse Publius content on web

Publius System Overview

r Publish

- m A publisher posts content across multiple servers in a source anonymous fashion

r Retrieve

- m A retriever gets content from multiple servers

r Delete

- m The original publisher of a document removes it from the Publius servers

r Update

- m The original publisher modifies a document

Publius Publishing

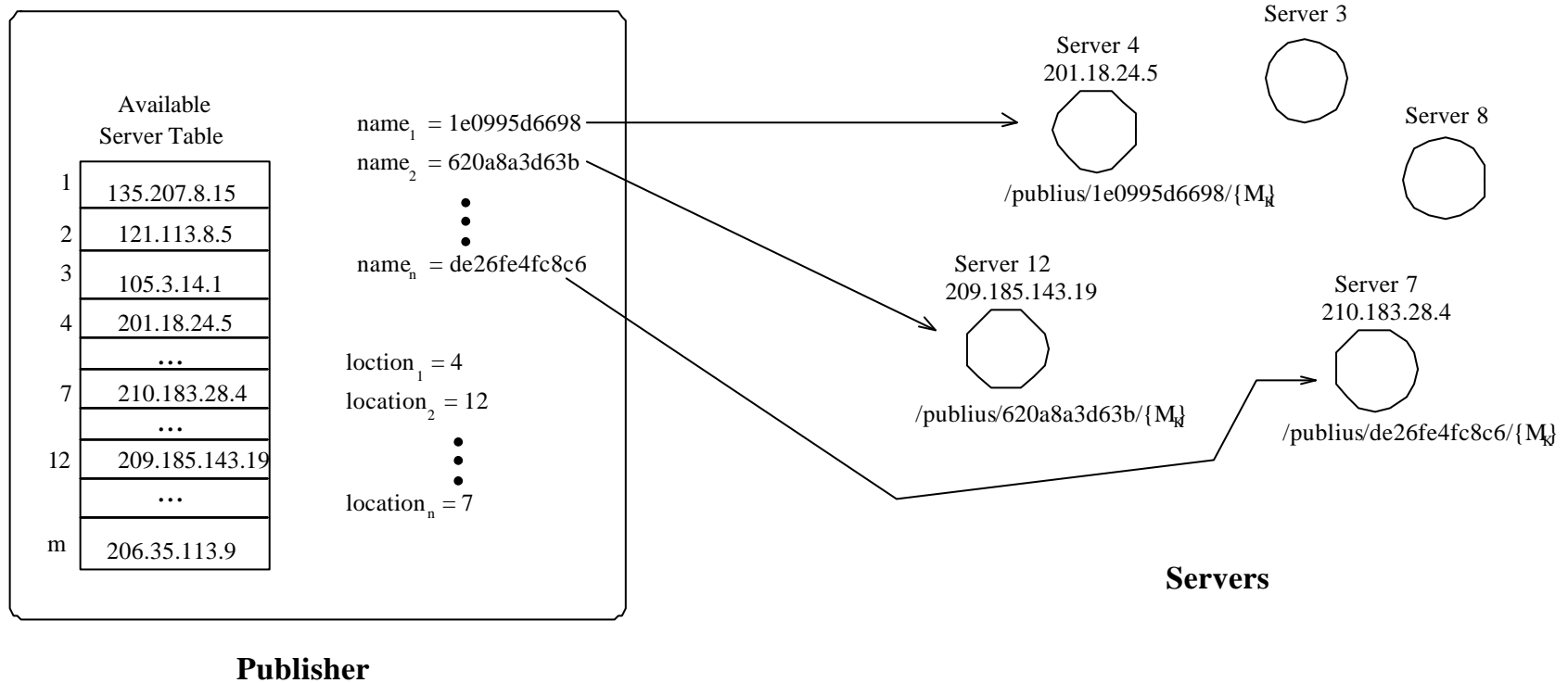
- r Alice generates a random symmetric key K
- r She encrypts message M with key K , producing $\{M\}_K$
- r She splits K into n shares, using Shamir secret sharing, such that any k can reproduce K
- r Each share is uniquely named:

$$name_i = wrap(H(M \cdot share_i))$$

Publius Publishing

- r A set of locations is chosen:
$$location_i = (name_i \text{ MOD } m) + 1$$
- r Each $location_i$ indexes into the list of m servers
- r If $d \geq k$ unique values are not obtained, start over
- r Alice publishes $\{M\}_K$ and $share_i$ into a directory $name_i$ on the server at $location_i$
- r A URL containing at least the d $name_i$ values is produced

Publius Publishing



Publius Retrieval

- r Bob parses out each $name_i$ from URL, and for each, computes:

$$location_i = (name_i \text{ MOD } m) + 1$$

- r Bob chooses k of these, and retrieves the encrypted file $\{M\}_K$ and $share_i$ at each server
- r Bob combines the shares to get K , and decrypts the file
- r Bob verifies that each name value is correct:

$$name_i = \text{wrap}(H(M \cdot share_i))$$

Publius Delete

- r Alice generates a password PW when publishing a file
- r Alice includes $H(server_domain_name \cdot PW)$ in server directory when publishing
 - m Note that each server has its own hash, to prevent a malicious server operator from deleting content on all servers
- r Alice deletes by sending $H(server_domain_name \cdot PW)$ and $name_i$ to each of the n servers hosting content

Publius Update

- r Idea: change content without changing original URL, as links to that URL may exist
- r In addition to the file, the share, and the password, there may be an update file in the *name_i* directory
- r This update file will not exist if Alice has not updated the content

Publius Update

- r To update, Alice specifies a new file, the original URL, the original password PW , and a new password
- r First, the new content is published, and a new URL is generated
- r Then, each of the n old files is deleted, and an update file, containing the new URL, is placed in each $name_i$ directory

Publius Update

- r When Bob retrieves updated content, the server returns the update file instead
- r Bob checks that all of the URLs are identical, then retrieves the content at the new URL

Linking Documents

- r Simple case: file A links to file B
 - m Solution: Publish B first, then rewrite URLs in A

- r Harder: files C and D link to each other
 - m Cannot use simple solution above
 - m Alice publishes C and D in any order
 - m She then rewrites the URLs in each file, and uses the Publius Update procedure on the new files

Other Features

- r Entire directories can be published by exploiting the updateability of Publius
- r Mechanism exists to encode MIME type into Publius content
- r Publius URLs include option fields and other flags, the value of k , and other relevant values
 - m Older browsers preclude URLs of length >255 characters
 - m Once this limitation is removed, URLs can include server list, making this list non-static

Limitations and Threats

- r Share deletion or corruption
 - m If all n copies of a file, or $n-k+1$ copies of the shares, are deleted, then the file is unreadable
 - m Increasing n , or decreasing k , makes this attack harder

Limitations and Threats

- r Update file deletion or corruption 1
 - m If there is no update file, malicious server operator Mallory could create one, pointing to bad content
 - m This requires the assistance of at least k other server operator, and motivates a higher value of k
 - m The Publius URL has several fields, among them a *no_update* flag, which will prevent this sort of attack

Limitations and Threats

- r Update file deletion or corruption 2
 - m If Publius content has already been updated, Mallory must corrupt update files on $n-k+1$ servers
 - m Of course, if Mallory can do this, she can censor any document
 - m Larger n and smaller k make this more difficult
- r Deciding upon good values for n and k is difficult
 - m No suggestions from Waldman et al.

Limitations and Threats

- r Publius, like all internet services, is subject to DoS attacks
 - m Flooding is less effective, as $n-k+1$ servers must be attacked
 - m A malicious user could attempt to fill disk space on servers
 - Some mechanisms in place to prevent this

Limitations and Threats

- r If the Publius content contains any identifying information, anonymity will be lost
- r Publius does not provide any connection based anonymity
 - m If you act as a publisher, you must anonymize your connections with the Publius servers

Questions

- r How do you publish Publius URLs anonymously?
 - m Freenet keys can be guessed at, but Publius URLs are entirely machine generated
 - m The first person to publish a Publius URL must have some connection with the publisher of the content
 - m If you have somewhere secure and anonymous to publish the Publius URLs, why do you need Publius?
 - One possible answer: censorship resistance
 - But server operators are then potentially liable

Questions

- r How deniable is Publius?
 - m Publius URLs are public
 - m With minimal effort, a Publius server operator could determine the content being served

Questions

- r How does Publius compare to Freenet?
 - m Both provide publisher anonymity, deniability, and censorship resistance
 - m Freenet provides anonymity for retrievers and servers, as well
 - Cost is high: data must be cached at many nodes
 - m Publius provides persistence of data
 - Freenet does not
 - Can any p2p system provide persistence?

Questions

- r Could Publius be made into a p2p service?
- r Would it be appropriate to do so?

6. Chord: A Scalable Peer-to-peer Lookup Service for Internet Applications

Ion Stoica, Robert Morris, David Karger,
M. Frans Kaashoek, Hari Balakrishnan

MIT and Berkeley

Now we see some CS in strength – Hash and Content based....for more
scaleble (distributed) directory lookup

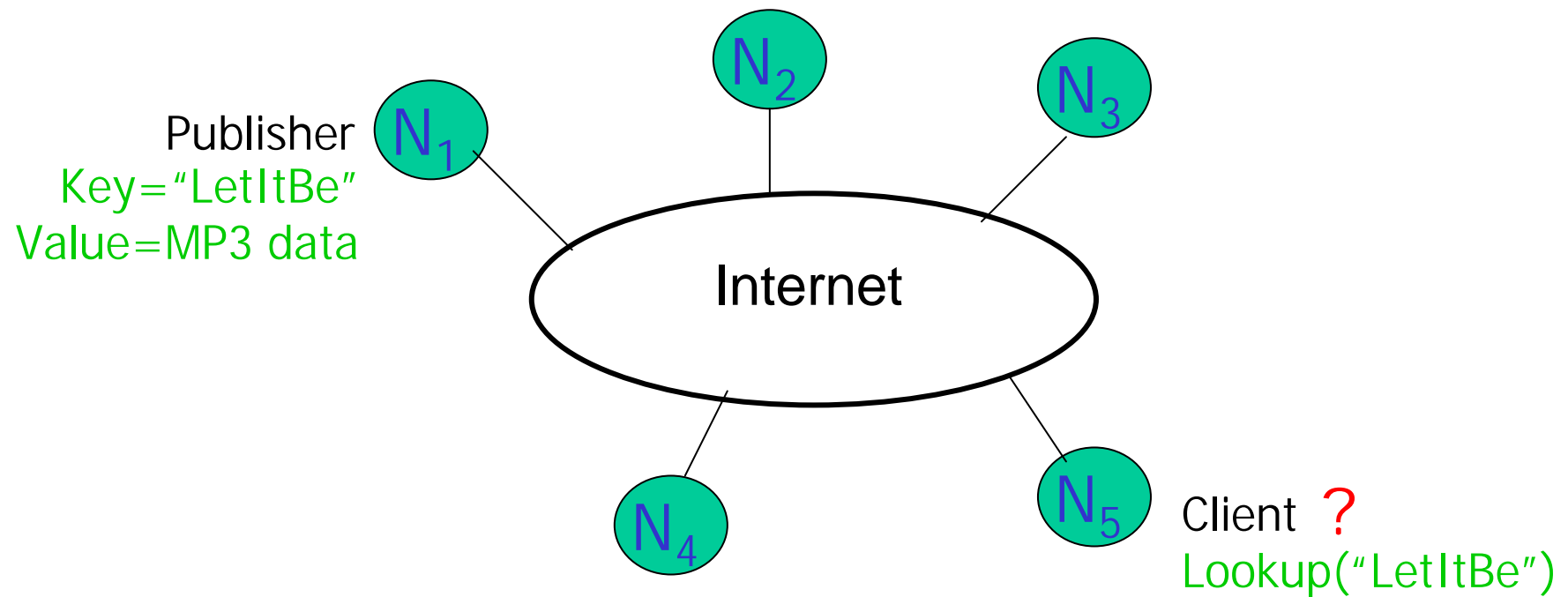
 presentation based on slides by Robert Morris (SI GCOMM'01)

Outline

- ✍ Motivation and background
- ✍ Consistency caching
- ✍ Chord
- ✍ Performance evaluation
- ✍ Conclusion and discussion

Motivation

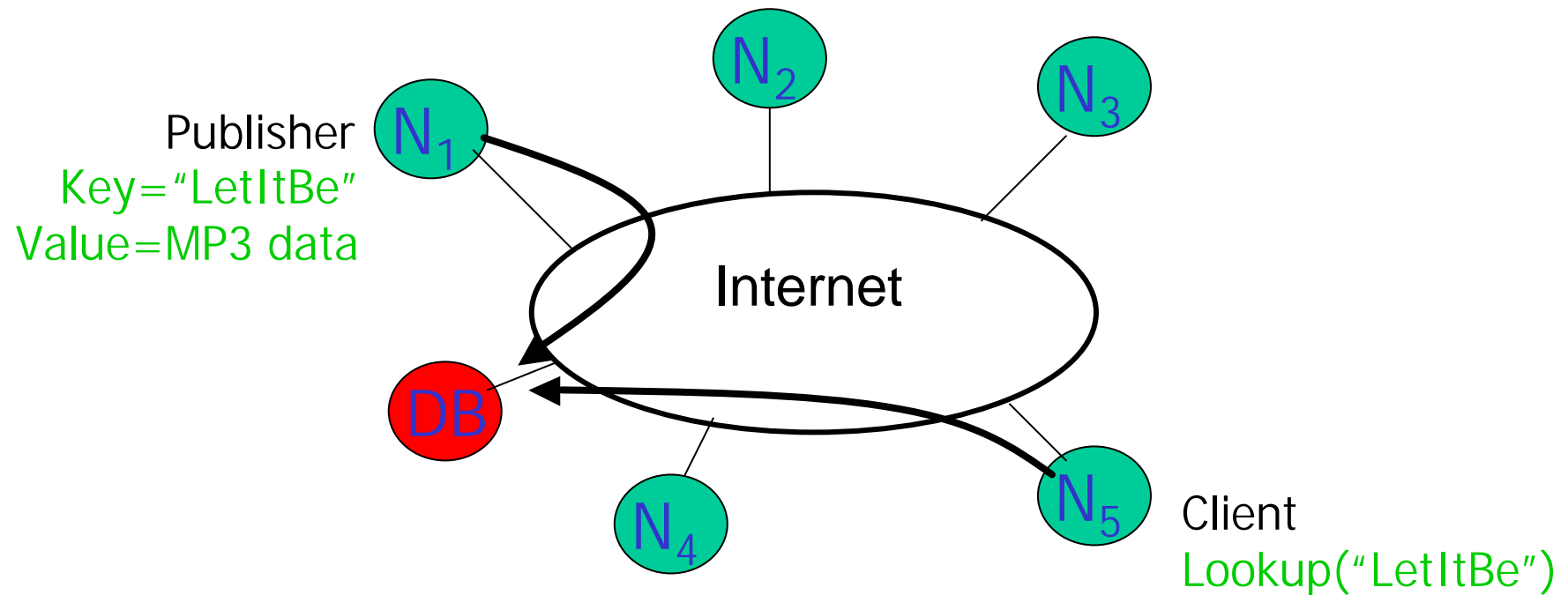
How to find data in a distributed file sharing system?



 Lookup is the key problem

Centralized Solution

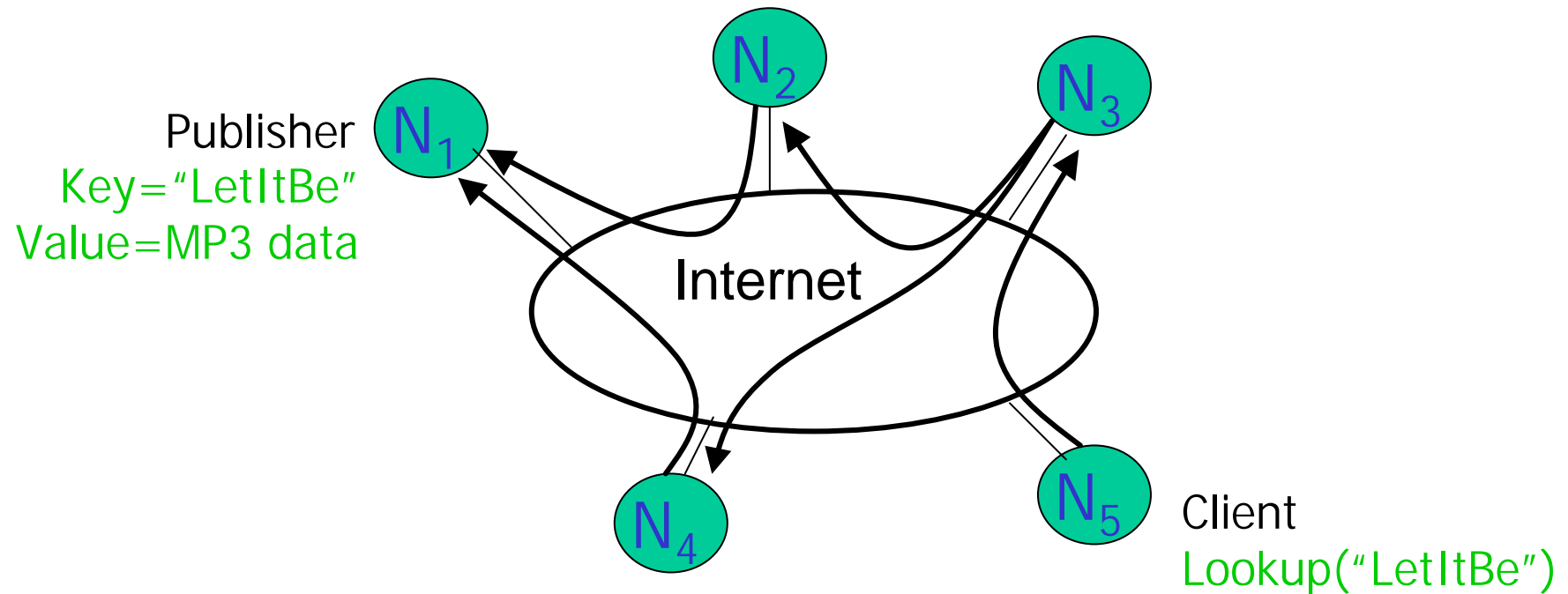
- ✗ Central server (Napster)



- ✗ Requires $O(M)$ state
- ✗ Single point of failure

Distributed Solution (1)

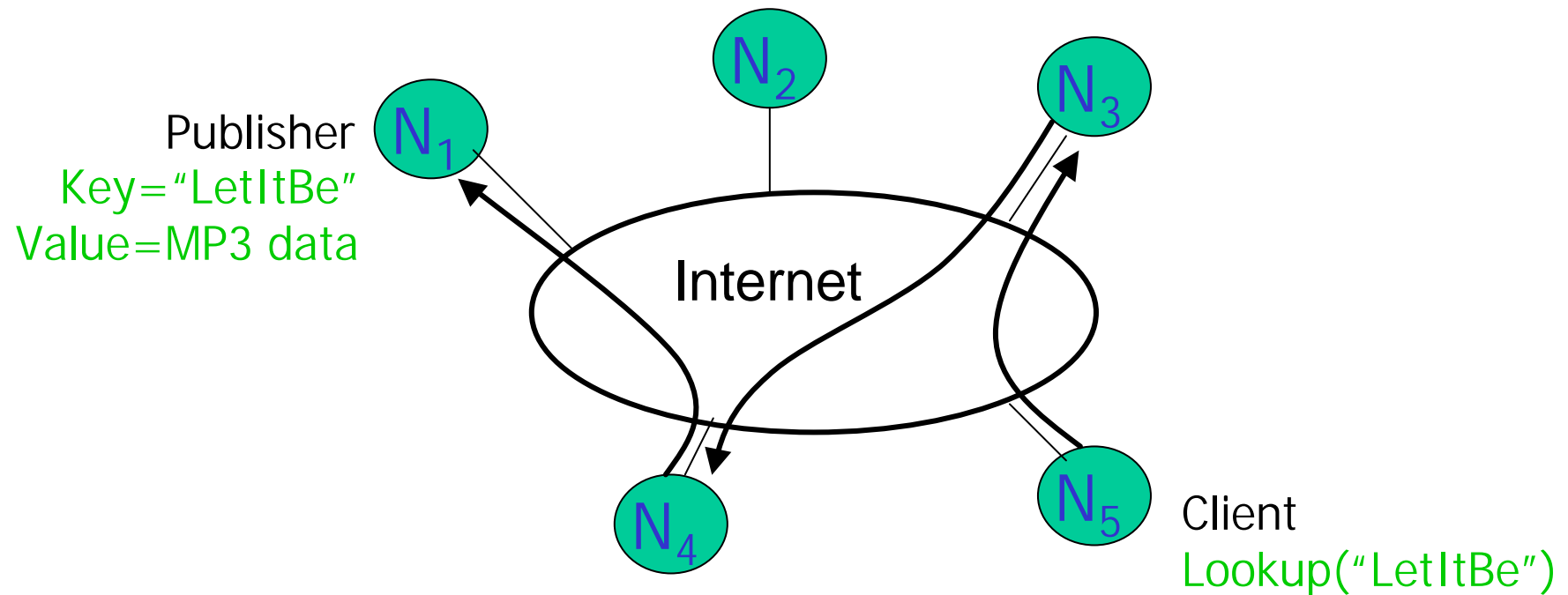
✍ Flooding (Gnutella, Morpheus, etc.)



✍ Worst case $O(N)$ messages per lookup

Distributed Solution (2)

✍ Routed messages (Freenet, Tapestry, Chord, CAN, etc.)



✍ Only exact matches

Routing Challenges

- ✍ Define a useful key nearness metric
- ✍ Keep the hop count small
- ✍ Keep the routing tables “right size”
- ✍ Stay robust despite rapid changes in membership

Authors claim:

- ✍ Chord: emphasizes efficiency and simplicity

Chord Overview

- ✍ Provides peer-to-peer hash lookup service:
 - ✍ Lookup(key) ? IP address
 - ✍ Chord does not store the data
- ✍ How does Chord locate a node?
- ✍ How does Chord maintain routing tables?
- ✍ How does Chord cope with changes in membership?

Chord properties

- ✍ Efficient: $O(\log N)$ messages per lookup
 - ✍ N is the total number of servers
- ✍ Scalable: $O(\log N)$ state per node
- ✍ Robust: survives massive changes in membership
- ✍ Proofs are in paper / tech report
 - ✍ Assuming no malicious participants

Chord I Ds

✍ m bit identifier space for both keys and nodes

✍ Key identifier = SHA-1(key)

Key="LetItBe" $\xrightarrow{\text{SHA-1}}$ ID=60

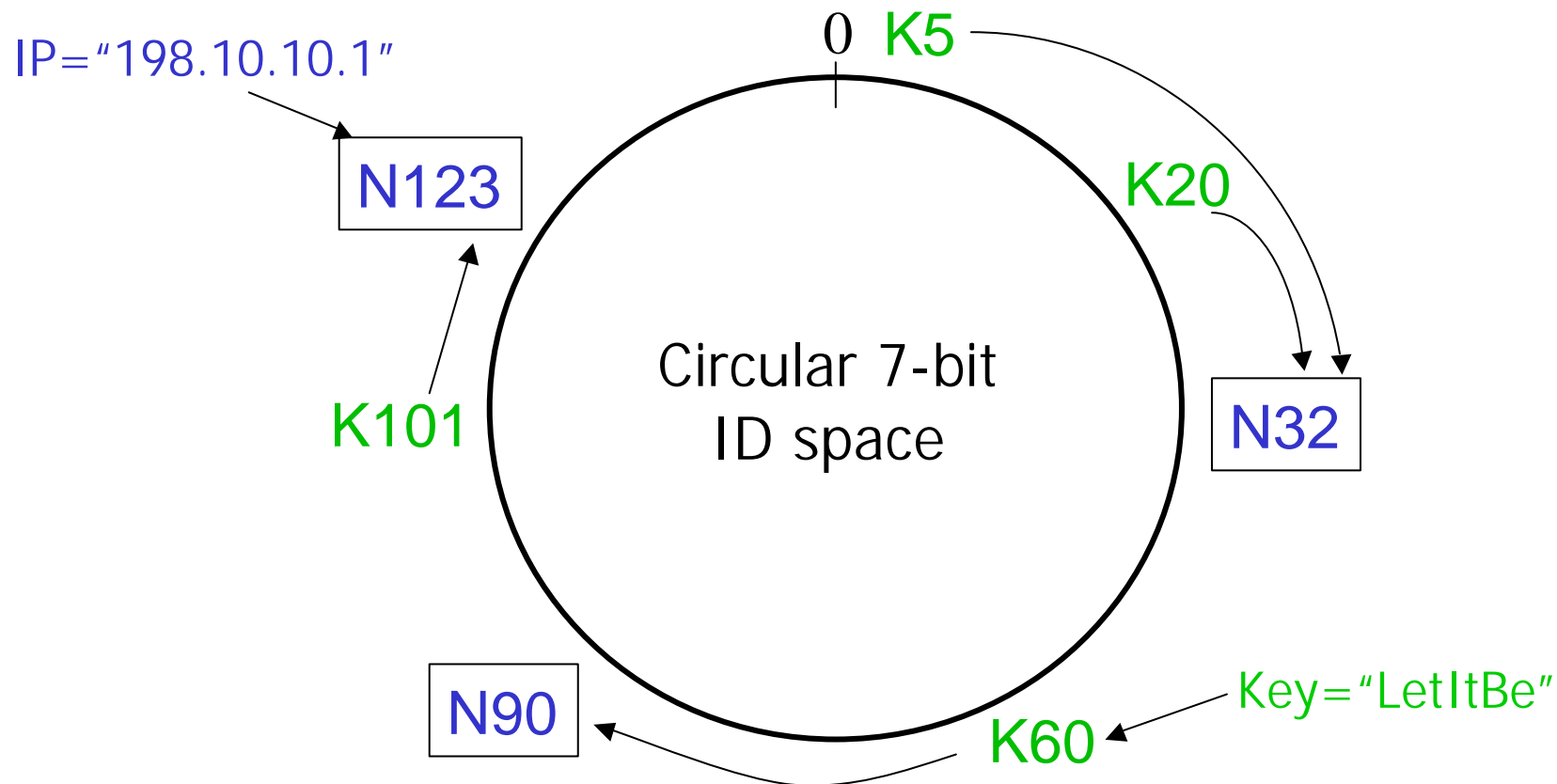
✍ Node identifier = SHA-1(IP address)

IP="198.10.10.1" $\xrightarrow{\text{SHA-1}}$ ID=123

✍ Both are uniformly distributed

✍ How to map key I Ds to node I Ds?

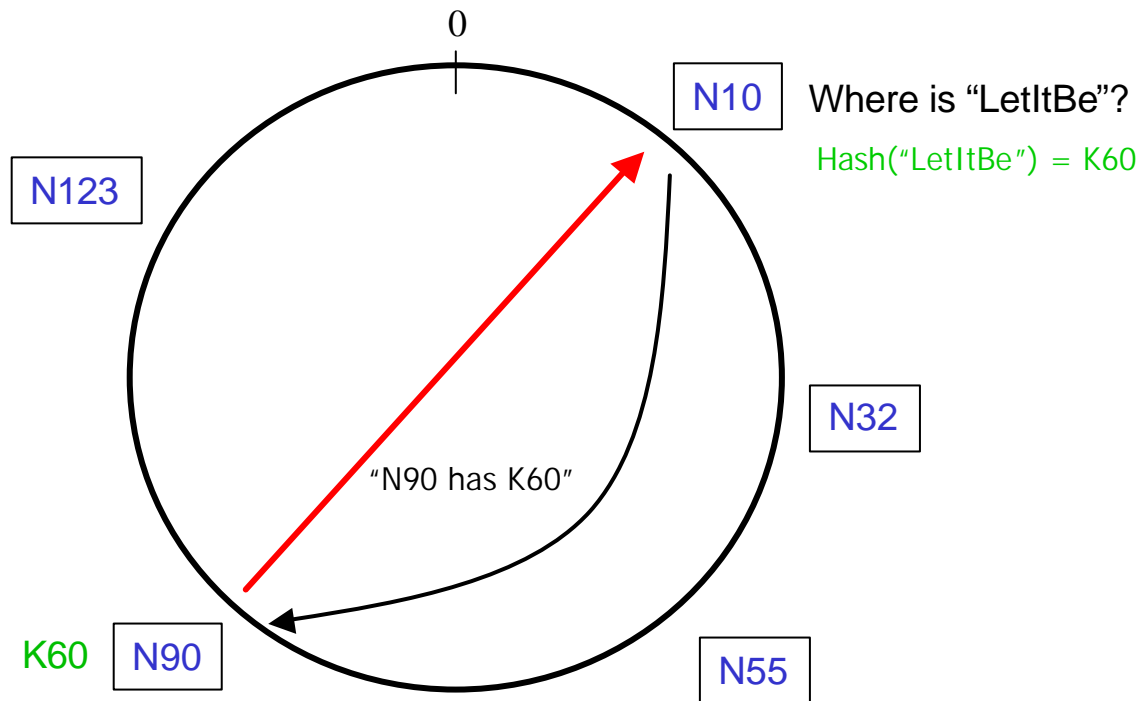
Consistent Hashing [Karger 97]



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

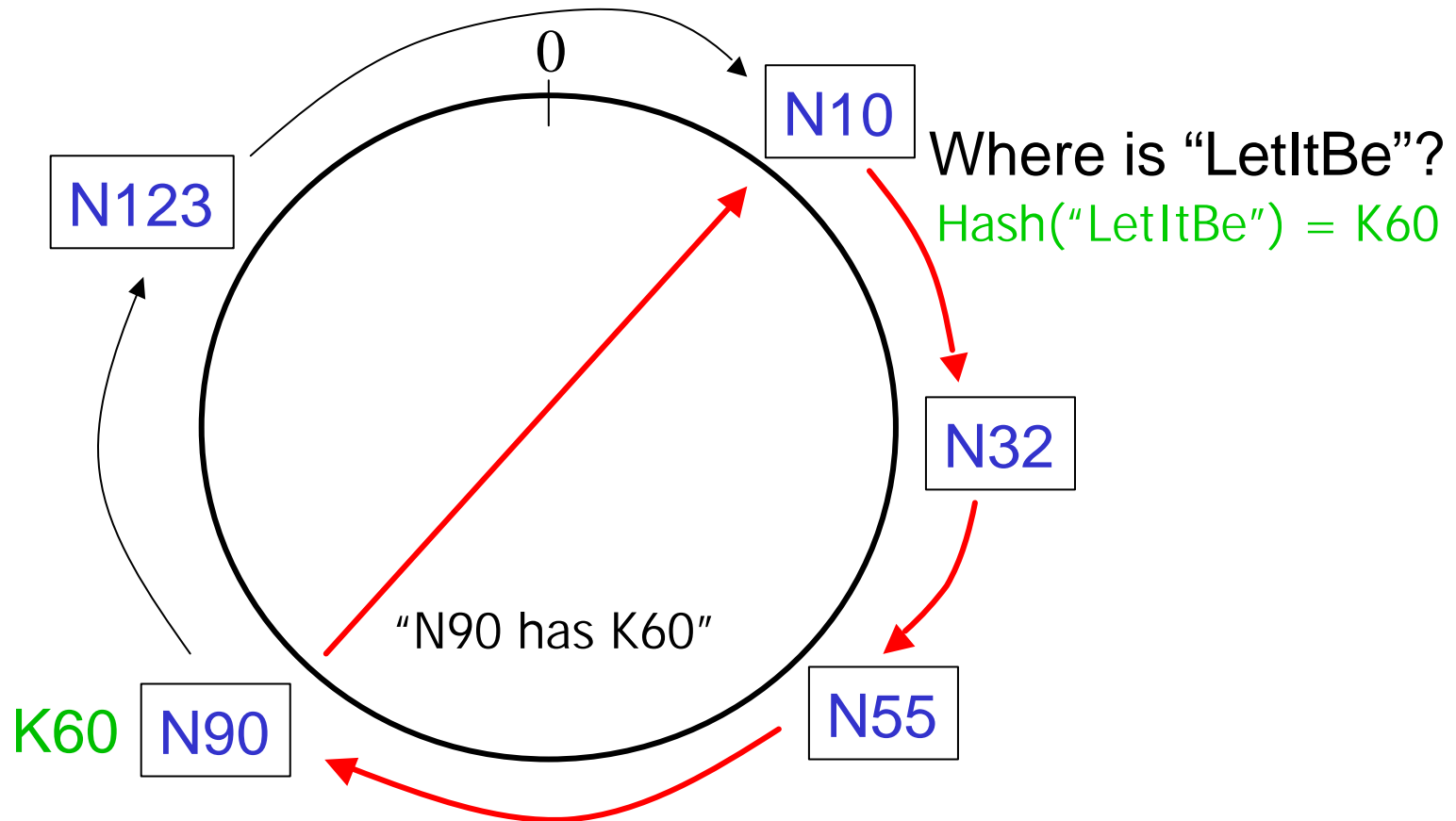
Consistent Hashing

- ✂ Every node knows of every other node
 - ✂ requires global information
- ✂ Routing tables are large $O(N)$
- ✂ Lookups are fast $O(1)$



Chord: Basic Lookup

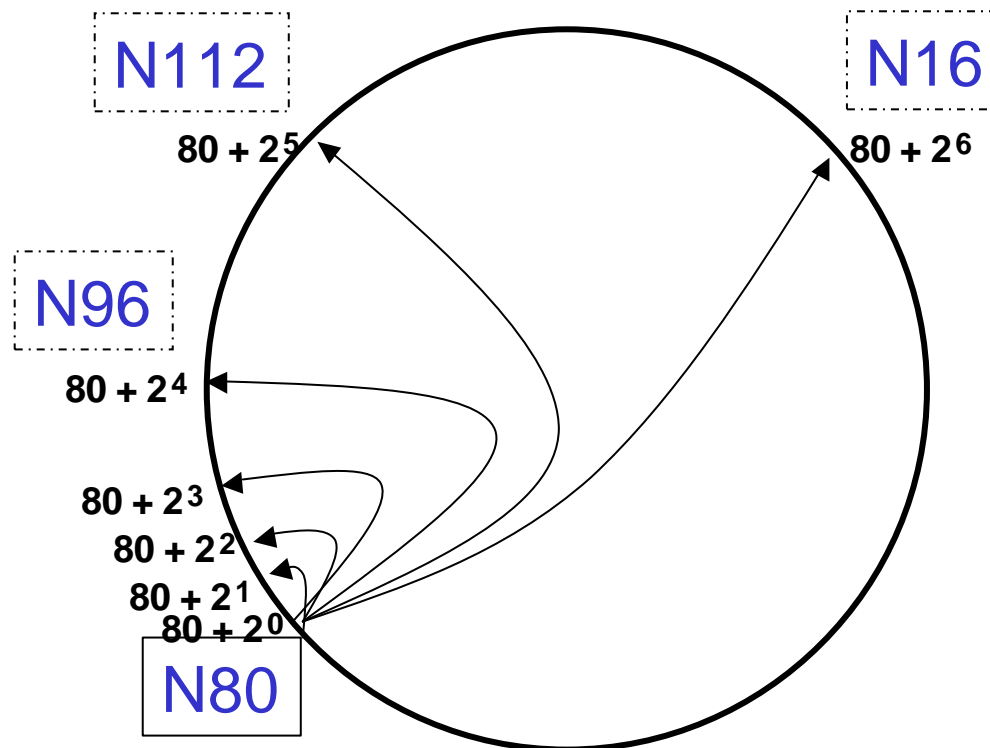
- ✗ Every node knows its successor in the ring



- ✗ requires $O(N)$ time

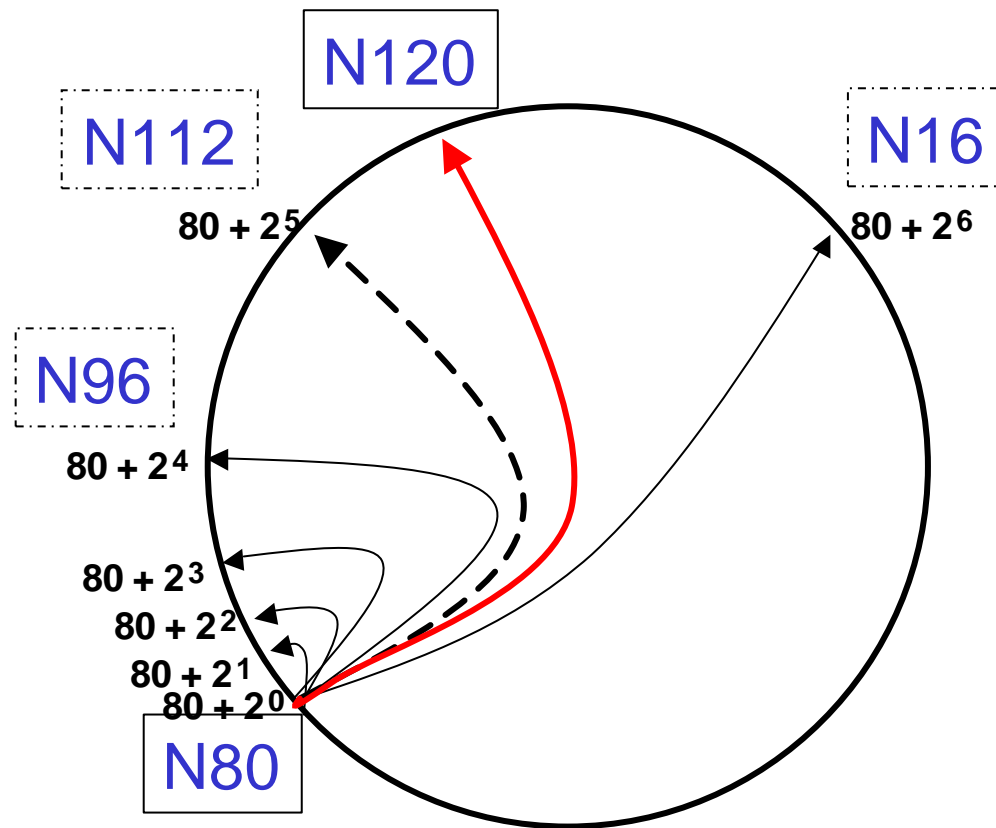
"Finger Tables"

- ✍ Every node knows m other nodes in the ring
- ✍ Increase distance exponentially



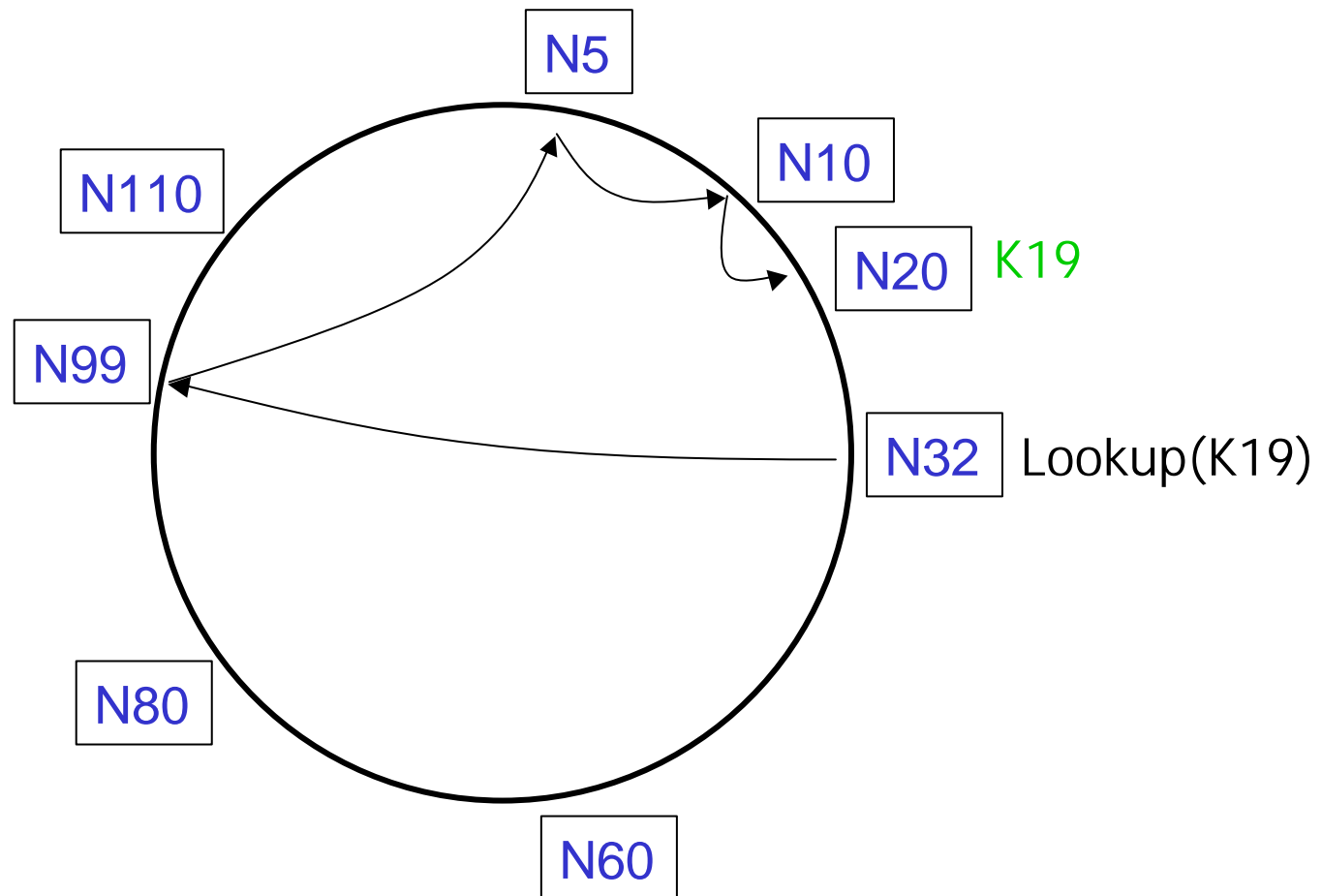
"Finger Tables"

✍ Finger i points to **successor** of $n+2^i$






Lookups are Faster

✂ Lookups take $O(\log N)$ hops






Joining the Ring

Three step process:

-  Initialize all fingers of new node
-  Update fingers of existing nodes
-  Transfer keys from successor to new node

Less aggressive mechanism (lazy finger update):

-  Initialize only the finger to successor node
-  Periodically verify immediate successor, predecessor
-  Periodically refresh finger table entries

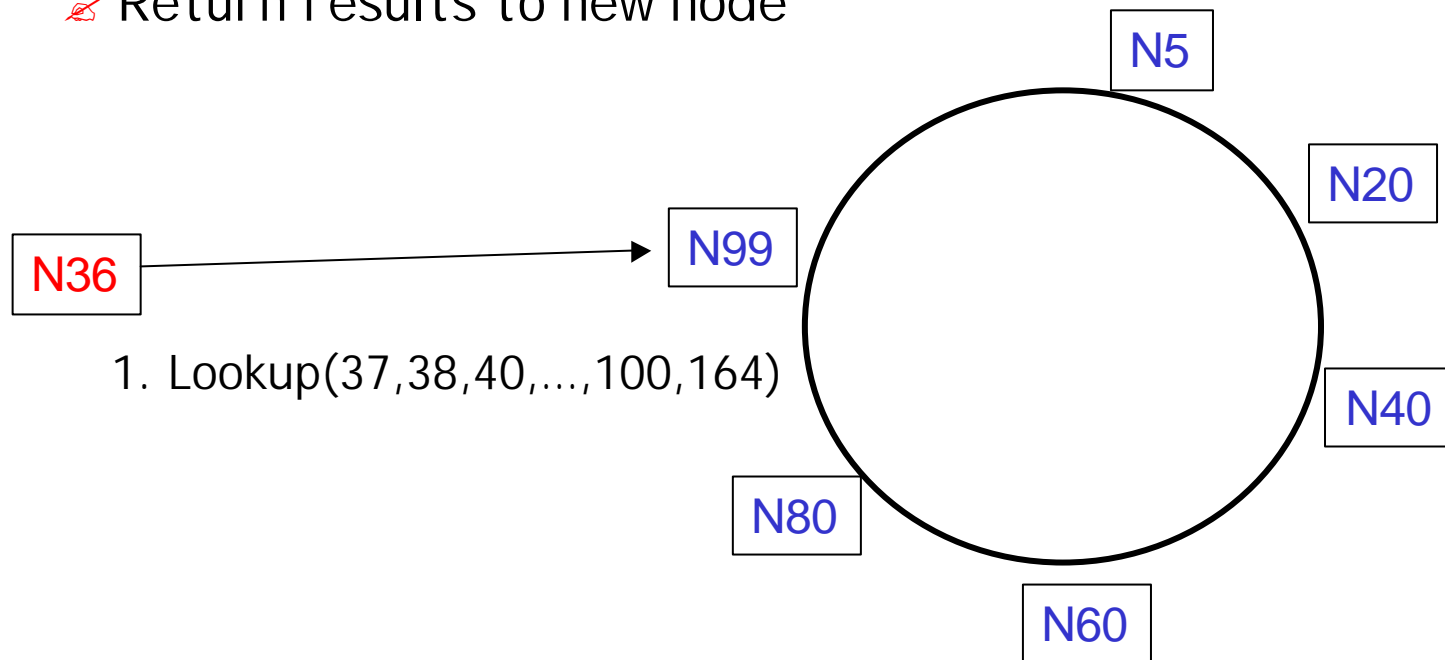
Joining the Ring - Step 1

✍ Initialize the new node finger table

✍ Locate any node p in the ring

✍ Ask node p to lookup fingers of new node N36

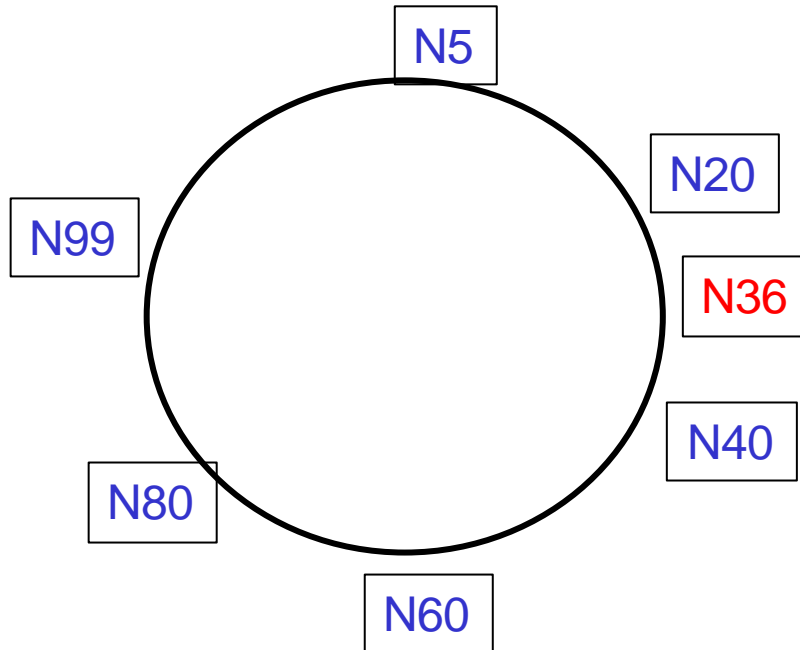
✍ Return results to new node



Joining the Ring - Step 2

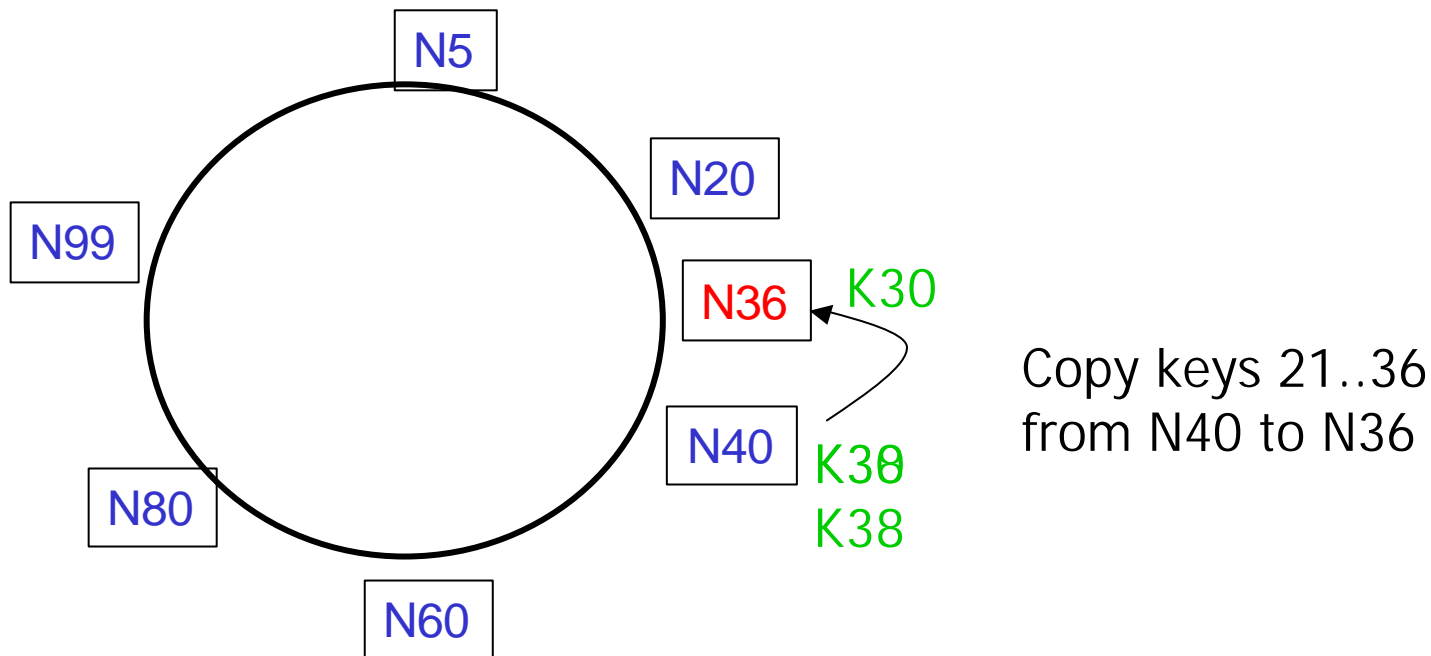
✍ Updating fingers of existing nodes

- ✍ new node calls update function on existing nodes
- ✍ existing nodes can recursively update fingers of other nodes



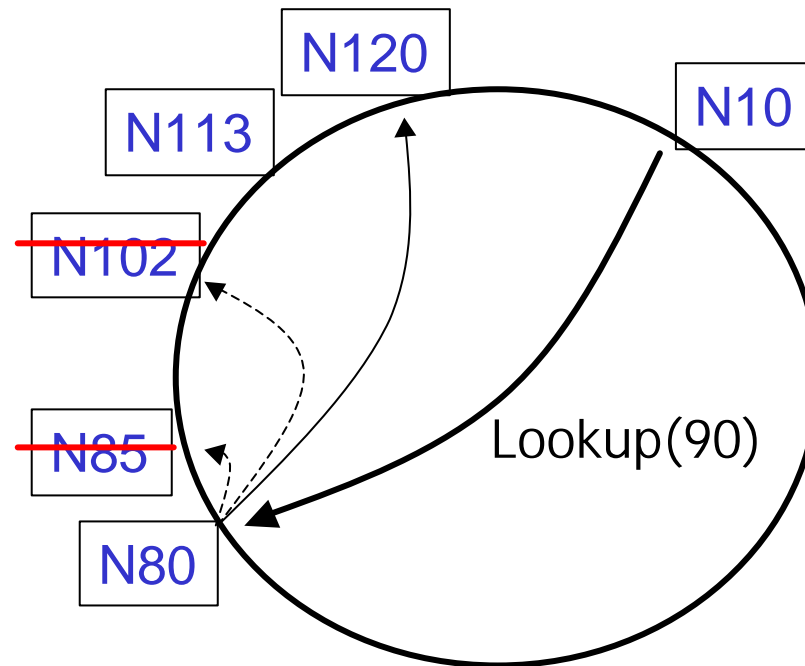
Joining the Ring - Step 3

- ✍ Transfer keys from successor node to new node
 - ✍ only keys in the range are transferred



Handling Failures




- ✗ Failure of nodes might cause incorrect lookup




- ✗ N80 doesn't know correct successor, so lookup fails
- ✗ Successor fingers are enough for correctness

Handling Failures

Use successor list

-  Each node knows r immediate successors
-  After failure, will know first live successor
-  Correct successors guarantee correct lookups

Guarantee is with some probability

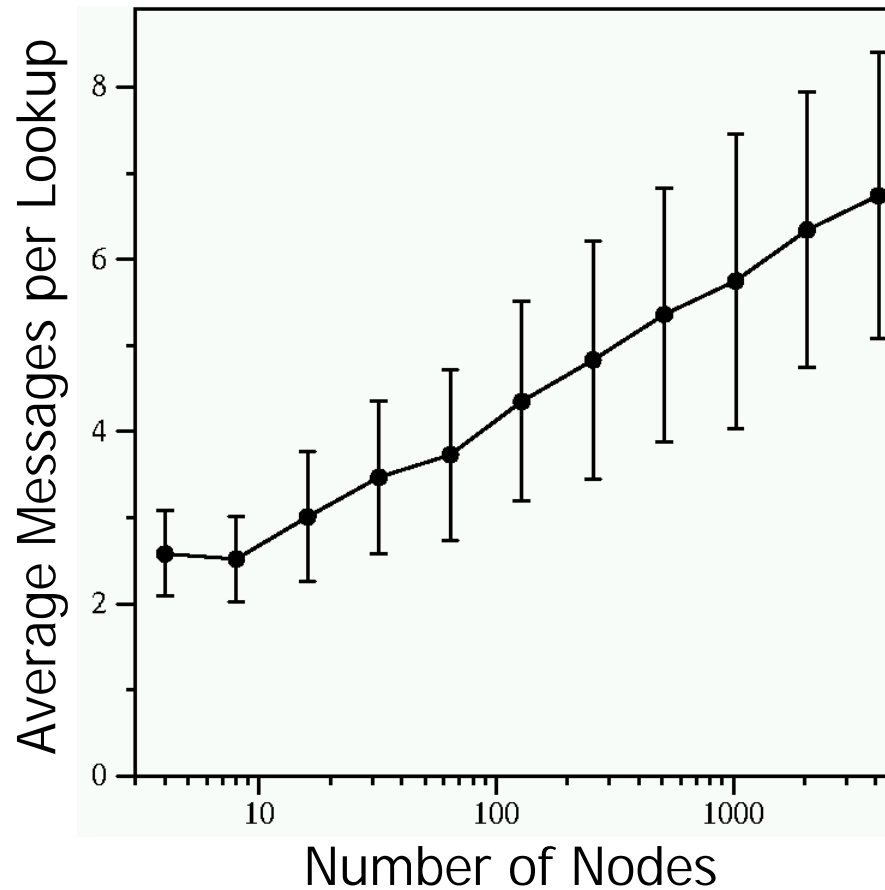
-  Can choose r to make probability of lookup failure arbitrarily small

Evaluation Overview

- ✍ Quick lookup in large systems
- ✍ Low variation in lookup costs
- ✍ Robust despite massive failure
- ✍ Experiments confirm theoretical results

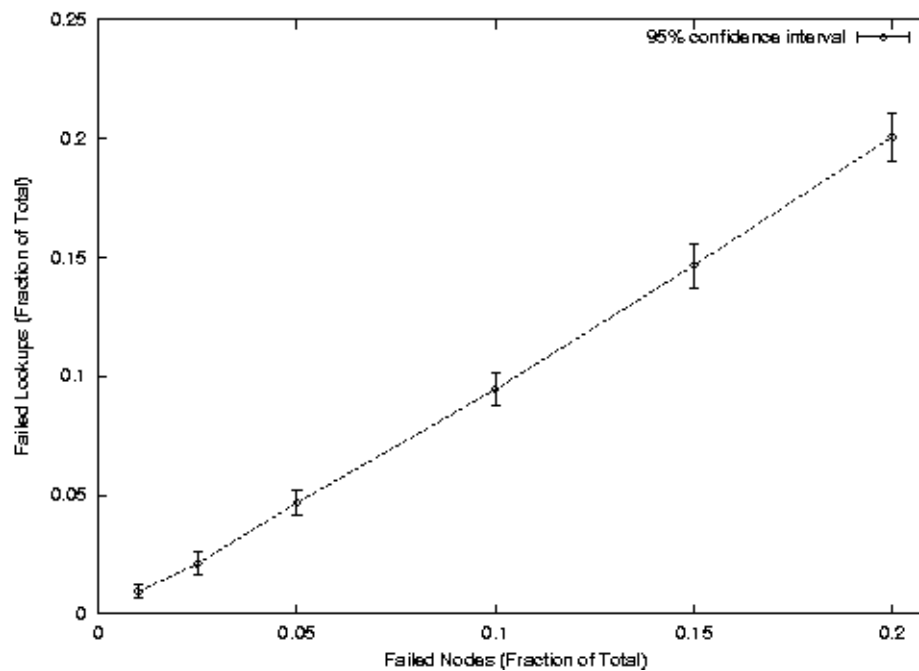
Cost of lookup

- ✗ Cost is $O(\log N)$ as predicted by theory
- ✗ constant is $1/2$



Robustness

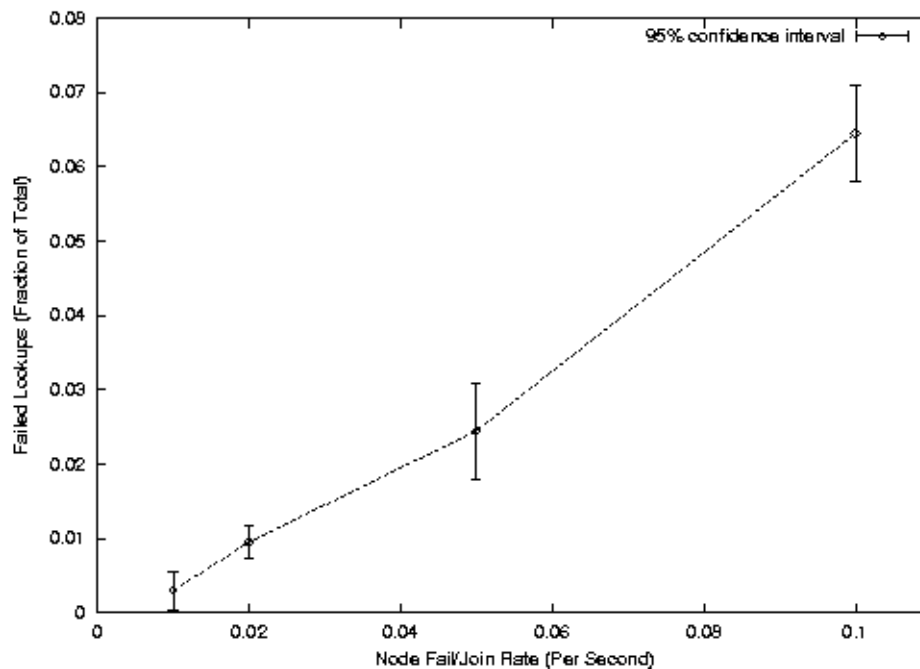
- ✍ Simulation results: static scenario
- ✍ Failed lookup means original node with key failed (no replica of keys)



- ✍ Result implies good balance of keys among nodes!

Robustness

- ✍ Simulation results: dynamic scenario
- ✍ Failed lookup means finger path has a failed node



- ✍ 500 nodes initially
- ✍ average *stabilize()* call 30s
- ✍ 1 lookup per second (Poisson)
- ✍ x join/fail per second (Poisson)

Current implementation

- ✍ Chord library: 3,000 lines of C++
- ✍ Deployed in small Internet testbed
- ✍ Includes:
 - ✍ Correct concurrent join/fail
 - ✍ Proximity-based routing for low delay (?)
 - ✍ Load control for heterogeneous nodes (?)
 - ✍ Resistance to spoofed node IDs (?)

Strengths

- ✍ Based on theoretical work (consistent hashing)
- ✍ Proven performance in many different aspects
 - ✍ “with high probability” proofs
- ✍ Robust (Is it?)

Weakness

- ✍ **NOT** that simple (compared to CAN)
- ✍ Member joining is complicated
 - ✍ aggressive mechanisms requires too many messages and updates
 - ✍ no analysis of convergence in lazy finger mechanism
- ✍ Key management mechanism mixed between layers
 - ✍ upper layer does insertion and handle node failures
 - ✍ Chord transfer keys when node joins (no leave mechanism!)
- ✍ Routing table grows with # of members in group
- ✍ Worst case lookup can be slow

Discussions

- ✍ Network proximity (consider latency?)
- ✍ Protocol security
 - ✍ Malicious data insertion
 - ✍ Malicious Chord table information
- ✍ Keyword search and indexing
- ✍ ...

7. Tapestry: Decentralized Routing and Location

Ben Y. Zhao
CS Division, U. C. Berkeley

Outline

- r Problems facing wide-area applications
- r Tapestry Overview
- r Mechanisms and protocols
- r Preliminary Evaluation
- r Related and future work

Motivation

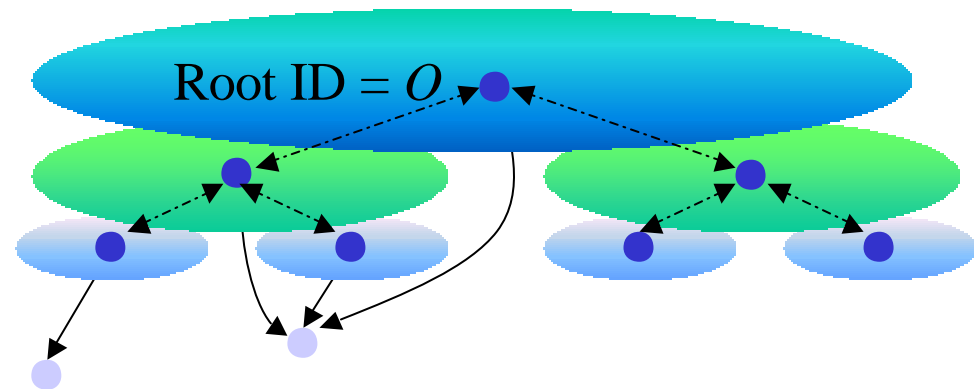
- r Shared Storage systems need an data location/routing mechanism
 - m Finding the peer in a scalable way is a difficult problem
 - m Efficient insertion and retrieval of content in a large distributed storage infrastructure
- r Existing solutions
 - m Centralized: expensive to scale, less fault tolerant, vulnerable to DoS attacks (e.g. Napster, DNS, SDS)
 - m Flooding: not scalable (e.g. Gnutella)

Key: Location and Routing

- r Hard problem:
 - m Locating and messaging to resources and data
- r Approach: wide-area *overlay infrastructure*:
 - m Scalable, Dynamic, Fault-tolerant, Load balancing

Decentralized Hierarchies

- r Centralized hierarchies
 - m Each higher level node responsible for locating objects in a greater domain
- r Decentralize: Create a tree for object O (really!)
 - m Object O has its own root and subtree
 - m Server on each level keeps pointer to **nearest** object in domain
 - m Queries search up in hierarchy



Directory servers tracking 2 replicas

What is Tapestry?

- r A prototype of a *decentralized, scalable, fault-tolerant, adaptive* location and routing infrastructure
(Zhao, Kubiatowicz, Joseph et al. U.C. Berkeley)
- r Network layer of **OceanStore** global storage system
Suffix-based hypercube routing
 - m Core system inspired by **Plaxton Algorithm** (Plaxton, Rajamaran, Richa (SPAA97))
- r Core API :
 - m *publishObject(ObjectID, [serverID])*
 - m *sendmsgToObject(ObjectID)*
 - m *sendmsgToNode(NodeID)*

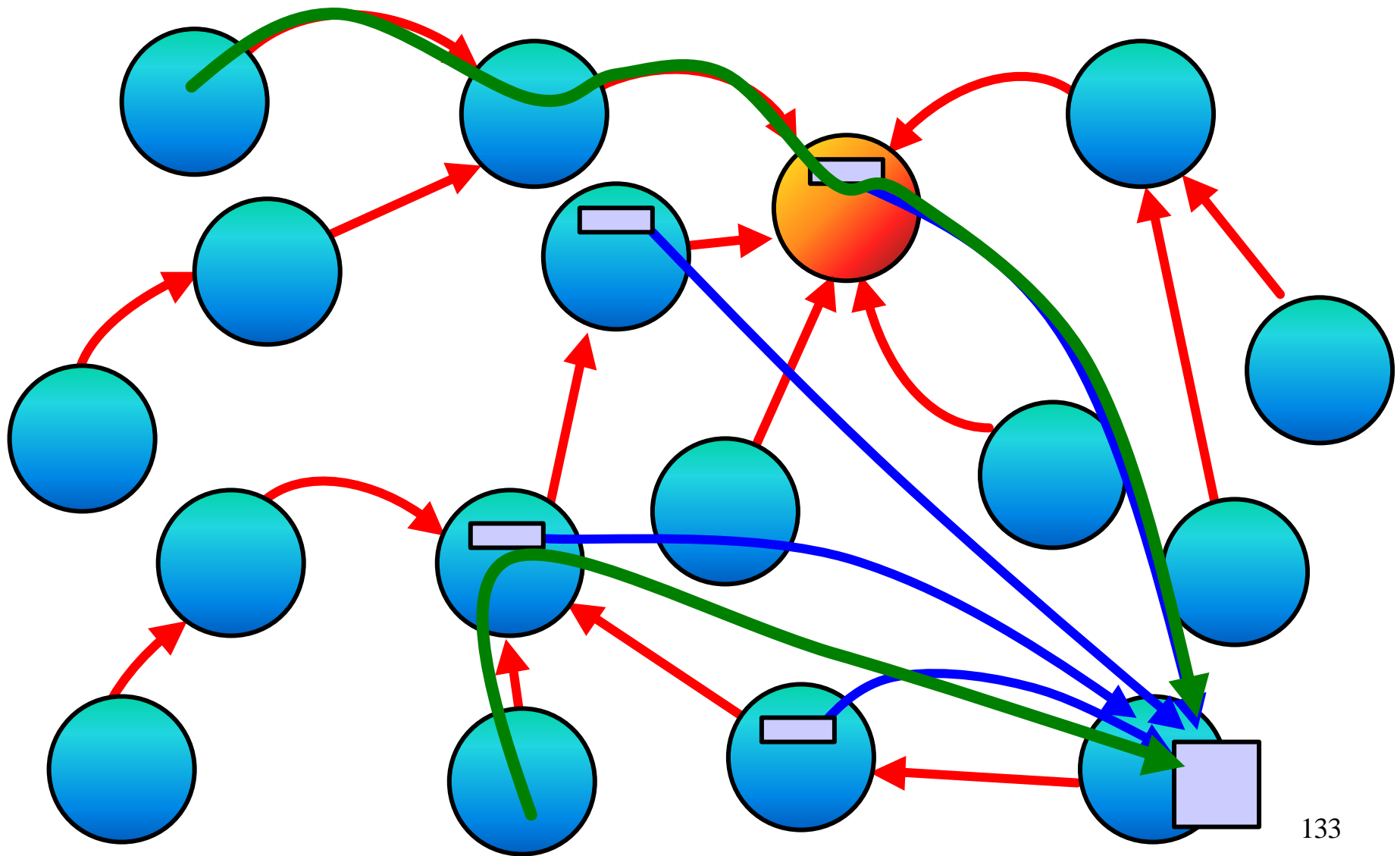
Incremental Suffix Routing

- r Namespace (nodes and objects)
 - m large enough to avoid collisions ($\sim 2^{160}$?)
(size N in $\log_2(N)$ bits)
- r Insert Object:
 - m Hash Object into namespace to get ObjectID
 - m For $(i=0, i < \log_2(N), i+j)$ { //Define hierarchy
 - j is base of digit size used, ($j = 4$ ~~hex~~ hex digits)
 - Insert entry into nearest node that matches on last i bits
 - When no matches found, then pick node matching $(i - n)$ bits with highest ID value, terminate

Routing to Object

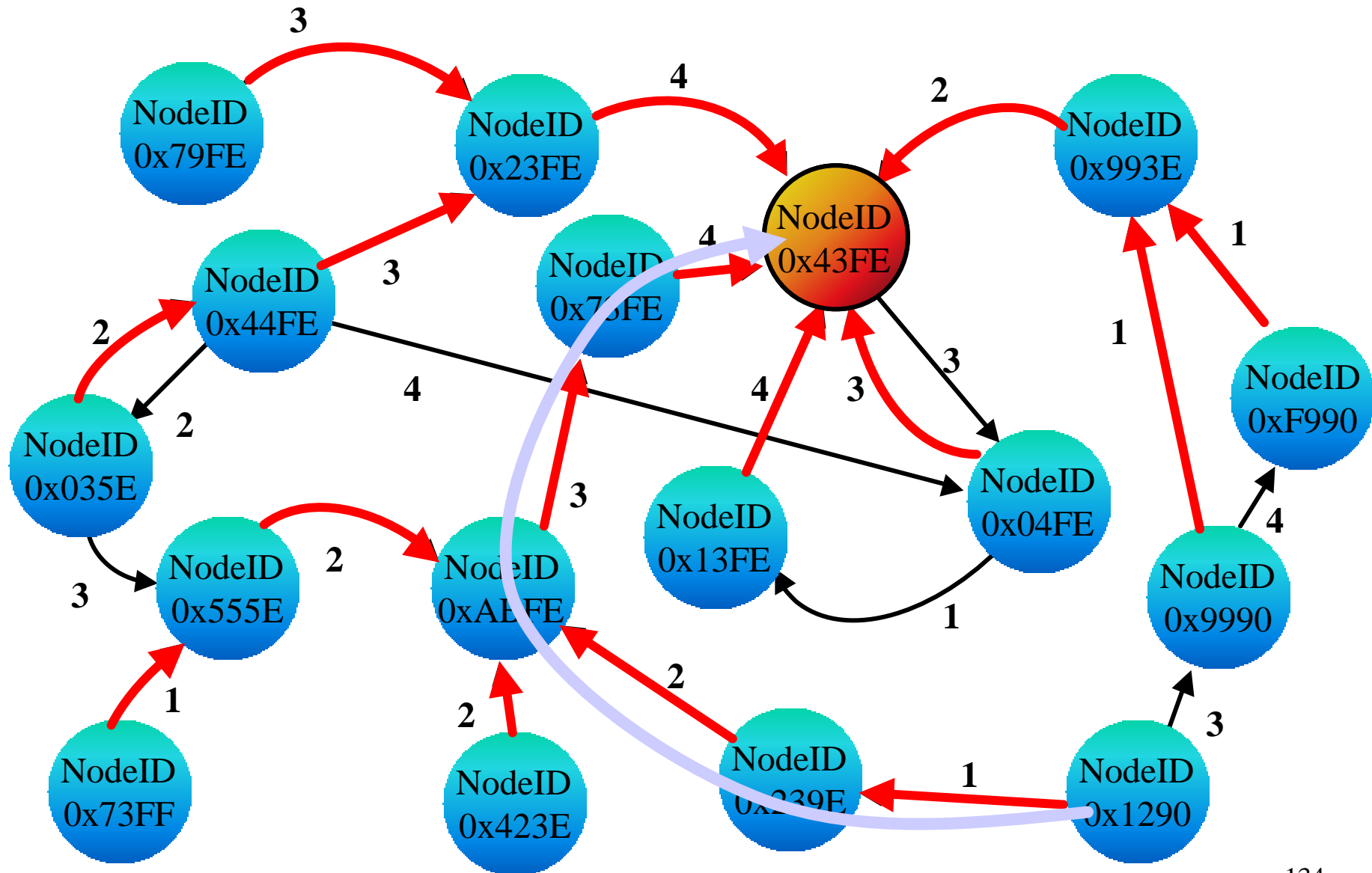
- r Lookup object
 - m Traverse same relative nodes as insert, **except searching for entry at each node**
 - m For $(i=0, i < \log_2(N), i+n)$
Search for entry in nearest node matching on last i bits
- r Each object maps to hierarchy defined by single root
 - m $f(\text{ObjectID}) = \text{RootID}$
- r Publish / search both route incrementally to root
- r Root node = $f(O)$, is responsible for “knowing” object’s location

Object Location Randomization and Locality



Tapestry Mesh

Incremental suffix-based routing



Contribution of this work

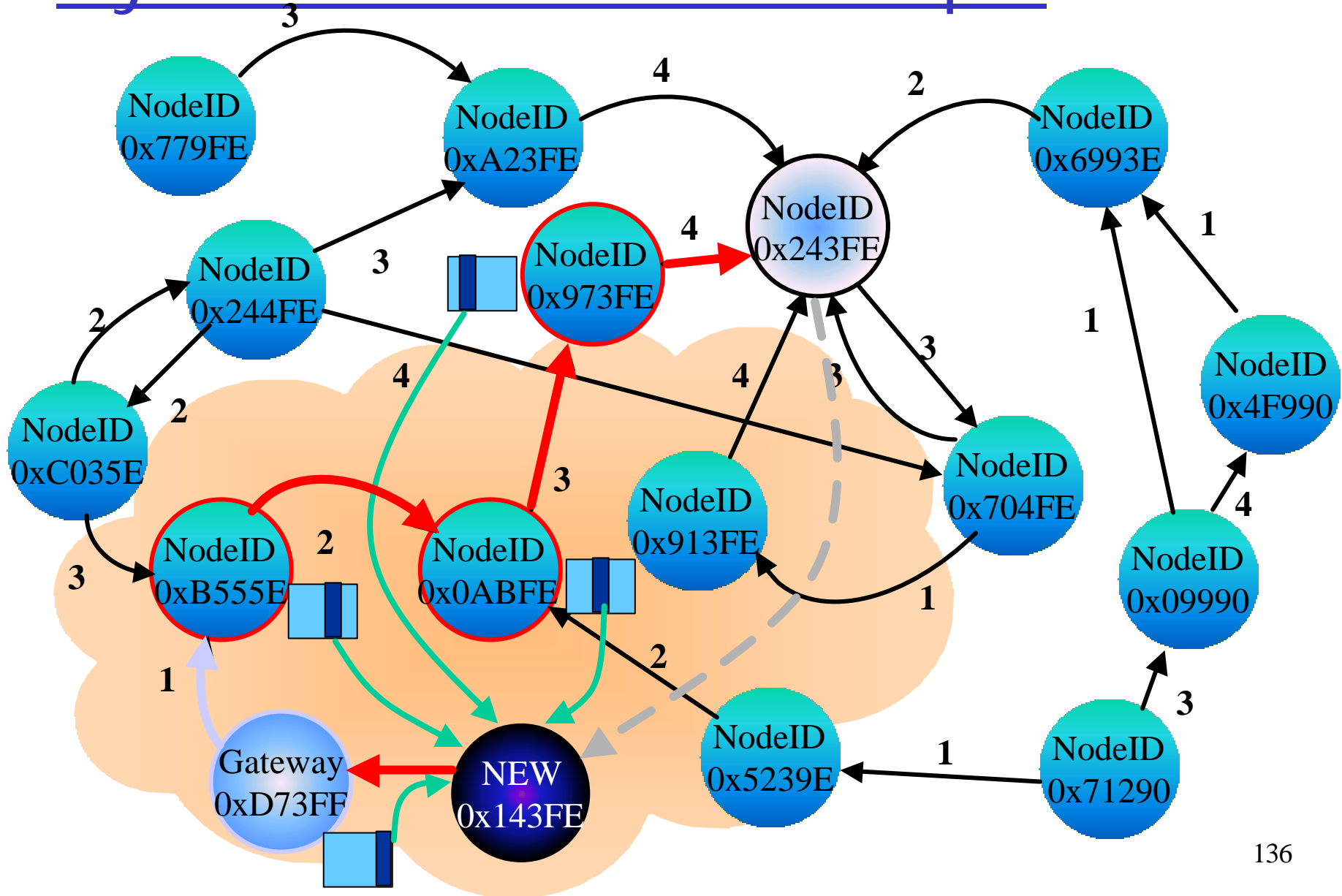
Plaxtor Algorithm

- r Limitations
 - m Global knowledge algorithms
 - m Root node vulnerability
 - m Lack of adaptability

Tapestry

- r Distributed algorithms
 - Dynamic node insertion
 - Dynamic root mapping
- m Redundancy in location and routing
- m Fault-tolerance protocols
- m *Self-configuring / adaptive*
- m *Support for mobile objects*
- r Application Infrastructure

Dynamic Insertion Example



Fault-tolerant Location

- r Minimized soft-state vs. explicit fault-recovery
- r Multiple roots
 - m Objects hashed w/ small salts ~~↗~~ multiple names/roots
 - m Queries and publishing utilize all roots in parallel
 - m $P(\text{finding Reference w/ partition}) = 1 - (1/2)^n$
where $n = \#$ of roots
- r Soft-state periodic republish
 - m 50 million files/node, daily republish,
 $b = 16$, $N = 2^{160}$, 40B/msg,
worst case update traffic: 156 kb/s,
 - m expected traffic w/ 2^{40} real nodes: 39 kb/s

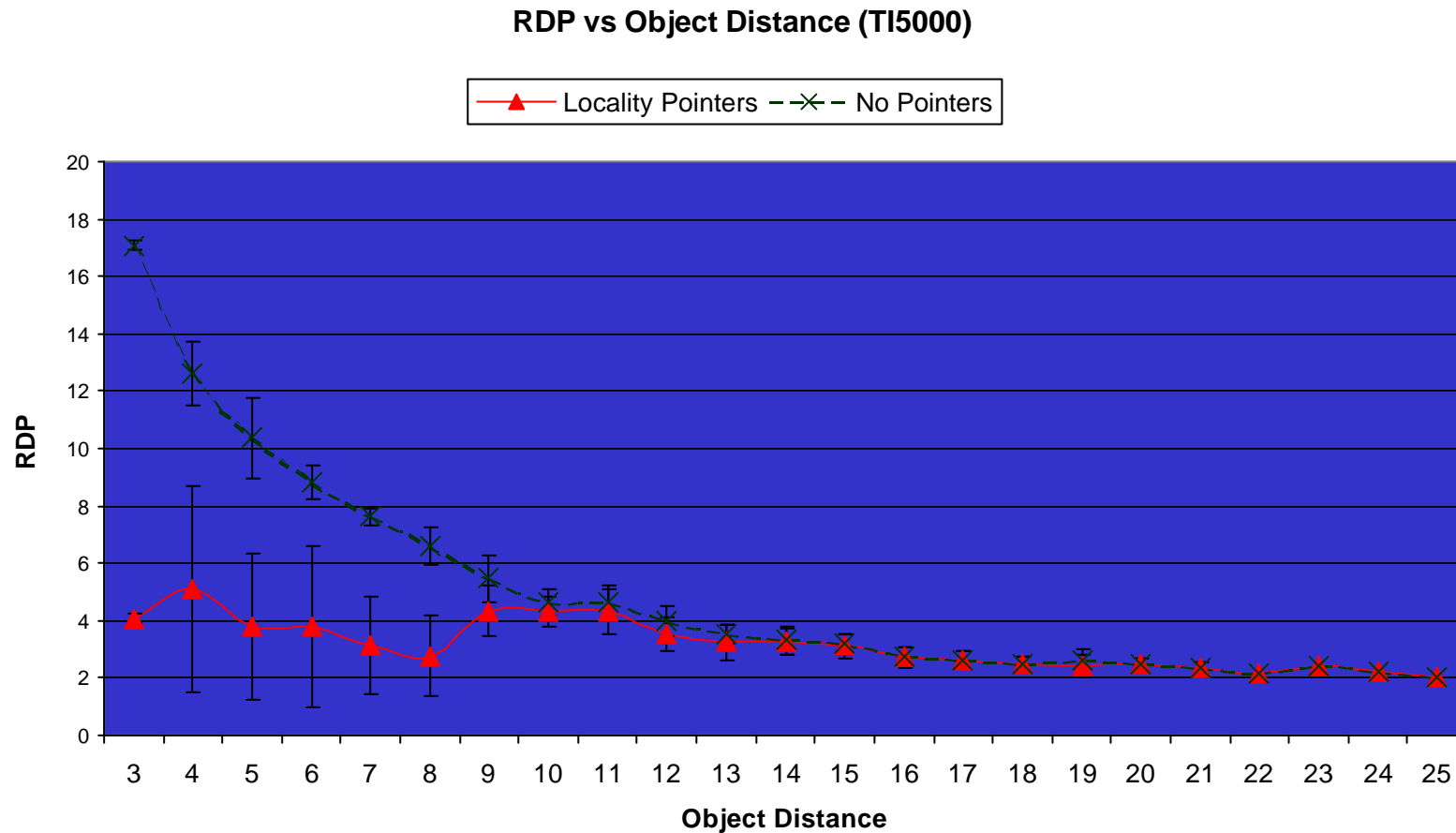
Fault-tolerant Routing

- r Detection:
 - m Periodic probe packets between neighbors
- r Handling:
 - m Each entry in routing map has 2 alternate nodes
 - m Second chance algorithm for intermittent failures
 - m Long term failures ✂ alternates found via routing tables
- r Protocols:
 - m First Reachable Link Selection
 - m Proactive Duplicate Packet Routing

Simulation Environment

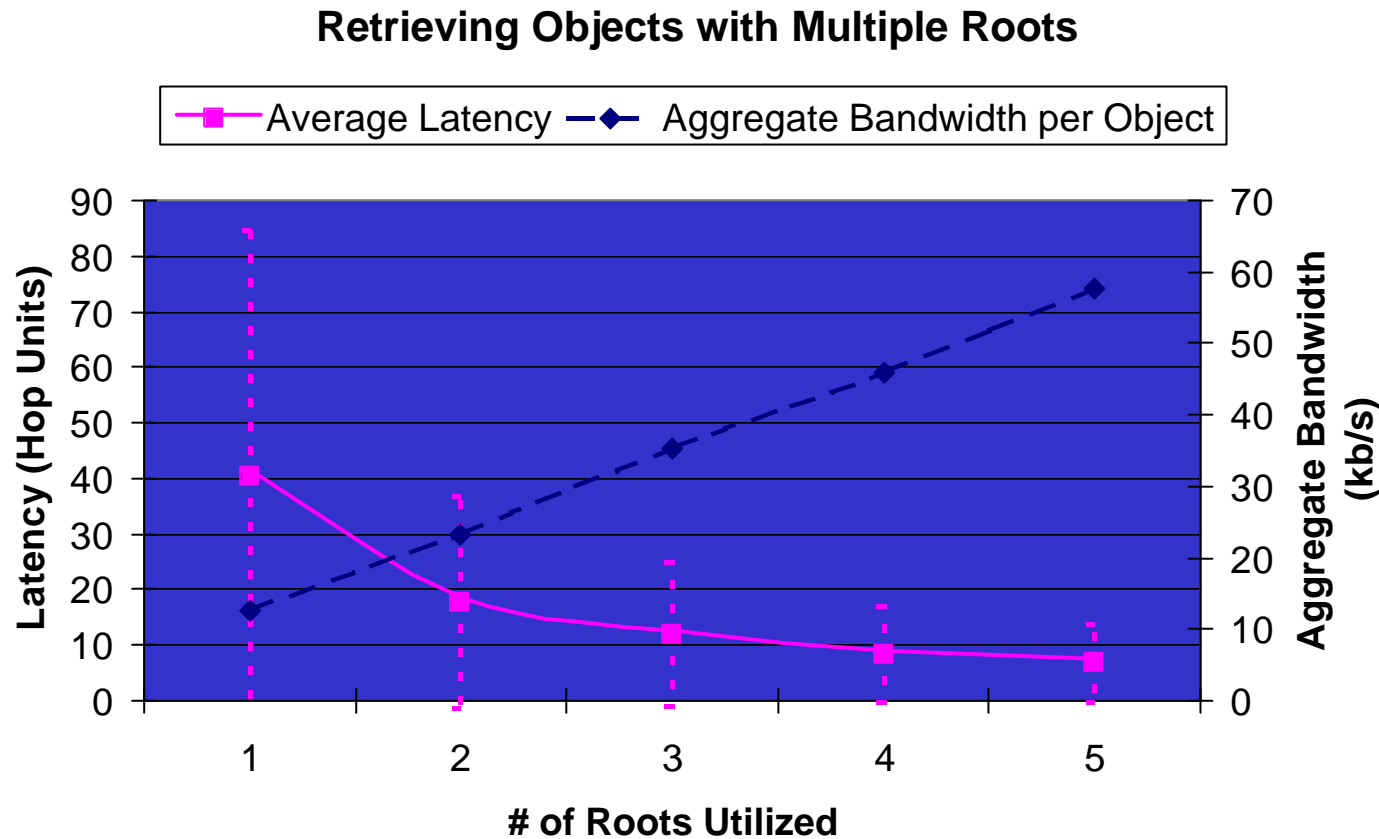
- r Implemented Tapestry routing as packet-level simulator
- r Delay is measured in terms of network hops
- r Do not model the effects of cross traffic or queuing delays
- r Four topologies: AS, MBone, GT-ITM, TIERS

Results: Location Locality



Measuring effectiveness of locality pointers (TI ERS 5000)

Results: Stability via Redundancy



Parallel queries on multiple roots. Aggregate bandwidth measures b/w used for soft-state republish 1/day and b/w used by requests at rate of 1/s.

Related Work

r Content Addressable Networks

m Ratnasamy et al.,
(ACI RI / UCB)

r Chord

m Stoica, Morris, Karger, Kaashoek,
Balakrishnan (MIT / UCB)

r Pastry

m Druschel and Rowstron
(Rice / Microsoft Research)

Strong Points

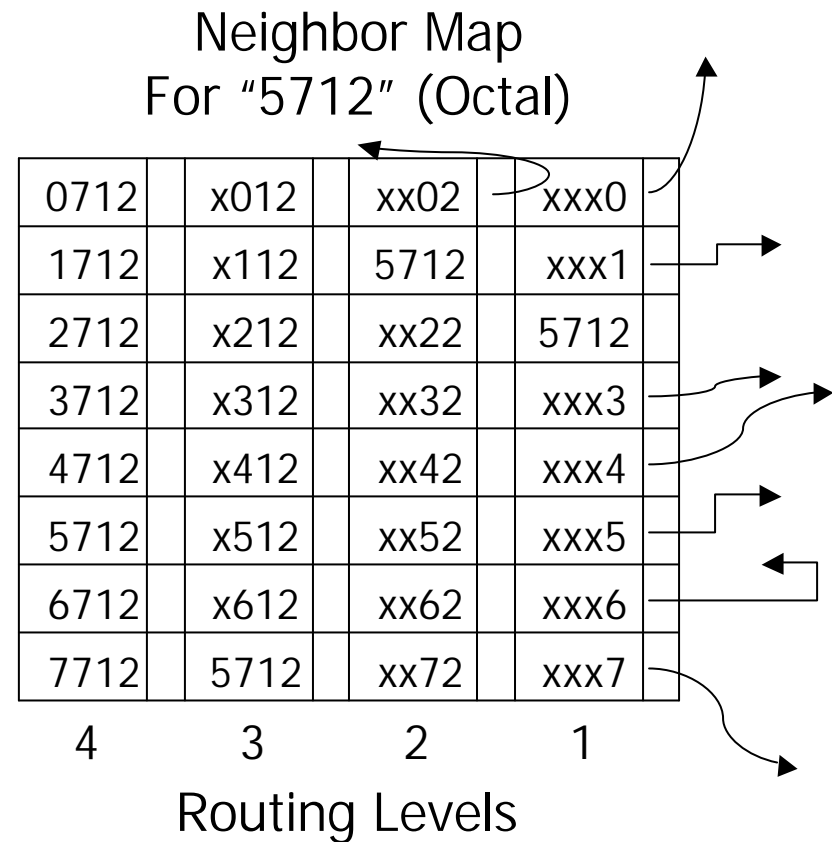
- r Designed system based on Theoretically proven idea (Plaxton Algorithm)
- r Fully decentralized and scalable solution for deterministic location and routing problem

Weaknesses/Improvements

- r Substantially complicated
 - m Esp, dynamic node insertion algorithm is non-trivial, and each insertion will take a non-negligible amount of time.
 - m Attempts to insert a lot of nodes at the same time
- r Where to put “root” node for a given object
 - m Needs universal hashing function
 - m Possible to put “root” to near expected clients dynamically?

Routing to Nodes

Example: Octal digits, 2^{18} namespace, 005712 ~~↗~~ 627510



Dynamic Insertion

Operations necessary for N to become fully integrated:

- r Step 1: Build up N 's routing maps
 - m Send messages to each hop along path from gateway to current node N' that best approximates N
 - m The i^{th} hop along the path sends its i^{th} level route table to N
 - m N optimizes those tables where necessary
- r Step 2: Send notify message via acked multicast to nodes with null entries for N 's ID, setup forwarding ptrs
- r Step 3: Each notified node issues republish message for relevant objects
- r Step 4: Remove forward ptrs after one republish period

Dynamic Root Mapping

- r Problem: choosing a root node for every object
 - m Deterministic over network changes
 - m Globally consistent
- r Assumptions
 - m All nodes with same matching suffix contains same null/non-null pattern in next level of routing map
 - m Requires: consistent knowledge of nodes across network

Plaxton Solution

- r Given desired I D N ,
 - m Find set S of nodes in existing network nodes n matching most # of suffix digits with N
 - m Choose S_i = node in S with highest valued I D
- r Issues:
 - m Mapping must be generated statically using global knowledge
 - m Must be kept as hard state in order to operate in changing environment
 - m Mapping is not well distributed, many nodes in n get no mappings

8. A Scalable, Content-Addressable Network

Sylvia Ratnasamy^{1,2}, Paul Francis,³ Mark Handley,¹

Richard Karp², Scott Shenker¹

¹
ACIRI

²
U.C. Berkeley

³
Tahoe
Networks

Outline

- r Introduction
- r Design
- r Evaluation
- r Strengths & Weaknesses
- r Ongoing Work

Internet-scale hash tables

- r Hash tables
 - m essential building block in software systems
- r Internet-scale distributed hash tables
 - m equally valuable to large-scale distributed systems?

Internet-scale hash tables

- r Hash tables

- m essential building block in software systems

- r Internet-scale distributed hash tables

- m equally valuable to large-scale distributed systems?

- peer-to-peer systems
 - Napster, Gnutella,, FreeNet, MojoNation...
 - large-scale storage management systems
 - Publius, OceanStore,, CFS ...
 - mirroring on the Web

Content-Addressable Network (CAN)

- r CAN: Internet-scale hash table
- r Interface
 - m insert(key,value)
 - m value = retrieve(key)







Content-Addressable Network (CAN)

- r CAN: Internet-scale hash table
- r Interface
 - m insert(key,value)
 - m value = retrieve(key)
- r Properties
 - m scalable
 - m operationally simple
 - m good performance (w/ improvement)

Content-Addressable Network (CAN)

- r CAN: Internet-scale hash table
- r Interface
 - m insert(key,value)
 - m value = retrieve(key)
- r Properties
 - m scalable
 - m operationally simple
 - m good performance
- r Related systems: Chord/Pastry/Tapestry/Buzz/Plaxton ...

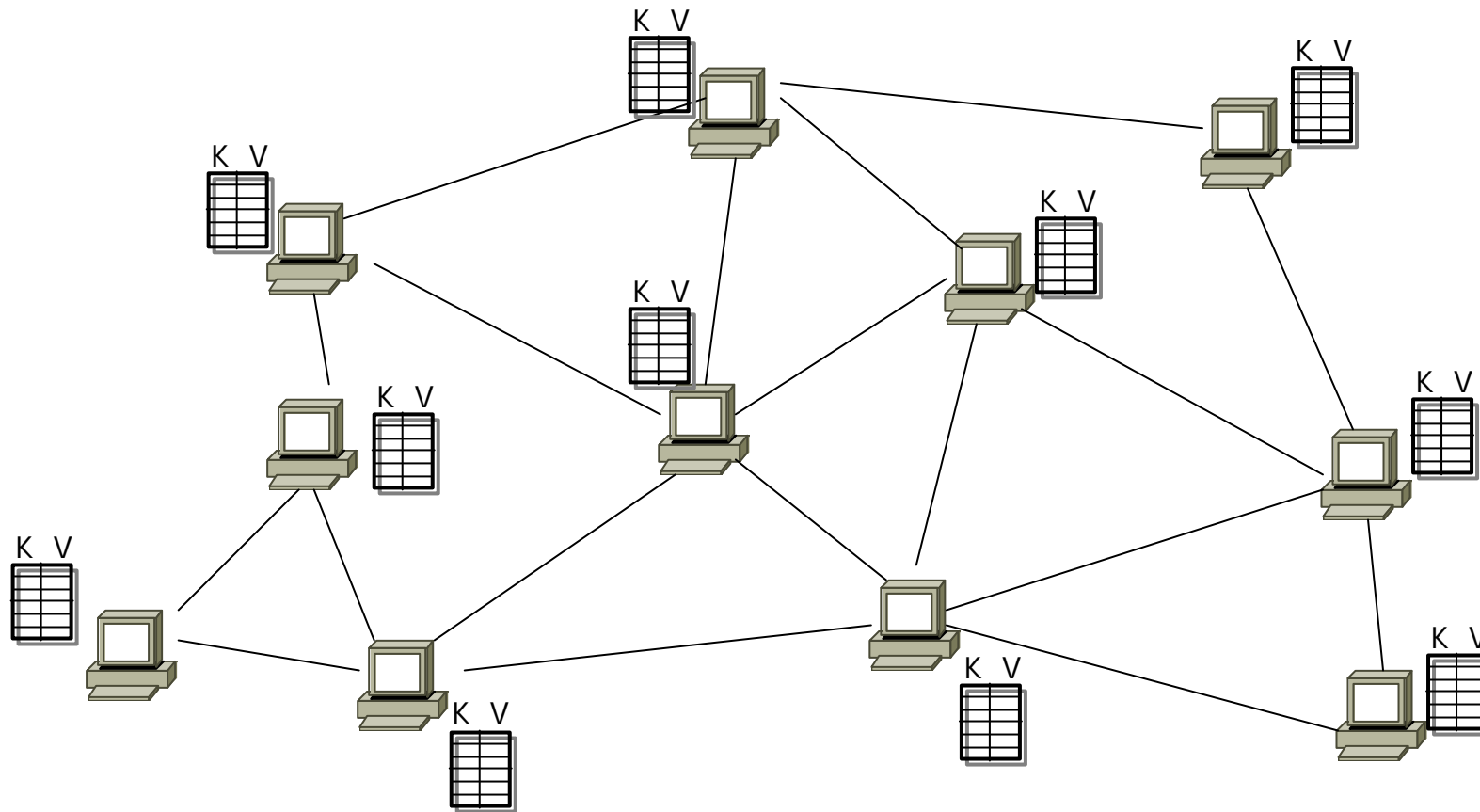
Problem Scope

- 4 Design a system that provides the interface
 -  scalability
 -  robustness
 -  performance
 - 5 security
- 6 Application-specific, higher level primitives
 -  keyword searching
 -  mutable content
 -  anonymity

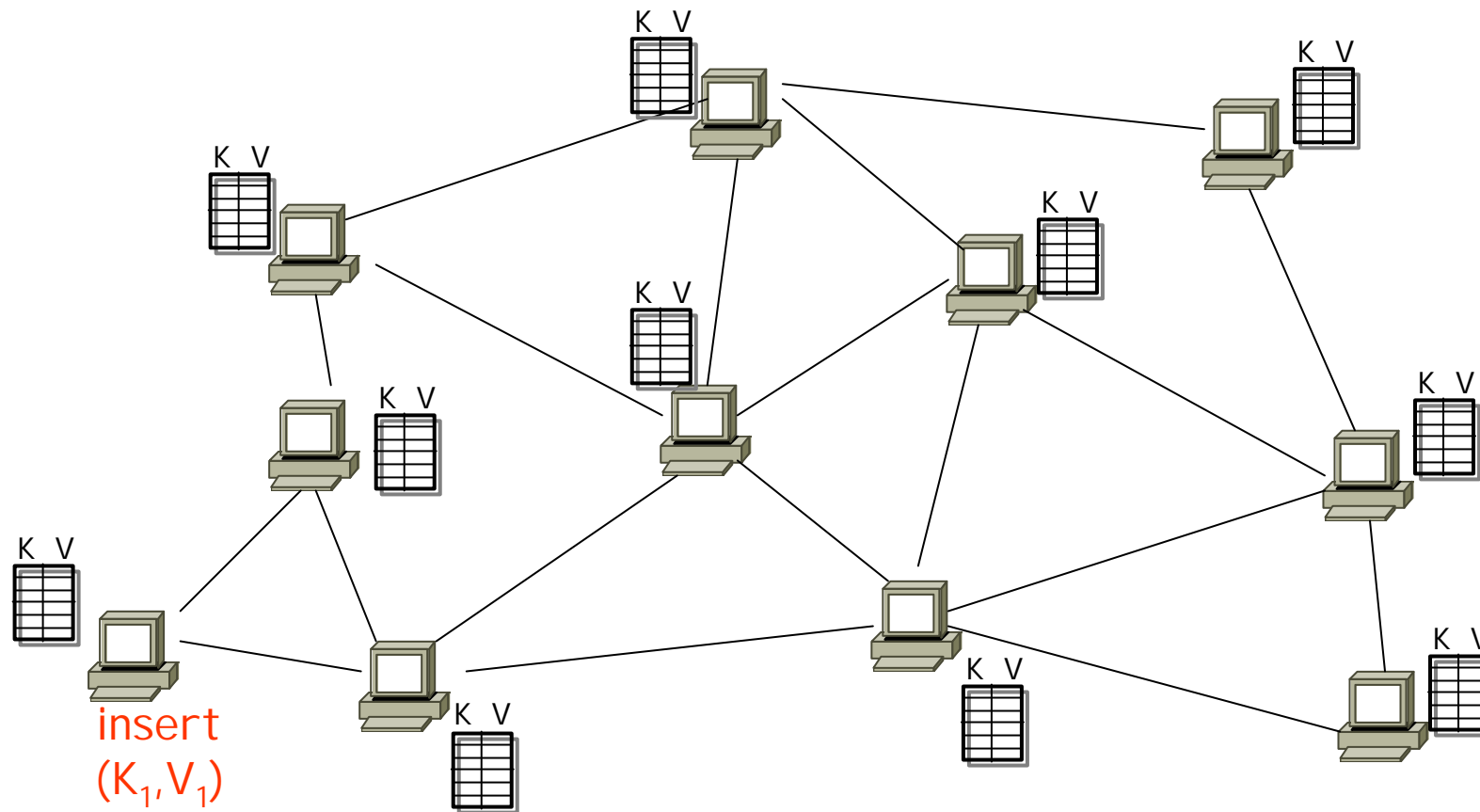
Outline

- r Introduction
- r **Design**
- r Evaluation
- r Strengths & Weaknesses
- r Ongoing Work

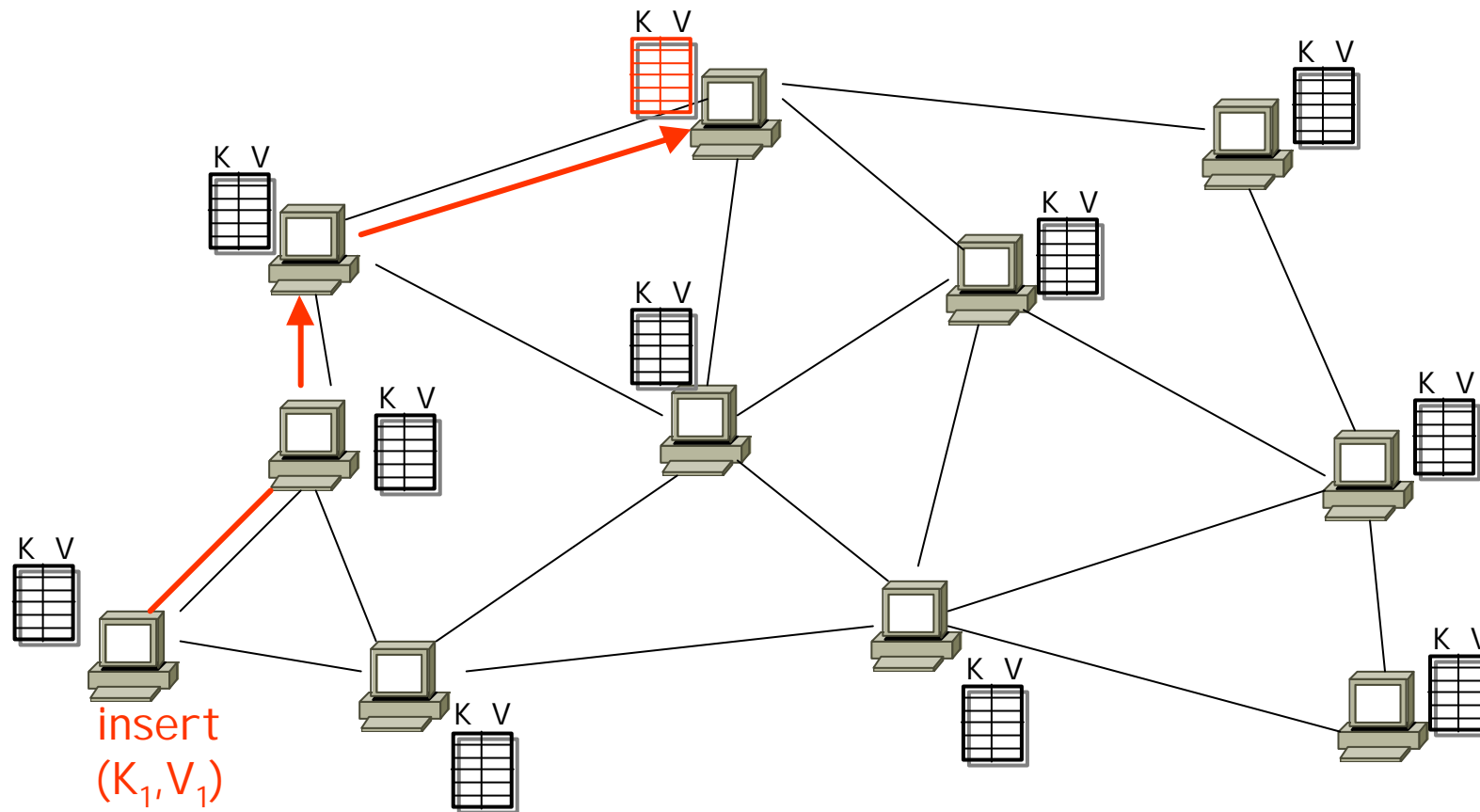
CAN: basic idea



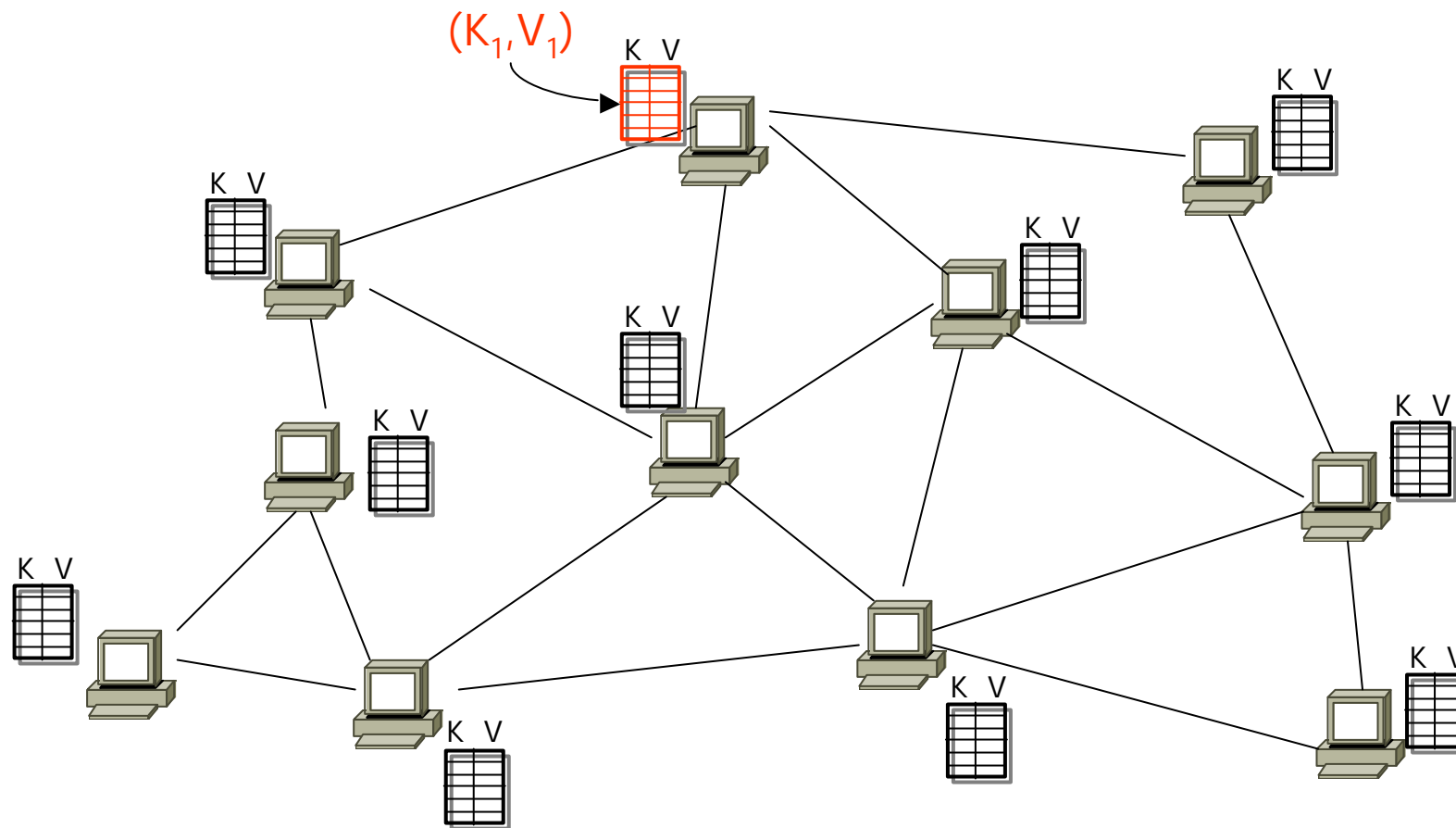
CAN: basic idea



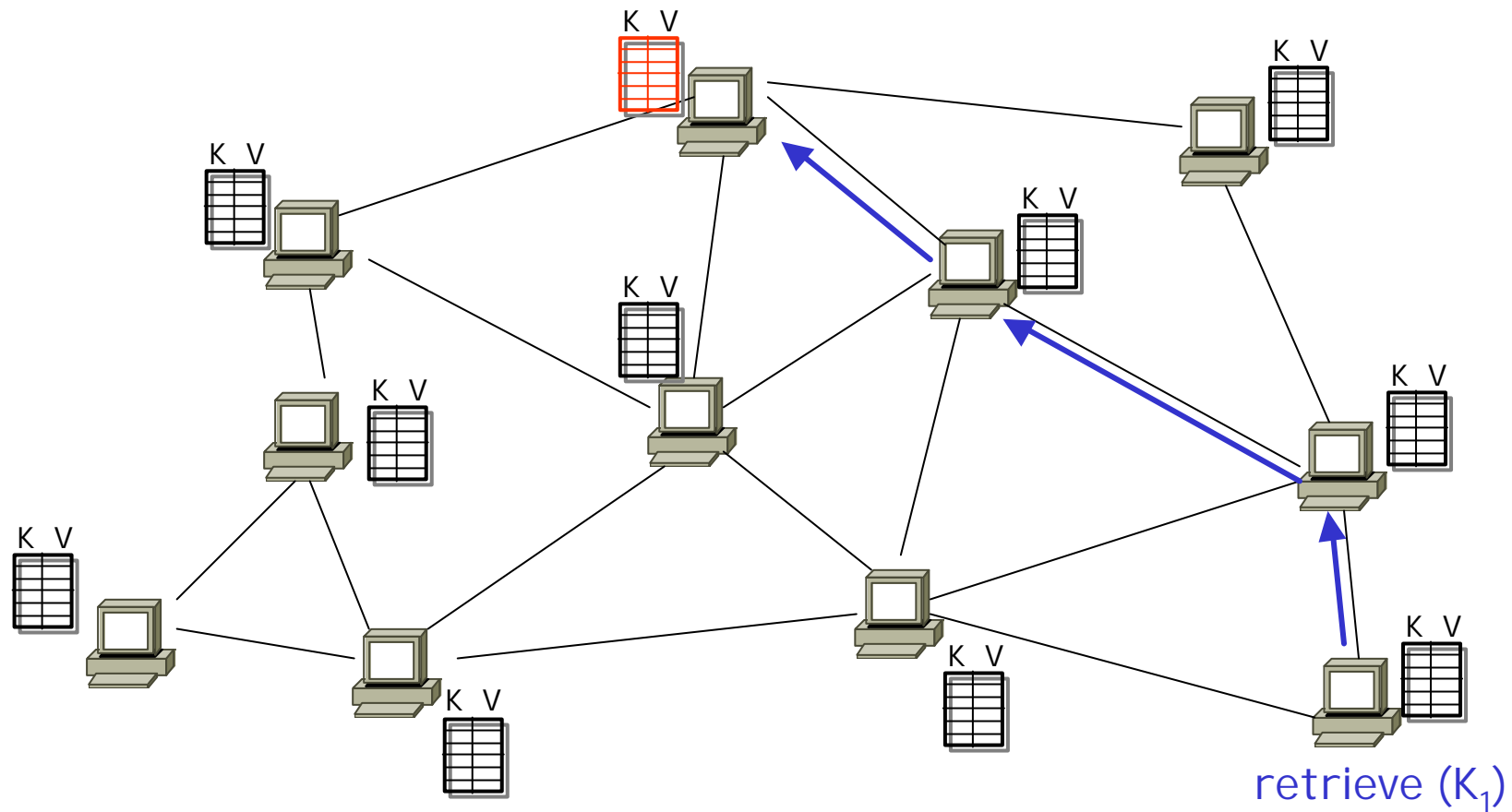
CAN: basic idea



CAN: basic idea



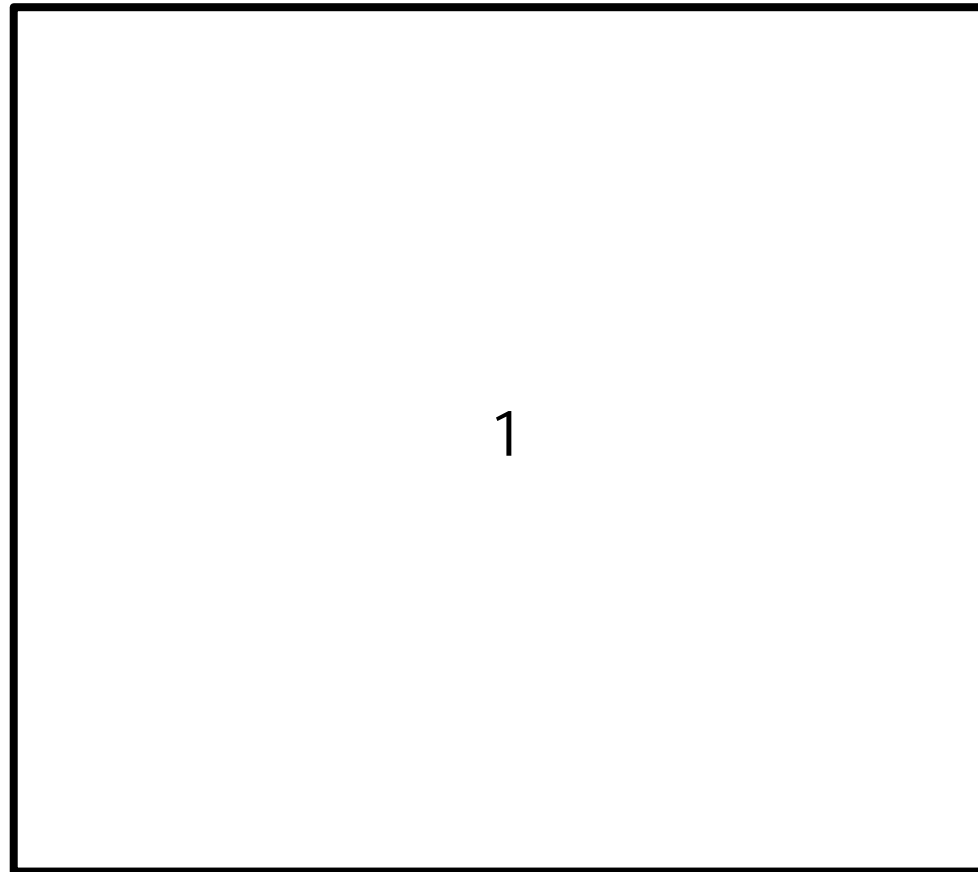
CAN: basic idea



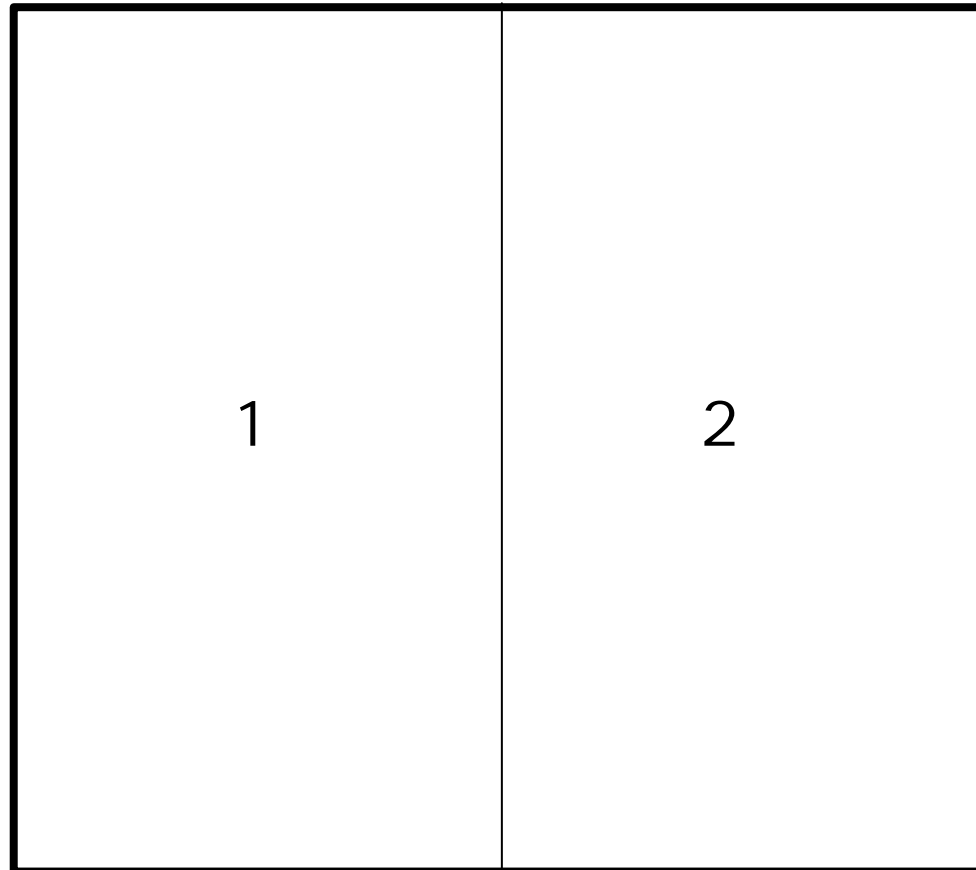
CAN: solution

- r virtual Cartesian coordinate space
- r entire space is partitioned amongst all the nodes
 - m every node “owns” a zone in the overall space
- r abstraction
 - m can store data at “points” in the space
 - m can route from one “point” to another
- r point = node that owns the enclosing zone

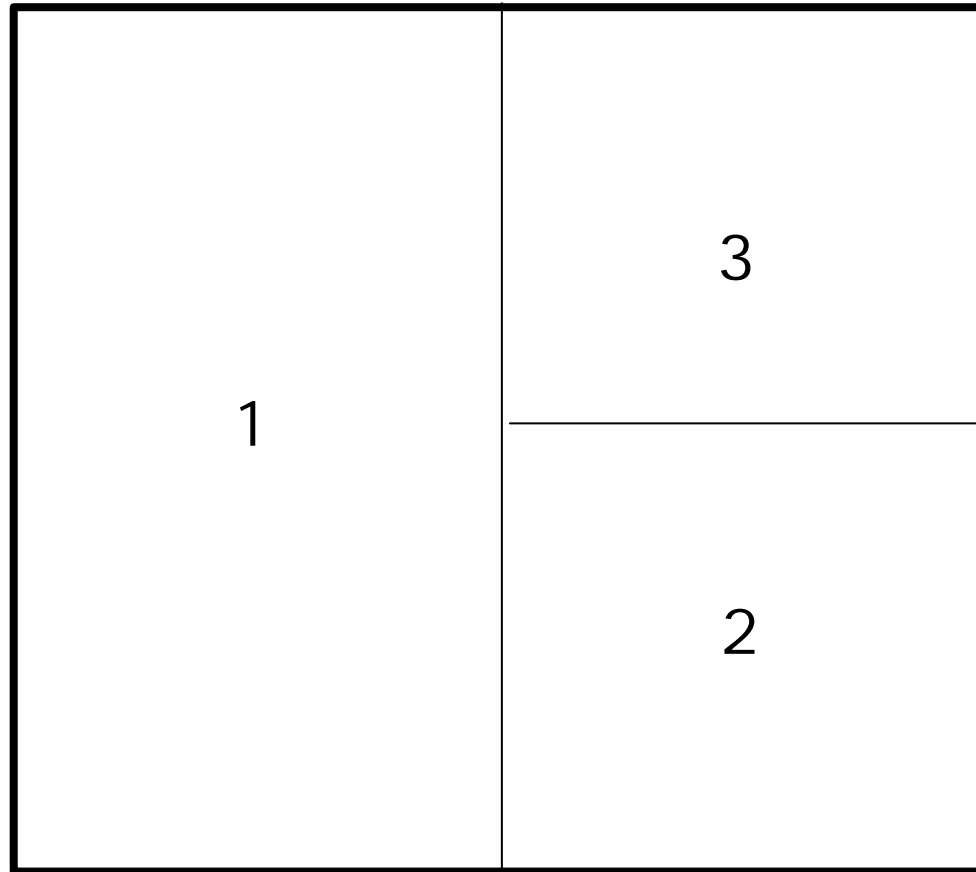
CAN: simple example



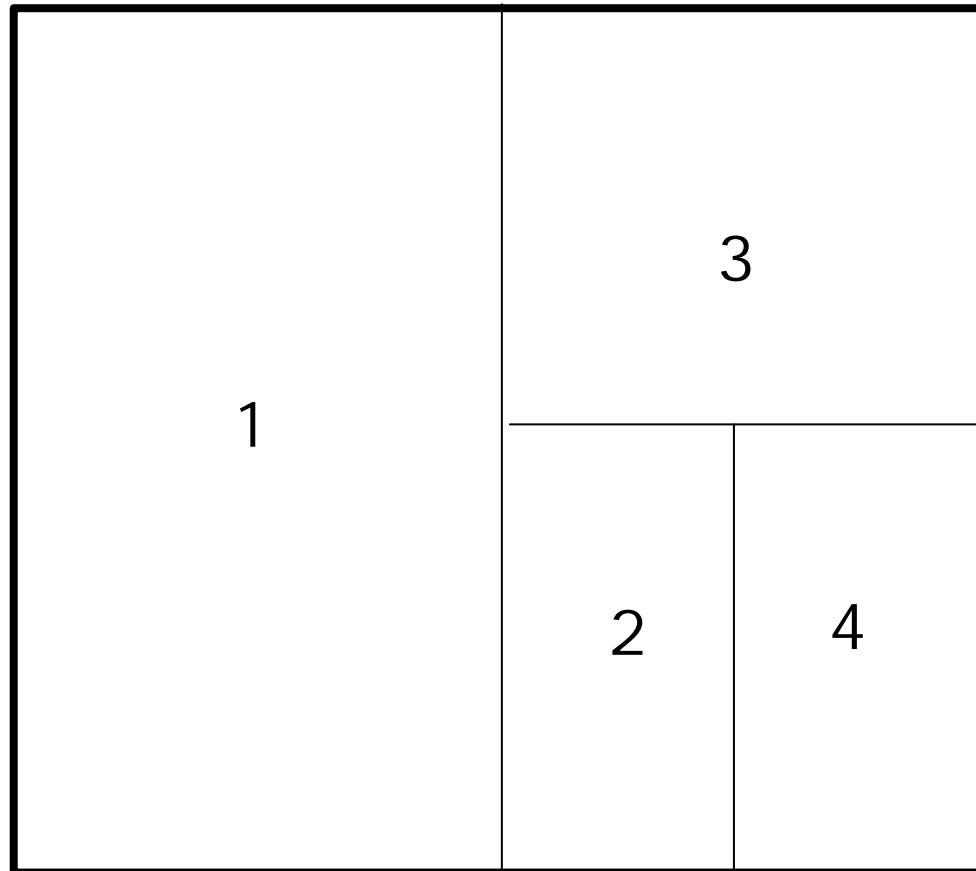
CAN: simple example



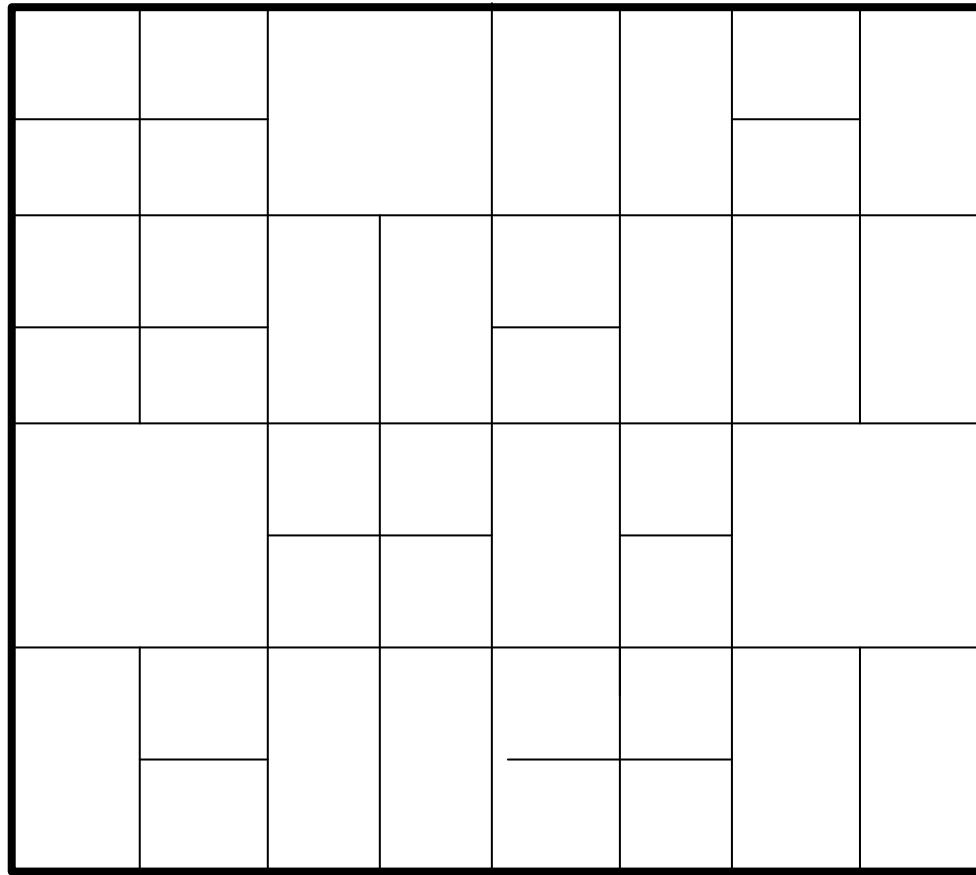
CAN: simple example



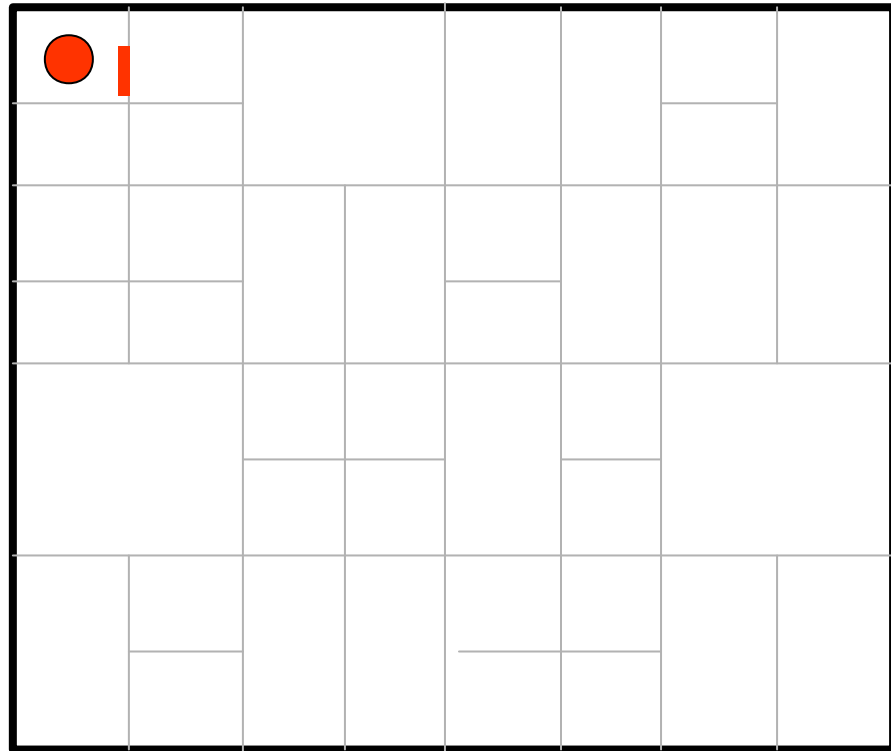
CAN: simple example



CAN: simple example

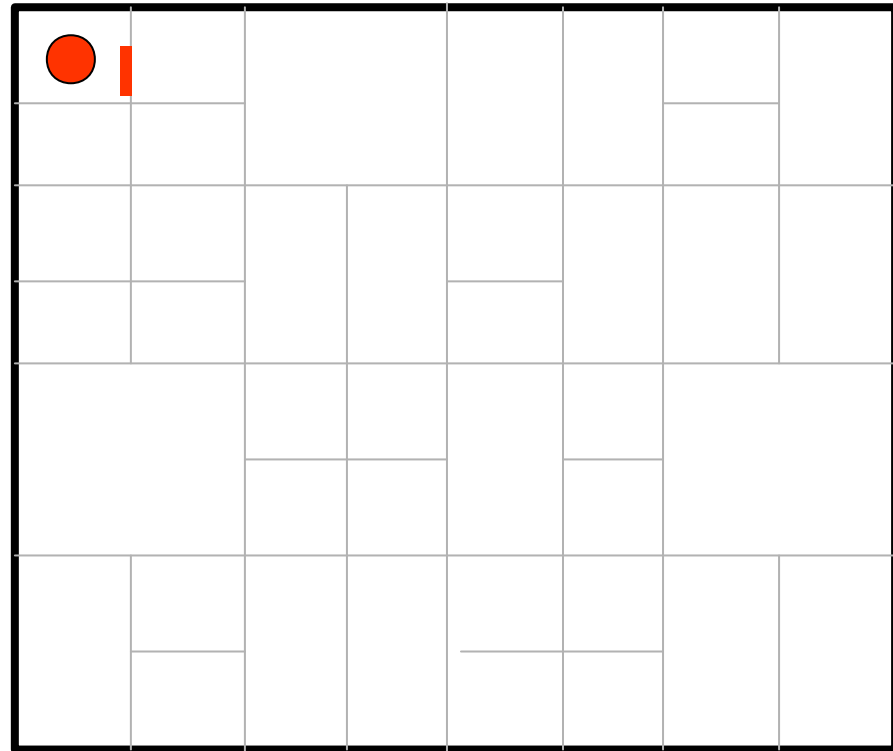


CAN: simple example



CAN: simple example

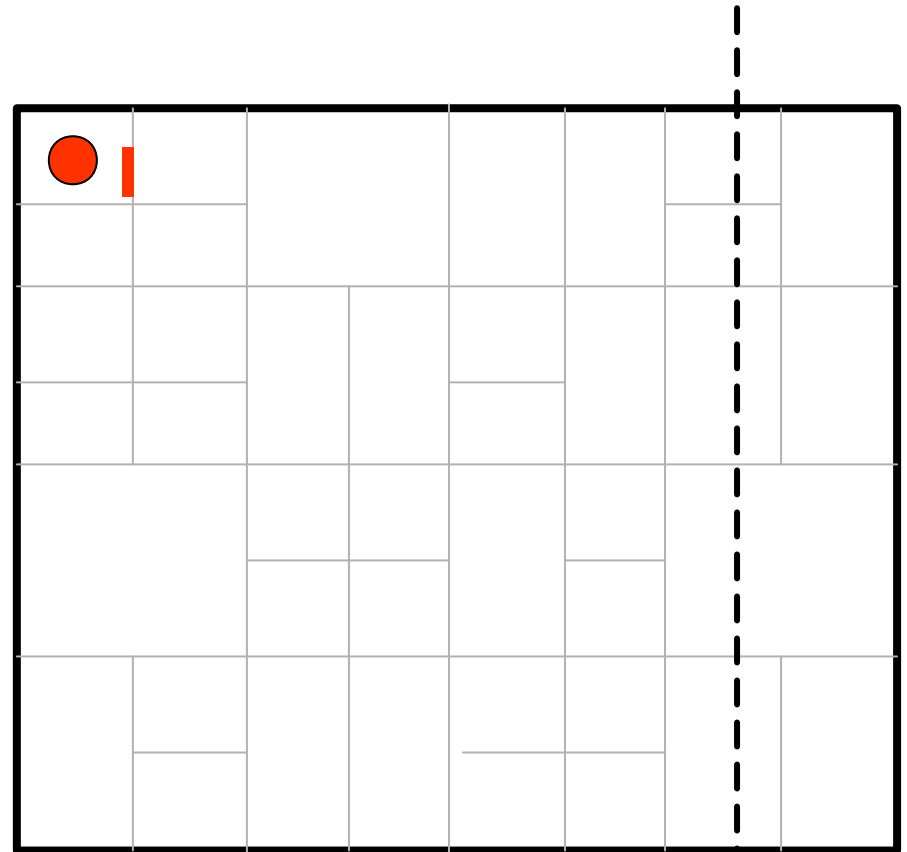
node I :: insert(K,V)



CAN: simple example

node I :: insert(K,V)

(1) $a = h_x(K)$

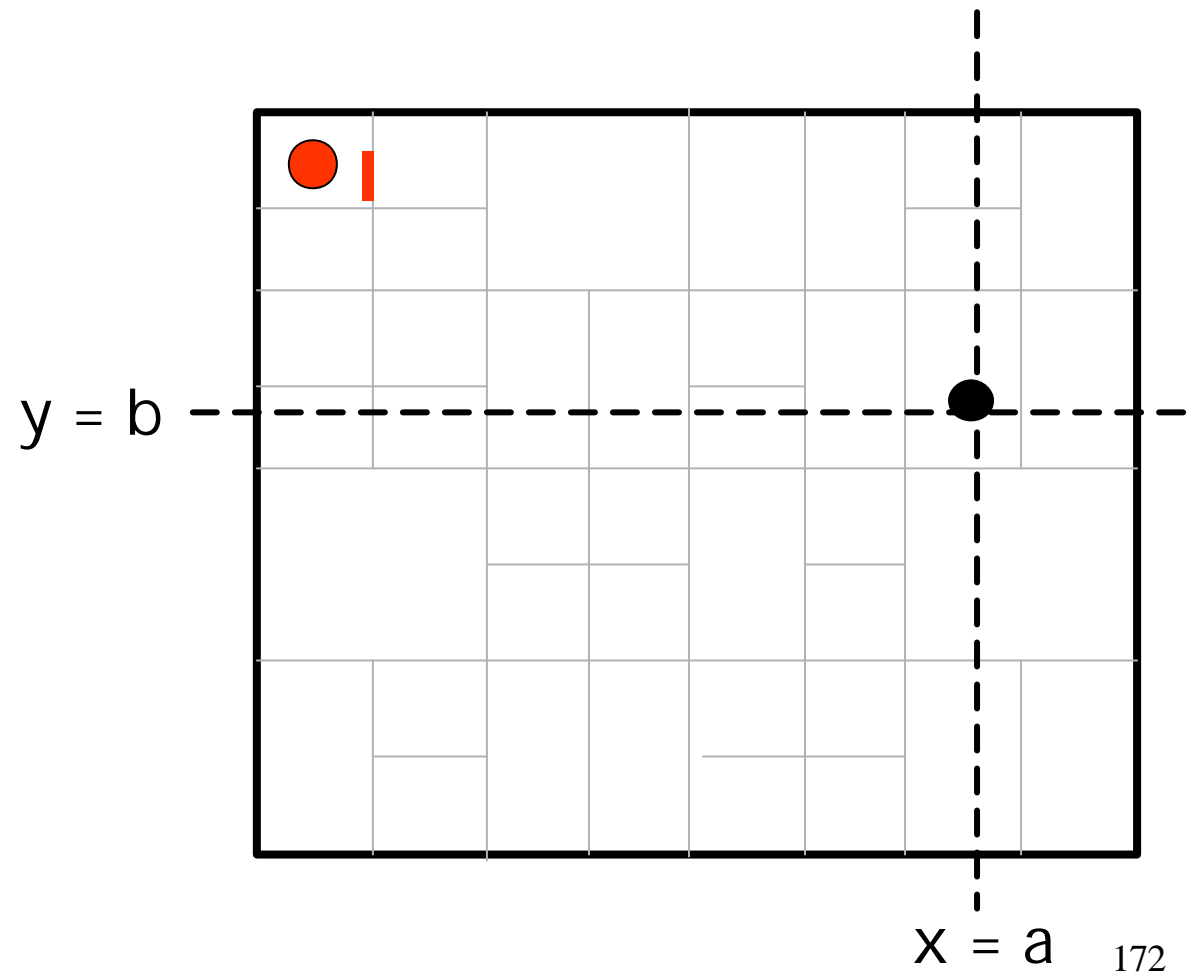


$x = a$ 171

CAN: simple example

node I :: insert(K,V)

(1) $a = h_x(K)$
 $b = h_y(K)$

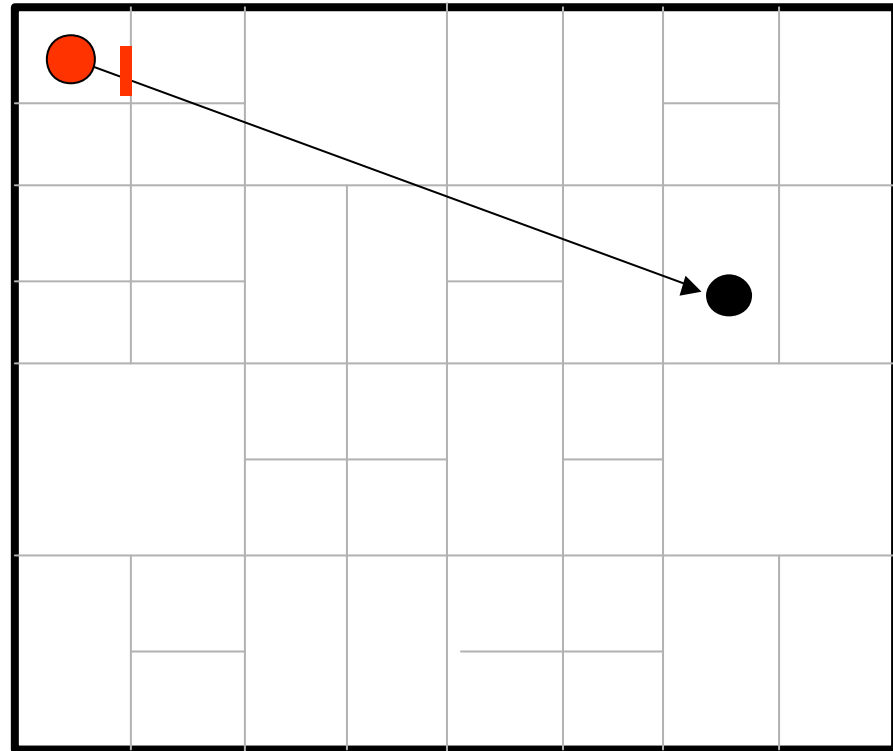


CAN: simple example

node I :: insert(K,V)

(1) $a = h_x(K)$
 $b = h_y(K)$

(2) route(K,V) \rightarrow (a,b)



CAN: simple example

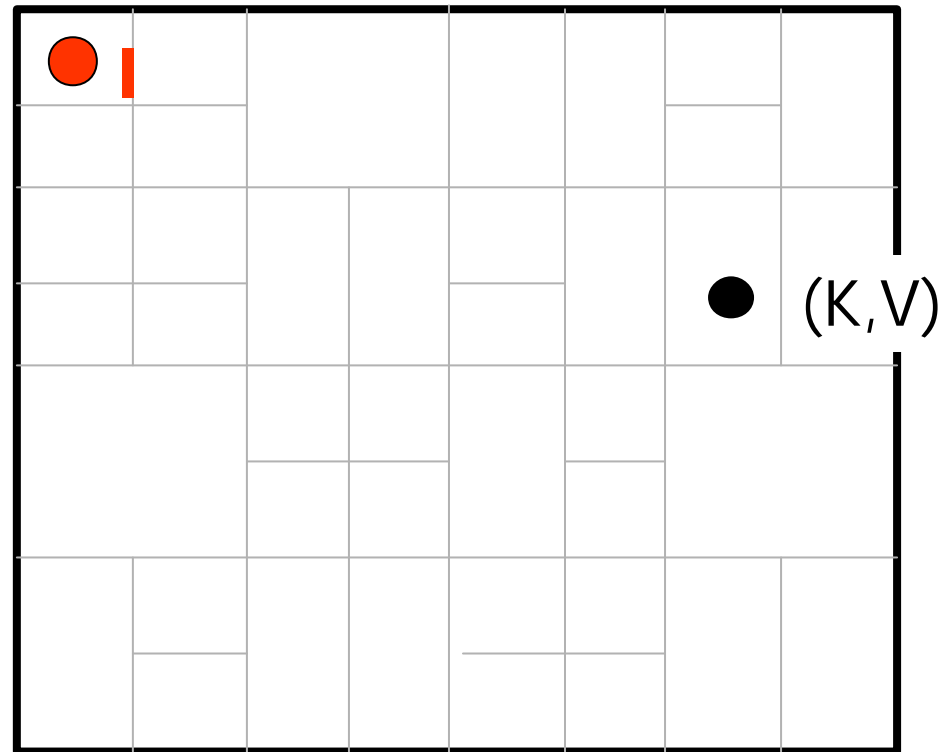
node I :: insert(K,V)

(1) $a = h_x(K)$

$b = h_y(K)$

(2) route(K,V) \rightarrow (a,b)

(3) (a,b) stores (K,V)



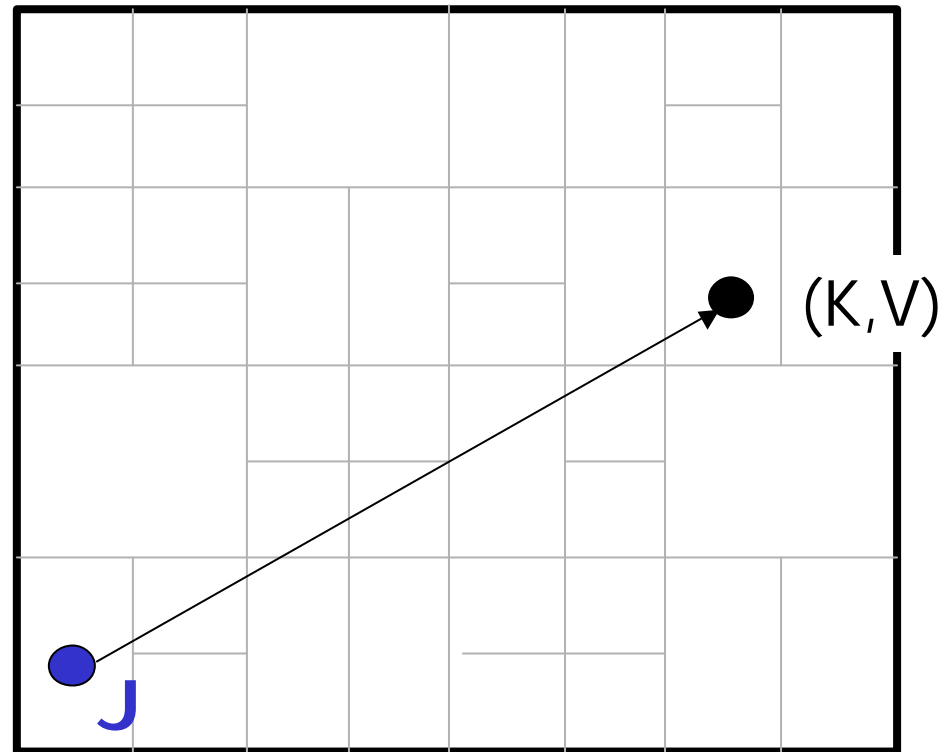
CAN: simple example

node J::retrieve(K)

(1) $a = h_x(K)$

$b = h_y(K)$

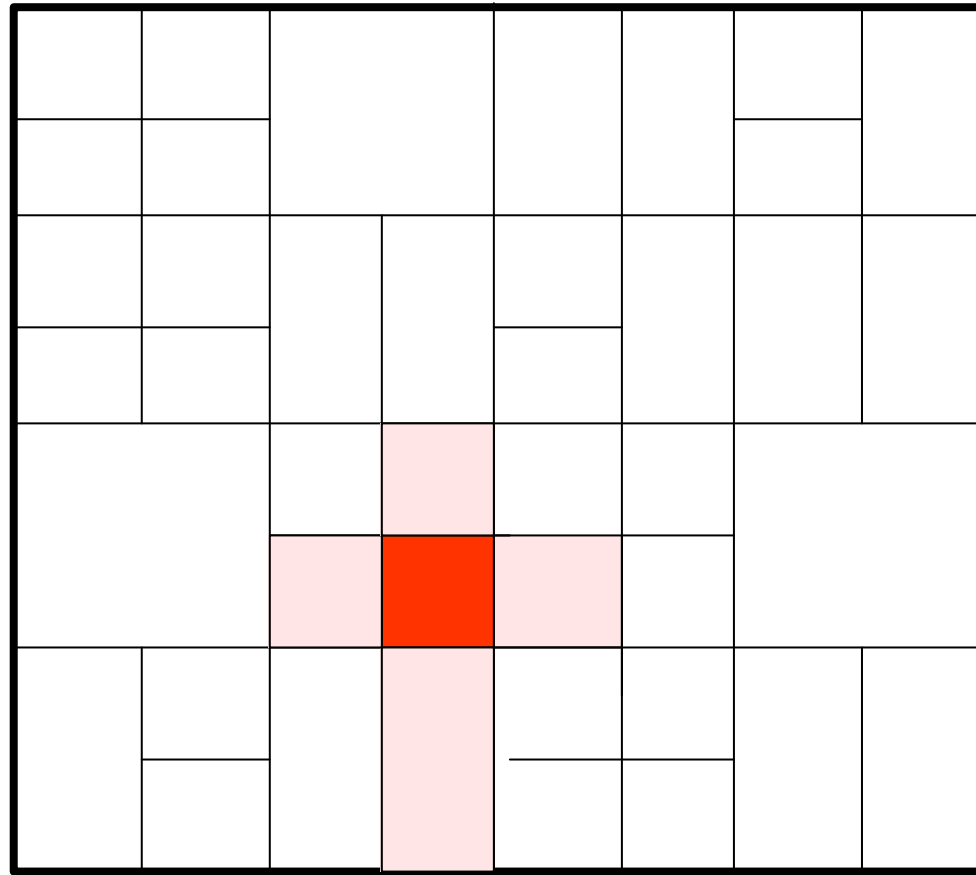
(2) route "retrieve(K)" to (a,b)



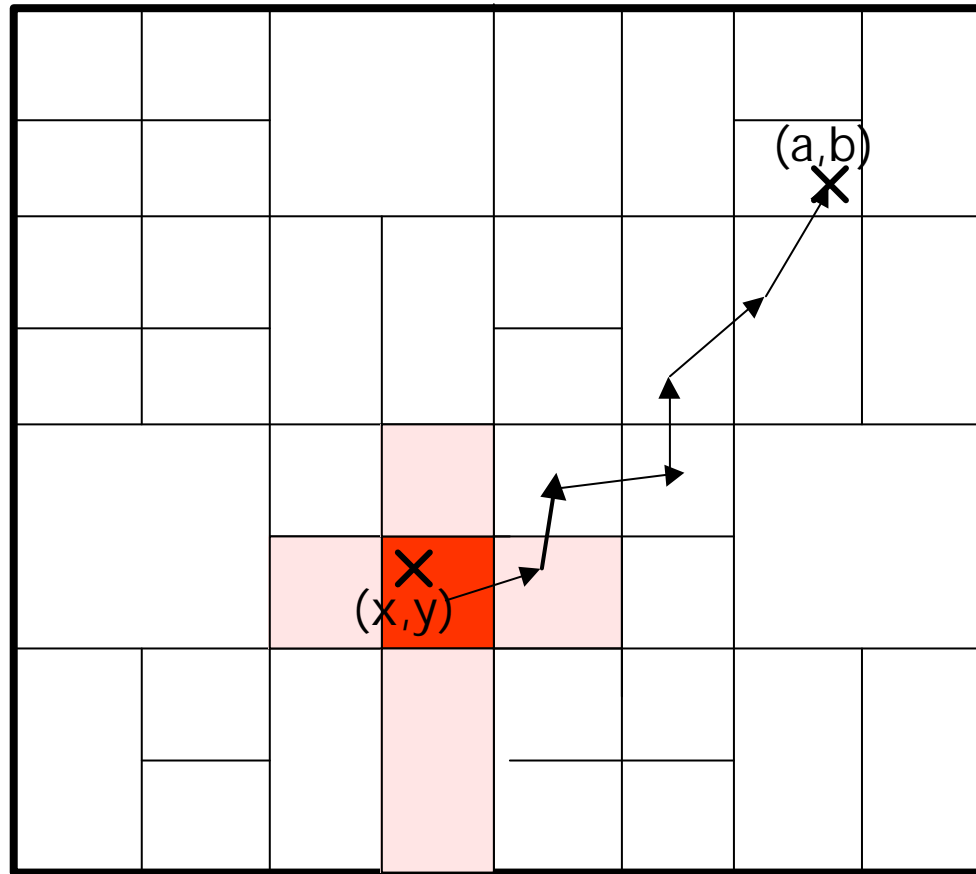
CAN

Data stored in the CAN is addressed by name (i.e. key), not location (i.e. IP address)

CAN: routing table



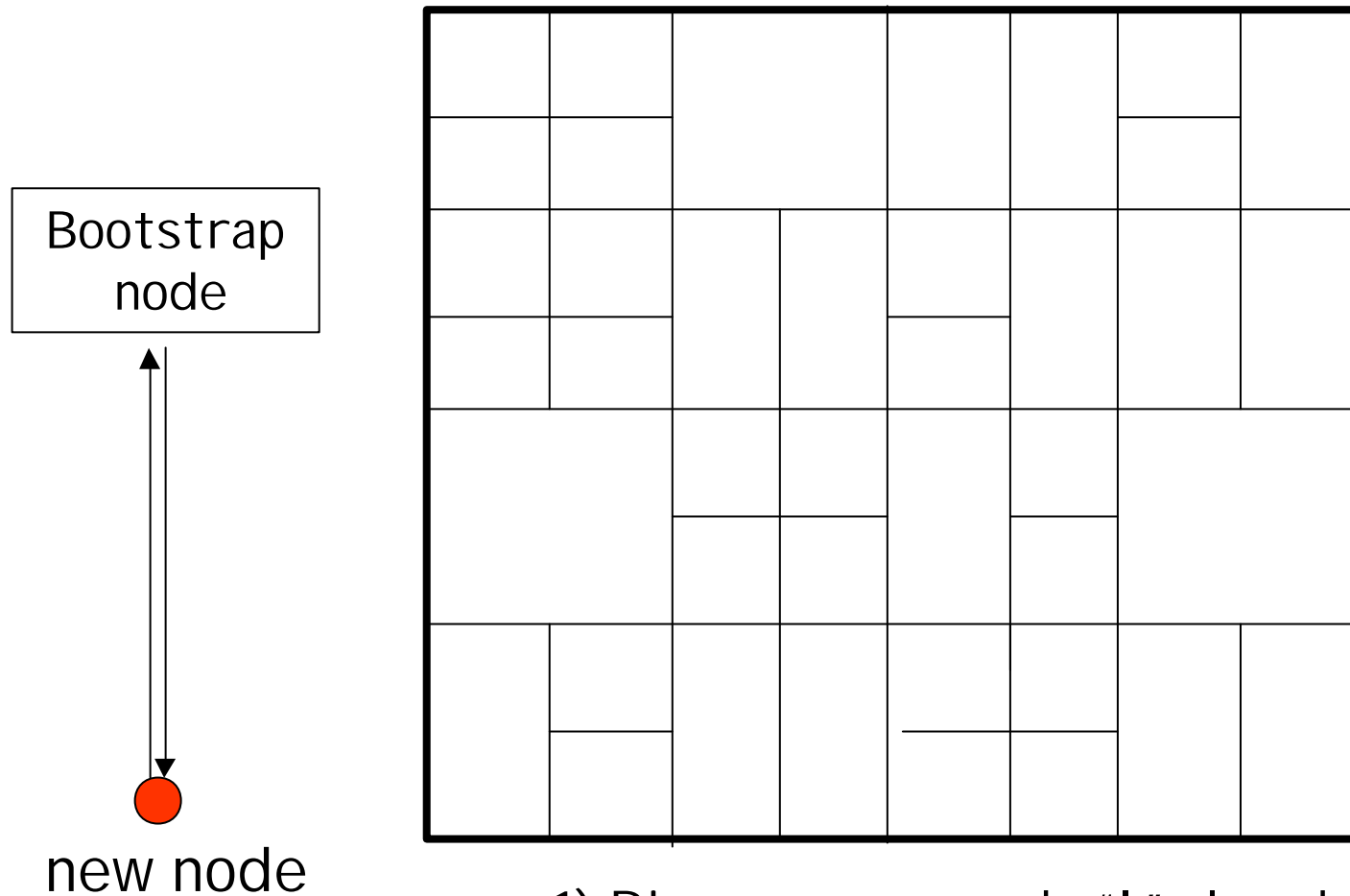
CAN: routing



CAN: routing

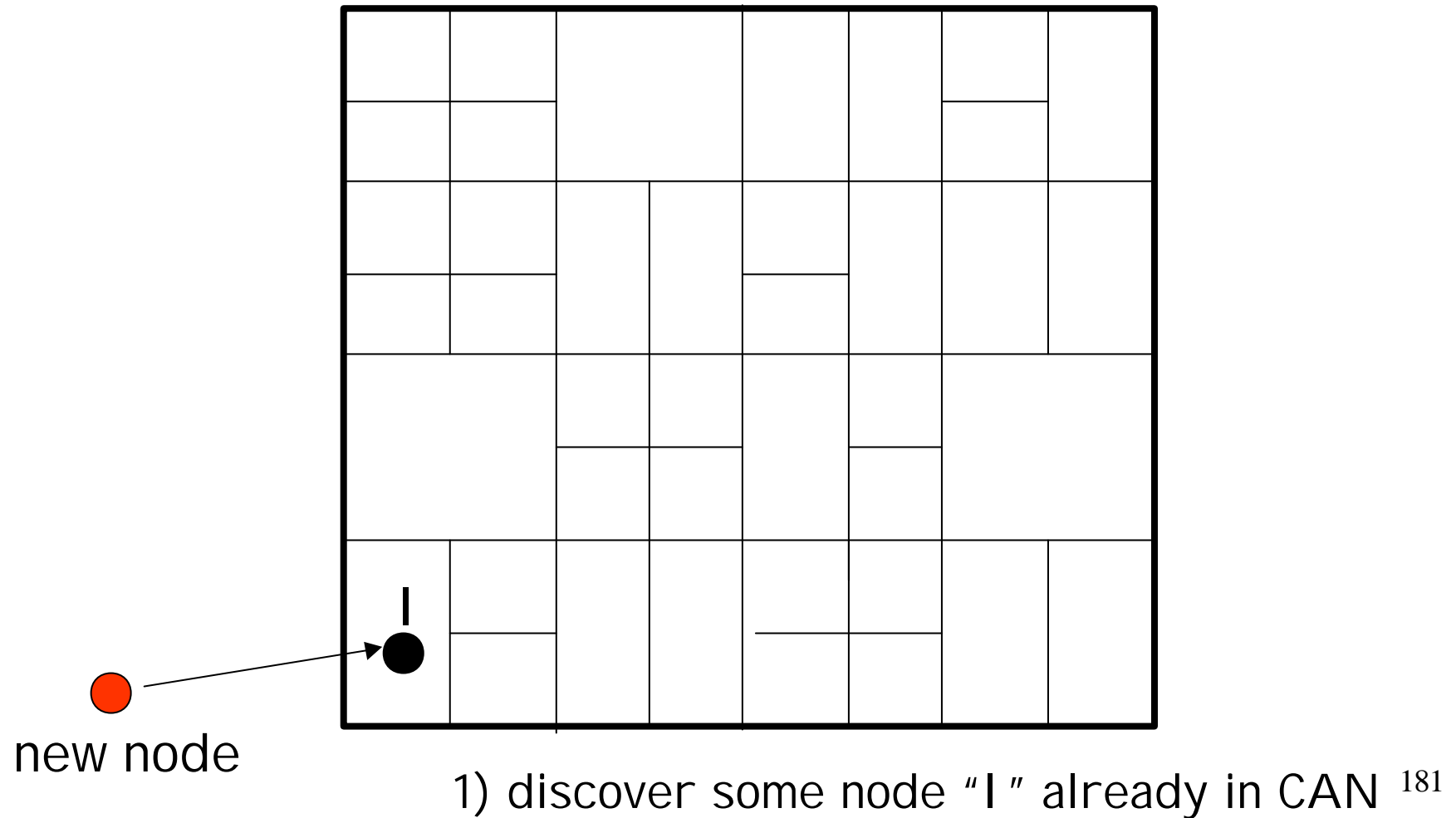
A node only maintains state for its immediate neighboring nodes

CAN: node insertion

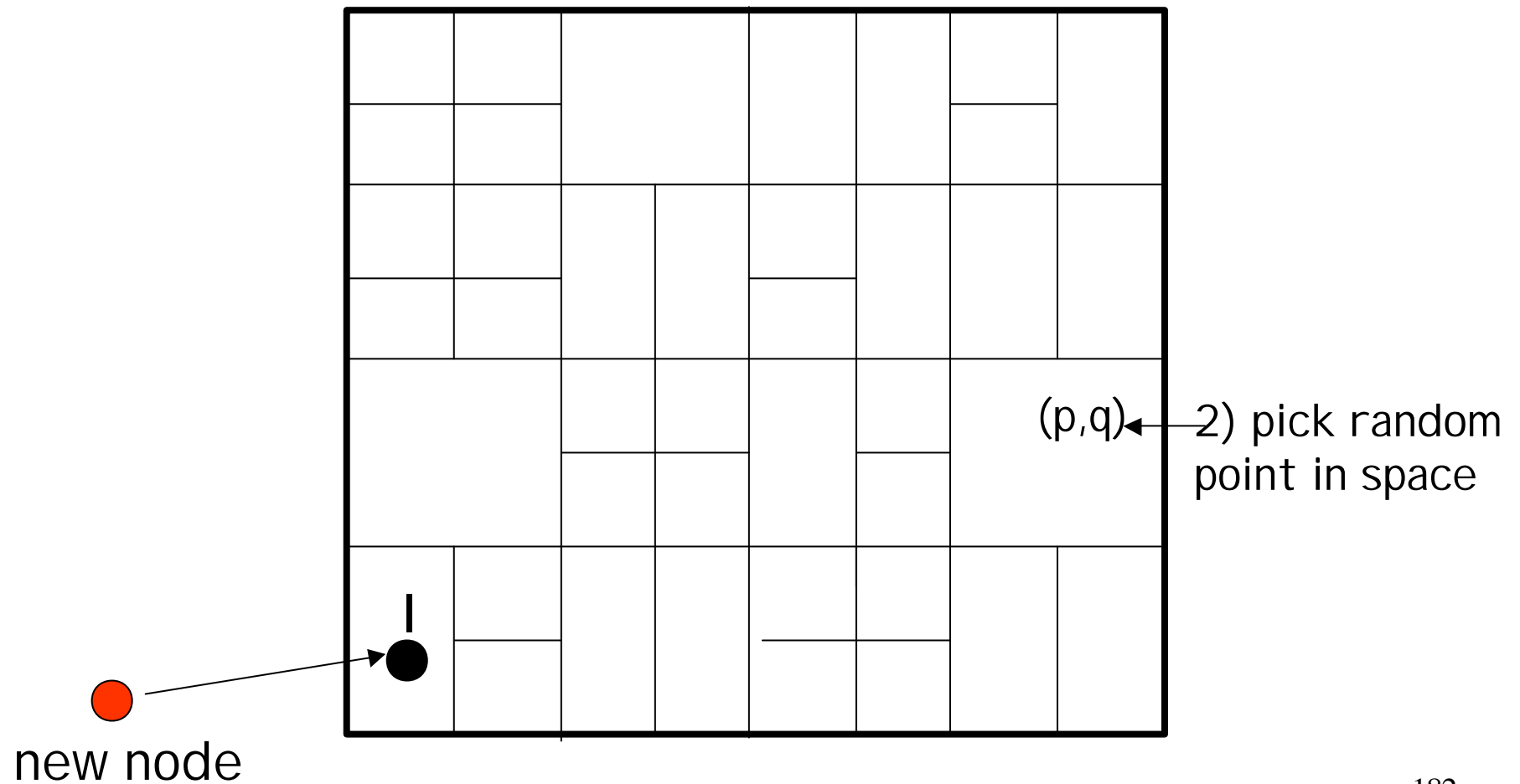


1) Discover some node "I" already in CAN

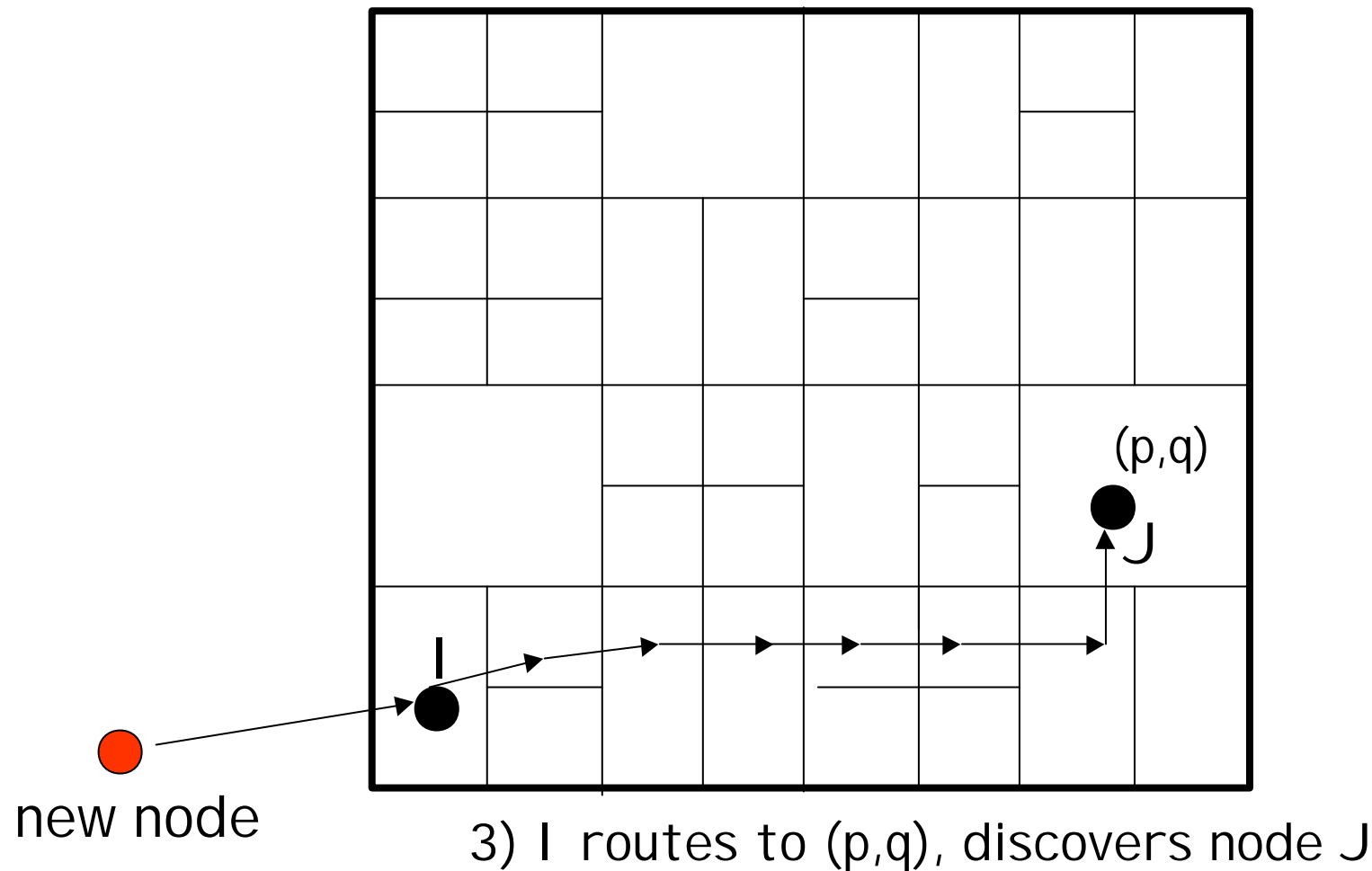
CAN: node insertion



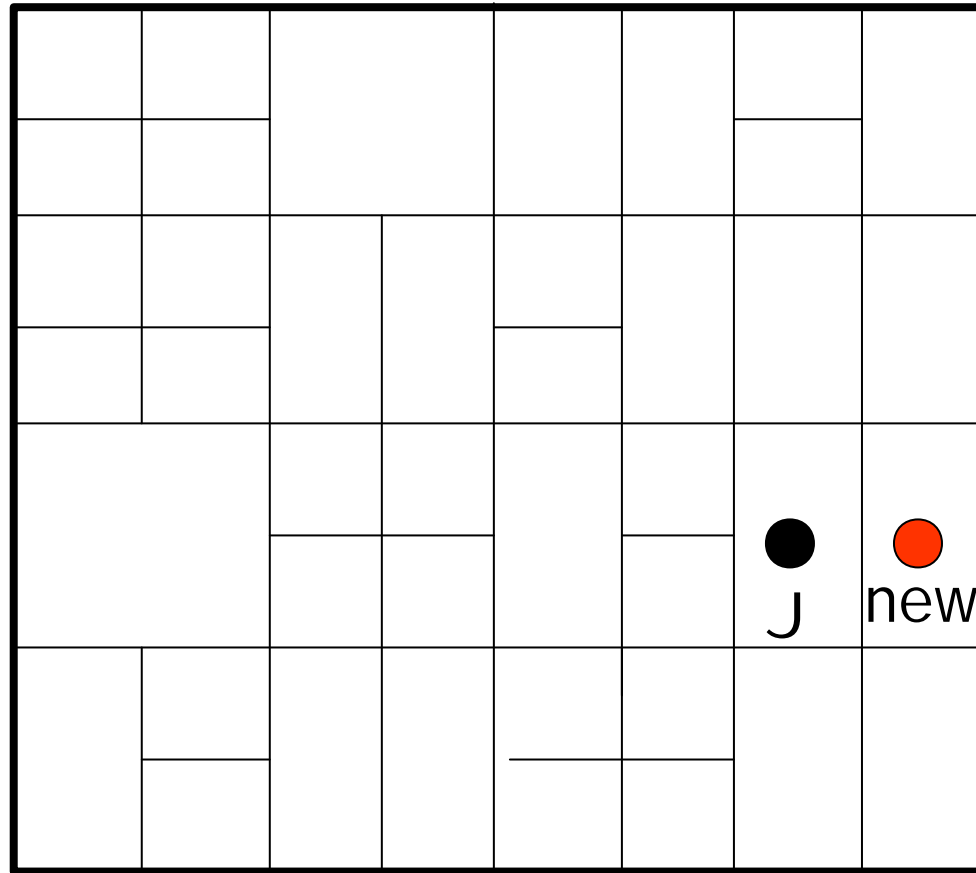
CAN: node insertion



CAN: node insertion



CAN: node insertion



4) split J's zone in half... new owns one half

CAN: node insertion

Inserting a new node affects only a single other node and its immediate neighbors

CAN: node failures

- r Need to repair the space
 - m recover database (weak point)
 - soft-state updates
 - use replication, rebuild database from replicas
 - m repair routing
 - takeover algorithm

CAN: takeover algorithm

- r Simple failures
 - m know your neighbor's neighbors
 - m when a node fails, one of its neighbors takes over its zone

- r More complex failure modes
 - m simultaneous failure of multiple adjacent nodes
 - m scoped flooding to discover neighbors
 - m hopefully, a rare event

CAN: node failures

Only the failed node's immediate neighbors are required for recovery

Design recap

- r Basic CAN
 - m completely distributed
 - m self-organizing
 - m nodes only maintain state for their immediate neighbors

- r Additional design features
 - m multiple, independent spaces (realities)
 - m background load balancing algorithm
 - m simple heuristics to improve performance

Outline

- r Introduction
- r Design
- r Evaluation
- r Strengths & Weaknesses
- r Ongoing Work

Evaluation

- r Scalability
- r Low-latency
- r Load balancing
- r Robustness

CAN: scalability

- r For a uniformly partitioned space with **n** nodes and **d** dimensions
 - m per node, number of neighbors is $2d$
 - m average routing path is $(dn^{1/d})/4$ hops
 - m simulations show that the above results hold in practice
- r Can scale the network without increasing per-node state
- r Chord/Plaxton/Tapestry/Buzz
 - m $\log(n)$ nbrs with $\log(n)$ hops

CAN: low-latency

r Problem

m latency stretch = $\frac{\text{(CAN routing delay)}}{\text{(IP routing delay)}}$

m application-level routing may lead to high stretch

r Solution

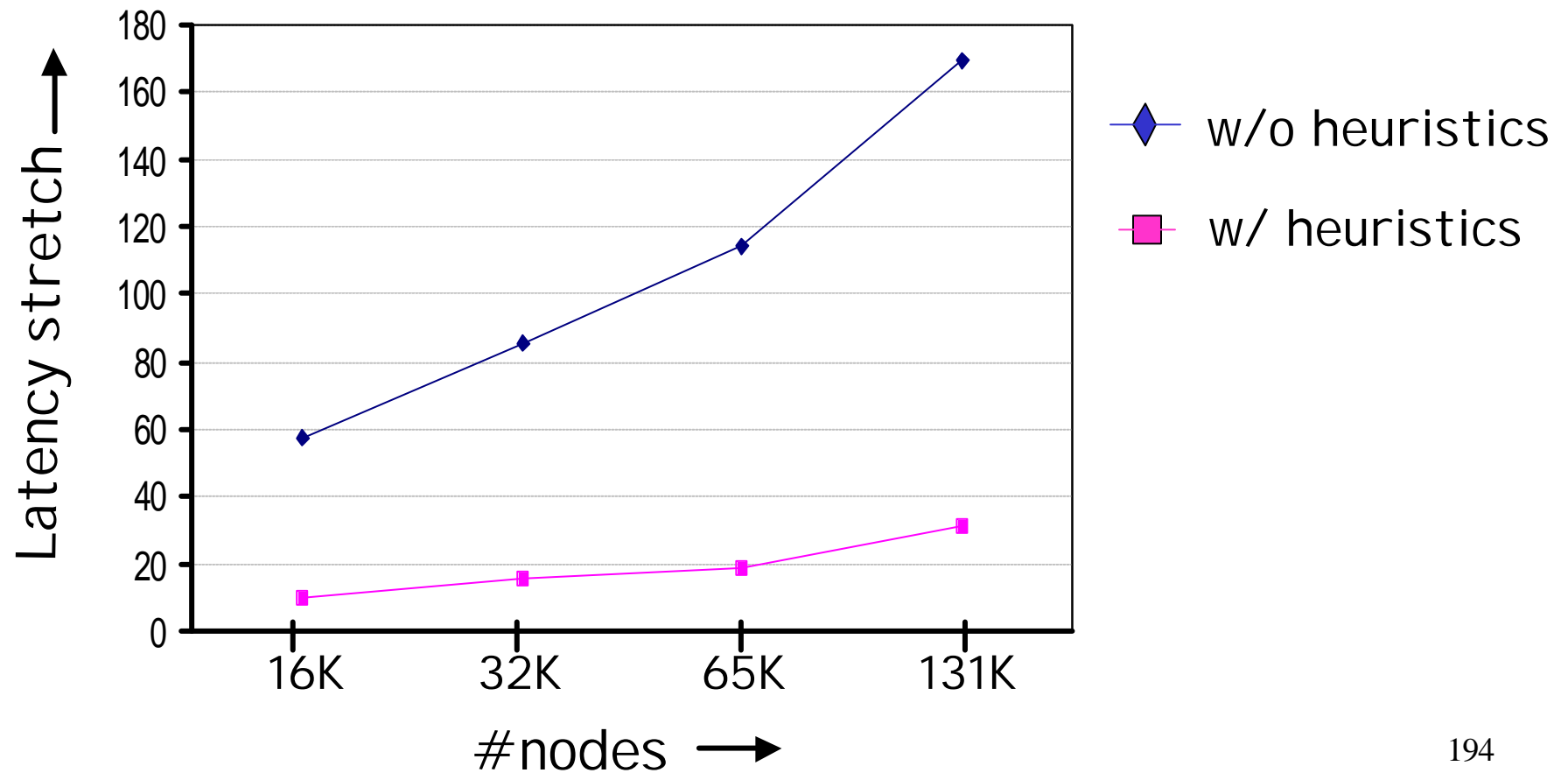
m increase dimensions, realities (reduce the path length)

m Heuristics (reduce the per-CAN-hop latency)

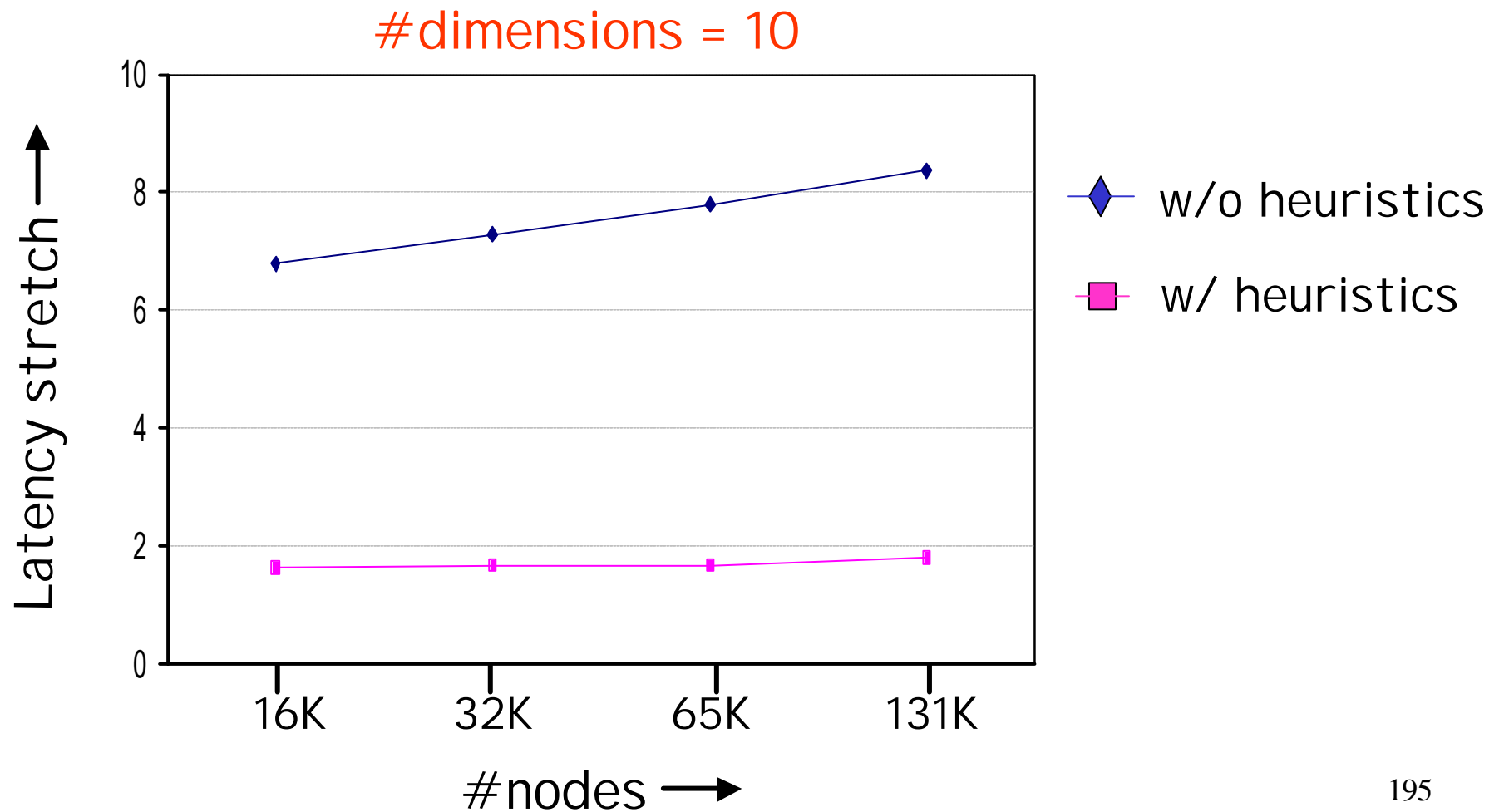
- RTT-weighted routing
- multiple nodes per zone (peer nodes)
- deterministically replicate entries

CAN: low-latency

#dimensions = 2



CAN: low-latency



CAN: load balancing

r Two pieces

m Dealing with hot-spots

- popular (key,value) pairs
- nodes cache recently requested entries
- overloaded node replicates popular entries at neighbors

m Uniform coordinate space partitioning

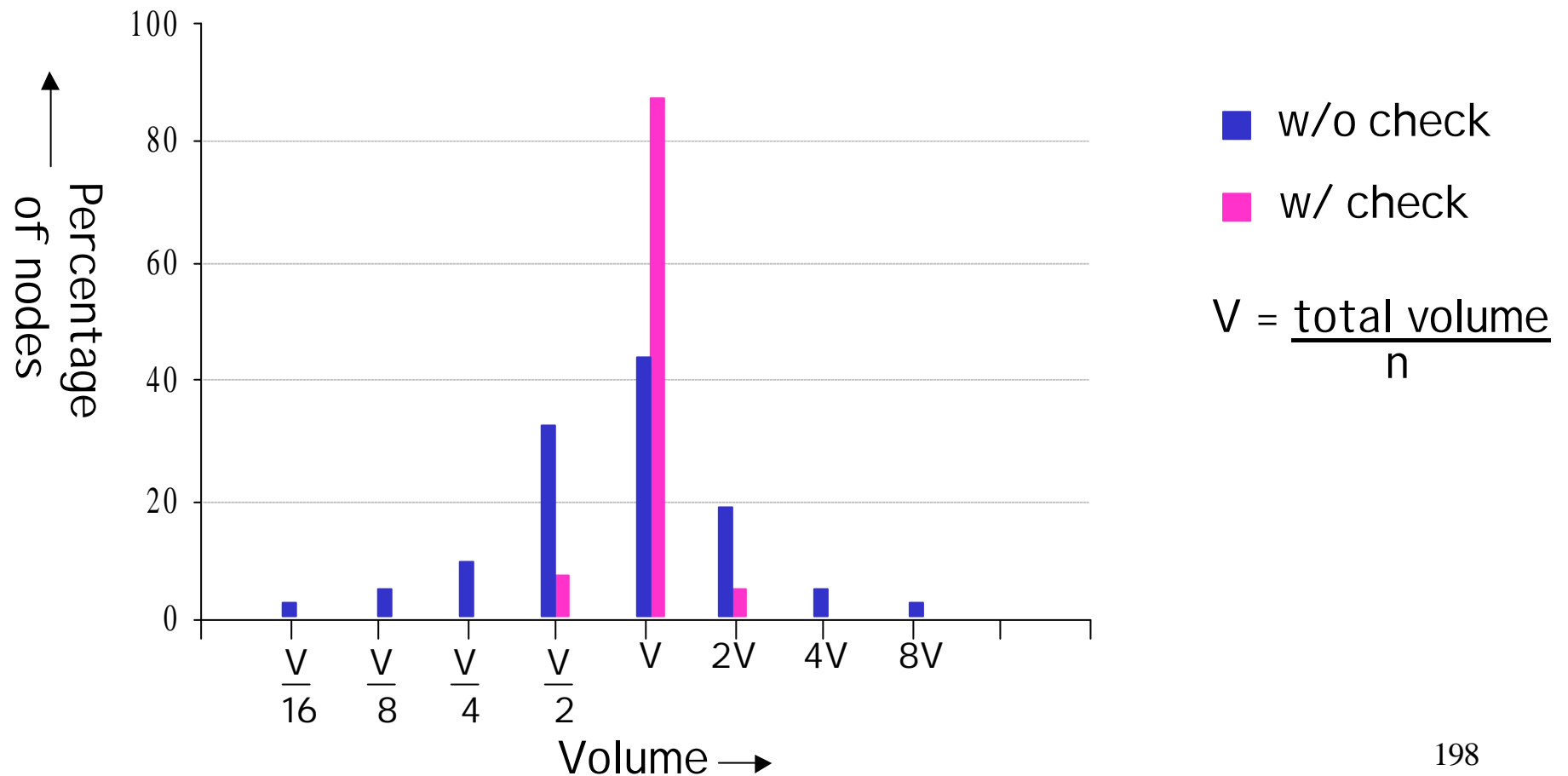
- uniformly spread (key,value) entries
- uniformly spread out routing load

Uniform Partitioning

- r Added check
 - m at join time, pick a zone
 - m check neighboring zones
 - m pick the largest zone and split that one

Uniform Partitioning

65,000 nodes, 3 dimensions



CAN: Robustness

- r Completely distributed
 - m no single point of failure (not applicable to pieces of database when node failure happens)
- r Not exploring database recovery (in case there are multiple copies of database)
- r Resilience of routing
 - m can route around trouble

Outline

- r Introduction
- r Design
- r Evaluation
- r Strengths & Weaknesses
- r Ongoing Work

Strengths

- r More resilient than flooding broadcast networks
- r Efficient at locating information
- r Fault tolerant routing
- r Node & Data High Availability (w/ improvement)
- r Manageable routing table size & network traffic

Weaknesses

- r Impossible to perform a fuzzy search
- r Susceptible to malicious activity
- r Maintain coherence of all the indexed data
(Network overhead, Efficient distribution)
- r Still relatively higher routing latency
- r Poor performance w/o improvement

Suggestions

- r Catalog and Meta indexes to perform search function
- r Extension to handle mutable content efficiently for web-hosting
- r Security mechanism to defense against attacks

Outline

- r Introduction
- r Design
- r Evaluation
- r Strengths & Weaknesses
- r Ongoing Work

Ongoing Work

- r Topologically-sensitive CAN construction
 - m distributed binning

Distributed Binning

r Goal

- m bin nodes such that co-located nodes land in same bin

r Idea

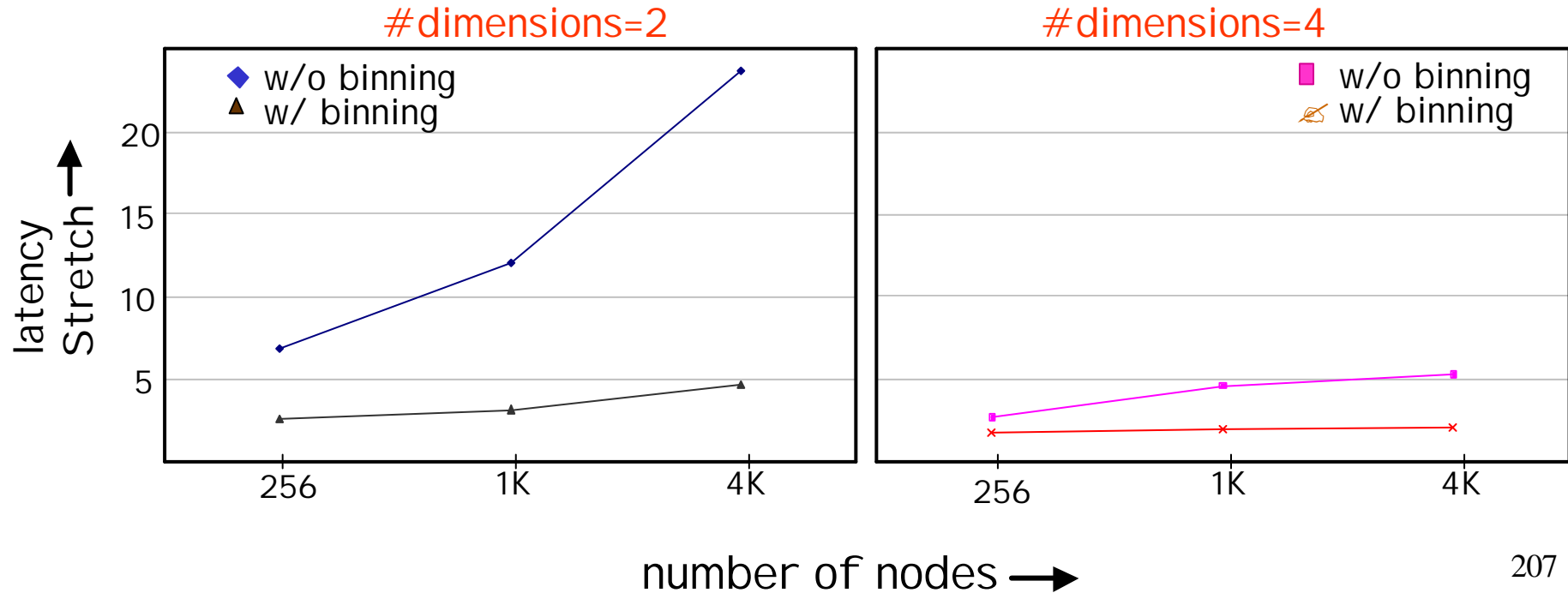
- m well known set of landmark machines
- m each CAN node, measures its RTT to each landmark
- m orders the landmarks in order of increasing RTT

r CAN construction

- m place nodes from the same bin close together on the CAN

Distributed Binning

- m 4 Landmarks (placed at 5 hops away from each other)
- m naïve partitioning



Ongoing Work (cont'd)

- r Topologically-sensitive CAN construction
 - m distributed binning

- r CAN Security (Petros Maniatis - Stanford)
 - m spectrum of attacks
 - m appropriate counter-measures

Ongoing Work (cont'd)

- r CAN Usage

- m Application-level Multicast (NGC 2001)

- m Grass-Roots Content Distribution

- m Distributed Databases using CANs

- (J.Hellerstein, S.Ratnasamy, S.Shenker, I.Stoica, S.Zhuang)

Summary

- r CAN

- m an Internet-scale hash table
 - m potential building block in Internet applications

- r Scalability

- m $O(d)$ per-node state

- r Low-latency routing

- m simple heuristics help a lot

- r Robust

- m decentralized, can route around trouble



9.Sun's Project JXTA

Technical Overview

Project JXTA Implementation Outline

- r Introduction - what is JXTA
- r Goal - what JXTA wants to be
- r Technology - what JXTA relies upon
- r Structure - how JXTA is built
- r Protocols - what protocols JXTA has
- r Security - whether JXTA is secure
- r Applications - what JXTA can be used for
- r Collaboration - how JXTA grows

Project JXTAduction

- r "JXTA" - pronounced as "**juxta**" as in "**juxtaposition**"
- r started by Sun's Chief Scientist **Bill Joy**
- r an effort to create a **common platform** for building distributed services and applications
- r Napster, Gnutella, and Freenet provide users with **limited ability** to share resources and are unable to share data with other, similar applications

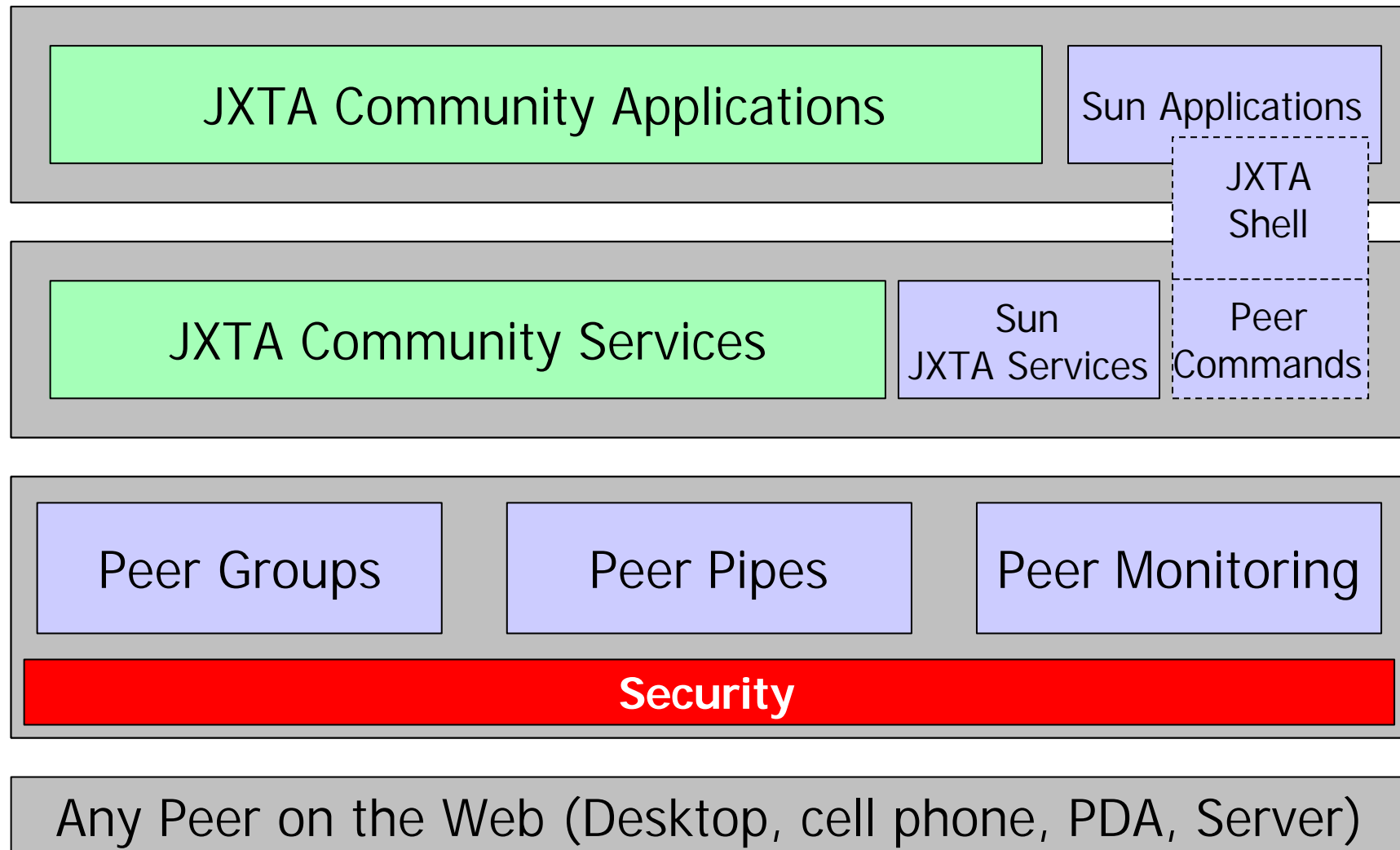
Project JXTA Purpose

- r enable a wide range of distributed computing applications by developing a common set of general purpose P2P protocols
- r achieve platform independence - any language, any OS, any hardware
- r overcome the limitations found in many today's P2P applications
- r enable new applications to run on any device that has a digital heartbeat (desktop computers, servers, PDAs, cell phones, and other connected devices)

Project JXTAology

- r JXTA technology is based on XML, Java technology, and key concepts of UNIX operating system
- r Transmitted information is packaged as messages. Messages define an XML envelop to transfer any kind of data.
- r The use of Java language is not required - JXTA protocols can be implemented in C, C++, Perl, or any other programming language

Project JXTA:ture



Project JXTA Multi-Layered Structure

r **JXTA Core**

- **Peer Groups** - mechanisms to/for create and delete, join, advertise, discover, communication, security, content sharing
- **Peer Pipes** - transfer of data, content, and code in a protocol-independent manner
- **Peer Monitoring** - including access control, priority setting, traffic metering and bandwidth balancing

r **JXTA Services**

- expand upon the capabilities of the core and facilitate application development
- mechanisms for searching, sharing, indexing, and caching code and content to enable cross-application bridging and translation of files

r **JXTA Shell** - much like UNIX OS

- facilitate access to core-level functions through a command line

r **JXTA Applications** - built using peer services as well as the core layer

Project JXTAcols

- r JXTA is a set of six protocols
- r Peer Discovery Protocol - find peers, groups, advertisements
- r Peer Resolver Protocol - send/receive search queries
- r Peer Information Protocol - learn peers' status/properties
- r Peer Membership Protocol - sign in, sign out, authentication
- r Pipe Binding Protocol - pipe advertisement to pipe endpoint
- r Endpoint Routing Protocol - available routes to destination

Project JXTAity

- r Confidentiality, integrity, availability - authentication, access control, encryption, secure communication, etc.
- r Developing more concrete and precise security architecture is an ongoing project
- r JXTA does not mandate certain security policies, encryption algorithms or particular implementations
- r JXTA 1.0 provides Security Primitives:
 - crypto library (MD5, RC4, RSA, etc.)
 - Pluggable Authentication Module (PAM)
 - password-based login
 - transport security mechanism modeled after SSL/TLS

Project JXTA Potential Applications

- r Search the entire web and all its connected devices (not just servers) for needed information
- r Save files and information to distributed locations on the network
- r Connect game systems so that multiple people in multiple locations
- r Participate in auctions among selected groups of individuals
- r Collaborate on projects from anywhere using any connected device
- r Share compute services, such as processor cycles or storage systems, regardless of where the systems or the users are located

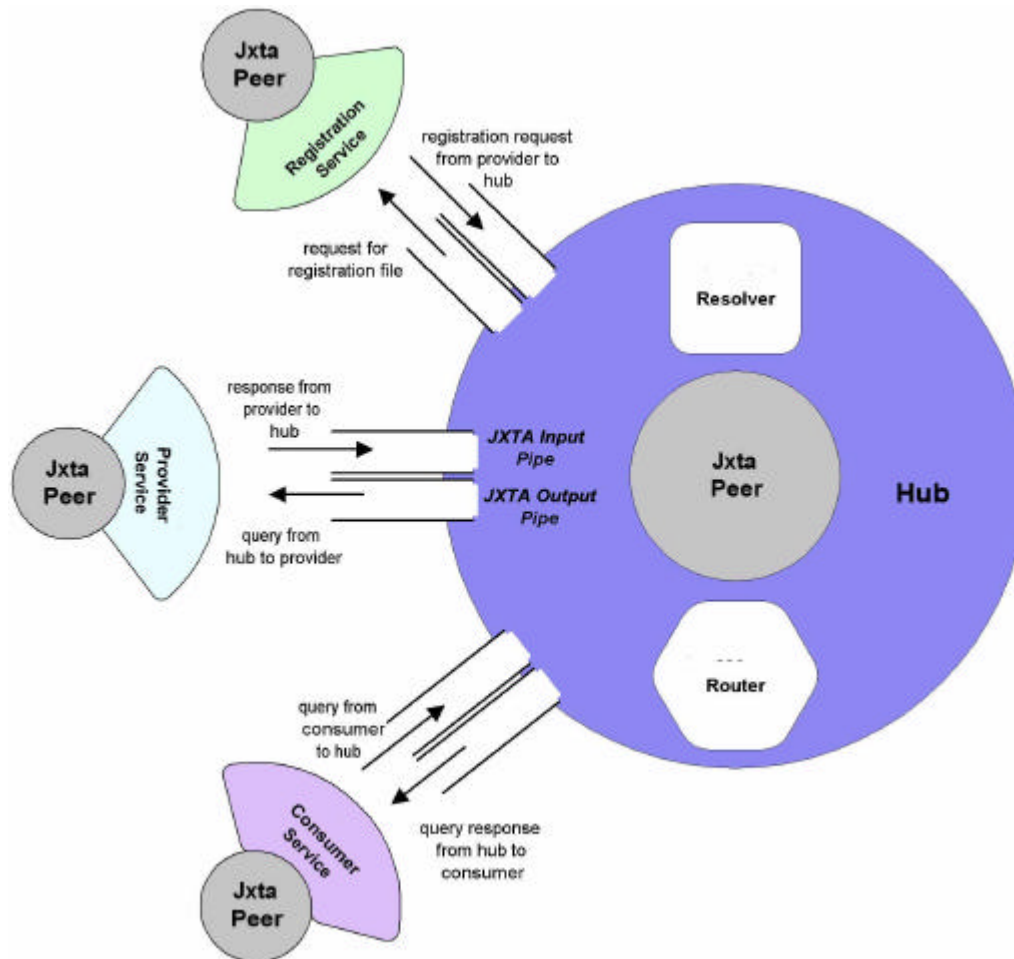
Project JXTA - A Search Overview

- r Started in June 2000 by Infrasearch as an idea to distribute queries to network peers best capable of answering them
- r Now it is the default searching methodology for the JXTA framework in the form of JXTA Search
- r Communication via an XML protocol called Query Routing Protocol (QRP)
- r Network components: Providers, Consumers, Hubs
- r Capable of providing both wide and deep search; deep search shows the most benefits
- r Design goals: Simplicity, Structure, Extensibility, Scalability

Project JXTA XTA Search Benefits

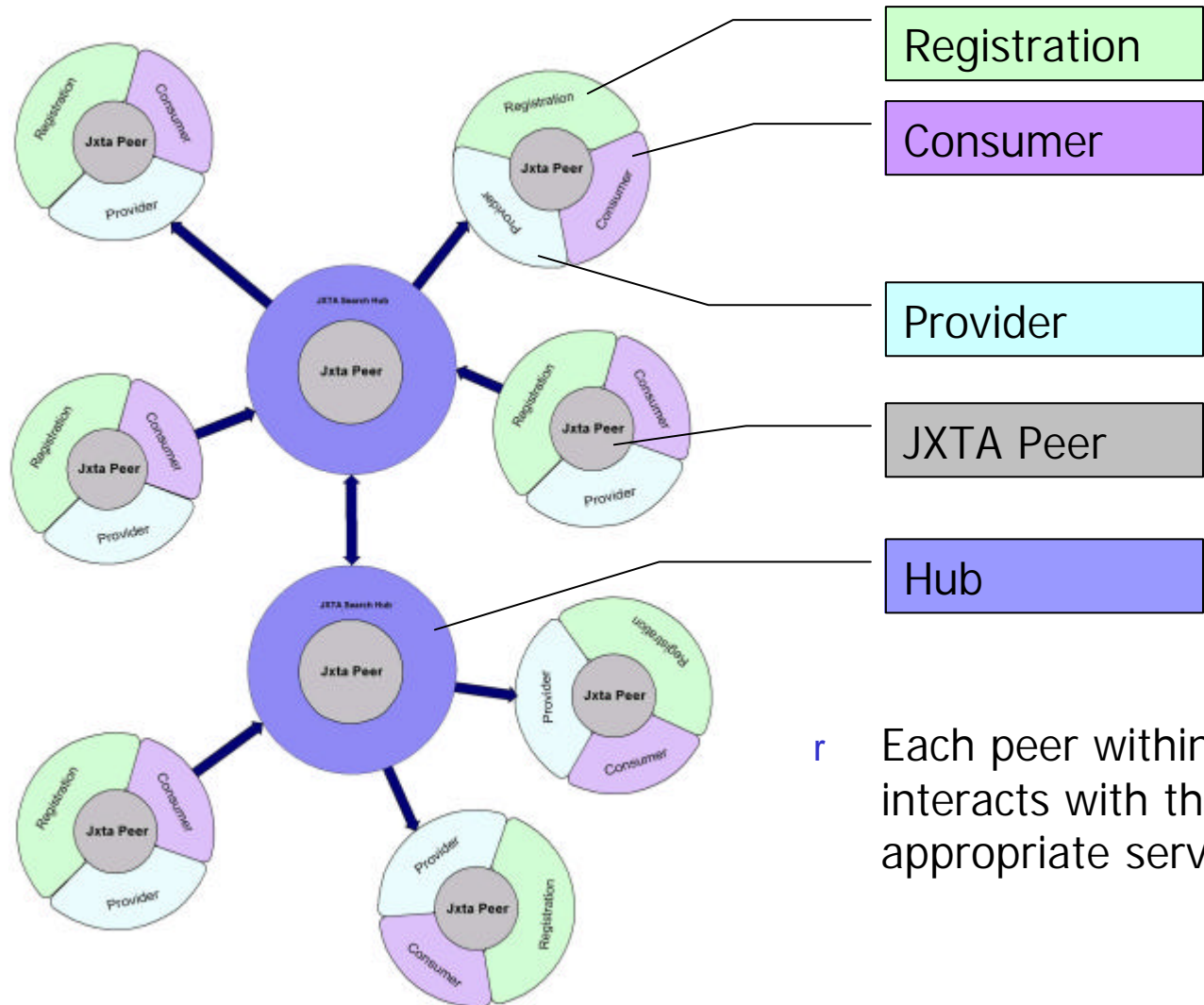
- r Speed of update - especially noticeable in deep search, where large data in databases are accessed directly without a need to create a central index.
- r Access - in crawling based approach many companies are resilient to grant access to web crawlers. In distributed approach the companies can serve the data as they feel appropriate.
- r Efficiency - no need to create a centrally placed and maintained index for the whole web.

Project JXTA JXTA Search Architecture



- r Each JXTA peer can run instances of Provider, Consumer, and Registration services on top of its JXTA core.
- r Each peer interacts with the JXTA Search Hub Services, which is also running on top of the JXTA core.

Project JXTA Search Architecture



Project **JXTA**collaboration

- r Currently over 25 companies are participating in developing JXTA projects.
- r Core (7 projects)
 - security, juxta-c, juxtaperl, pocketjxta
- r Services (20 projects)
 - search, juxtaspaces, p2p-email, juxta-grid, payment, monitoring
- r Applications (12 projects)
 - shell, jnushare, dfwbase, brando
- r Other projects (5) - demos, tutorials, etc.

r Future Work

- C/C++ implementation
- KVM based implementation (PDAs, cell phones)
- Naming and binding services
- Security services (authentication, access control)
- Solutions for firewalls and NAT gateways

r Is this the right structure?

r Do JXTA protocols dictate too much or too little?

10. Application Layer Anycasting: A Server Selection Architecture and Use in Replicated Web Service

Ellen W. Zegura

Mostafa H. Amamr

Zongming Fei

Networking and Telecommunications Group
Georgia Tech, Atlanta, GA

Samrat Battacharjee

Department of Computer Science
University of Maryland, College Park, MD

Agenda

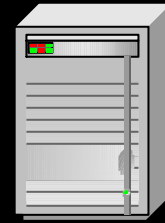
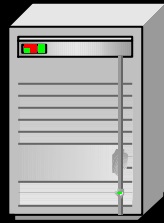
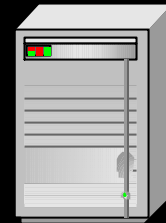
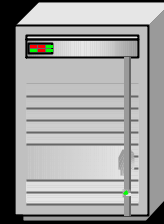
- r Problem Statement
- r The Anycasting Communication Paradigm
- r Some Related Work
- r Application Layer Anycasting
- r Experimental Results
- r Conclusions

Problem Statement

- r Efficient service provision in wide area networks
- r Replicated services
- r Applications want access to the best server
- r Best may depend on time, performance, policy

Server Replication

- r Standard technique to improve scalability of a service
- r Issues in server replication
 - m Location of servers
 - m Consistency across servers
 - m Server selection
- r Server selection problem
 - m How does a client determine which of the replicated servers to access ?



Server Selection

r Alternatives

- m Designated (e.g. nearest) server
- m Round robin assignment (e.g. DNS rotator)
- m Explicit list with user selection
- m Selection architecture (e.g. Cisco DistributedDirector)

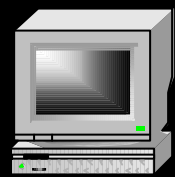
r Application-Layer Anycasting:

- m Client requests connection to anycast group
- m Anycast group consists of replicated (equivalent) servers
- m System connects client to any good server

Anycasting Communication Paradigm

- r Anycast identifier specifies a group of equivalent hosts
- r Requests are sent to best host in the group

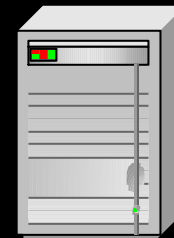
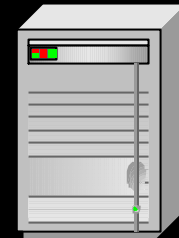
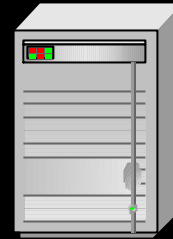
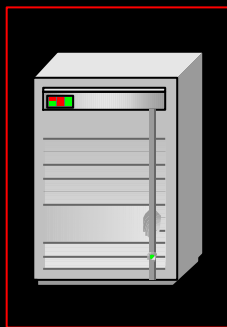
Anycast mechanism
resolves to one of
possible many



Request



Reply



Existing Anycast Solutions and Limitations

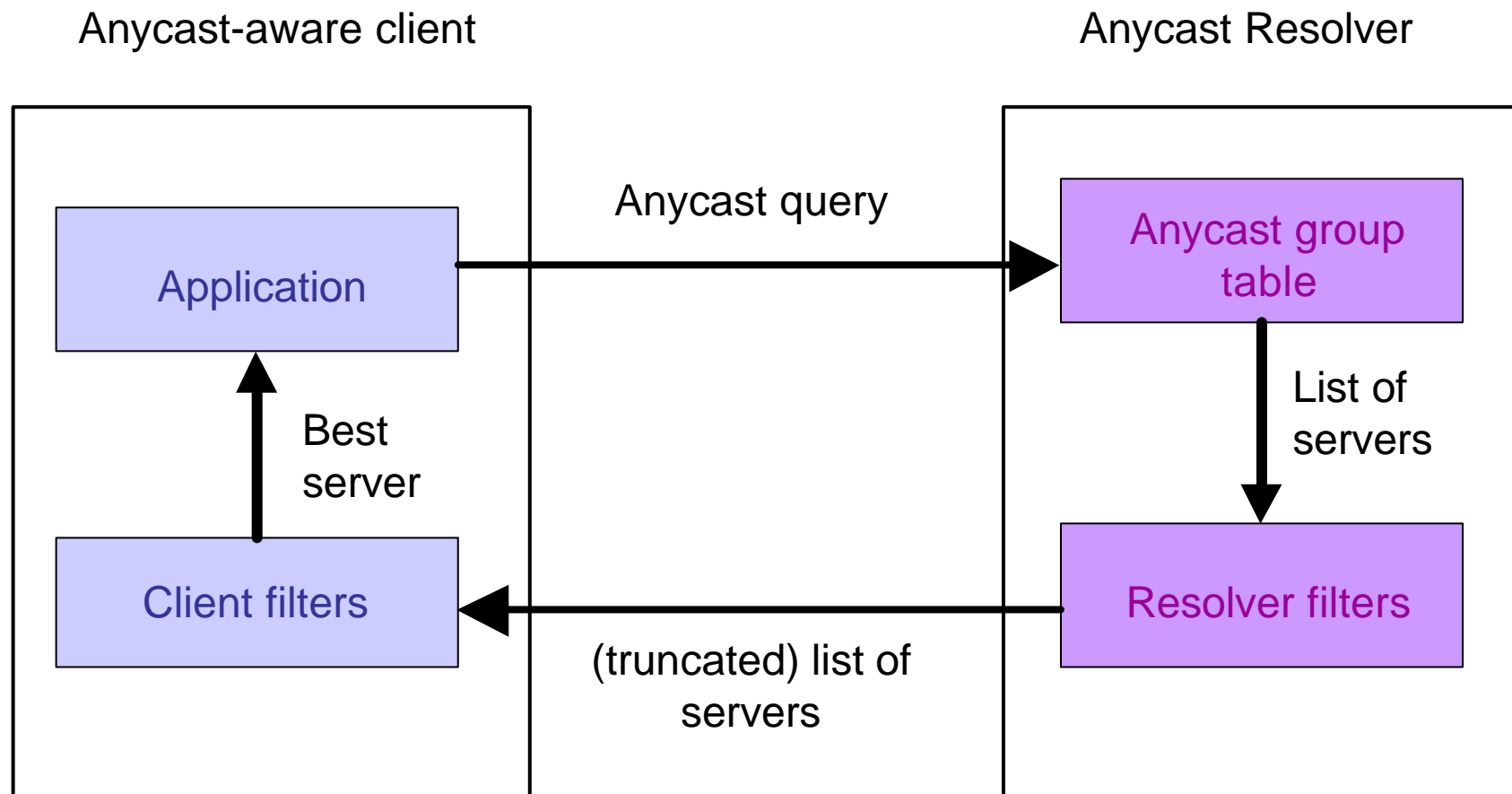
r Existing Solutions:

- m RFC 1546 Host Anycasting Service
 - Definition of anycasting communication paradigm
 - Implementation suggestions for the network layer
- m IETF Server Location Protocol – Still existing ?
- m AKAMAI and other commercial cache/CDNs
- m Cisco DistributedDirector

r Limitations

- m Global router support
- m Per diagram destination selection
- m Limited set of metrics
- m No option for user input in server selection
- m Allocation of IP address space for anycast address

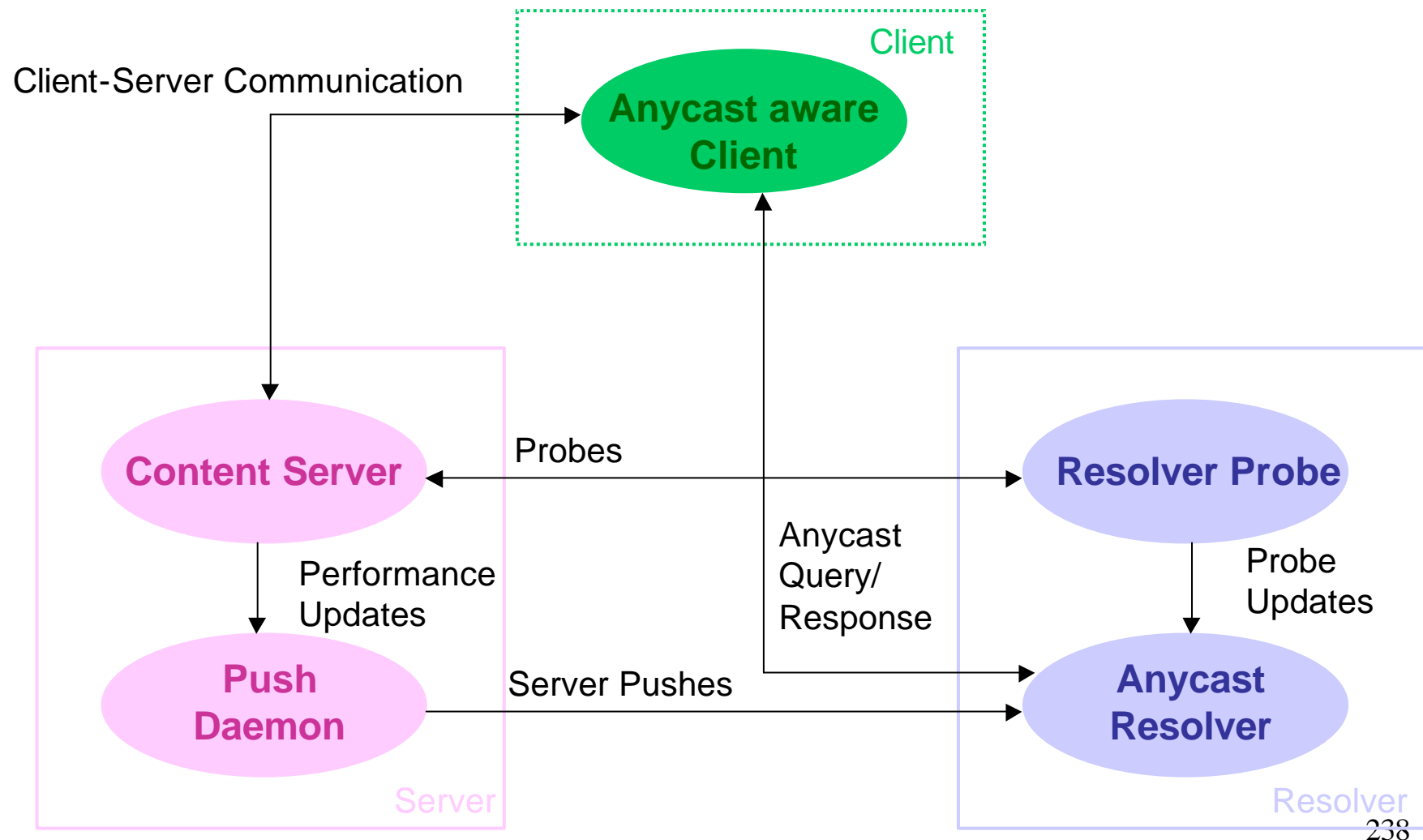
Application-Layer Anycasting



Filters

- r Content-independent filters
 - m E.g. Round-robin
- r Metric-based filters
 - m E.g. Minimum response time
- r Policy-based filters
 - m E.g. Minimum cost
- r Filter Specification – Metric Qualified ADN:
 - m <Metric_Service>%<Domain Name>
 - m ServerLoad.Daily_News%cc.gatech.edu

Application-layer Anycasting Architecture



Anycast Groups

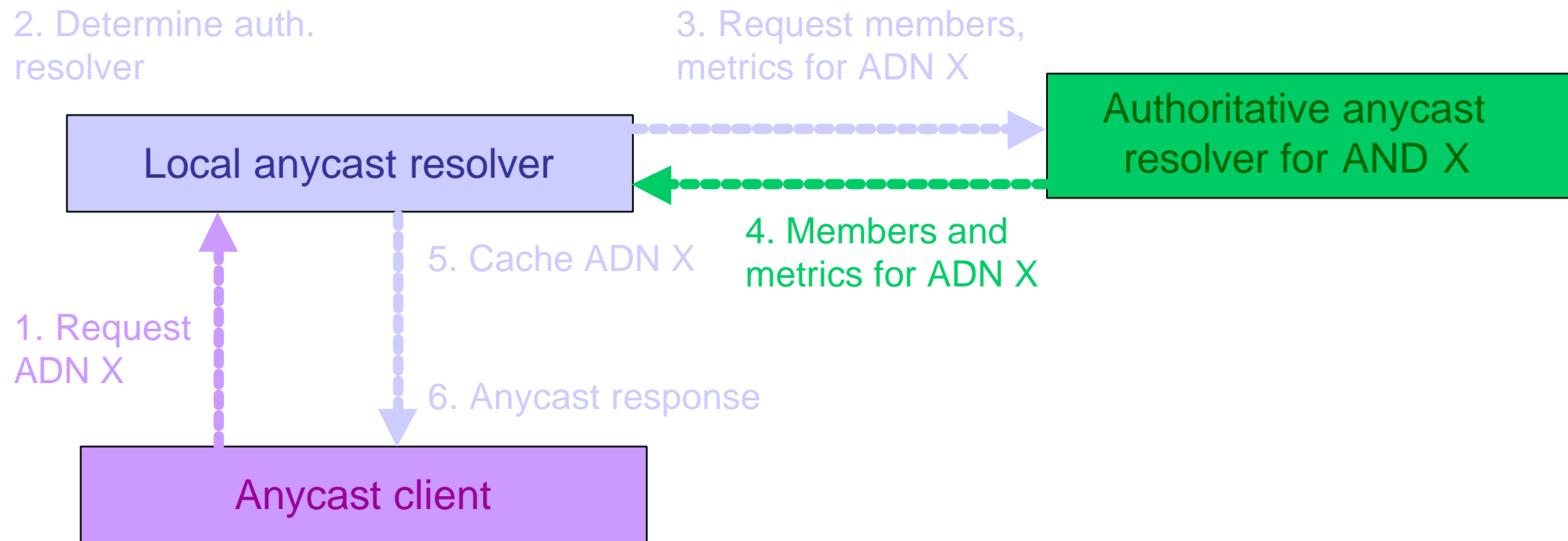
- r Anycast groups consist of collection of IP addresses, domain names or aliases
- r Group members provide equivalent service, e.g., mirrored FTP servers or web search engines
- r Anycast groups identified by Anycast Domain Names
- r Group membership an orthogonal issue architecture aliases

Anycast Domain Names

r Structure: <Service>%<Domain Name>

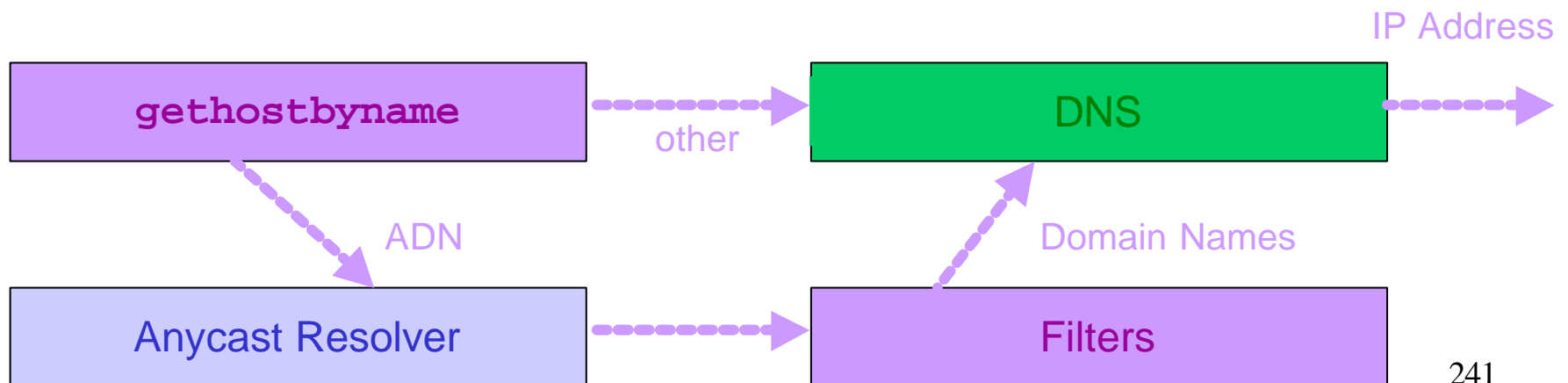
r Example:

m Daily-News%cc.gatech.edu



Implementation

- r Implementation using Metric Qualified ADNs
- r Intercept calls to **gethostbyname**
- r Transparent access to anycasting without modifying existing applications



Response Time Determination for Web Servers

- r Response time:
 - m Measured from time client issues request until receives last byte of file of network
 - m Round trip path delays + server processing delays
- r Overview of technique:
 - m Resolver probes for path-dependent response time
 - m Server measures and pushes path-independent processing time
 - m Lighter-weight push more frequent than heavier-weight probe
 - m Probe result used to calibrate pushed value

Performance Metric Determination

- r Metric collection techniques
 - m Server push algorithm
 - m Agent probe mechanism
 - m Hybrid push/probe technique

Server Push Process

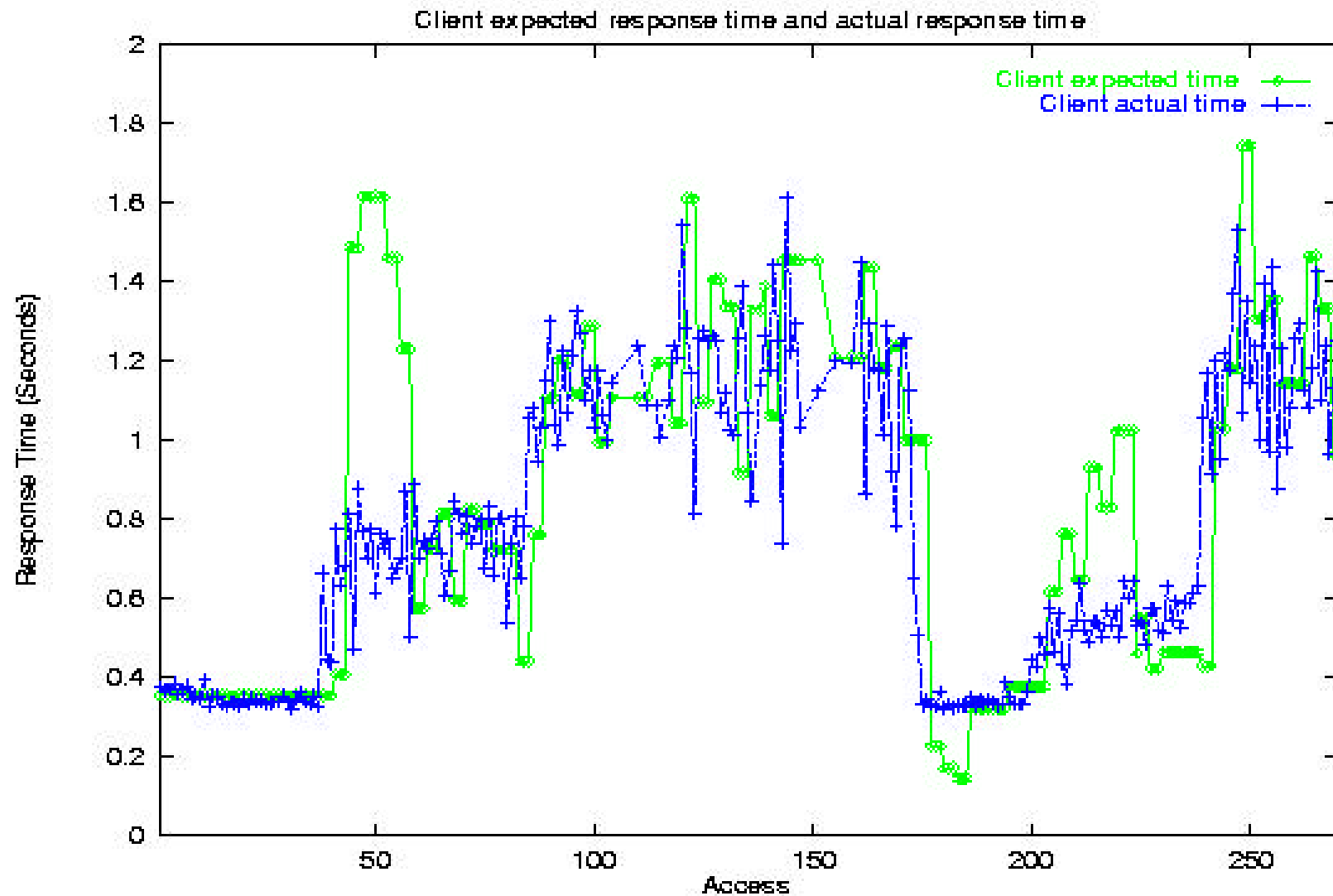
- r Typical server response cycle:
 - assign process to handle query
 - parse query
 - locate requested file
 - repeat until file is written
 - read from file
 - write to network

- r Process:
 - m Measure and smooth time until first read (TUFR)
 - m Push if significant change

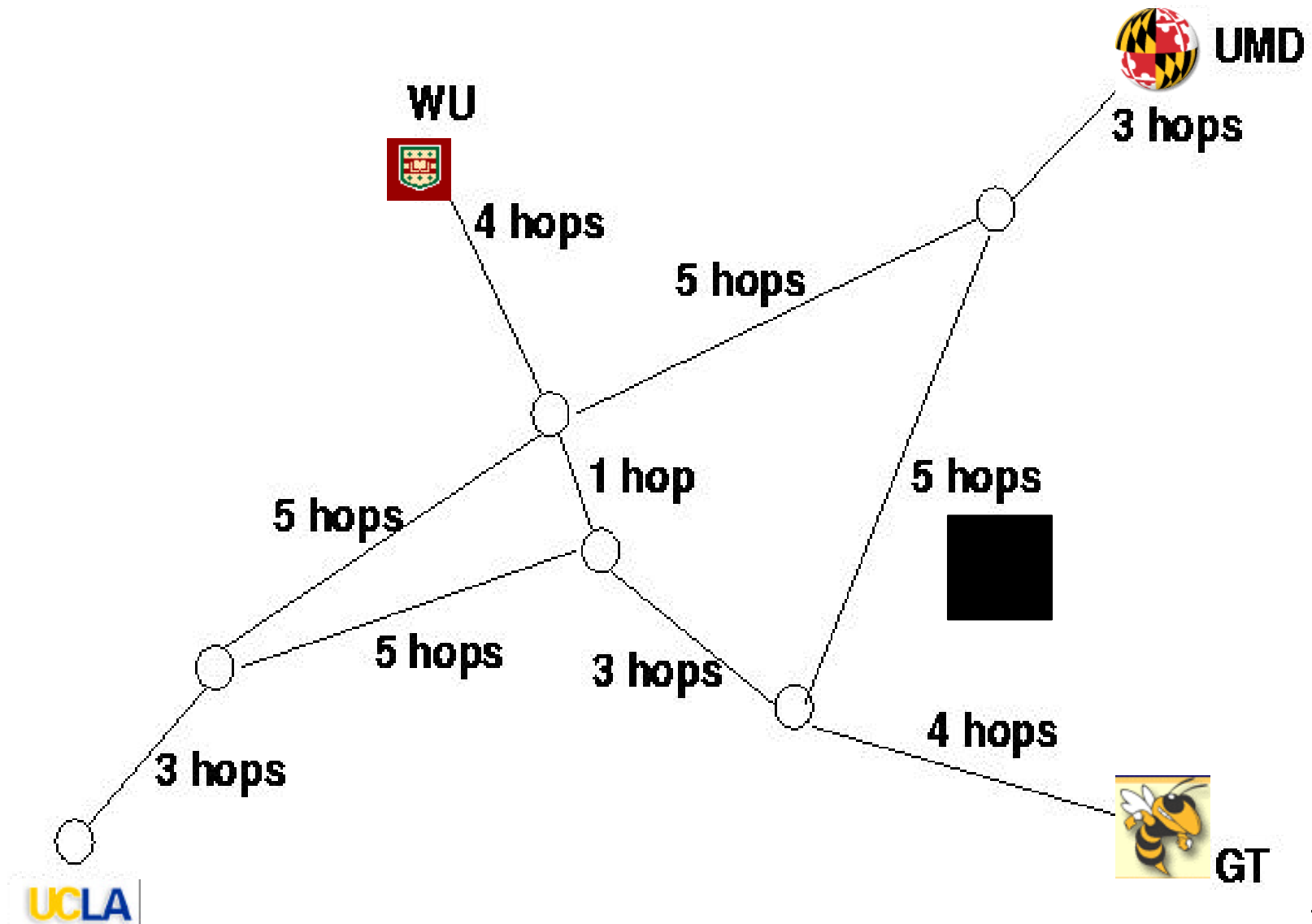
Resolver Process and Hybrid Technique

- r Resolver probe process:
 - m Request dummy file from server
 - m Measure response time (RT)
- r Hybrid push-probe technique
 - m Dummy file contains most recent TUFR
 - m Each probe: compute scaling factor $SF = RT/TUFR$
 - m Each push: estimate $RT = SF \times TUFR$

Performance of Hybrid Algorithm



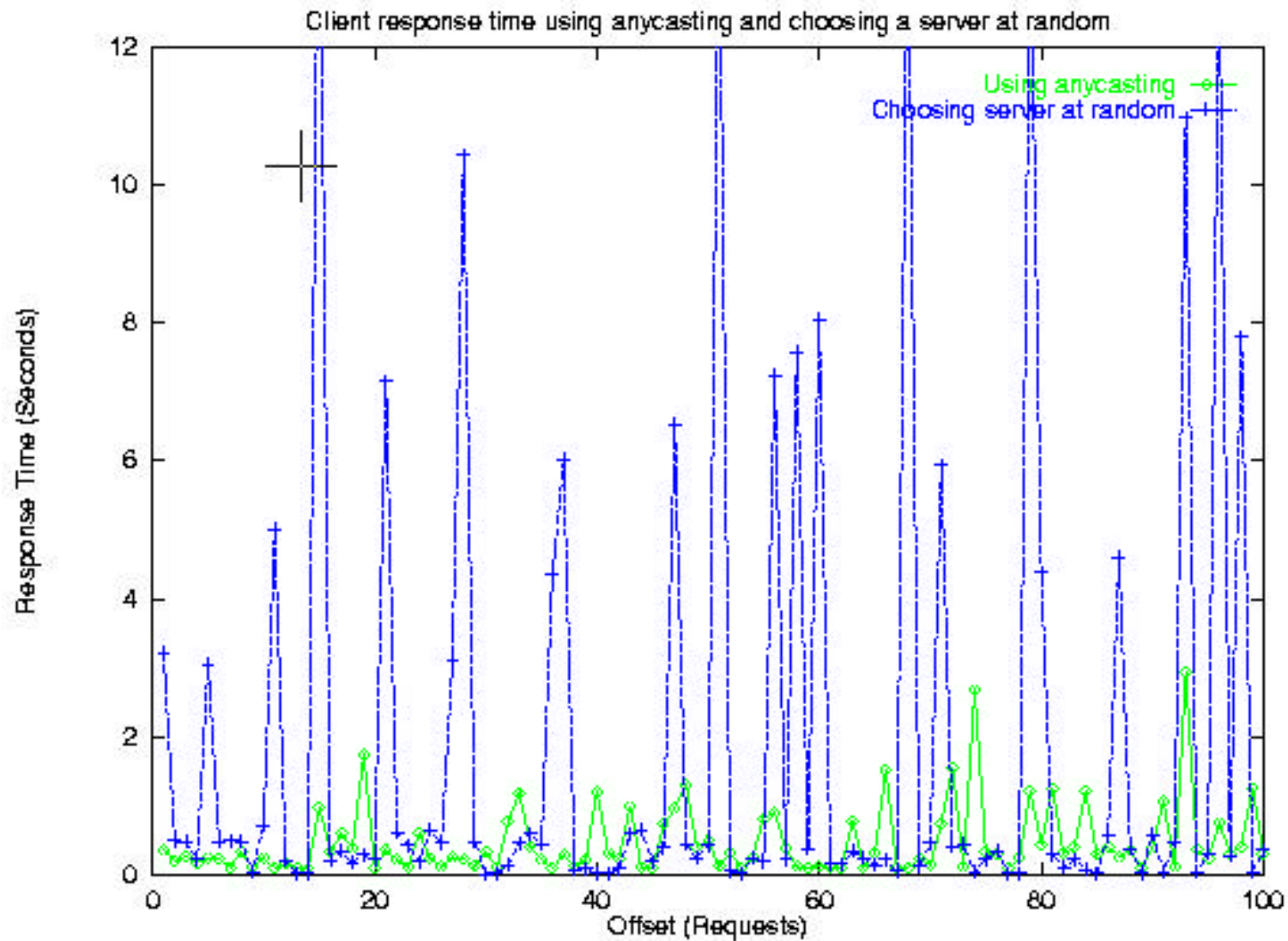
Wide-Area Experiment



Refinement

- r Problem of oscillation among servers
- r Set of Equivalent Servers (ES)
 - m Subset of the replicated servers whose measured performance is within a threshold of best performance

Performance of Anycast vs Random Selection



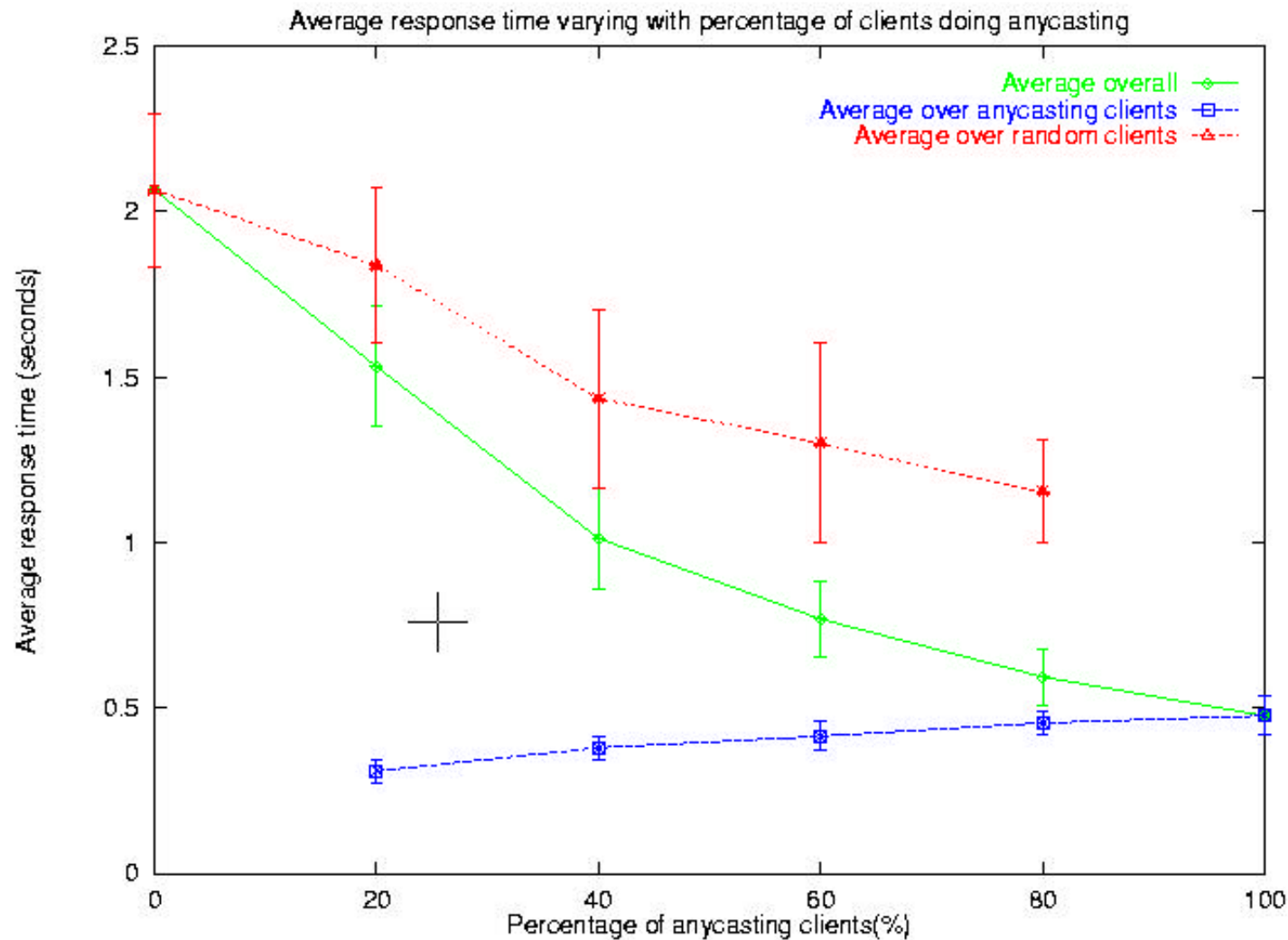
Performance of Server Location Schemes

Server Location Algorithm	Average Response Time (sec.)	Standard Deviation (sec.)
Anycasting	0.49	0.69
Nearest Server	1.12	2.47
Random	2.13	6.96

- 50% improvement using Nearest Server
- Another 50% improvement using Anycasting
- More predictable service

Performance as More Clients

Anycast



Avoid Oscillations Among Servers

- r Basic technique:
 - m Estimate response time for each server
 - m Indicate the best server when queried
- r Identifying one best server can result in oscillations
- r Use set of equivalent servers
- r Choose randomly among equivalent servers

Effect of Technique on Server Load

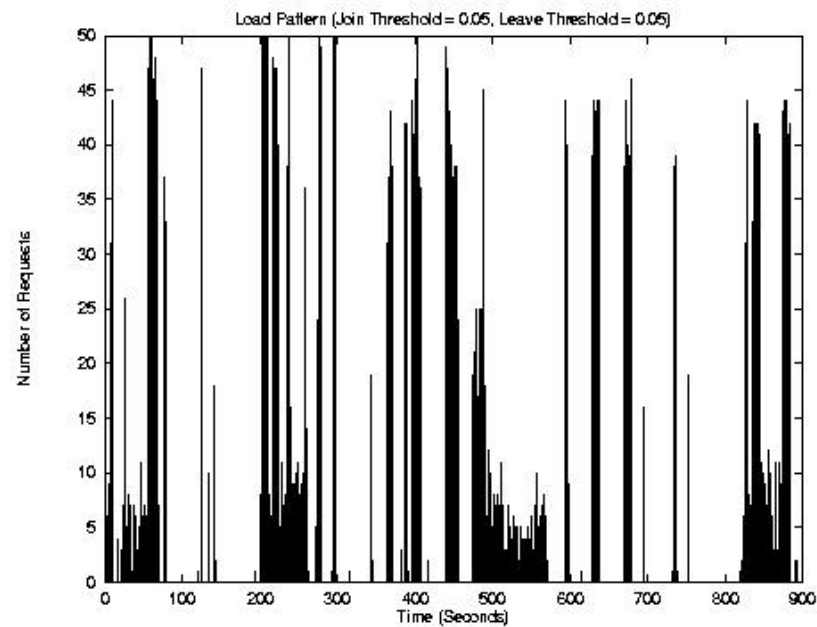


Figure 1. Low Threshold Values

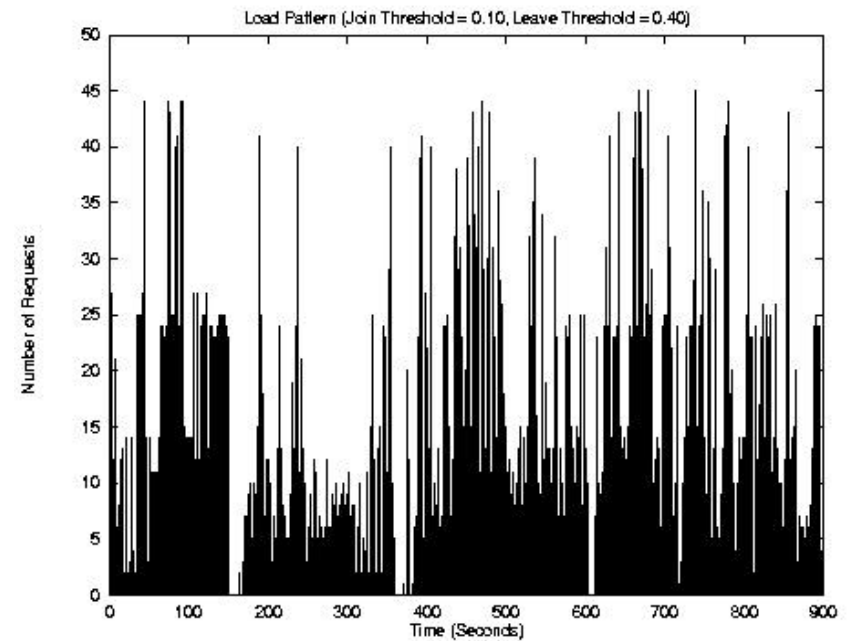


Figure 2: Higher Threshold Values

Scalability Techniques

- r Server can multicast pushed data
- r Server and resolver can control overhead
- r System can limit number of anycast groups
- r Resolver can track “most promising” servers
- r Users can pay premium for service

Conclusions

r Summary

- m Server replication increasingly important – web services etc.
- m Application layer architecture that is scalable using replicated resolvers organized in a DNS like hierarchy
- m Web server performance can be tracked with reasonable relative accuracy. Techniques used can be generalized to other servers
- m A hybrid push-probe technique provides scalable monitoring. May be useful in other contexts
- m Application-layer anycasting gives significant improvement over other server selection techniques

Any problems (truth in advertising)

r Discussions

m Was the study extensive enough ?

- 4 Anycast-aware servers – UCLA (1), WUSTL(1), Gatech (2)
- Anycast resolvers – UMD, Gatech
- 20 Anycast-aware clients – UMD (4), Gatech (16)

r Study of anycast vs random selection

m Experiments done one after another – any performance difference due to cached content ?

r Would performance improve if network-support for path performance metrics included ?

m Global IP-anycast [SIGCOMM00]

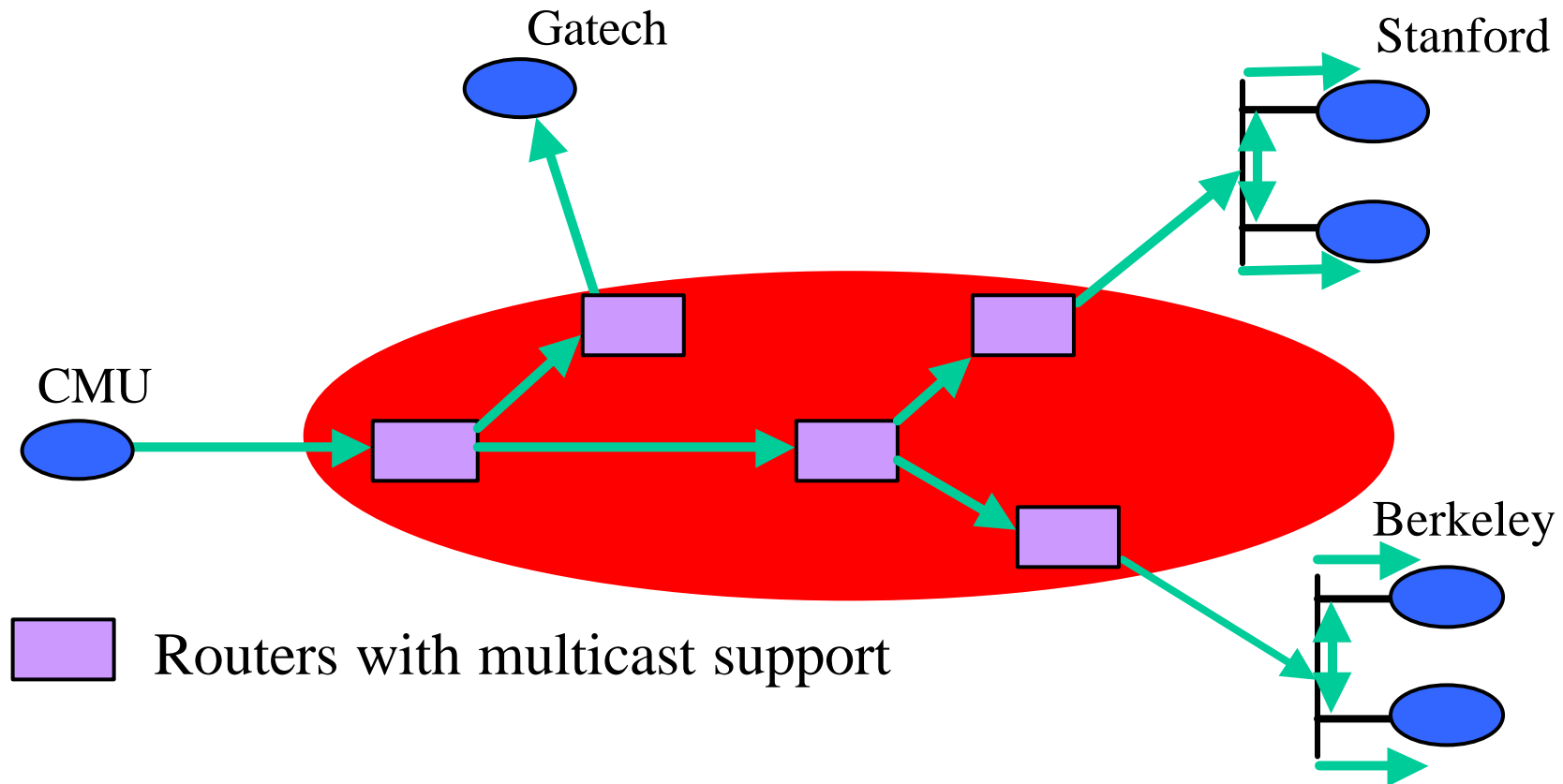
11. A Case For End System Multicast

Yang-hua Chu, Sanjay Rao and Hui Zhang
Carnegie Mellon University

Talk Outline

- r *Definition, potential benefits & problems*
- r *Look at Narada approach specifically*
- r *Performance in simulation*
- r *Performance in network experiment*
- r *Related work*
- r *Discussion*

IP Multicast



- No duplicate packets
- Highly efficient bandwidth usage

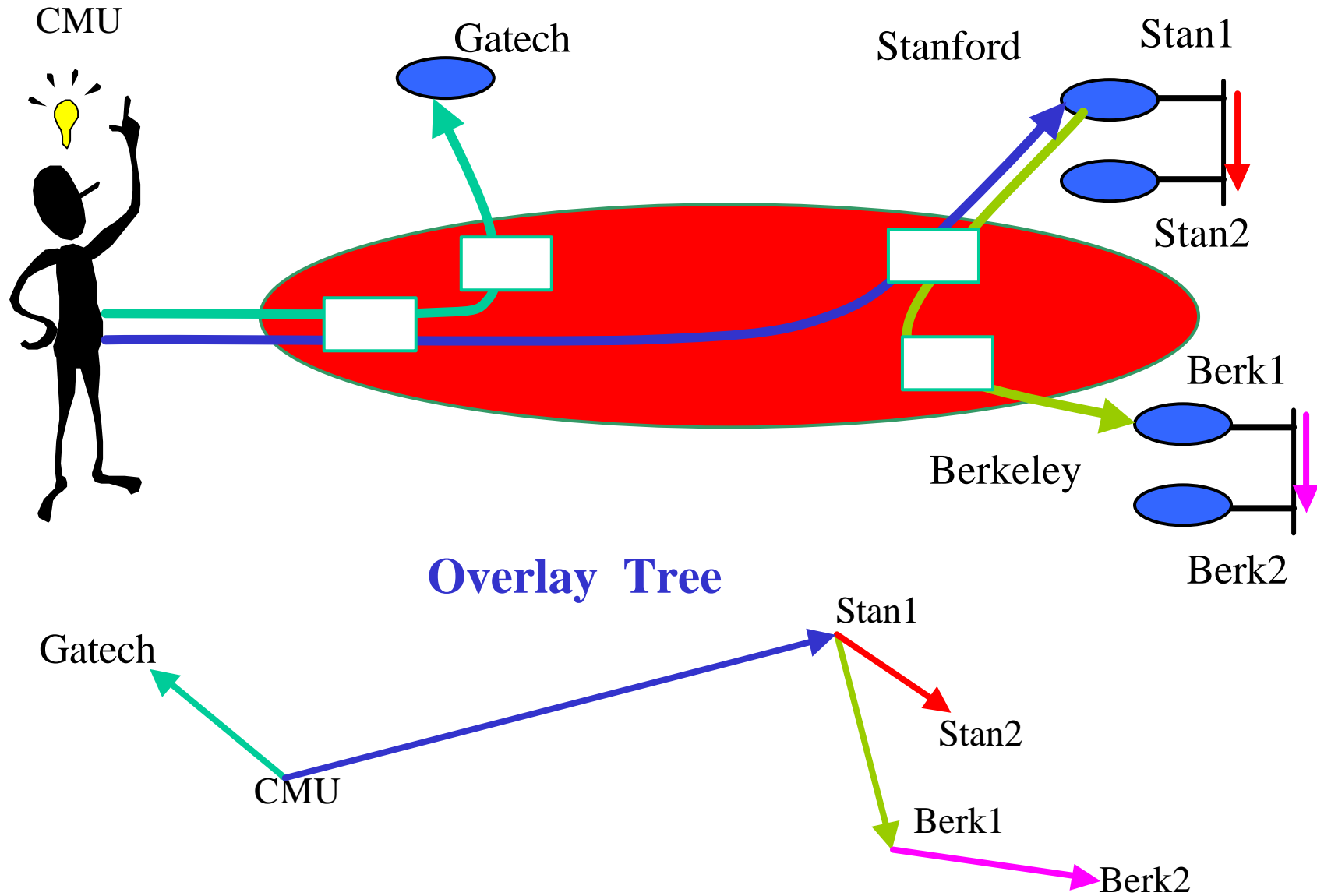
Key Architectural Decision: Add support for multicast in IP layer

Key Concerns with I P

Multicast

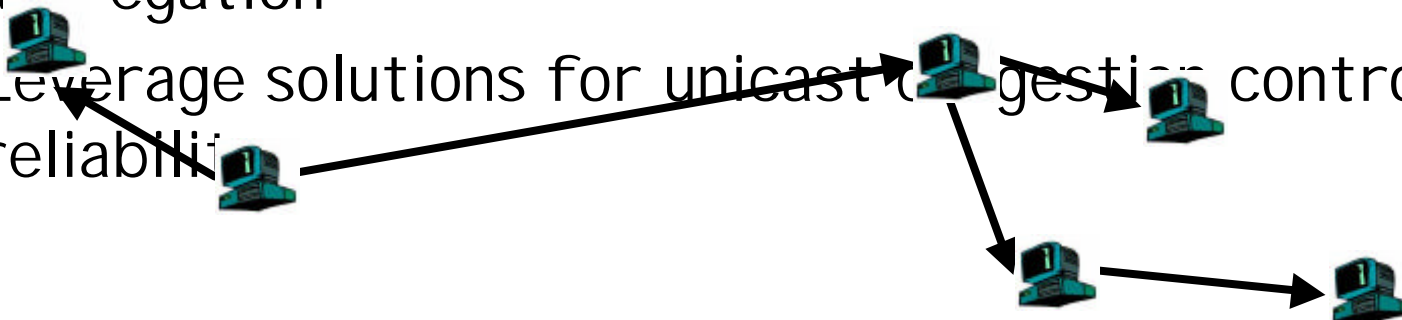
- r Scalability with number of groups
 - m Routers maintain **per-group state**
 - m Analogous to per-flow state for QoS guarantees
 - m Aggregation of multicast addresses is complicated
- r Supporting higher level functionality is difficult
 - m I P Multicast: **best-effort multi-point delivery** service
 - m End systems responsible for handling higher level functionality
 - m Reliability and congestion control for I P Multicast complicated
- r *Inter-domain routing is hard.*
- r *No management of flat address space.*
- r Deployment is difficult and slow
 - m I SP's reluctant to turn on I P Multicast

End System Multicast

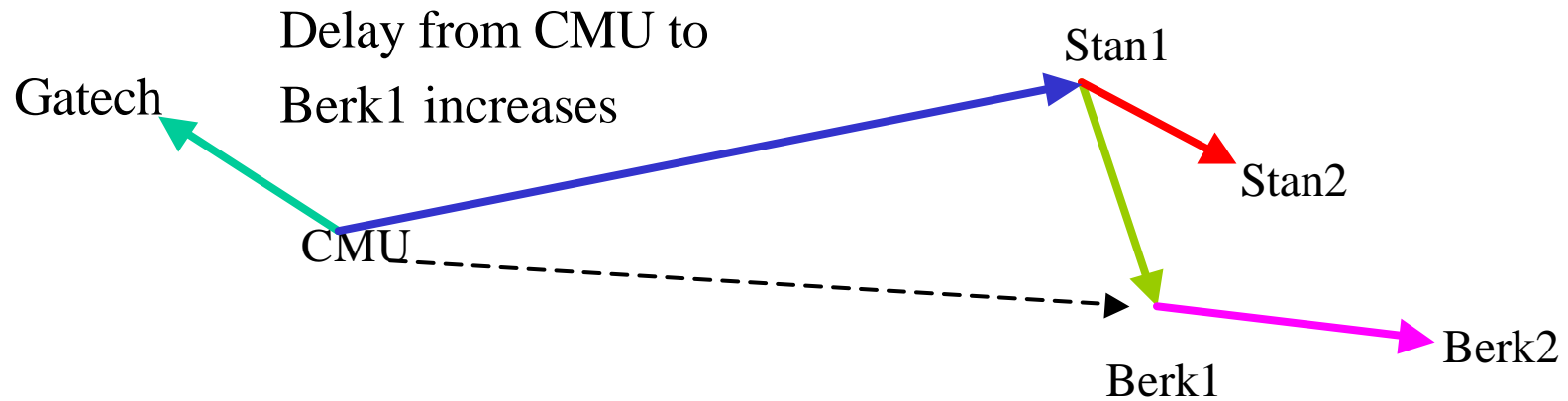


Potential Benefits

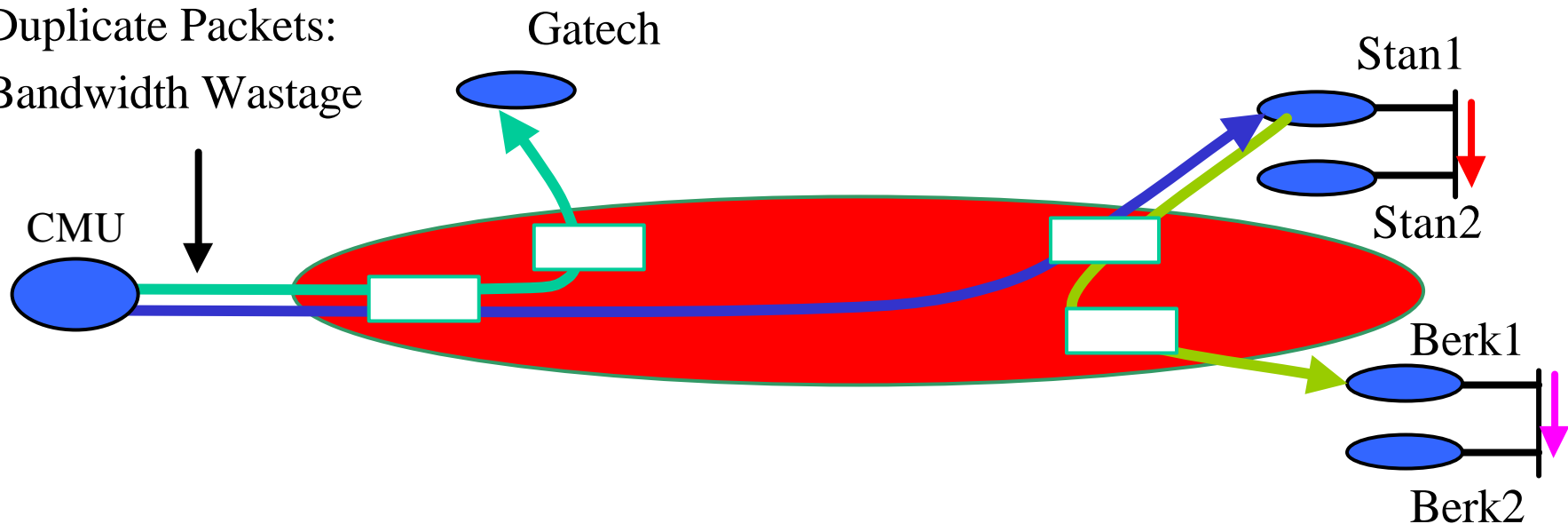
- r Scalability (*number of sessions in the network*)
 - m Routers do not maintain per-group state
 - m End systems do, but they participate in very few groups
- r Easier to deploy
- r Potentially simplifies support for higher level functionality
 - m Leverage computation and storage of end systems
 - m For example, for buffering packets, transcoding, ACK aggregation
 - m Leverage solutions for unicast congestion control and reliability



Performance Concerns



Duplicate Packets:
Bandwidth Wastage

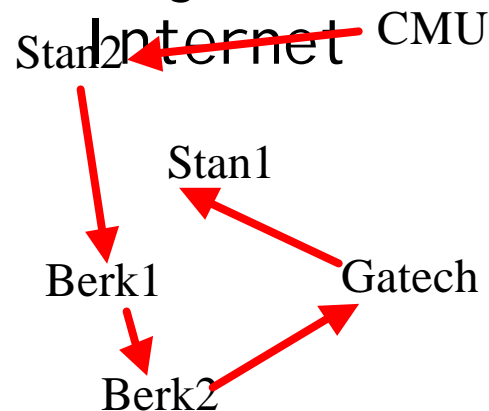


What is an efficient overlay

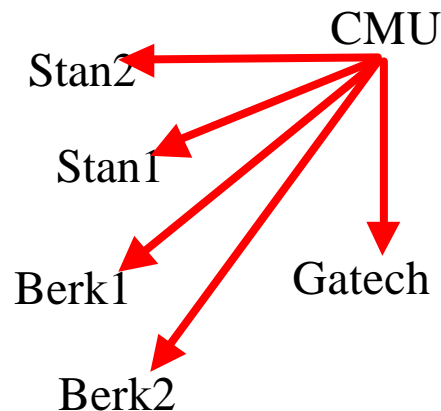
- r The delay between the source and receivers is small
- r Ideally,
 - m The number of redundant packets on any physical link is low

Heuristic we use:

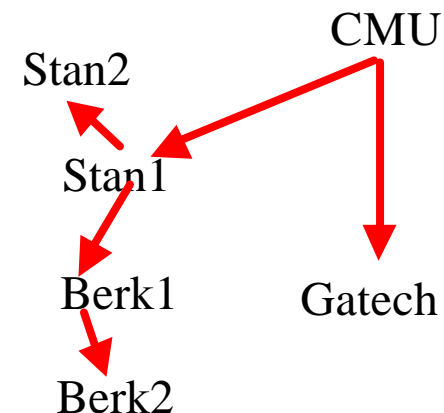
- m Every member in the tree has a small degree
- m Degree chosen to reflect bandwidth of connection to



High latency



High degree (unicast)



“Efficient” overlay

Why is self-organization hard?

- r Dynamic changes in group membership
 - m Members join and leave dynamically
 - m Members may die
- r Limited knowledge of network conditions
 - m Members do not know delay to each other when they join
 - m Members probe each other to learn network related information
 - m Overlay must **self-improve** as more information available
- r Dynamic changes in network conditions
 - m Delay between members may vary over time due to congestion

Narada Design (1)

Step 0

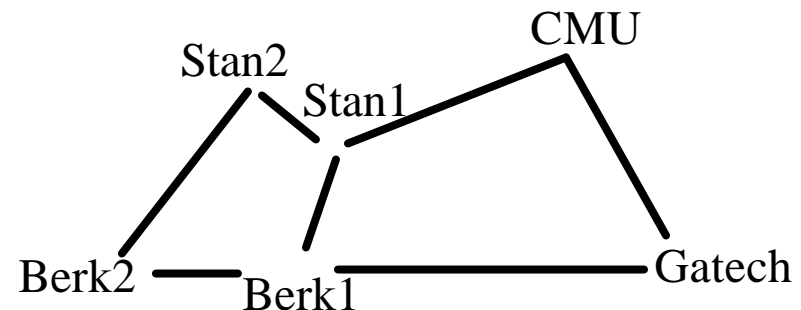
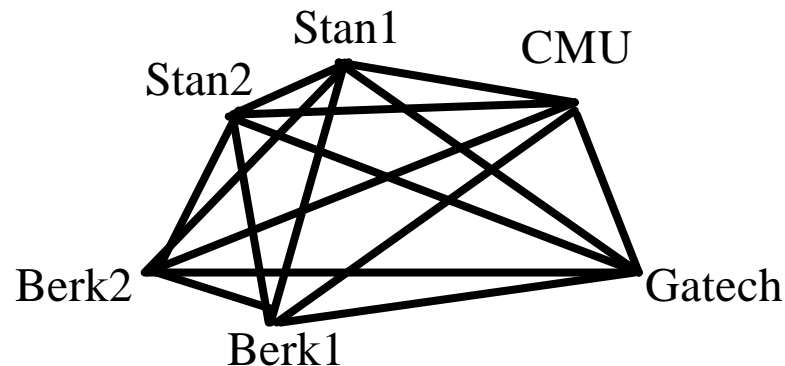
Maintain a complete overlay graph of all group members

- Links correspond to unicast paths
- Link costs maintained by polling

Step 1

“Mesh”: Subset of complete graph may have cycles and includes all group members

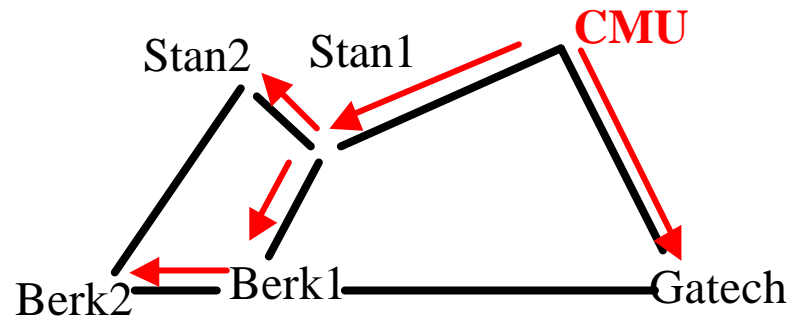
- Members have low degrees
- Shortest path delay between any pair of members along mesh is small



Narada Design (2)

Step 2

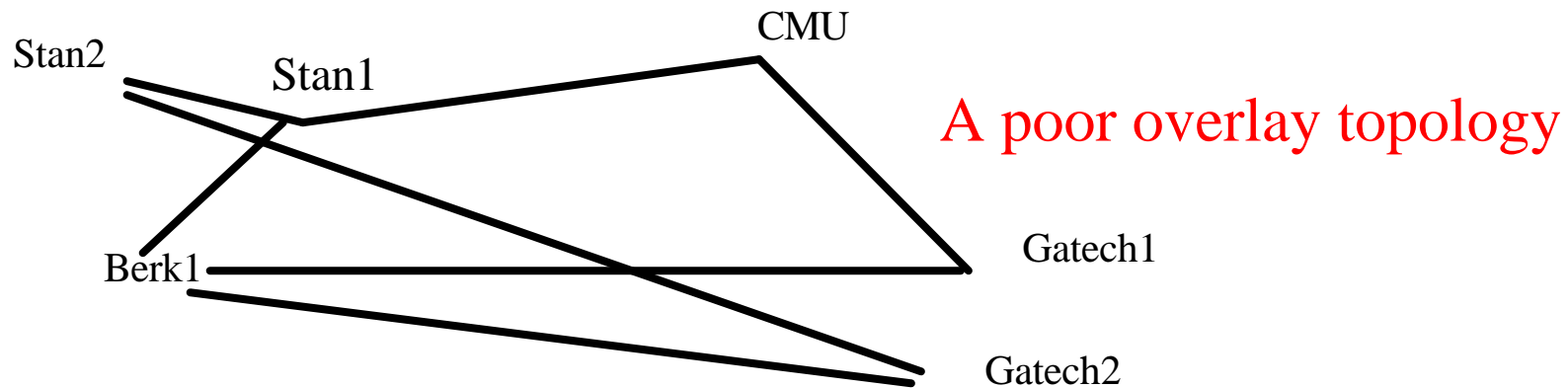
- Source rooted shortest delay spanning trees of mesh
- Constructed using well known routing algorithms
 - Members have low degrees
 - Small delay from source to receivers



Narada Components

- r Mesh Management:
 - m Ensures mesh remains connected in face of membership changes
- r Mesh Optimization:
 - m Distributed heuristics for ensuring shortest path delay between members along the mesh is small
- r Spanning tree construction:
 - m Routing algorithms for constructing data-delivery trees
 - m Distance vector routing, and reverse path forwarding

Optimizing Mesh Quality



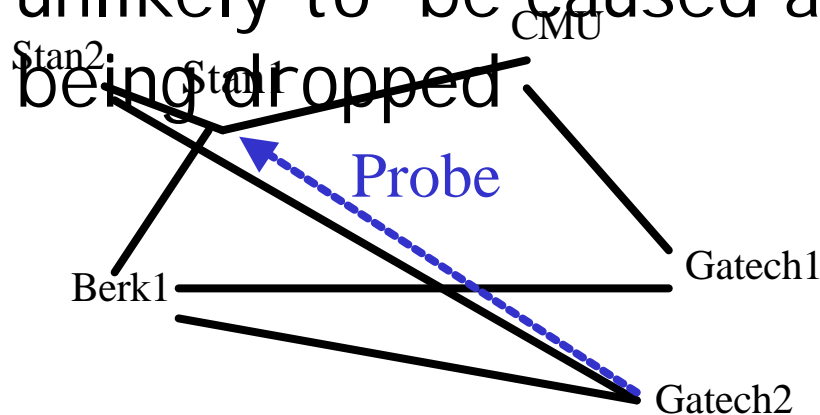
- r Members periodically probe other members at random
- r New Link added if
Utility Gain of adding link > Add Threshold
- r Members periodically monitor existing links
- r Existing Link dropped if
Cost of dropping link < Drop Threshold

The terms defined

- r **Utility gain of adding a link** based on
 - m The number of members to which routing delay improves
 - m How significant the improvement in delay to each member is
- r **Cost of dropping a link** based on
 - m The number of members to which routing delay increases, for either neighbor
- r **Add/Drop Thresholds** are functions of:
 - m Member's estimation of group size
 - m Current and maximum degree of member in the mesh

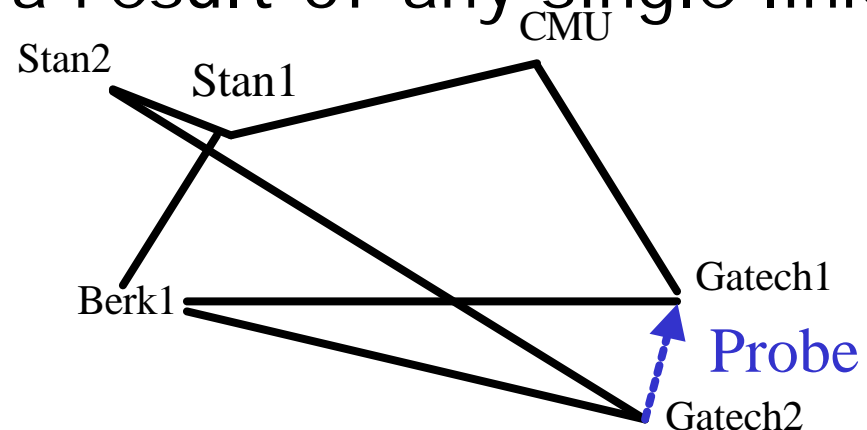
Desirable properties of heuristics

- r **Stability:** A dropped link will not be immediately readded
- r **Partition Avoidance:** A partition of the mesh is unlikely to be caused as a result of any single link



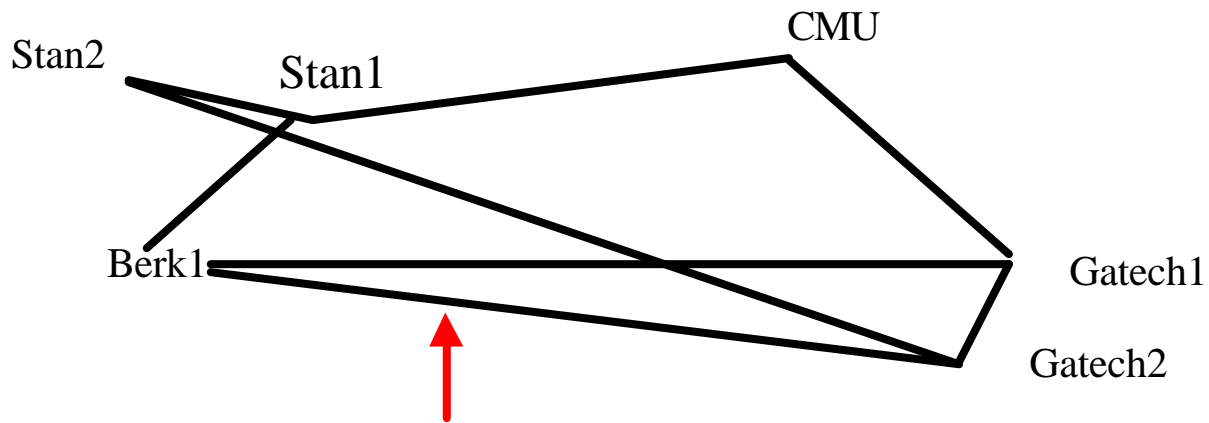
Delay improves to Stan1, CMU but marginally.

Do not add link!



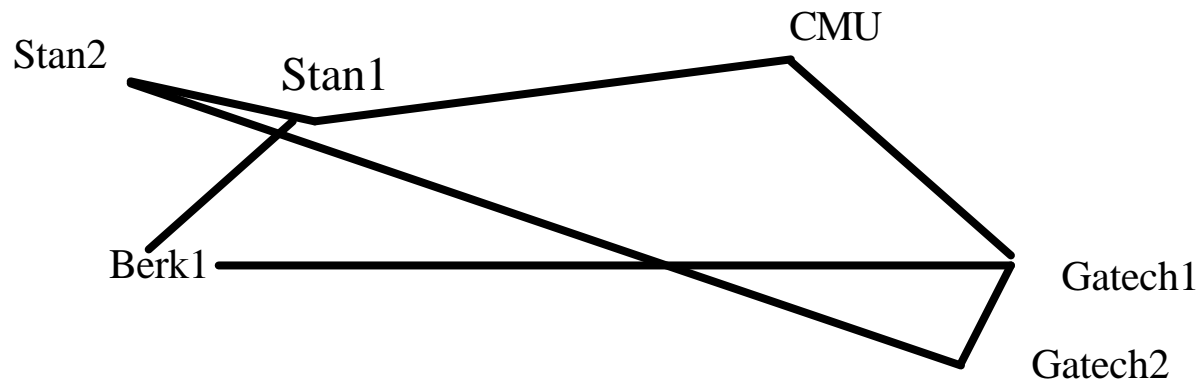
Delay improves to CMU, Gatech1 and significantly.

Add link!



Used by Berk1 to reach only Gatech2 and vice versa.

Drop!!



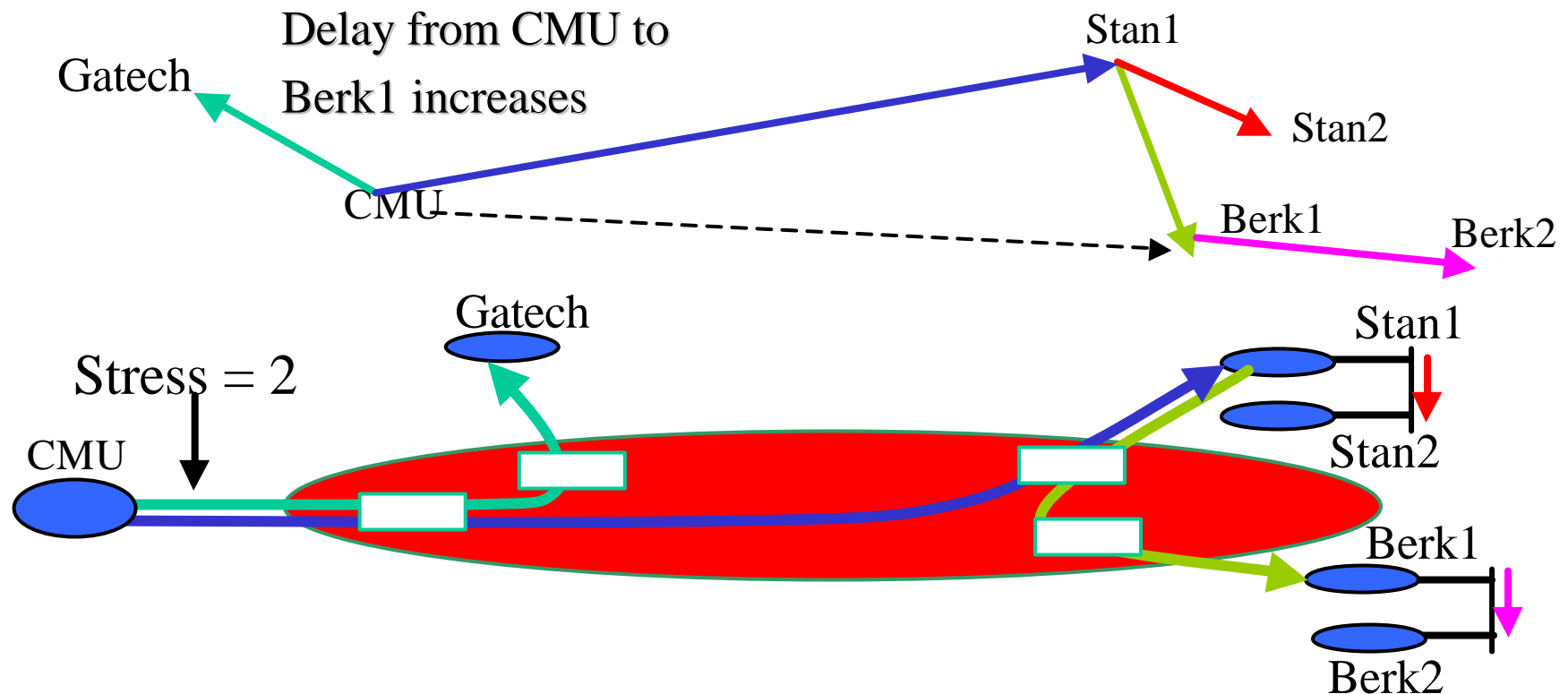
An improved mesh !!

Narada Evaluation

- r Simulation experiments
- r Evaluation of an implementation on the Internet

Performance Metrics

- r **Delay** between members using Narada
- r **Stress**, defined as the number of identical copies of a packet that traverse a physical link



Factors affecting performance

r **Topology Model**

- m Waxman Variant
- m Mapnet: Connectivity modeled after several ISP backbones
- m ASMap: Based on inter-domain Internet connectivity

r **Topology Size**

- m Between 64 and 1024 routers

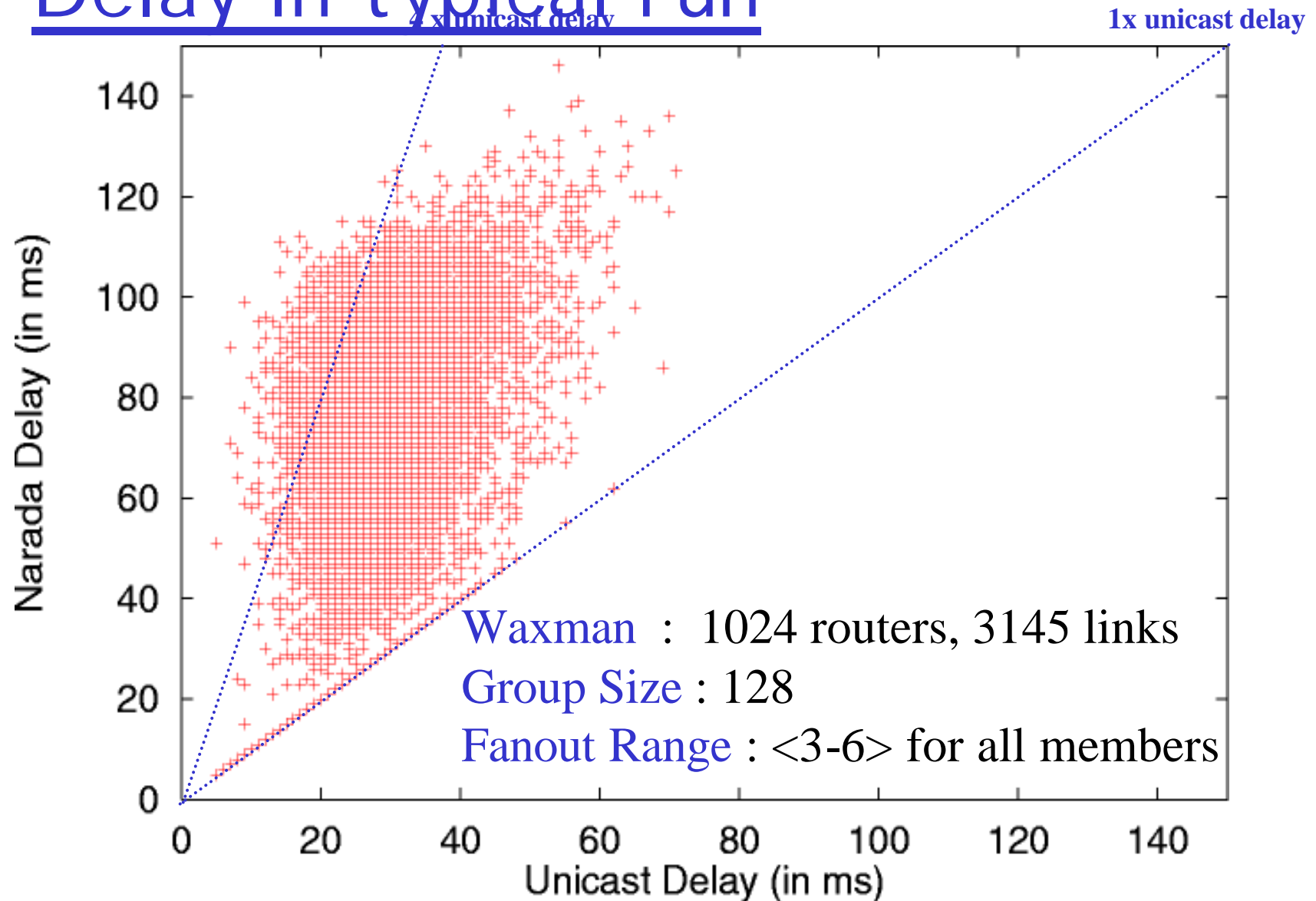
r **Group Size**

- m Between 16 and 256

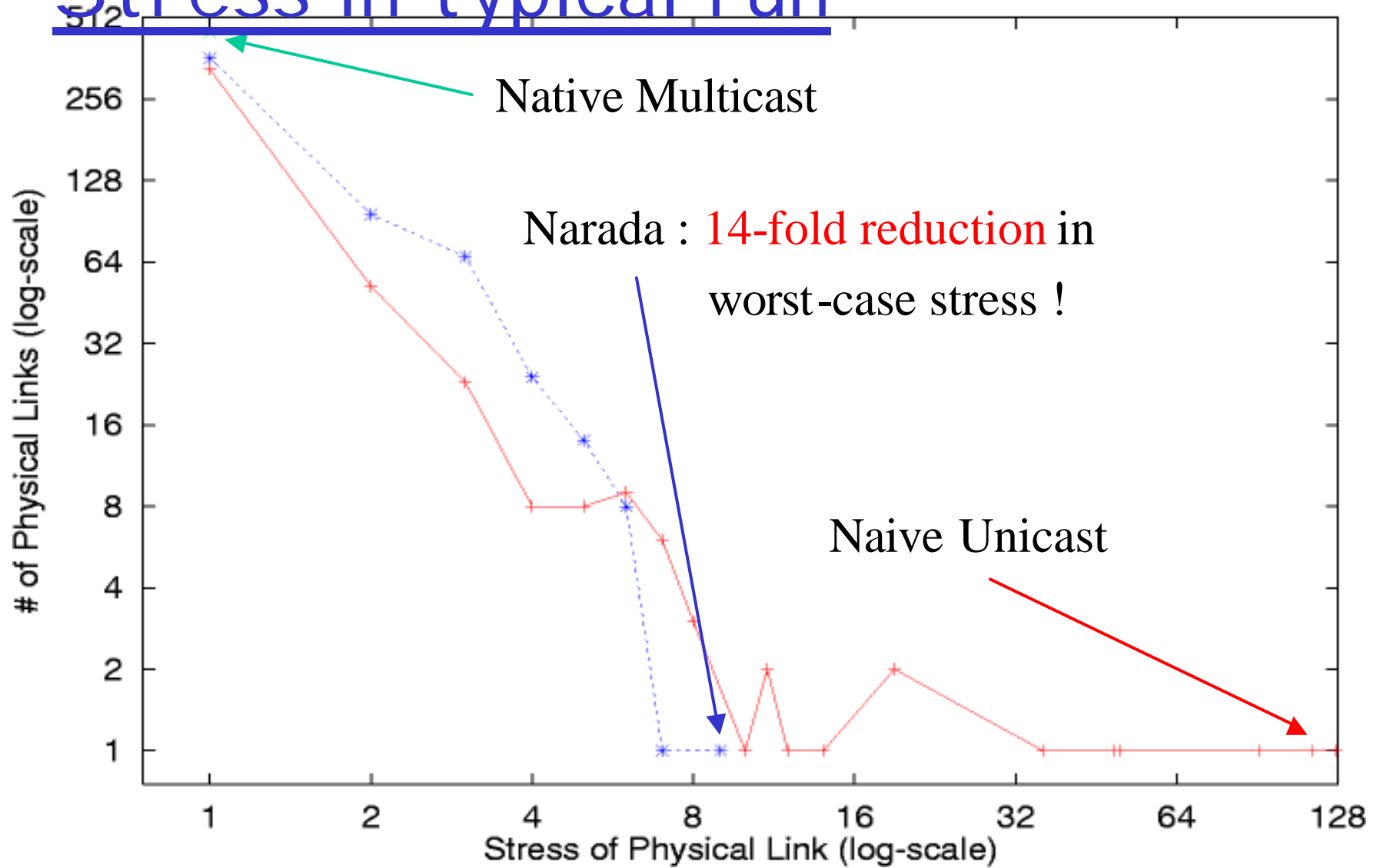
r **Fanout range**

- m Number of neighbors each member tries to maintain in the mesh

Delay in typical run



Stress in typical run



Overhead

- r *Two sources*
 - m *pairwise exchange of routing and control information*
 - m *polling for mesh maintenance.*
- r *Claim: Ratio of non-data to data traffic grows linearly with group size.*
- r *Narada is targeted at small groups.*

Related Work

- r Yoid (Paul Francis, ACI RI)
 - m More emphasis on architectural aspects, less on performance
 - m Uses a shared tree among participating members
 - More susceptible to a central point of failure
 - Distributed heuristics for managing and optimizing a tree are more complicated as cycles must be avoided
- r Scattercast (Chawathe et al, UC Berkeley)
 - m Emphasis on infrastructural support and proxy-based multicast
 - To us, an end system includes the notion of proxies
 - m Also uses a mesh, but differences in protocol details

Conclusions

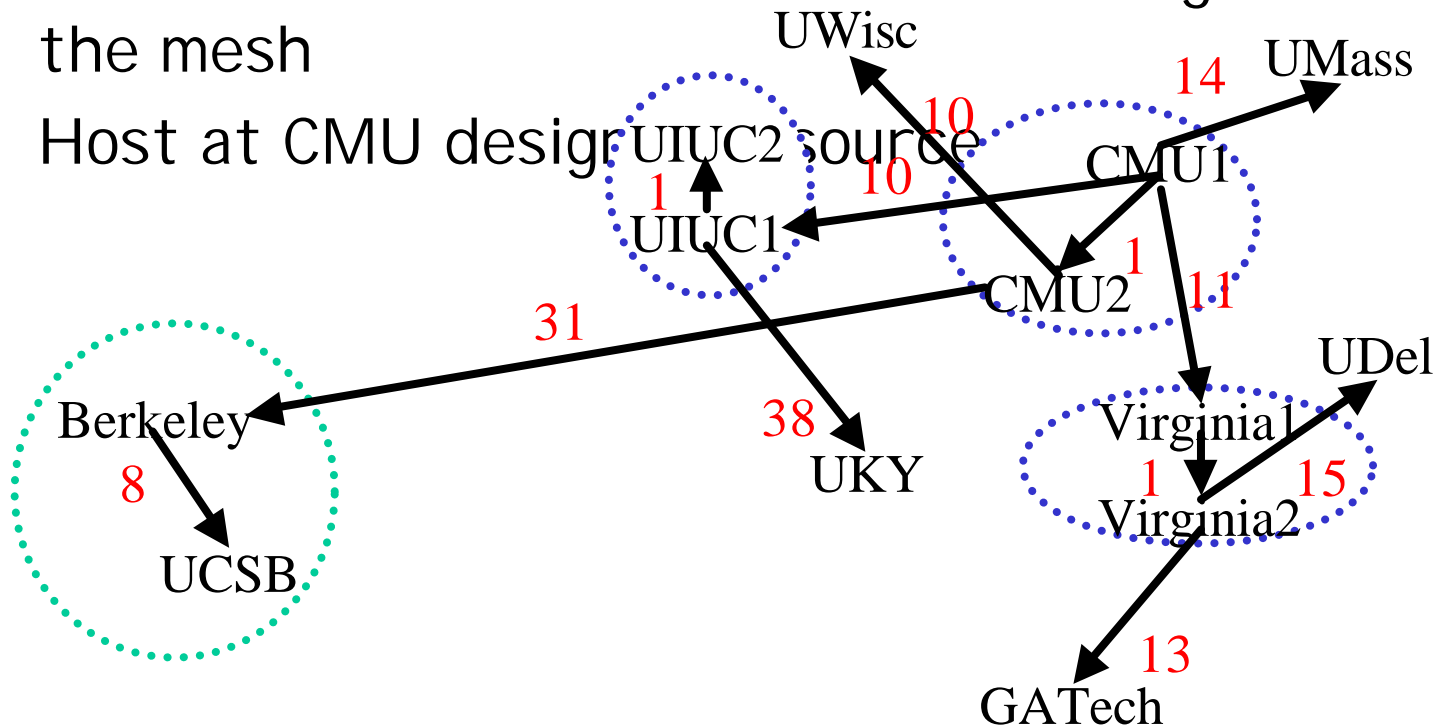
- r Proposed in 1989, IP Multicast is not yet widely deployed
 - m Per-group state, control state complexity and scaling concerns
 - m Difficult to support higher layer functionality
 - m Difficult to deploy, and get ISP's to turn on IP Multicast
- r Is IP the right layer for supporting multicast functionality?
- r For small-sized groups, an end-system overlay approach
 - m is **feasible**
 - m has a **low performance penalty** compared to IP Multicast
 - m has the potential to simplify support for higher layer functionality
 - m allows for application-specific customizations

Open Questions

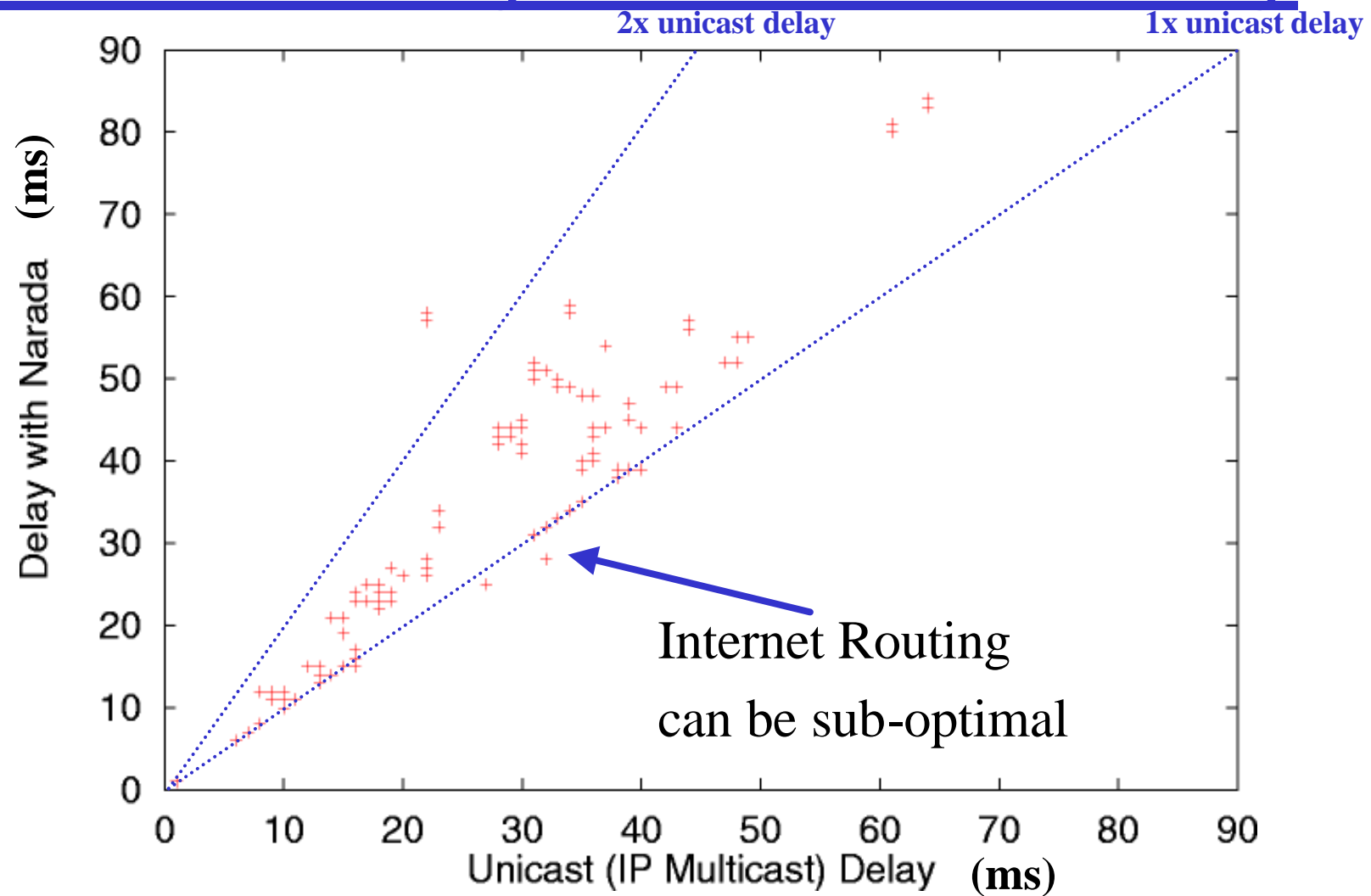
- r Theoretical bounds on how close an ESM tree can come to IP multicast performance.*
- r Alternate approach: Work with complete graph but modify multicast routing protocol.*
- r Leveraging unicast reliability and congestion control.*
- r Performance improvements: Reduce polling overhead.*

Internet Evaluation

- r 13 hosts, all join the group at about the same time
- r No further change in group membership
- r Each member tries to maintain 2-4 neighbors in the mesh
- r Host at CMU design



Narada Delay Vs. Unicast Delay



12. Overcast: Reliable Multicasting with an Overlay Network

Paper authors: Jannotti, Gifford,
Johnson, Kaashoek, O'Toole Jr.

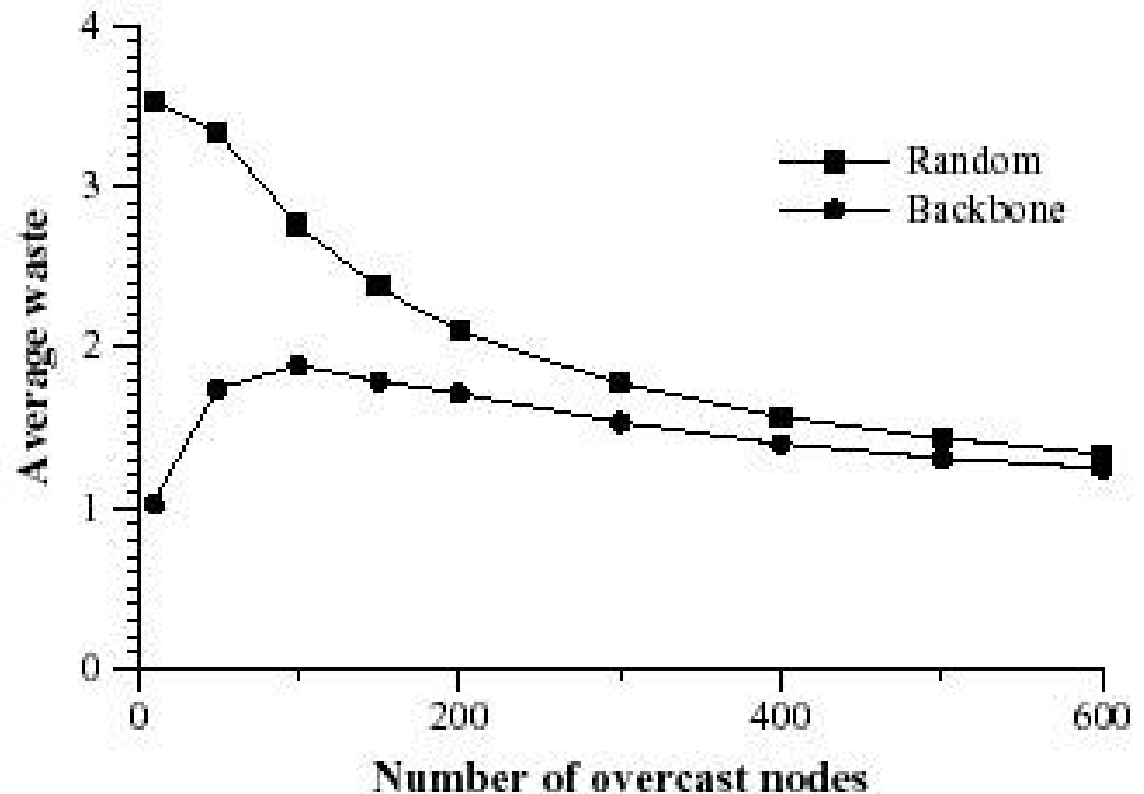
Design Goals

- r Provide application-level multicasting using already existing technology via **overcasting**
- r Scalable, efficient, and reliable distribution of high quality video
- r Compete well against **IP Multicasting**



Overcasting vs IP Multicast

Overcasting imposes about twice as much **network load** as IP multicast (for large networks). This difference is in $O(n)$. (*Figure 4*)



Architecture of an Overcast system

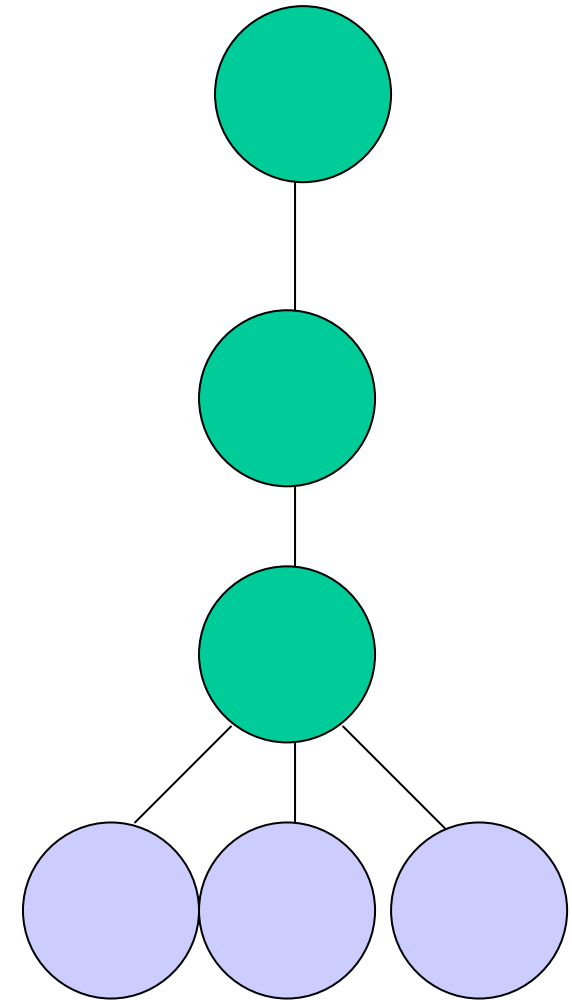
- r The entities which form the architecture of the Overcast system are called **nodes**.
- r Nodes are connected via an organizational scheme **distribution tree**.
- r Each **group** provides **content**, which is replicated at each of the nodes. A node can conceivably participate in more than one group.
- r **Clients** are end-system consumers of content.

Root nodes (1)

- r Content originates at the root node and is streamed down the distribution tree.
- r The root is the administrative center of the group.
- r When clients join the group, they go to a webpage at the root, which redirects them to the “best” overcast node
- r Is the root a single point of failure?

Root nodes (2)

- **Linear roots** can alleviate this problem. For the source to fail, all roots must fail.
- Using **round robin IP resolution**, you can stop your content serving cluster from being ‘slashdotted’ (overloaded by sheer popularity) by client requests.



Distribution tree

- r The distribution tree is built and maintained using a **self-organizing** algorithm.
- r The primary heuristic of this algorithm is to maximize **bandwidth** from the root to an overcast node.
- r **Backbone nodes** are nodes which are located on or close to a network backbone. Overcast performs better when these backbone nodes are located at the top of the tree (ie, they are switched on first)

Tree building protocol (1)

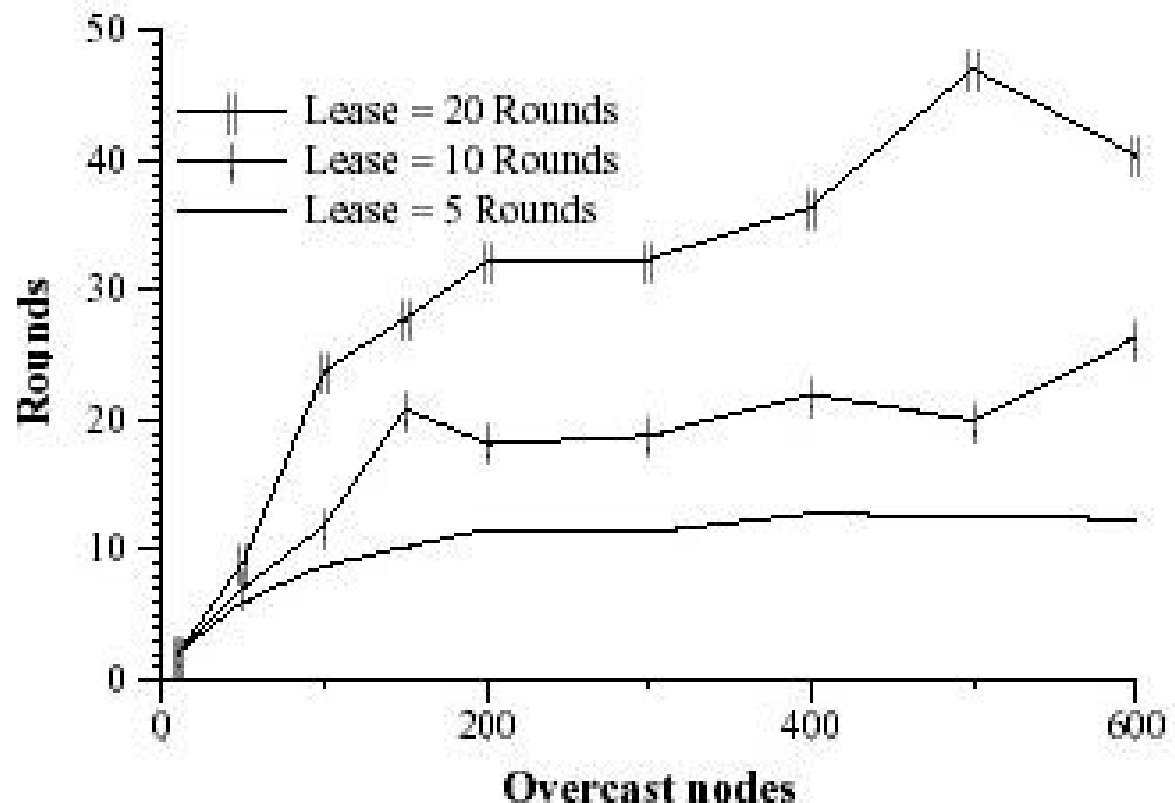
- r A node **initializes** by booting up, obtaining its IP, and contacts a “global, well-known registry” (Possible point of failure?) with a unique **serial number**.
- r **Registry** provides a list of groups to join.
- r This node initially chooses the root as its **parent**. A series of **rounds** will begin in which the node decides where on the tree it should be.

Tree building protocol (2)

- r For each round, evaluate the bandwidth we have to our parent. Also consider the bandwidth to the rest of our parent's children.
- r If there is a tie (bandwidth differences between 2 or more nodes within 10%), break it by the number of hops reported by **traceroute**.
- r The child selects the best of its parents children as its new parent.
- r Nodes maintain an **ancestor list** and can rejoin further up the tree when its ancestors fail.

Reaching a stable state

Overcast nodes can still receive content from the root, even when the tree is not stabilized. A typical round period is about 1 to 2 seconds.



Tree building protocol (3)

- r By having consecutive rounds of tree building, the distribution tree can overcome conditions occurring on the underlying network.
- r The **up/down** protocol is used to maintain information about the status of nodes.
- r Each node in the network maintains a table of information about all of its descendants.
- r After the tree stabilizes, nodes will continue to consider relocating after a **reevaluation period**.

Up/down protocol

- r A parent that gets a new child gives its parent a **birth certificate**, which propagates to the root. The node's **sequence number** is incremented by 1. Sequence numbers are used to prevent a race condition.
- r When a parent's child fails to report its status (called **checkin**), after a length of time called a **lease period**, the parent propagates a **death certificate** for its child up the tree.
- r A child also presents any certificates or changes it has accumulated from its last checkin.

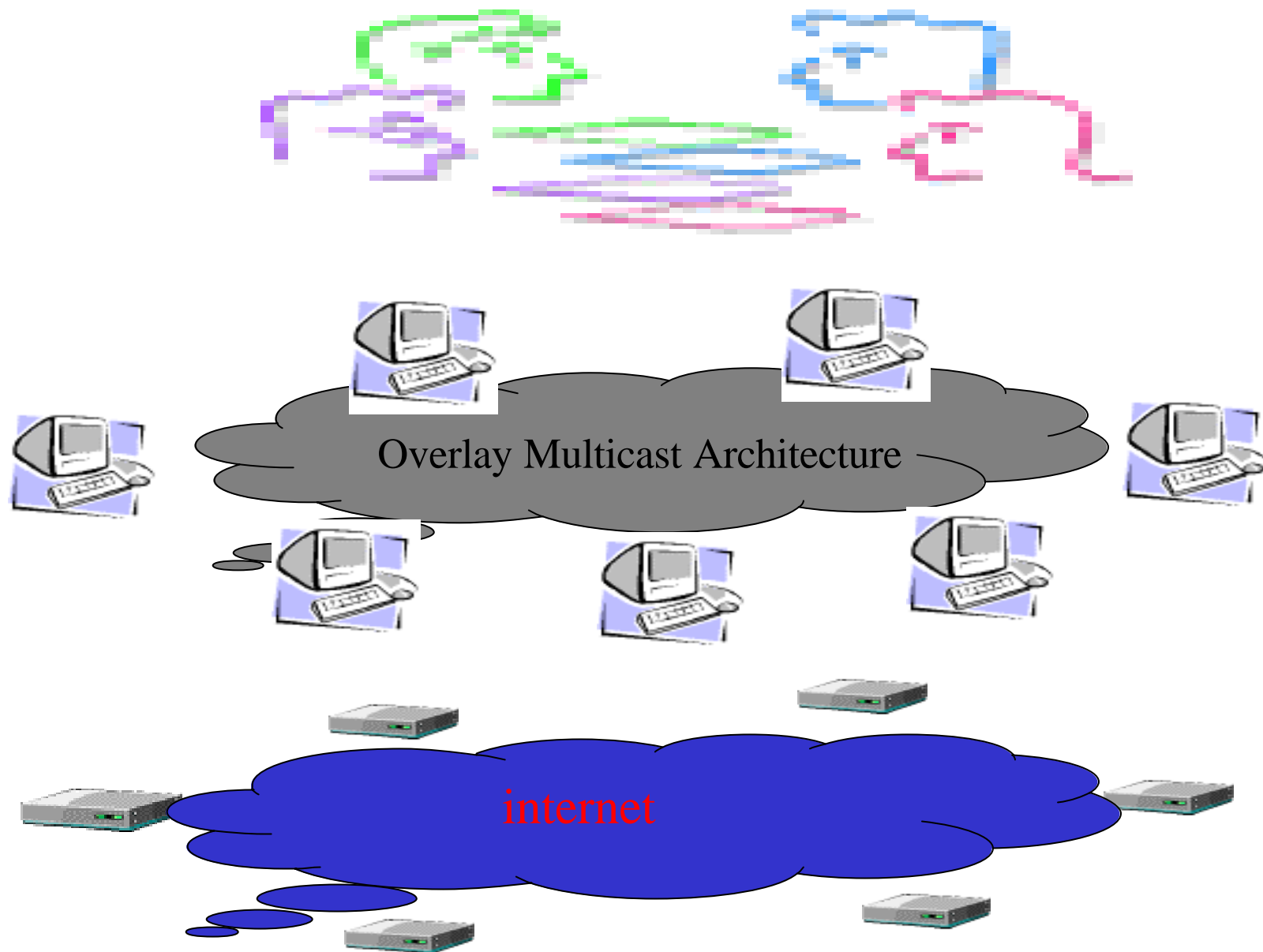
Problems

- r A simple approach leads to a simple solution:
- r Not appropriate for software (same as mirroring a file!)
- r Not appropriate for games (latency is too high!)
- r What about teleconferencing? Authors suggest that if a non-root node wishes to send to other nodes, the node should first **unicast** (send via normal TCP/IP to the root, which will then overcast it as normal). Is this a good solution?

Closing thoughts

- r Simulated on a virtual network topology (Georgia Tech Internetwork Topology Models). Take results with a grain of salt (or two).
- r Might work out well commercially. Consider a high demand for high definition video over the net, and corporations/entities willing to deliver it (CNN, Hollywood studios, Olympics). Overcast could be a premium service for ISP subscribers [like newsgroups, www hosting].

13. Enabling Conferencing
Applications
on the Internet using an
Overlay Multicast Architecture
Yanghua Chu et al. (CMU)
SIGCOMM'01



Past Work

- r Self-organizing protocols
 - m Yoid (ACI RI), Narada (CMU), Scattercast (Berkeley), Overcast (CI SCO), Bayeux (Berkeley), ...
 - m Construct overlay trees in distributed fashion
 - m Self-improve with more network information

- r Performance results showed promise, but...
 - m Evaluation conducted in simulation
 - m Did not consider impact of network dynamics on overlay performance

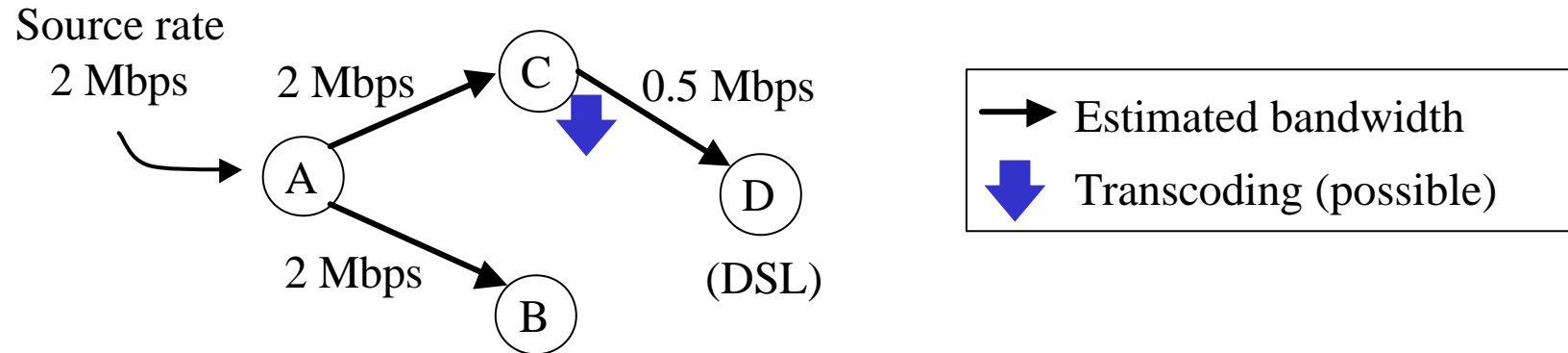
Focus of This Paper

- r Can End System Multicast support real-world applications on the Internet?
 - m Study in context of conferencing applications
 - m Show performance acceptable even in a dynamic and heterogeneous Internet environment
- r First detailed Internet evaluation to show the feasibility of End System Multicast

Why Conferencing?

- r Important and well-studied
 - m Early goal and use of multicast (vic, vat)
- r Stringent performance requirements
 - m High bandwidth, low latency
- r Representative of interactive applications
 - m E.g., distance learning, on-line games

Supporting Conferencing in ESM



r Framework

- m Bandwidth estimation
- m Adapt data rate to bandwidth est. by packet dropping

r Objective

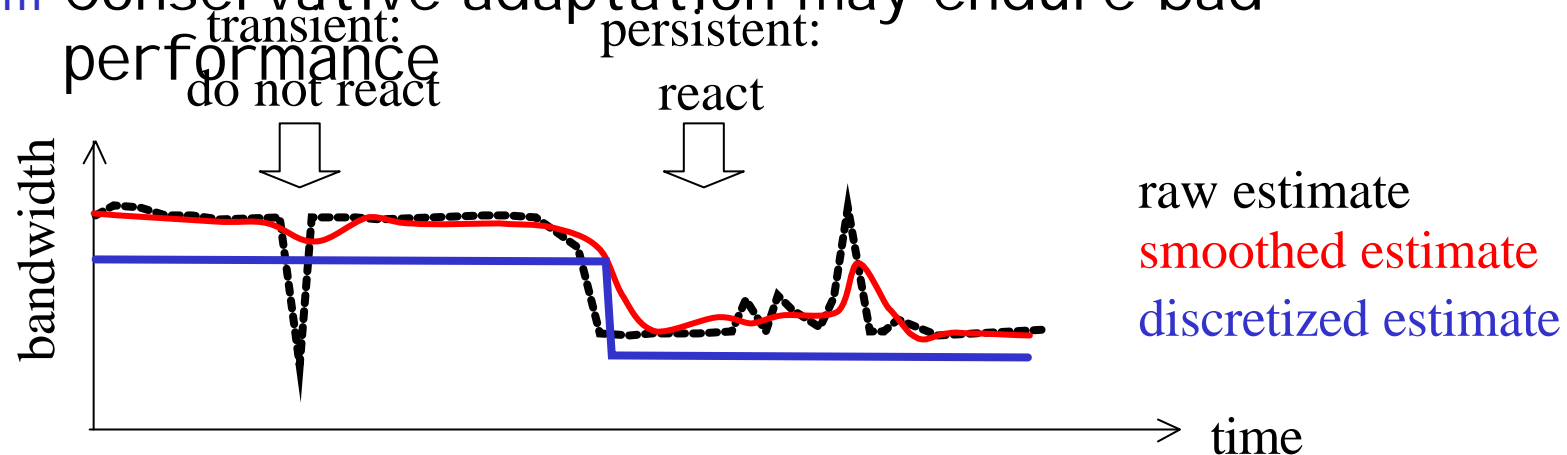
- m High bandwidth and low latency to all receivers along the overlay

Enhancements of Overlay Design

- r Two new issues addressed
 - m Dynamically adapt to changes in network conditions
 - m Optimize overlays for multiple metrics
 - Latency and bandwidth
- r Study in the context of the Narada protocol (Sigmetrics 2000)
 - m Use Narada to define logical topology
 - m Techniques presented apply to all self-organizing protocols

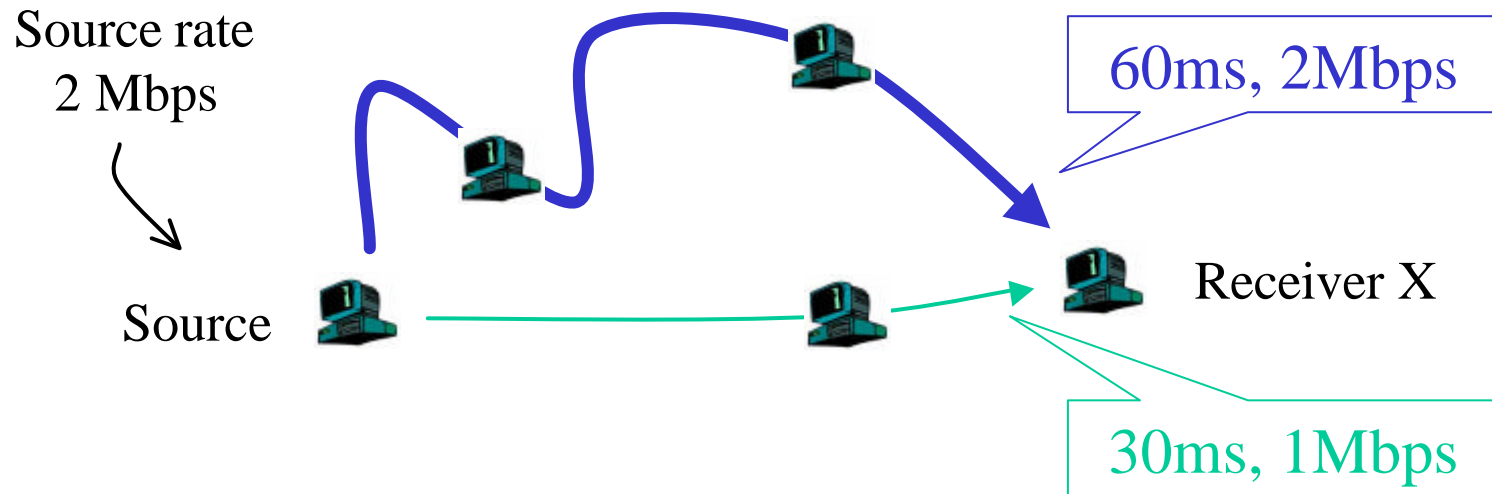
Adapt to Dynamic Metrics

- r Adapt overlay trees to changes in network condition
 - m Monitor bandwidth and latency of overlay links
- r Link measurements can be noisy
 - m Aggressive adaptation may cause overlay instability
 - m Conservative adaptation may endure bad



- Capture the long term performance of a link
 - Exponential smoothing, Metric discretization

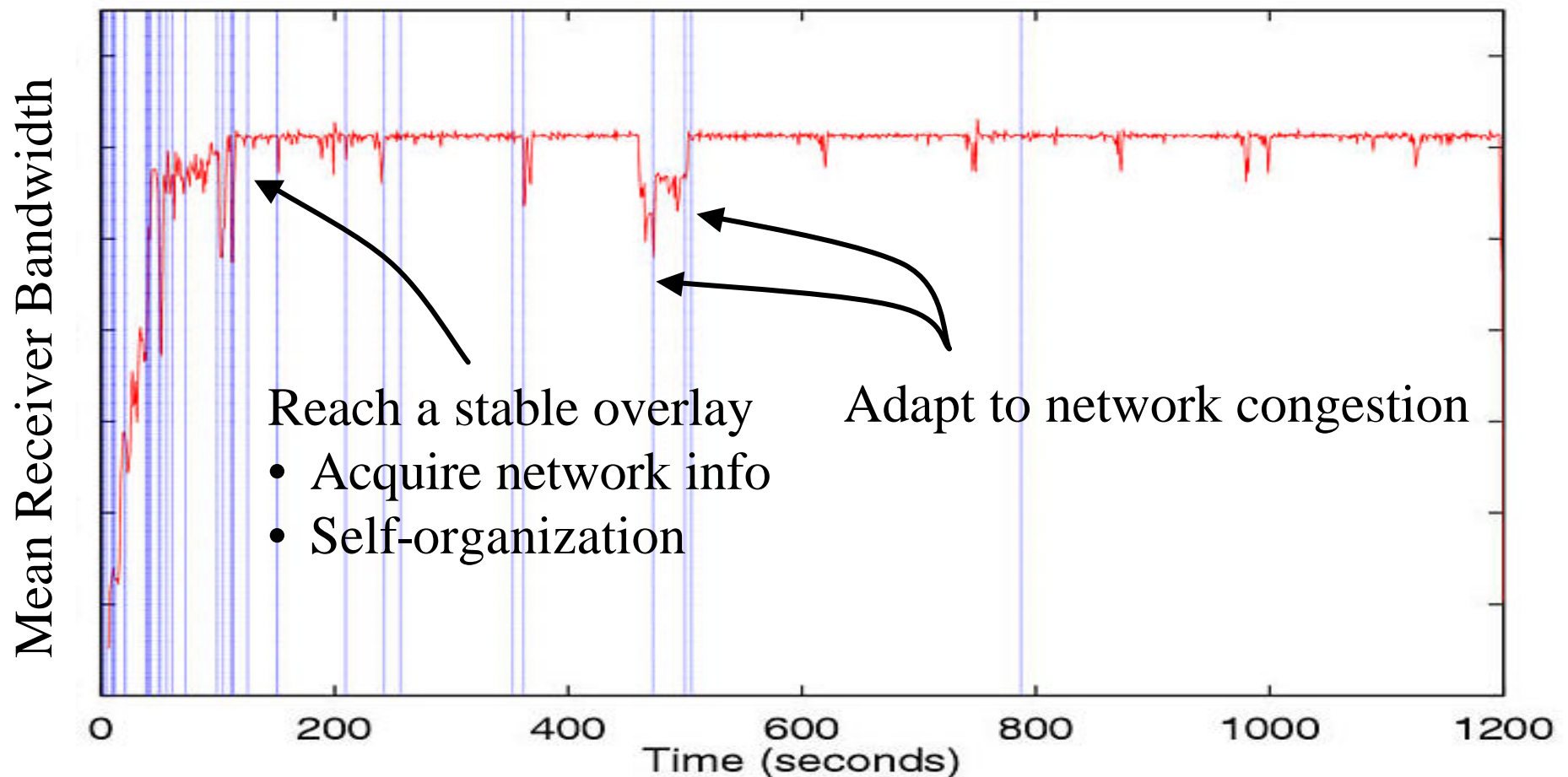
Optimize Overlays for Dual Metrics



- r Prioritize bandwidth over latency
- r Break tie with shorter latency

Example of Protocol Behavior

- r All members join at time 0
- r Single sender, CBR traffic



Evaluation Goals

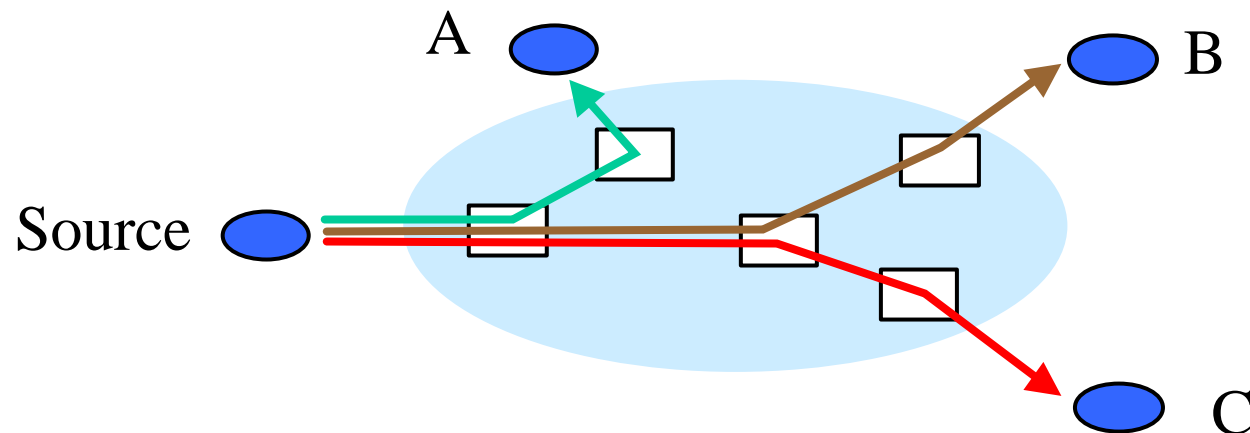
- r Can ESM provide application level performance comparable to IP Multicast?
- r What network metrics must be considered while constructing overlays?
- r What is the network cost and overhead?

Evaluation Overview

- r Compare performance of our scheme with
 - m Benchmark (sequential unicast, mimicing I P Multicast)
 - m Other overlay schemes that consider fewer network metrics
- r Evaluate schemes in different scenarios
 - m Vary host set, source rate
- r Performance metrics
 - m Application perspective: latency, bandwidth
 - m Network perspective: resource usage, overhead

Benchmark Scheme

- r IP Multicast not deployed
- r **Sequential Unicast**: an approximation
 - m Bandwidth and latency of unicast path from source to each receiver
 - m Performance similar to IP Multicast with ubiquitous deployment



Overlay Schemes

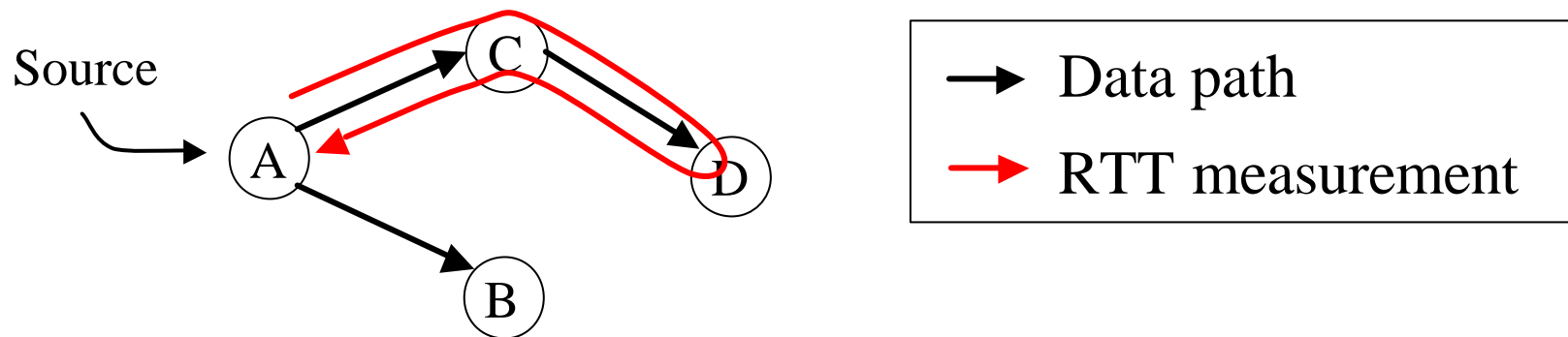
Overlay Scheme	Choice of Metrics	
	Bandwidth	Latency
Bandwidth-Latency	✓	✓
Bandwidth-Only	✓	✗
Latency-Only	✗	✓
Random	✗	✗

Experiment Methodology

- r Compare different schemes on the Internet
 - m Ideally: run different schemes concurrently
 - m Interleave experiments of schemes
 - m Repeat same experiments at different time of day
 - m Average results over 10 experiments
- r For each experiment
 - m All members join at the same time
 - m Single source, CBR traffic
 - m Each experiment lasts for 20 minutes

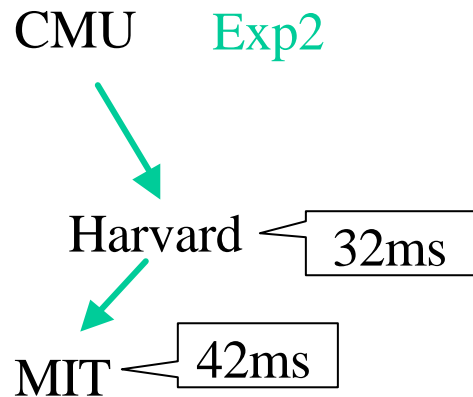
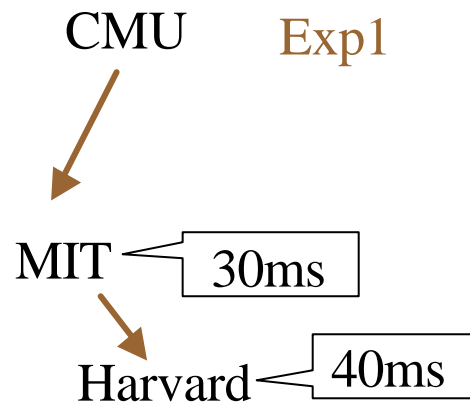
Application Level Metrics

- r **Bandwidth** (throughput) observed by each receiver
- r **RTT** between source and each receiver along overlay

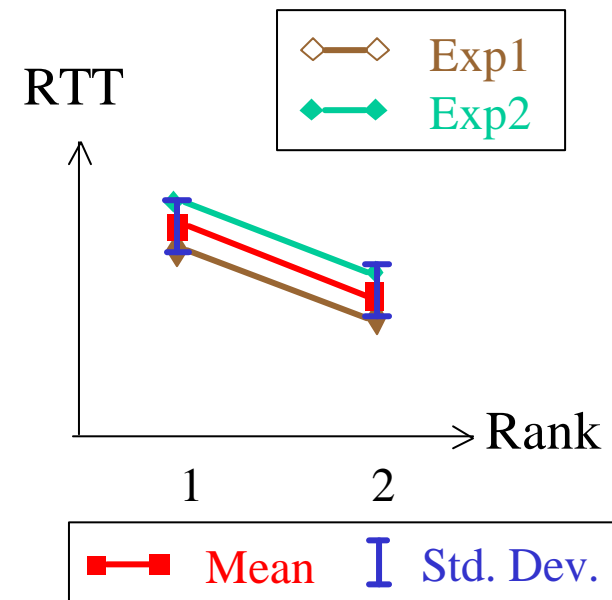


These measurements include queueing and processing delays at end systems

Performance of Overlay Scheme



Different runs of the same scheme may produce different but “similar quality” trees



“Quality” of overlay tree produced by a scheme

- r Sort (“rank”) receivers based on performance
- r Take mean and std. dev. on performance of same rank across multiple experiments
- r Std. dev. shows variability of tree quality

Factors Affecting Performance

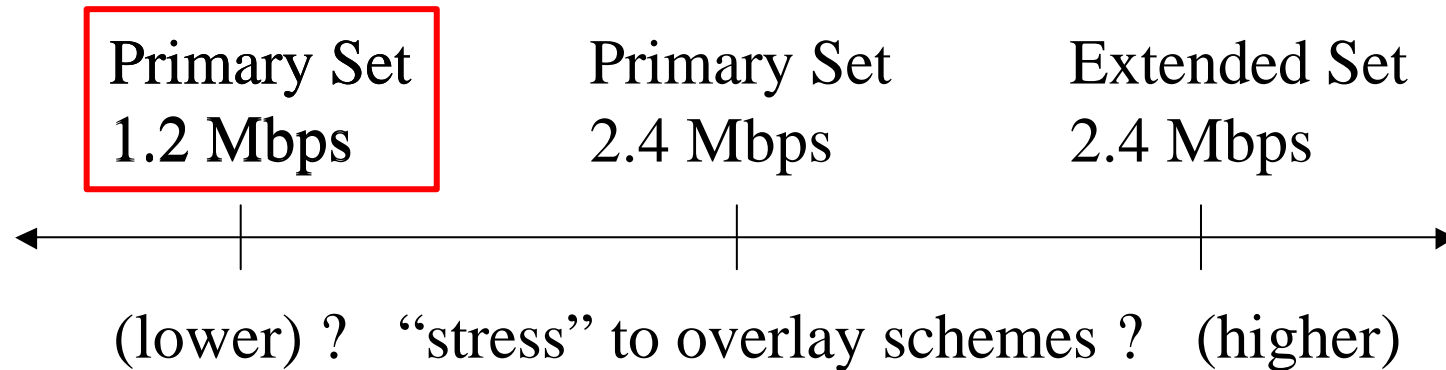
r Heterogeneity of host set

- m *Primary Set*: 13 university hosts in U.S. and Canada
- m *Extended Set*: 20 hosts, which includes hosts in *Primary Set*, Europe, `Asia, and behind ADSL

r Source rate

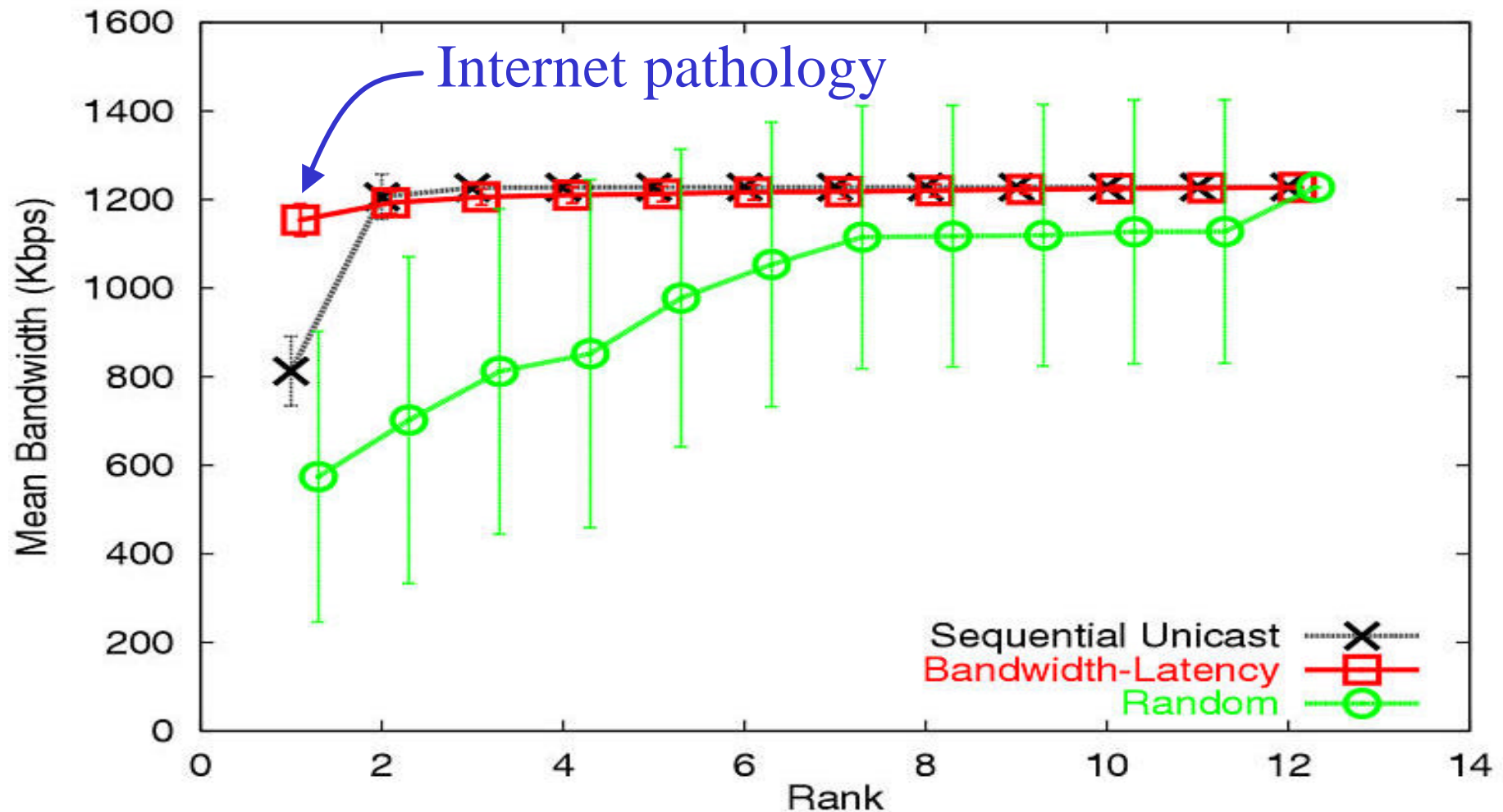
- m Fewer Internet paths can sustain higher source rate

Three Scenarios Considered



- r Does ESM work in different scenarios?
- r How do different schemes perform under various scenarios?

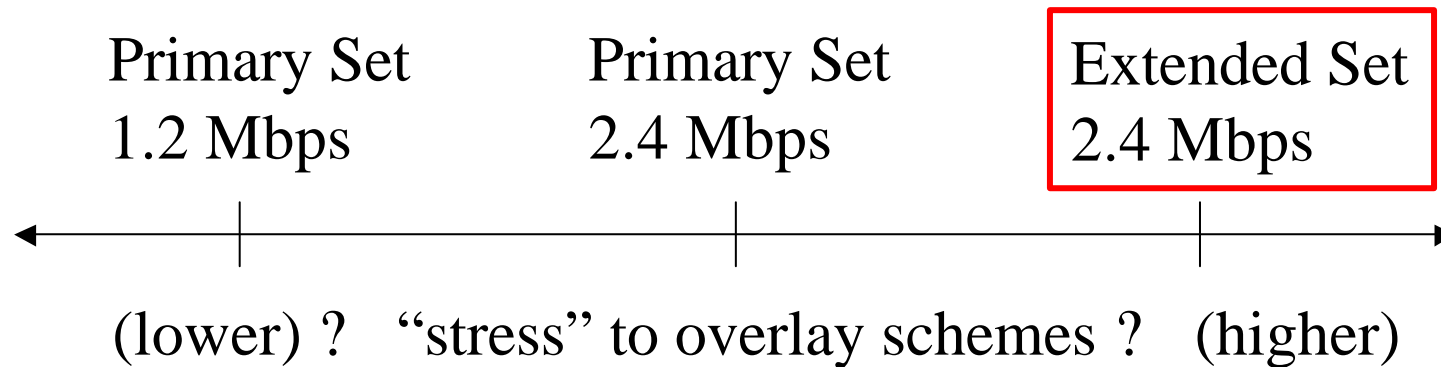
BW, Primary Set, 1.2 Mbps



Naïve scheme performs poorly even in a less “stressful” scenario

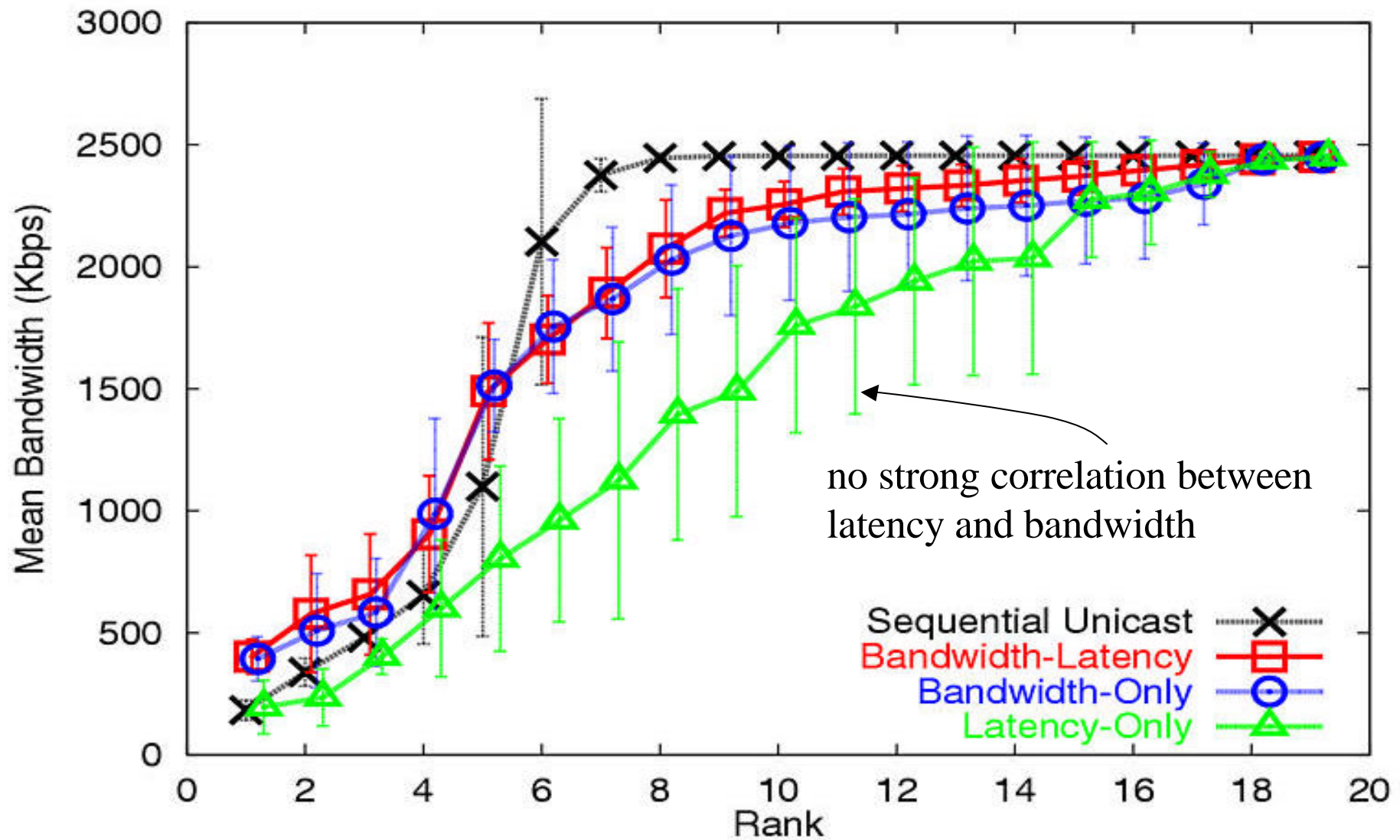
RTT results show similar trend

Scenarios Considered



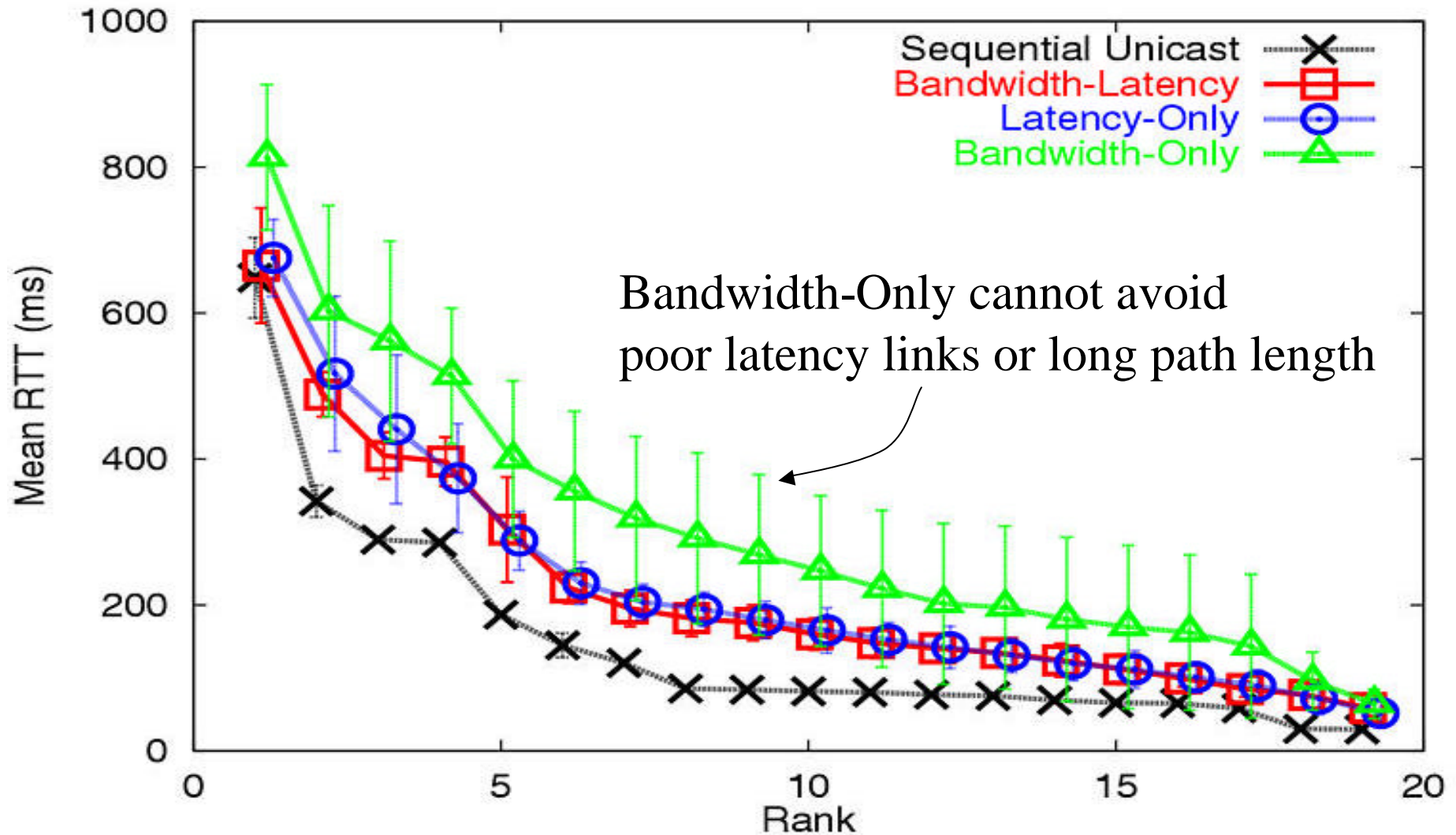
- r Does an overlay approach continue to work under a more “stressful” scenario?
- r Is it sufficient to consider just a single metric?
 - m *Bandwidth-Only, Latency-Only*

BW, Extended Set, 2.4 Mbps



Optimizing only for latency has poor bandwidth performance₃₉

RTT, Extended Set, 2.4Mbps



Optimizing only for bandwidth has poor latency performance

Summary so far...

- r For best application performance: adapt dynamically to both latency and bandwidth metrics
- r *Bandwidth-Latency* performs comparably to IP Multicast (*Sequential-Unicast*)
- r What is the network cost and overhead?

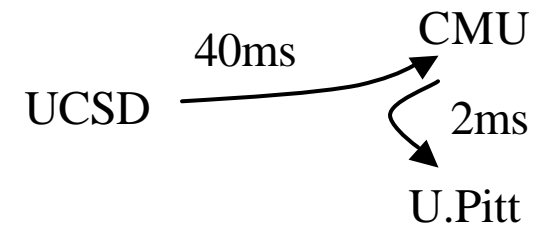
Resource Usage (RU)

Captures consumption of network resource of overlay tree

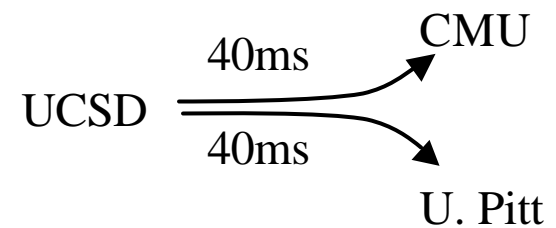
- Overlay link RU = propagation delay
- Tree RU = sum of link RU

Scenario: Primary Set, 1.2 Mbps
(normalized to IP Multicast RU)

IP Multicast	1.0
Bandwidth-Latency	1.49
Random	2.24
Naïve Unicast	2.62



Efficient (RU = 42ms)



Inefficient (RU = 80ms)

Protocol Overhead

$$\text{Protocol overhead} = \frac{\text{total non-data traffic (in bytes)}}{\text{total data traffic (in bytes)}}$$

- r Results: Primary Set, 1.2 Mbps
 - m Average overhead = 10.8%
 - m 92.2% of overhead is due to bandwidth probe
- r Current scheme employs active probing for available bandwidth
 - m Simple heuristics to eliminate unnecessary probes
 - m Focus of our current research

Contribution

- r First detailed **Internet evaluation** to show the feasibility of **End System Multicast** architecture
 - m Study in context of a/v conferencing
 - m Performance comparable to IP Multicast

- r Impact of metrics on overlay performance
 - m For best performance: **use both latency and bandwidth**

- r More info: <http://www.cs.cmu.edu/~narada>

Discussion (1)

- r Peer-to-peer versus proxy based architecture

Peer-to-peer	Proxy based
Distributed	Share network information , history across groups
Scalable to number of groups	Long term connection, more stable
	Manage resource allocation among groups

Discussion (2)

- r Multipath framework where each recipient gets data from the source along multiple paths, with a fraction of the data flowing along any given path
 - m Any individual path doesn't radically affect overall performance
 - m Receive data while monitoring

Discussion (3)

- r Rigorous change detection algorithm
 - m On the Constancy of Internet Path Properties
 - www.aciri.org/vern/imw2001-papers/38.ps.gz
 - m Idea of more formally identifying changes
 - www.variation.com/cpa/tech/changepoint.html

14. Fault-tolerant replication management in large-scale distributed storage systems

Richard Golding

*Storage Systems Program, Hewlett Packard
Labs*

golding@hpl.hp.com

Elizabeth Borowsky

*Computer Science Dept., Boston
College*

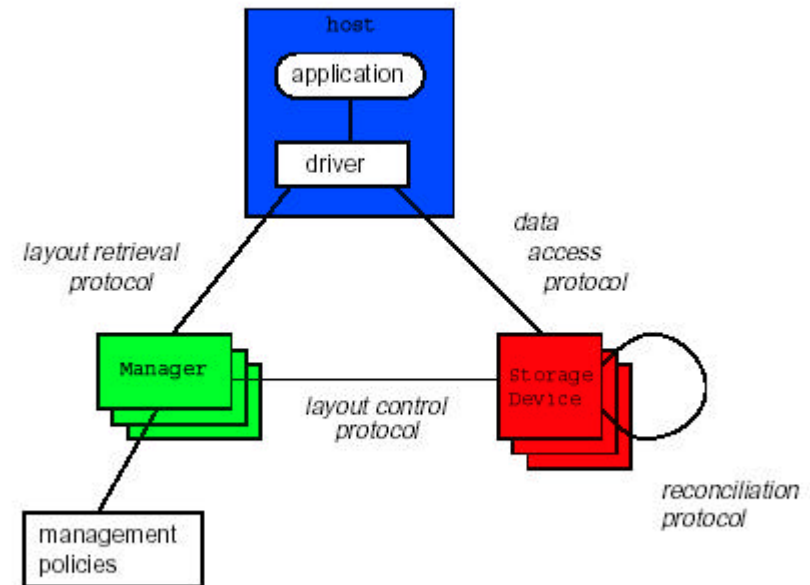
borowsky@cs.bc.edu

Introduction

- r **Palladio** - solution for detecting, handling, and recovering from both small- and large-scale failures in a distributed storage system.
- r **Palladio** - provides virtualized data storage services to applications via set of **virtual stores**, which are structured as a logical array of bytes into which applications can write and read data. The store's **layout** maps each byte in its address space to an address on one or more devices.
- r **Palladio** - **storage devices** take an active role in the recovery of the stores they are part of. **Managers** keep track of the virtual stores in the system, coordinating changes to their layout and handling recovery from failure.

- **Provide robust read and write access to data in virtual stores.**

- Atomic and serialized read and write access.
- Detect and recover from failure.
- Accommodate layout changes.



Palladio implementation structure

Entities

Hosts
Stores
Managers
Management policies

Protocols

Layout Retrieval protocol
Data Access protocol
Reconciliation protocol
Layout Control protocol

Protocols

Access protocol allows hosts to read and write data on a storage device as long as there are no failures or layout changes for the virtual store. It must provide serialized, atomic writes that can span multiple devices.

Layout retrieval protocol allows hosts to obtain the current layout of a virtual store — the mapping from the virtual store's address space onto the devices that store parts of it.

Reconciliation protocol runs between pairs of devices to bring them back to consistency after a failure.

Layout control protocol runs between managers and devices — maintains consensus about the layout and failure status of the devices, and in doing so coordinates the other three protocols.

Layout Control Protocol

The layout control protocol tries to maintain agreement between a store's manager and the storage devices that hold the store.

- The layout of data onto storage devices
- The identity of the store's *active manager*.

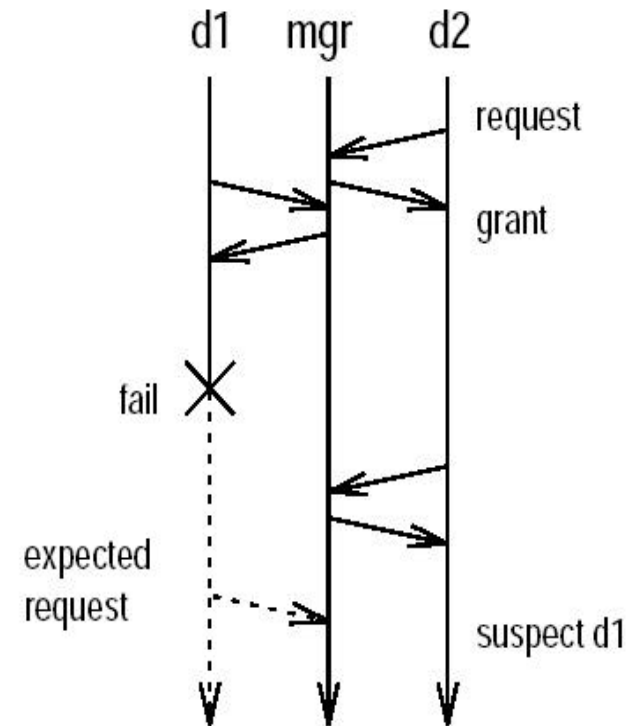


The notion of ***epochs***

- The layout and manager are fixed during each epoch
- Epochs are numbered
- Epoch transitions
- Device leases acquisition and renewal
- Device leases used to detect possible failure.

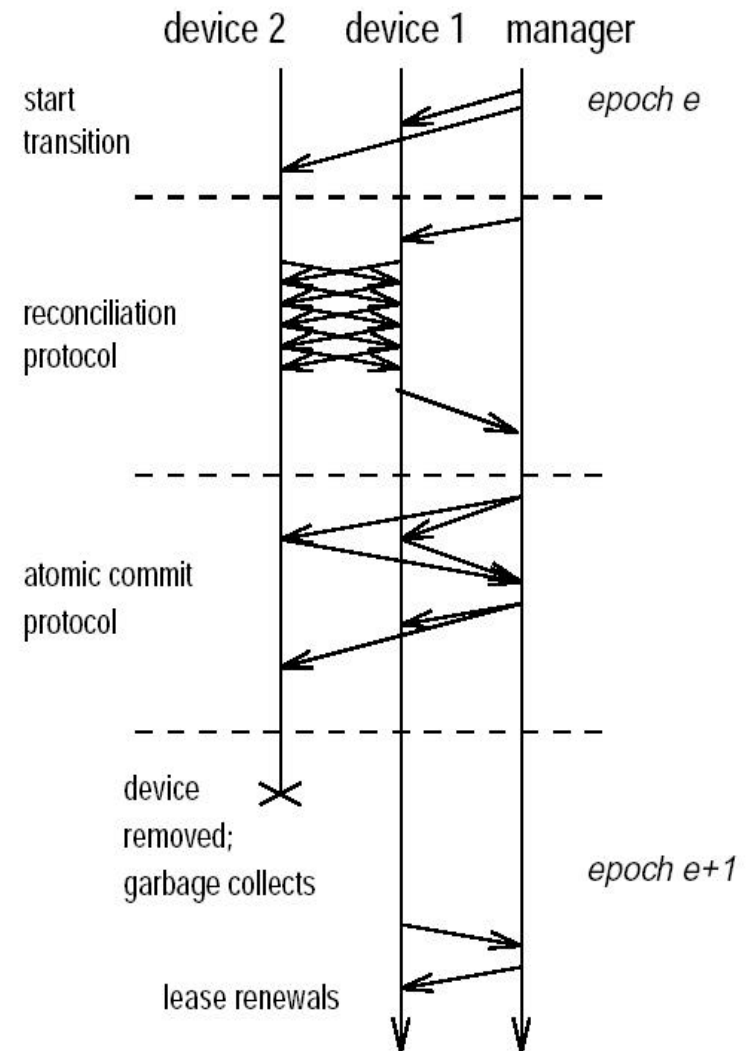
Operation during an epoch

- The manager has quorum and coverage of devices.
- Periodic lease renewal
 - »In case a device fails to report and try to renew its lease, the manager considers it failed
 - »In case the manager fails to renew the lease, the device considers the manager failed and starts a *manager recovery sequence*
- When the manager loses quorum or coverage the epoch ends and a state of epoch transition is entered.



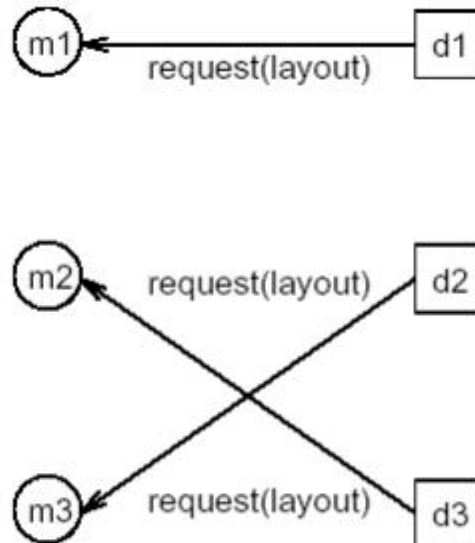
Epoch transition

- Transaction initiation
- Reconciliation
- Transaction commitment
- Garbage collection

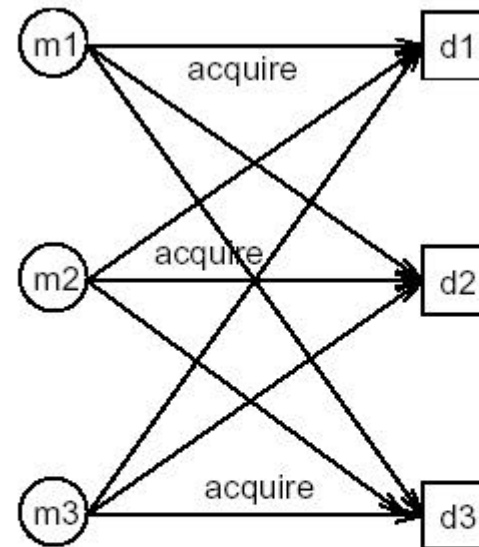


The recovery sequence

- Initiation - querying a recovery manager with the current layout and epoch number



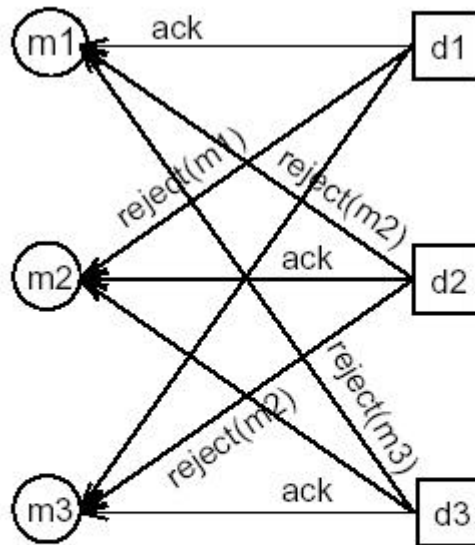
Step 1: devices request recovery



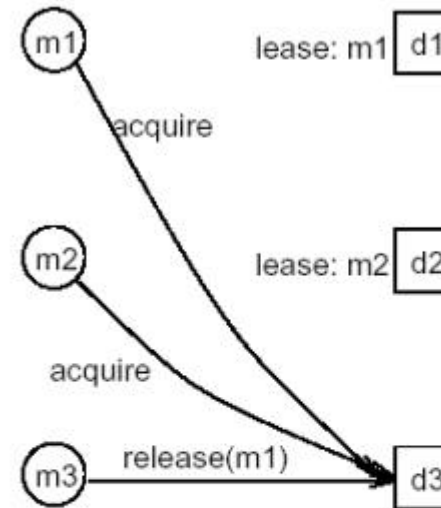
Step 2: managers try to acquire devices

The recovery sequence (continued)

- Contention - managers struggle to obtain quorum and coverage and to become active managers for the store - (recovery leases, acks and rejections)



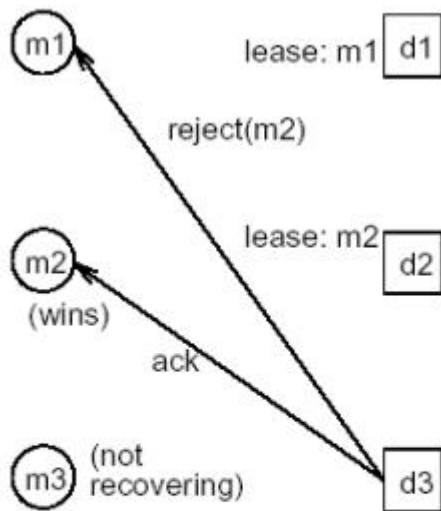
Step 3: devices accept first lease, reject others



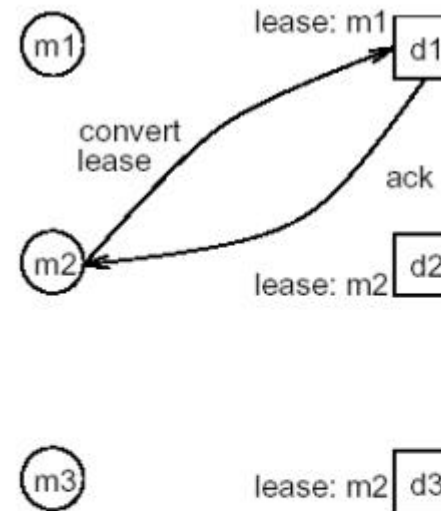
Step 4: least precedence manager gives up, others try again

The recovery sequence (continued)

- Completion - setting correct recovery leases & starting epoch transition
- Failure - failure of devices and managers during recovery



Step 5: remaining device accepts one lease, rejects other



Step 6: manager 2 wins contention, converts leases on other devices

Extensions

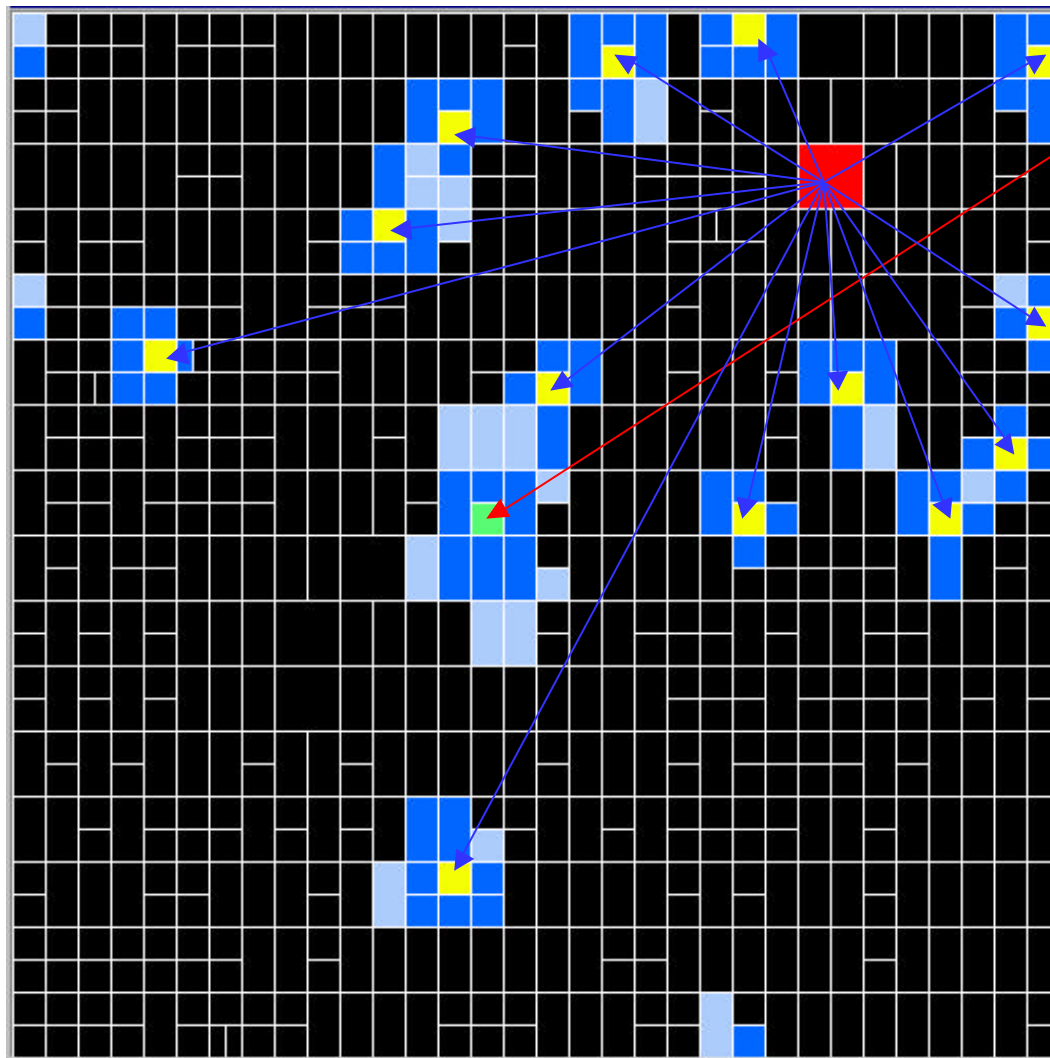
- *Single manager v.s. Multiple managers*
- *Whole devices v.s. Device parts (chunks)*
- *Reintegrating devices*
- *Synchrony model (future)*
- *Failure suspects (future)*

Conclusions & recap

Palladio - Replication management system featuring

- » *Modular protocol design*
- » *Active device participation*
- » *Distributed management function*
- » *Coverage and quorum condition*

Application example



Very popular content

Popularity indicator

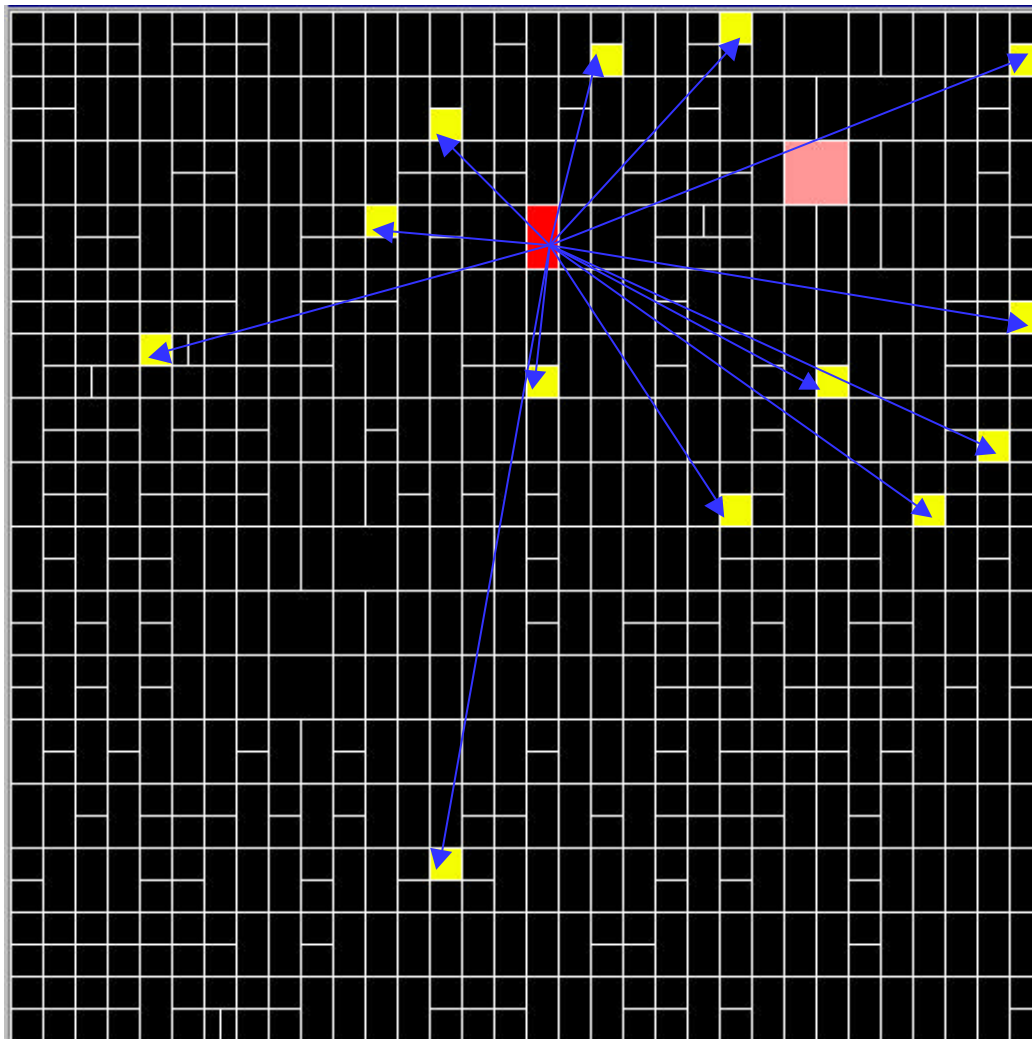
Manager node

$ID = \text{hash}\langle \text{FileID}, \text{MGR} \rangle$

Storage nodes

$ID = \text{hash}\langle \text{FileID}, \text{STR}, n \rangle$

Application example - benefits



■ Stable manager node

■ Stable storage nodes

- Self-manageable storage
- Increased availability
- Popularity is hard to fake
- Less per node load
- Could be applied recursively (?)

Wrapup discussion questions (1):

- r What is a peer-peer network (what is not a peer-to-peer network?). Necessary:
 - m every node is designed to (but may not by user choice) *provide some service* that helps other nodes in the network get service
 - m no 1-N service providing
 - m *each node potentially has the same responsibility, functionality (maybe nodes can be polymorphic)*
 - corollary: by design, nothing (functionally) prevents two nodes from communicating directly
 - m some applications (e.g., Napster) are a mix of peer-peer and centralized (lookup is centralized, file service is peer-peer) [recursive def. of peer-peer]
 - m *(logical connectivity rather than physical connectivity) routing will depend on service and data*

Overlays?

- r What is the relationship between peer-peer and application overlay networks?
 - m Peer-peer and application overlays are different things. It is possible for an application level overlay to be built using peer-peer (or vice versa) but not always necessary
 - m Overlay: in a wired net: if two nodes can communicate in the overlay using a path that is *not* the path the network level routing would define for them. Logical network on top of underlying network
 - source routing?
 - m Wireless ad hoc nets – what commonality is there REALLY?

Wrapup discussion questions (2):

- r What were the best p2p idea
- r Vote now (and should it be a secret ballot
usign Eternity✍)

Wrapup discussion questions (3):

- r Is ad hoc networking a peer-peer application?
 - m Yes (30-1)
- r Why peer-peer over client-server?
 - m A well-deigned p2p provides better “scaability”
- r Why client-server of peer-peer
 - m peer-peer is harder to make reliable
 - m availability different from client-server (p2p is more often at least partially “up”)
 - m more trust is required
- r If all music were free in the future (and organized), would we have peer-peer.
 - m Is there another app: ad hoc networking, any copyrighted data, peer-peer sensor data gathering and retrieval, simulation
- r Evolution #101 – what can we learn about systems?