

## Lecture Coverage

1. Survey of Approaches
2. Napster/Gnutella
3. Distributed Hash Table Based Systems
  1. Chord
  2. Can
  3. Freenet
4. Semantic (keyword search)
5. Topology Aware Routing and Search

### **Assumed Knowledge**

**Public Key Encryption**

**Hashing Algorithms**

# The Search/Discovery Problem

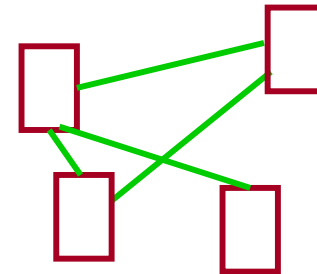
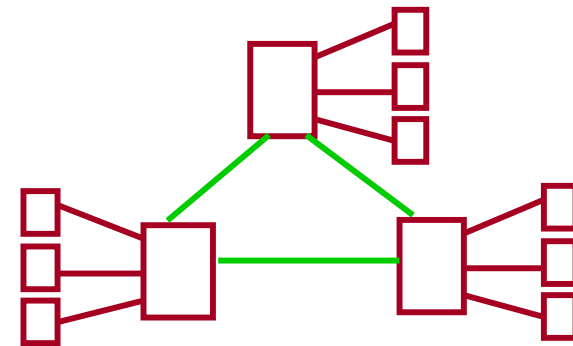
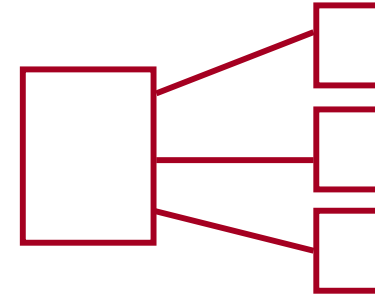
- “Location Resolution”  
Given an object (might be name, attribute, or even content)  
Return a channel to a node (peer) that has that object
- Approaches:
  - Centralized Index (Napster)
  - Broadcast information to be resolved (Gnutella)
  - Search Strategies
  - **Distributed Hashing**

# Design Goals

- Scalability
- **Low latency (efficient resolution)**
- Load balancing
- Completely distributed/self-organizing
- Robust
- Deployable
- Simple

## Spectrum of “Purity”

- Hybrid
  - Centralized index, P2P file storage and transfer
- Super-peer
  - A “pure” network of “hybrid” clusters
- Pure
  - functionality completely distributed



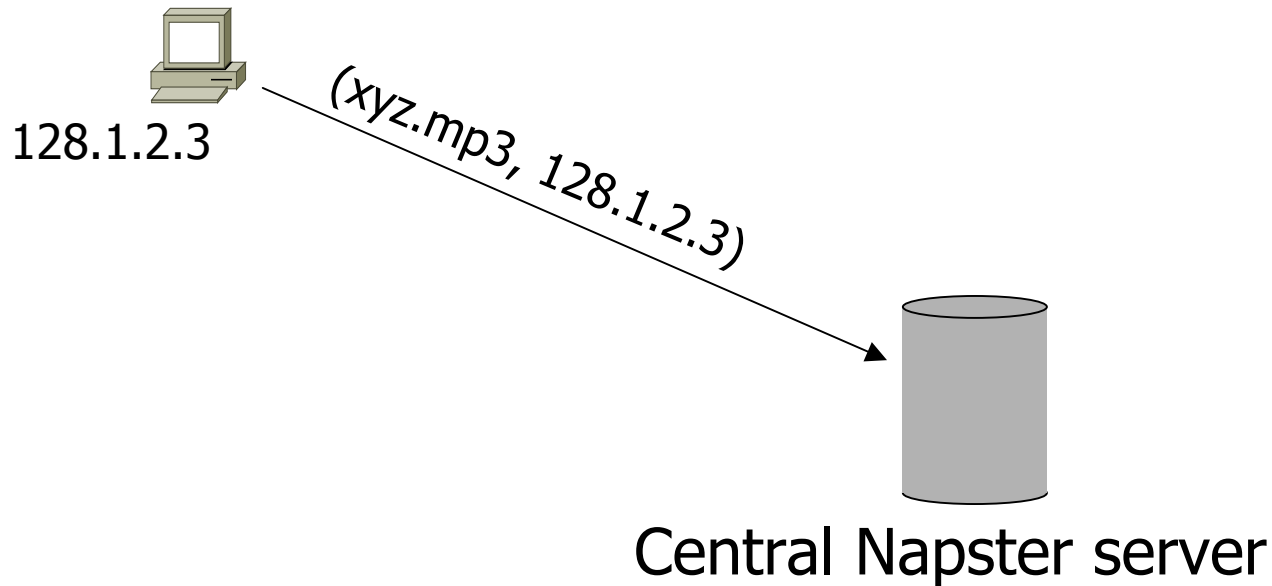
## Metrics

- Cost (aggregate)
  - Bandwidth
  - Processing Power
- Quality of Results
  - Number of results
  - Satisfaction (true if # results  $\geq X$ , false otherwise)
  - Time to satisfaction

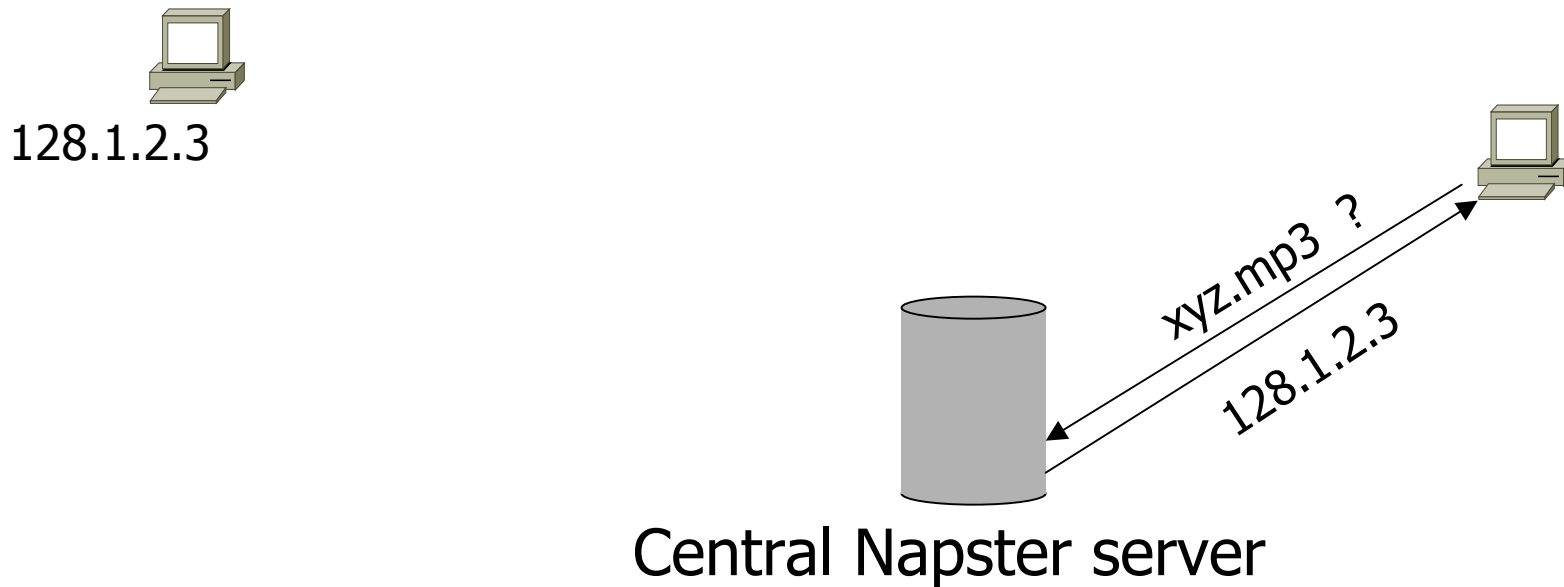
# P2P File-sharing

- Napster
  - Decentralized storage of actual content
    - transfer content directly from one peer (client) to another
  - Centralized index and search
- Gnutella
  - Like Napster, with decentralized indexing
  - Search via flooding
  - Direct download

# Napster

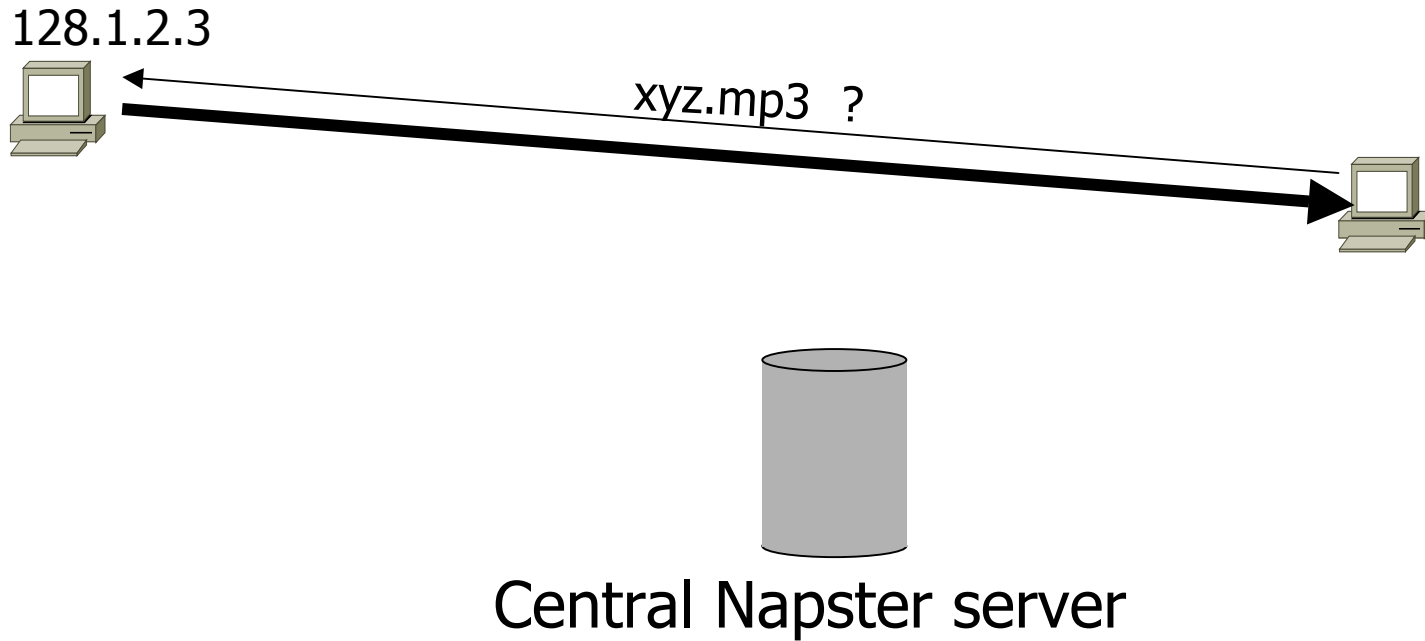


# Napster

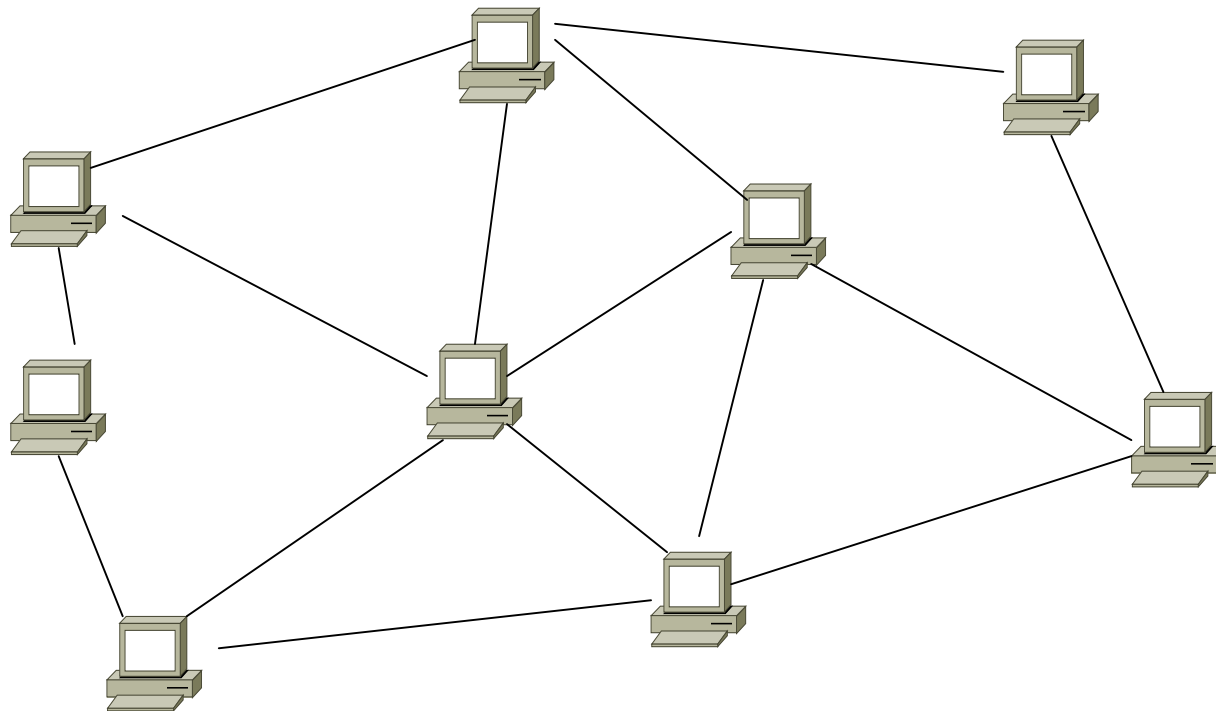




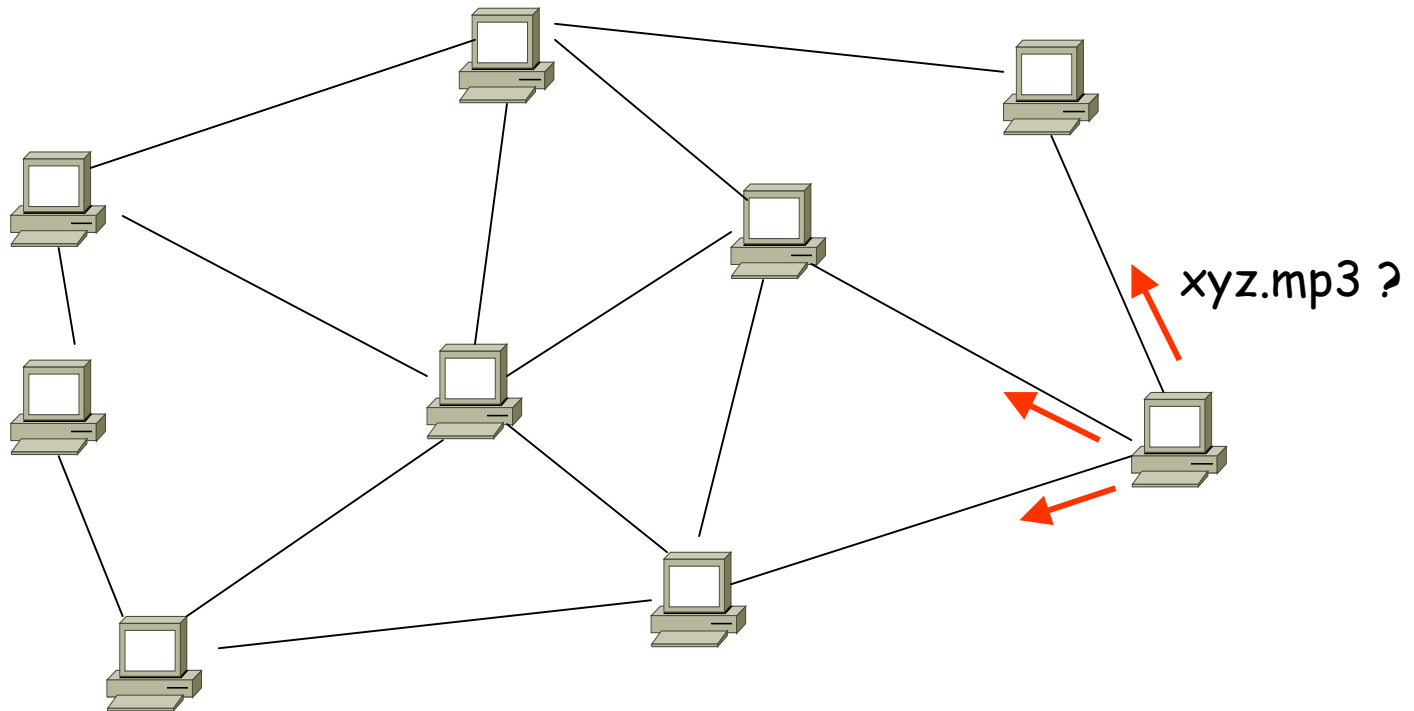
# Napster



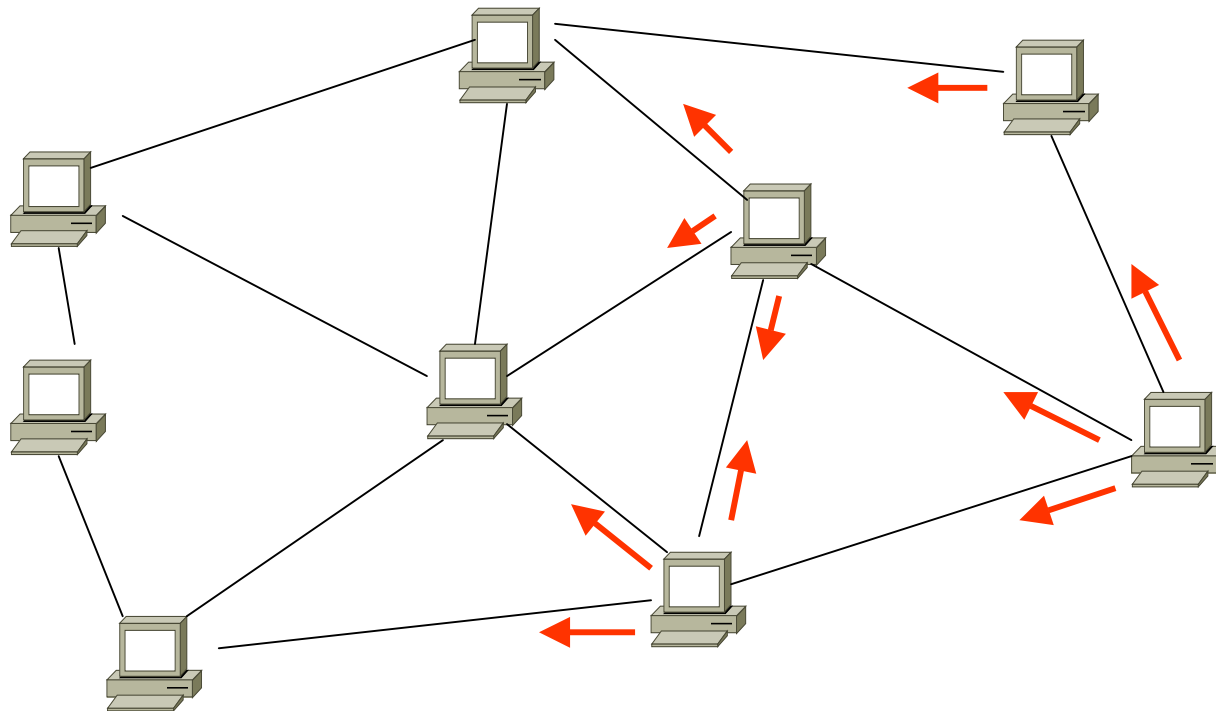
# Gnutella



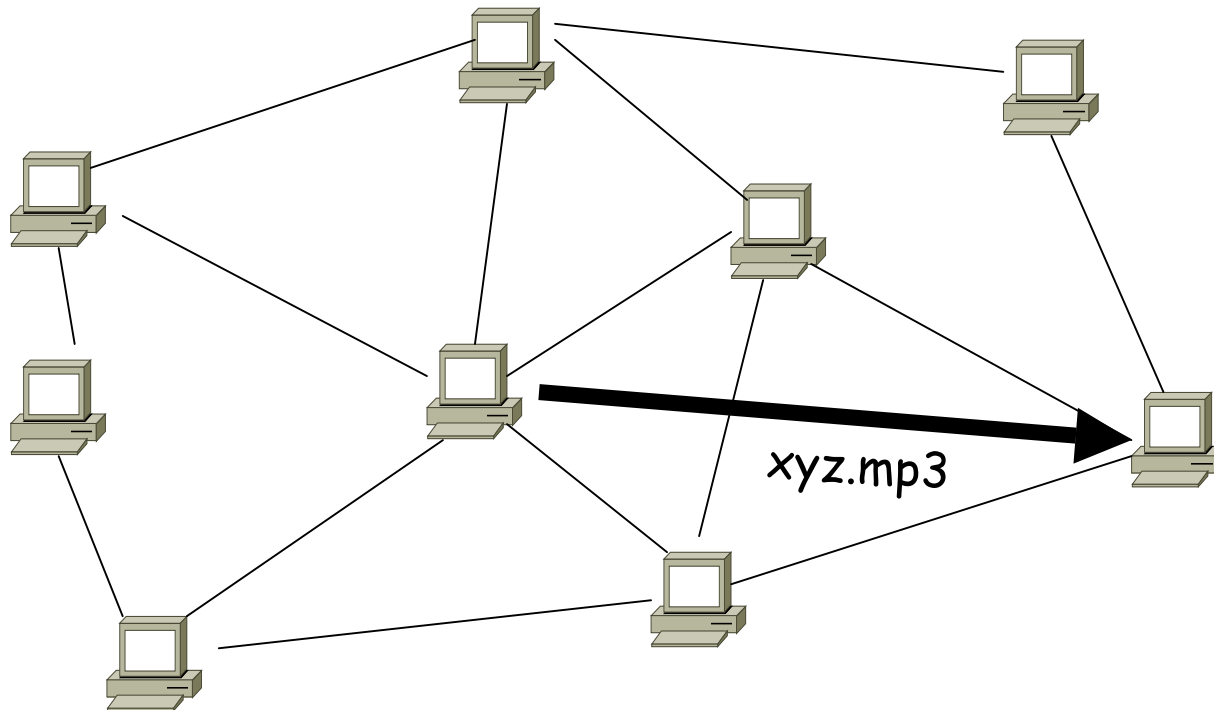
# Gnutella



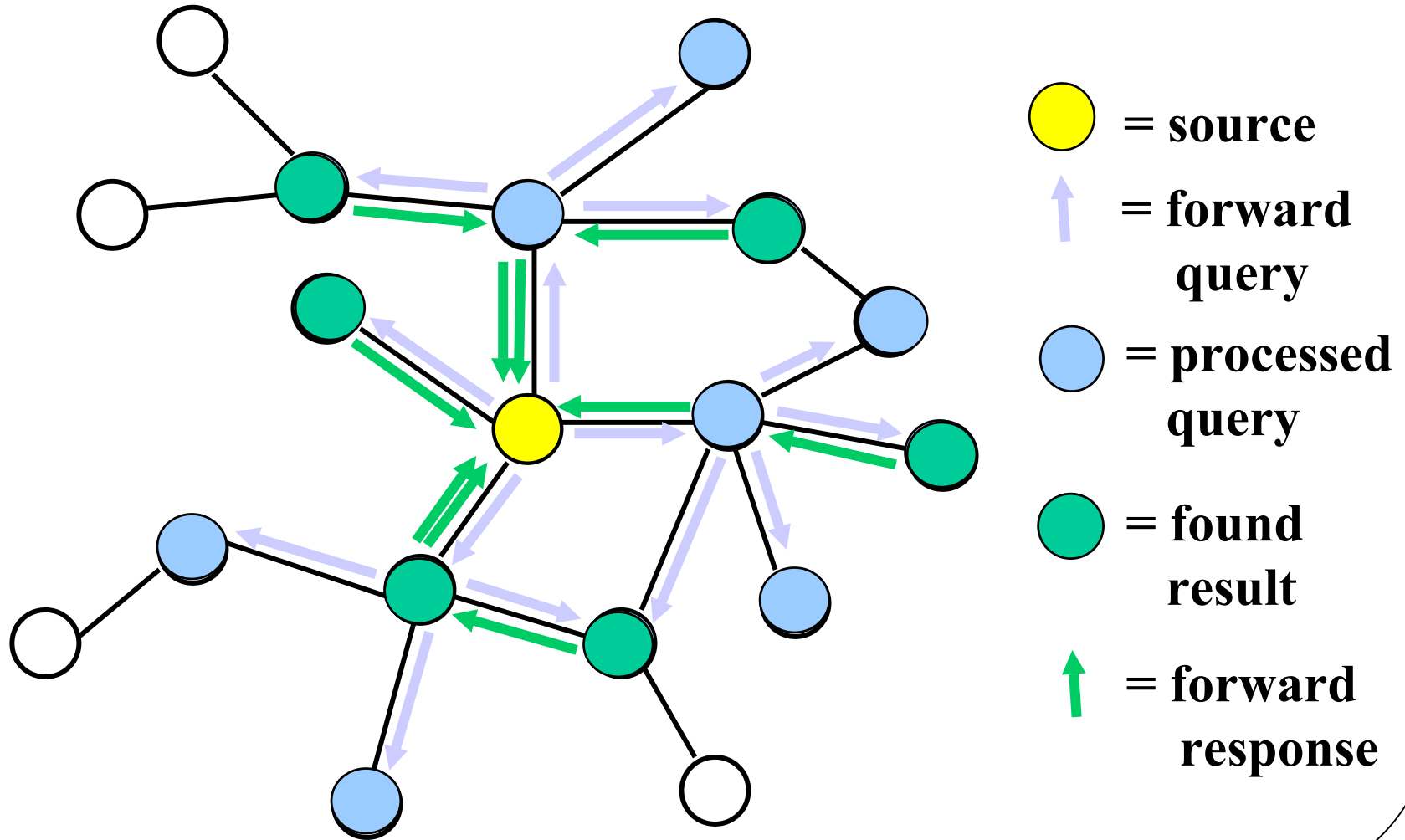
# Gnutella



# Gnutella

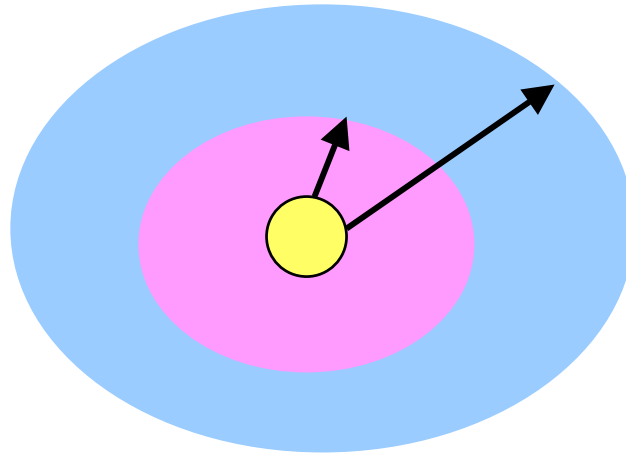


## Current Techniques: Gnutella: **Breadth-First Search (BFS)**

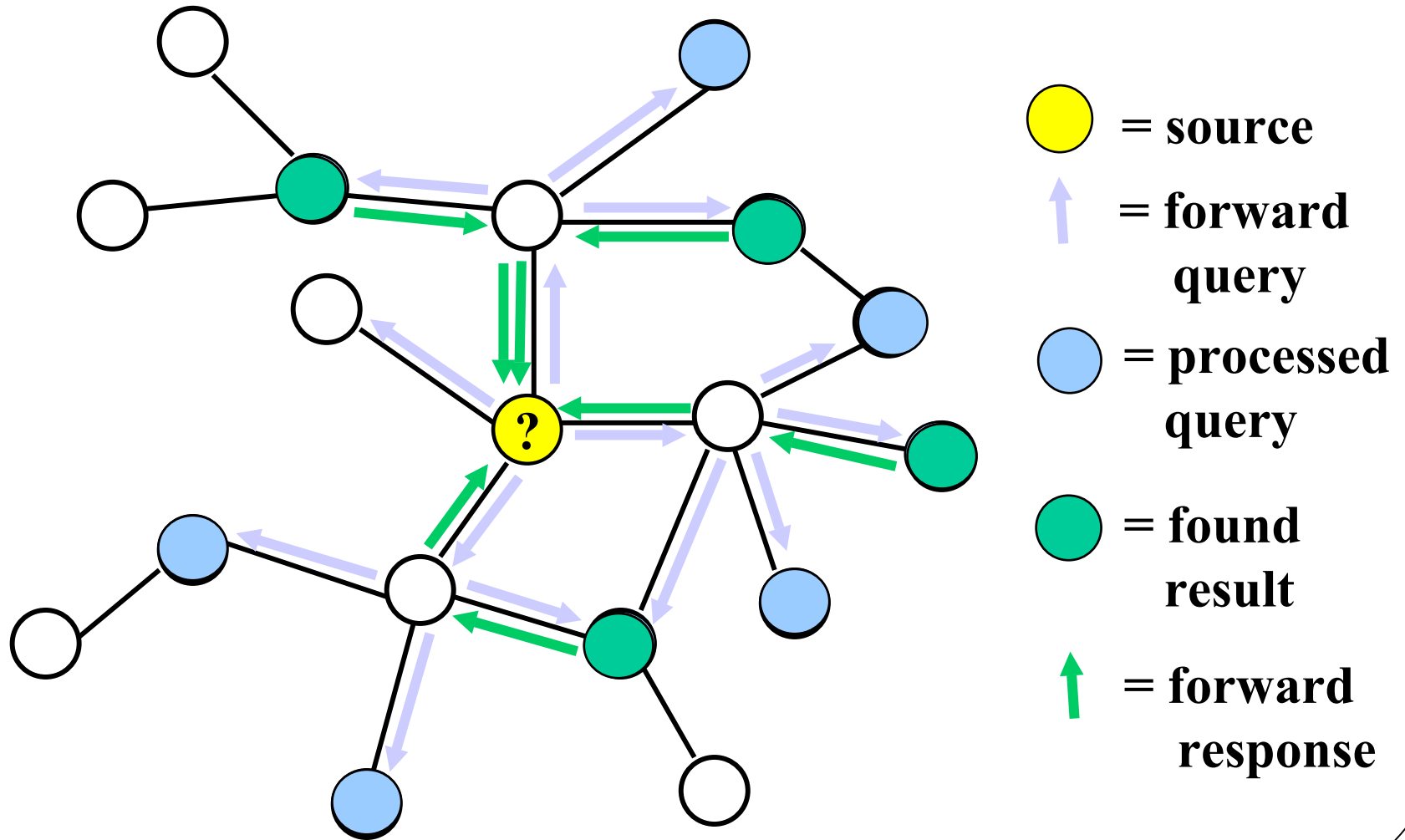


## Iterative Deepening

- Interested in **satisfaction**, not # of results
- BFS returns “too many” results expensive
- Iterative Deepening: common technique to reduce the cost of BFS
  - Intuition: A search at a small depth is much cheaper than at a larger depth



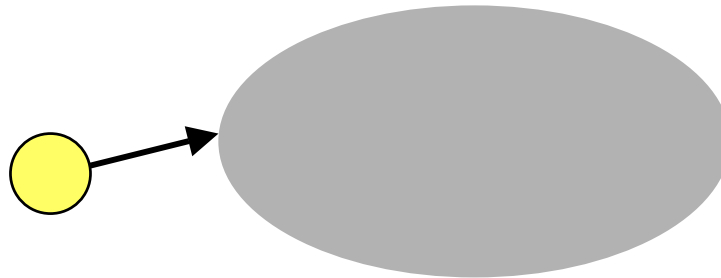
## Iterative Deepening





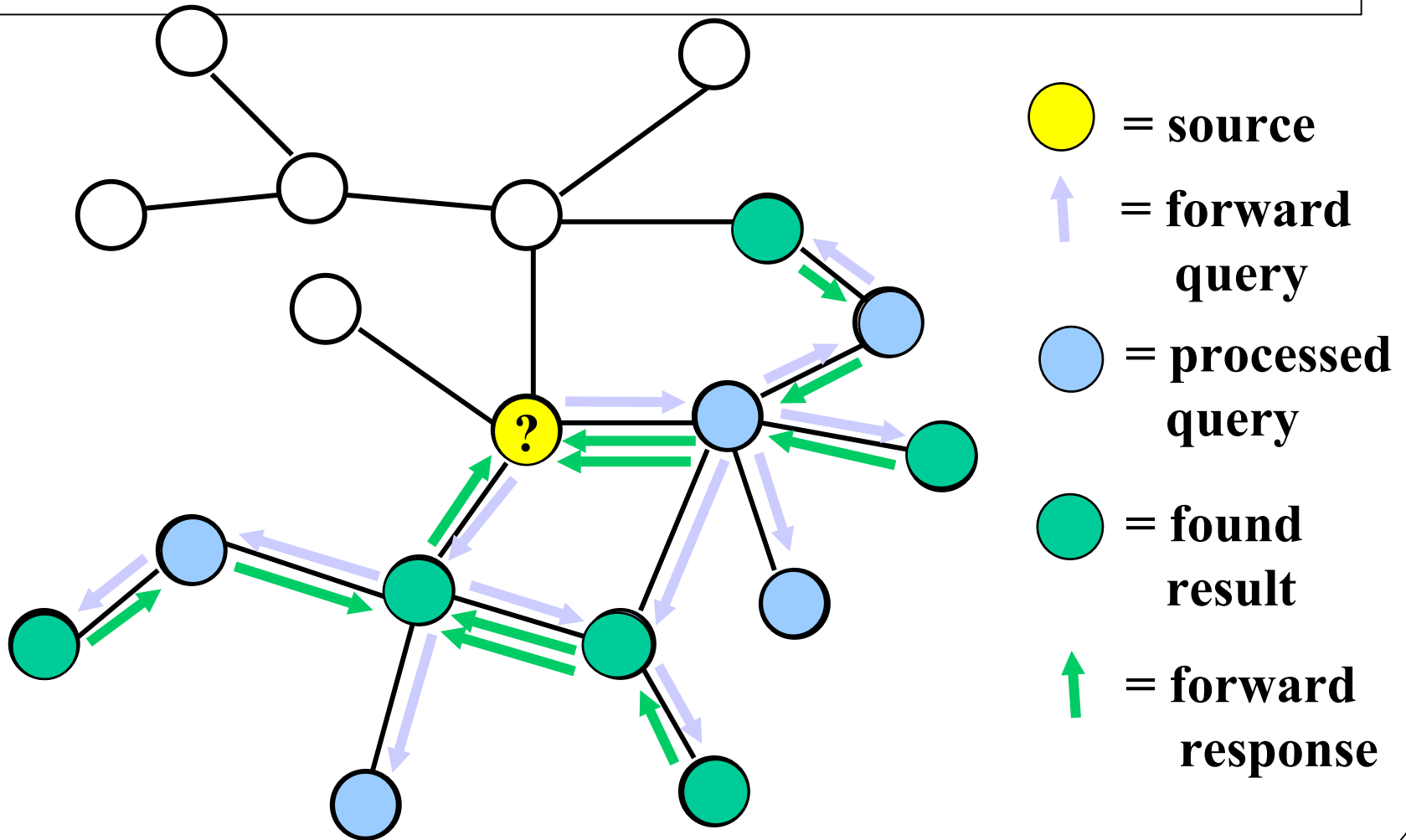
## Directed BFS

- Sends query to a subset of neighbors
- Maintains statistics on neighbors
  - E.g., ping latency, history of number of results
- Chooses subset intelligently (via heuristics), to maximize quality of results
  - E.g., Neighbors with shortest message queue, since long message queue implies neighbor is saturated/dead



# Search/Discovery (and Insertion) in Distributed Systems

## Directed BFS

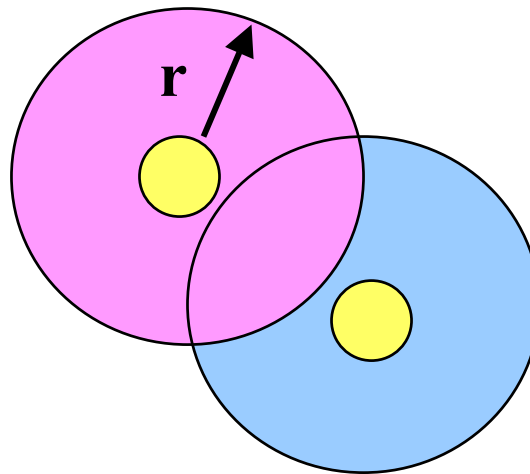


# Directed BFS: Heuristics

RAND	(Random)
RES	Returned greatest # results in past
TIME	Had shortest avg. time to satisfaction in past
HOPS	Had smallest avg. # hops for response messages in past
MSG	Sent our client greatest # of messages
QLEN	Shortest message queue
DEG	Highest degree

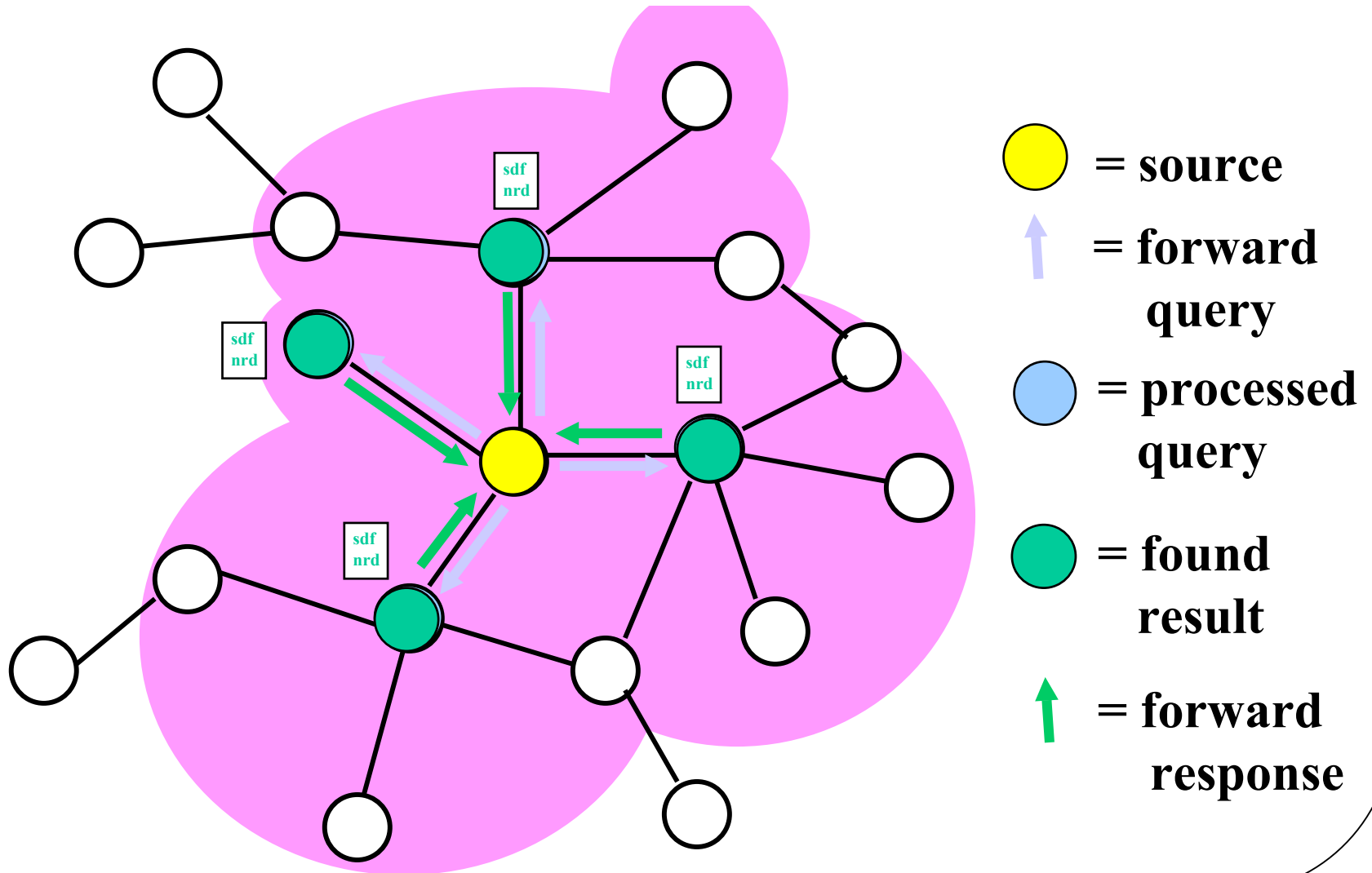
## Local Indices

- Each node maintains index over other nodes' collections
  - $r$  is the **radius** of the index
  - Index covers all nodes within  $r$  hops away
- Can process query at fewer nodes, but get just as many results back

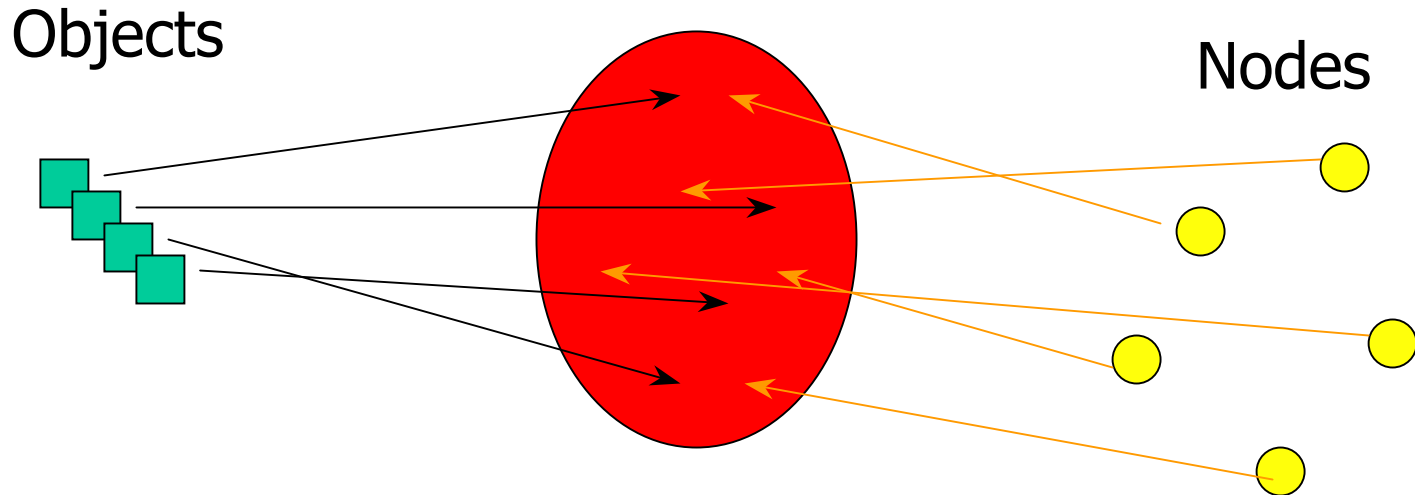


# Search/Discovery (and Insertion) in Distributed Systems

## Local Indices (r=1)

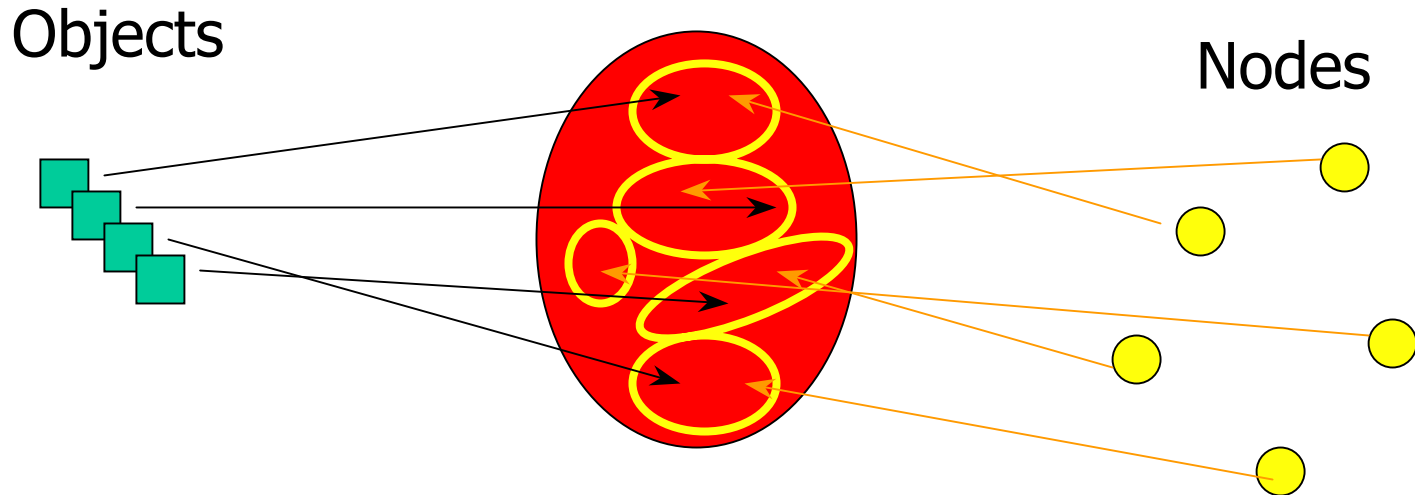


## Distributed Hashing — General Approach



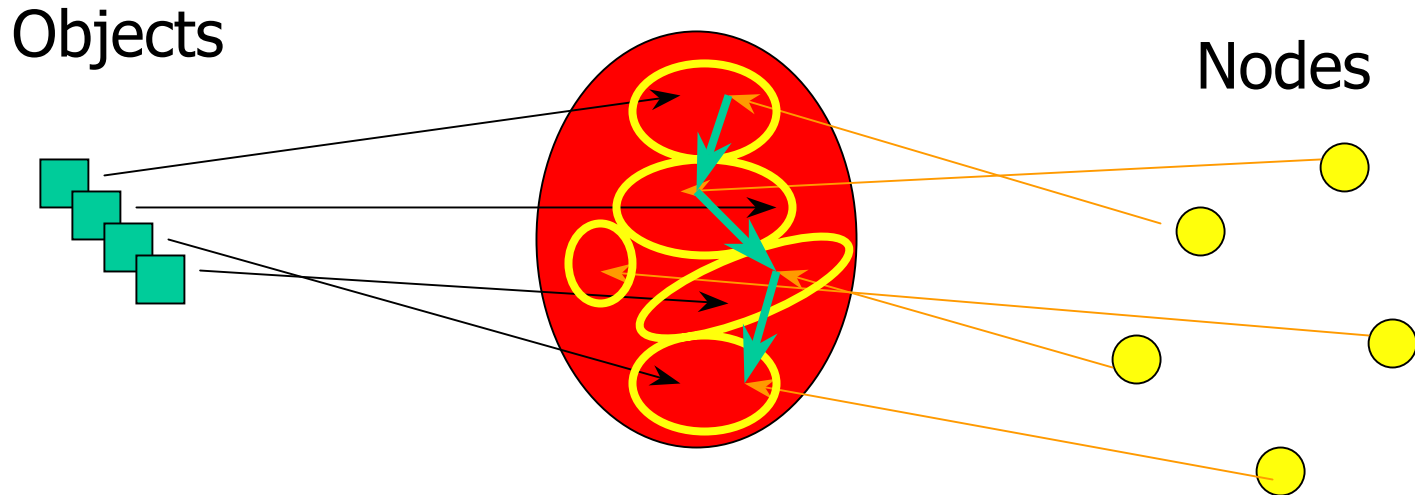
1. Map both objects and nodes into some topology (“id space”)

## Distributed Hashing — General Approach



1. Map both objects and nodes into some topology (“id space”)
2. Each node “owns” some neighborhood in the topology, has channel to some neighbors

## Distributed Hashing — General Approach

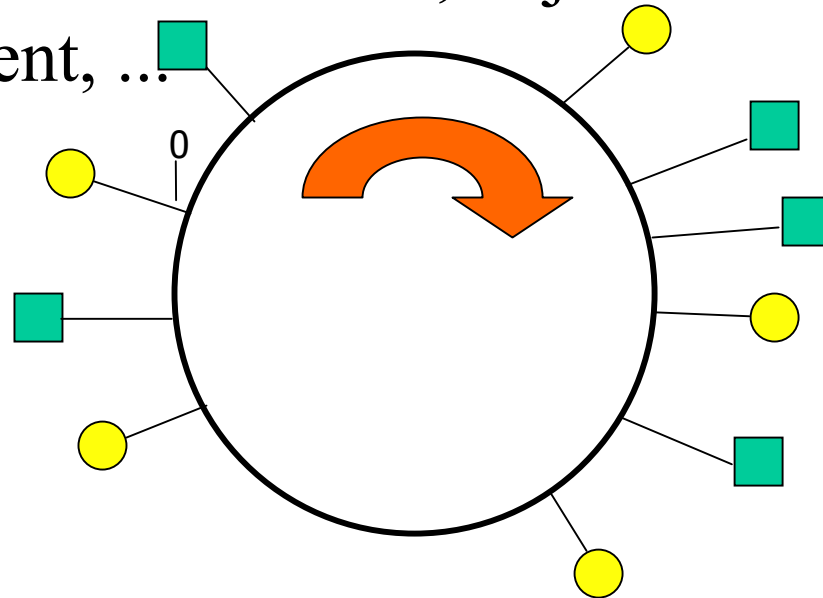


1. Map both objects and nodes into some topology (“id space”)
2. Each node “owns” some neighborhood in the topology, has channel to some neighbors
3. Topological structure lets query be routed to the “owner” of a given point



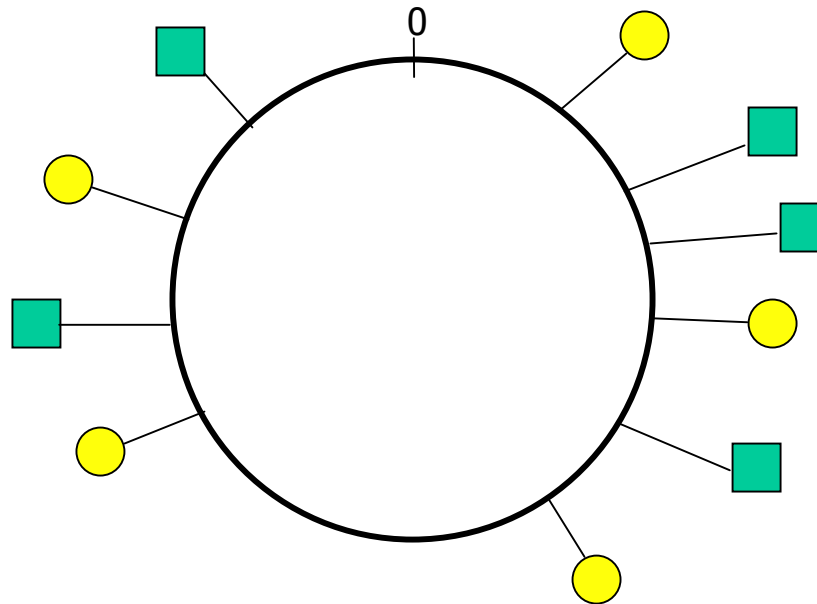
## Chord - Basic Idea

- Topology is a ring of ordered, fixed-size IDs (say 32 bits)
  - Node ID based on IP address, object ID based on name, content, ...



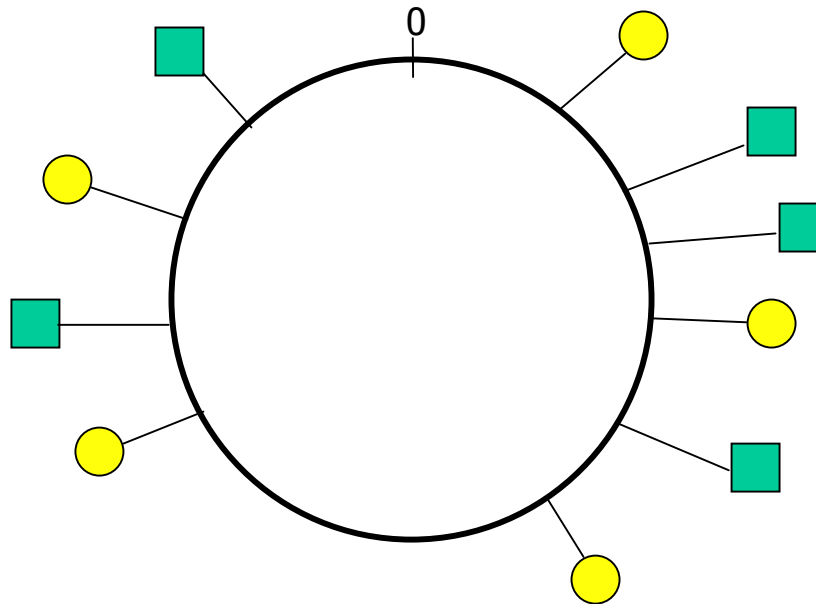
## Chord - Basic Idea

- Nodes “own” the part of the ID space between their ID and their predecessor’s ID.



## Chord - Basic Idea

- Each node has a channel to its successors at distances  $1, 2, 4, 8, 16, \dots, 2^{(m-1)}$ 
  - where  $m = \log_2$  of the ring size (32 in this case)



## Chord: Resolution

- Get ID of desired object
- Find the last node whose ID is LESS than the desired ID
  - Look in finger table to find farthest-away neighbor whose ID is LESS than the desired ID
  - Ask it for somebody closer
- That node's successor is the “owner” of the object

## Chord: Performance

- Resolution:  $O(\log N)$
- Joining:  $O(\log^2 N)$  = find all your neighbors
  - Doesn't count cost of “moving” objects that have a new owner
- Stability: provable

## CAN: Basic Idea

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

## CAN: Basic Idea

- Topology is an N-dimensional torus
  - N=2 for simple examples
- Each node is responsible for a subrange in each dimension
  - Space is partitioned among all nodes
- Route via neighbors -- move in direction of destination

## CAN: simple example



1

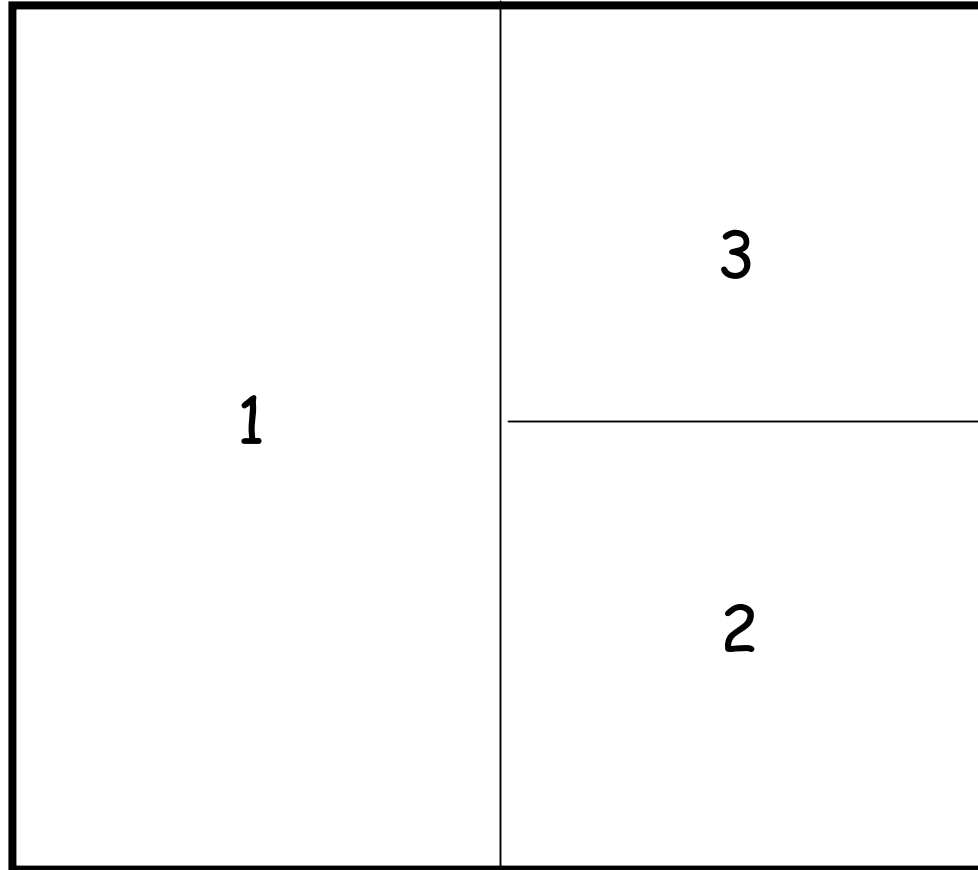


## CAN: simple example

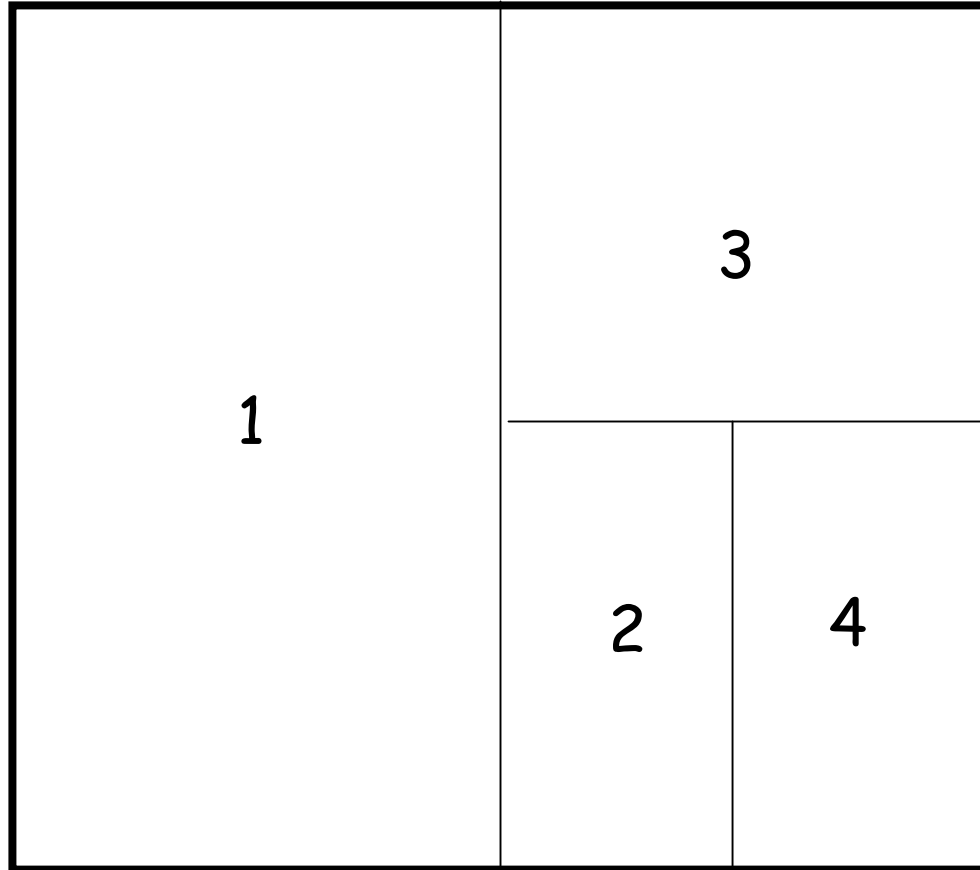
1

2

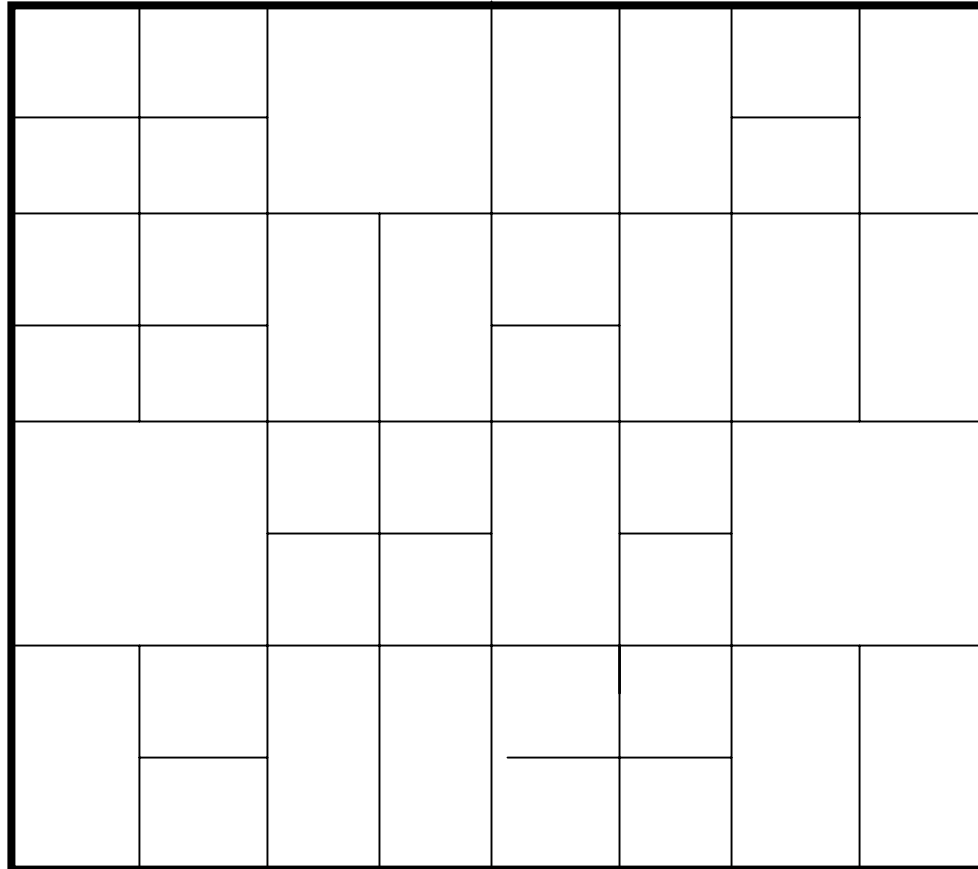
## CAN: simple example



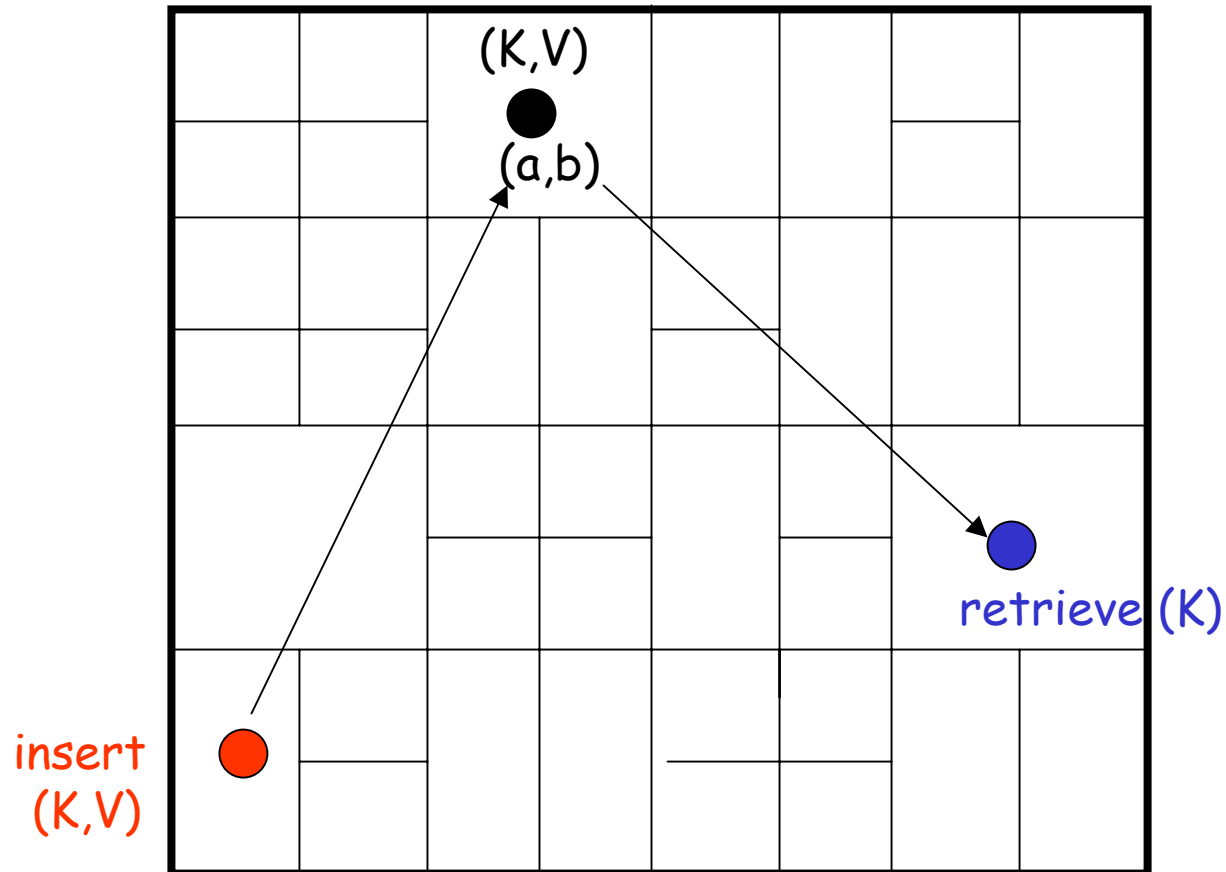
## CAN: simple example



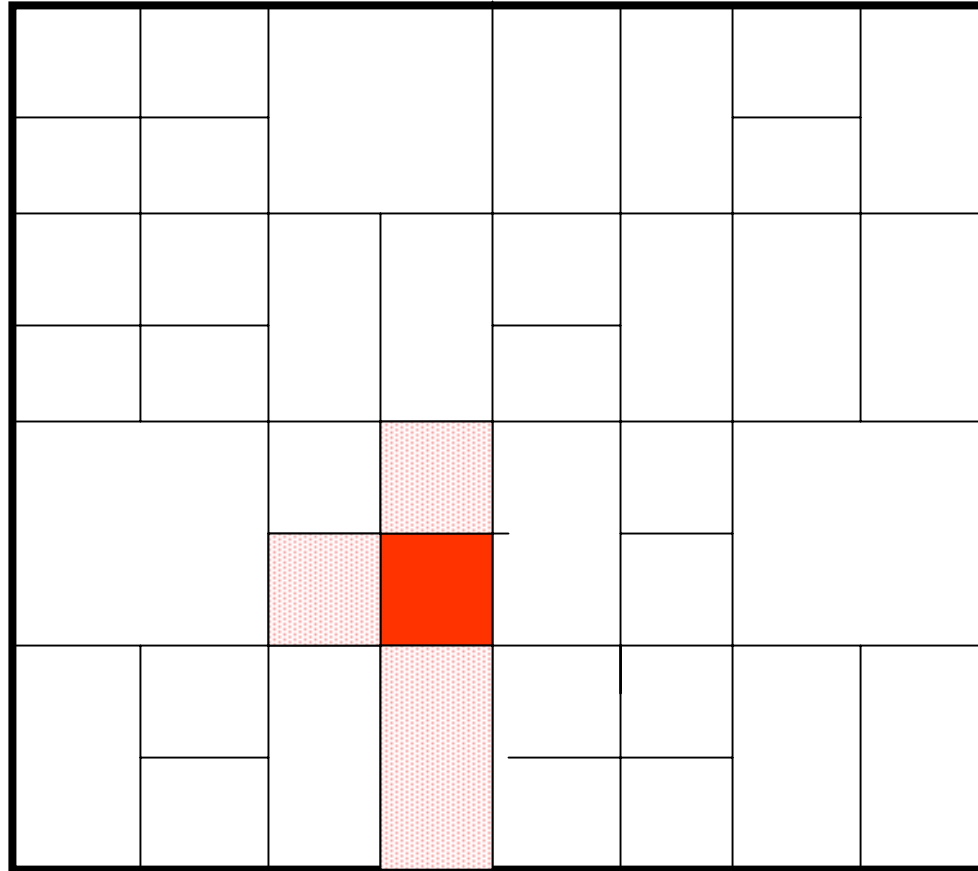
## CAN: simple example



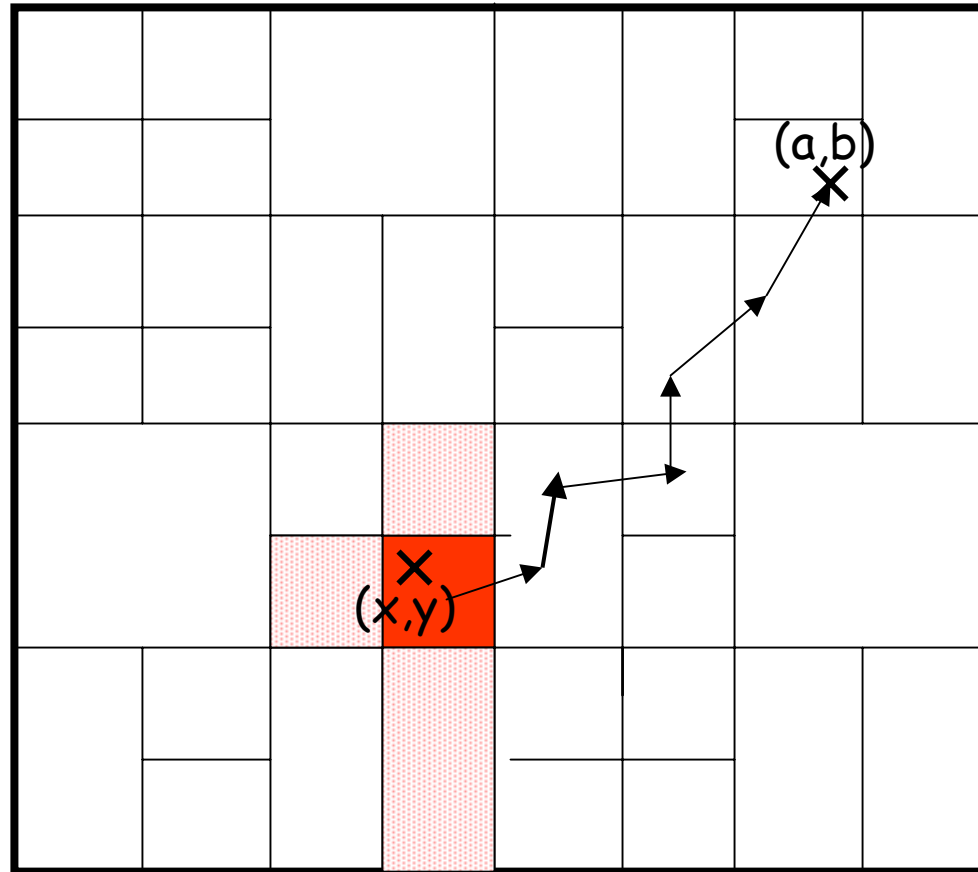
## CAN: simple example



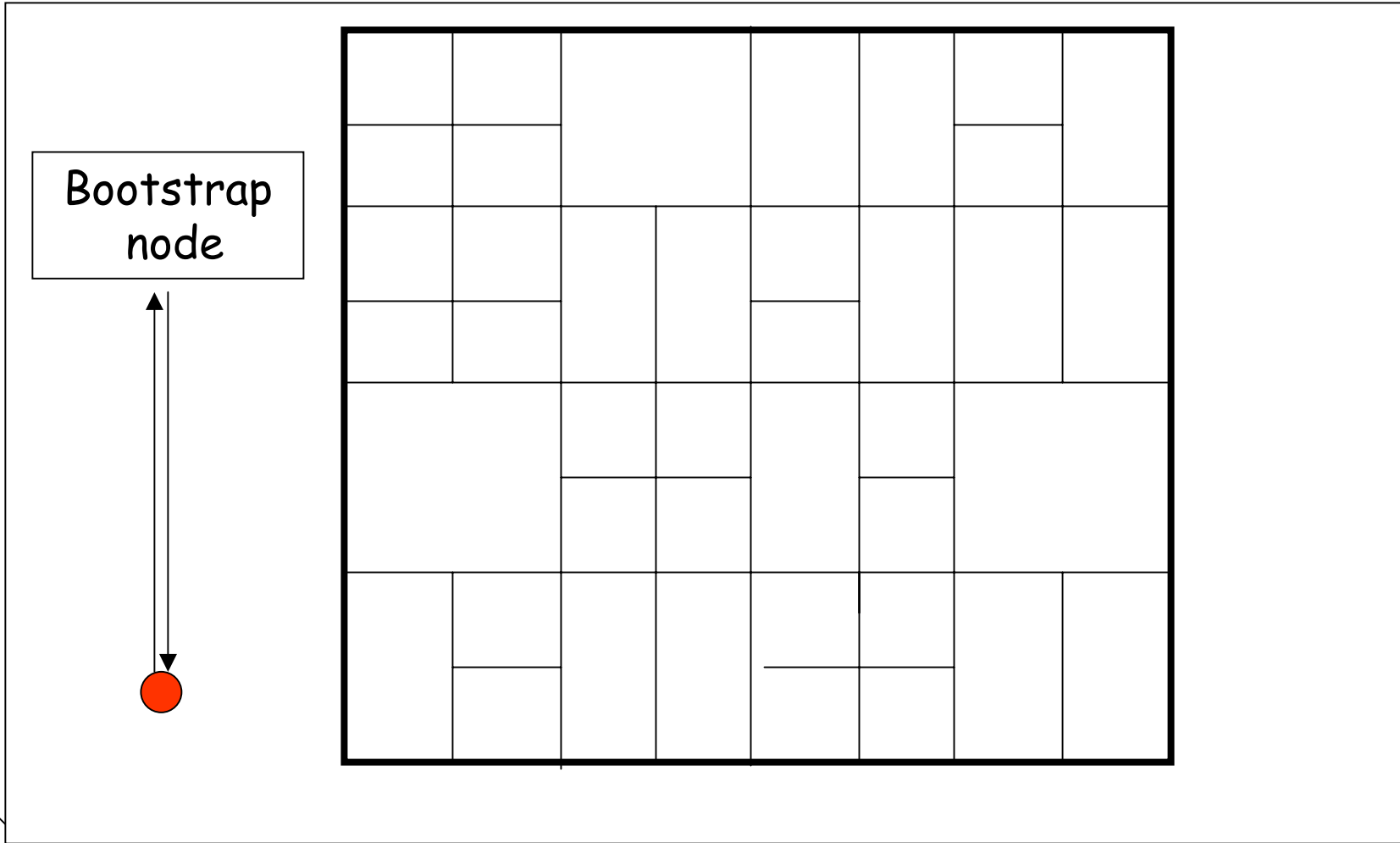
## CAN: routing table



## CAN: routing



## CAN: node insertion



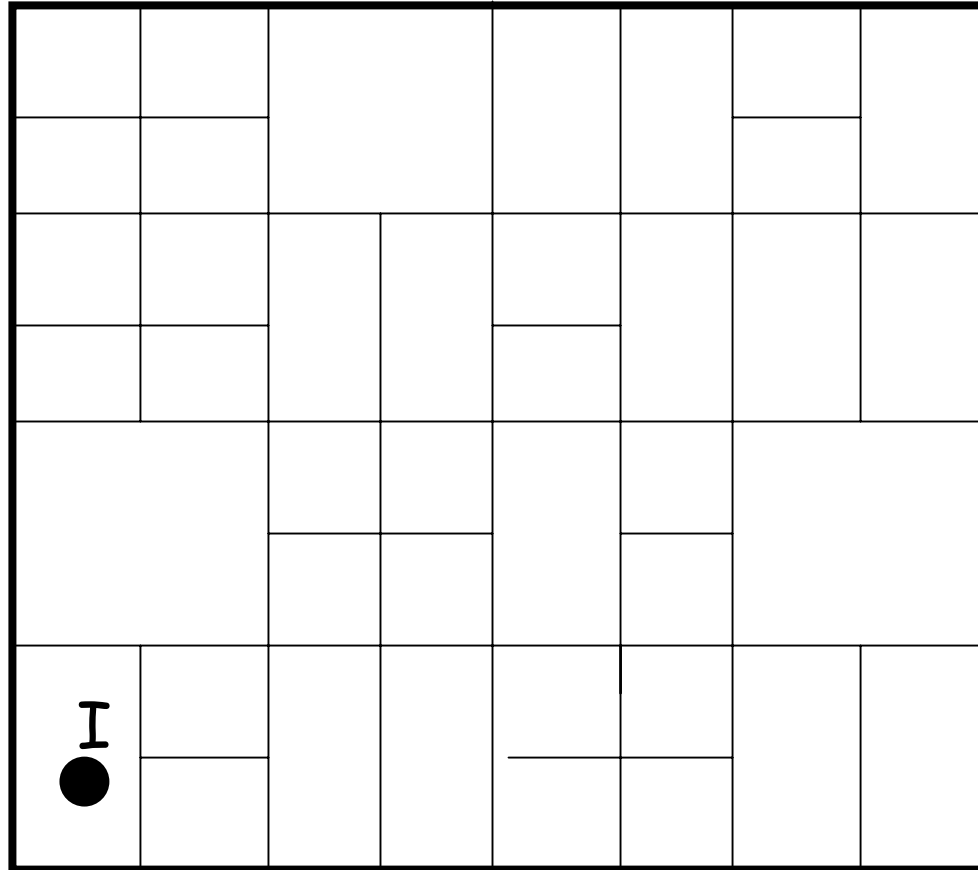
new node


Internet and Grid Computing - Fall  
1) Discover some node "I" already in CAN



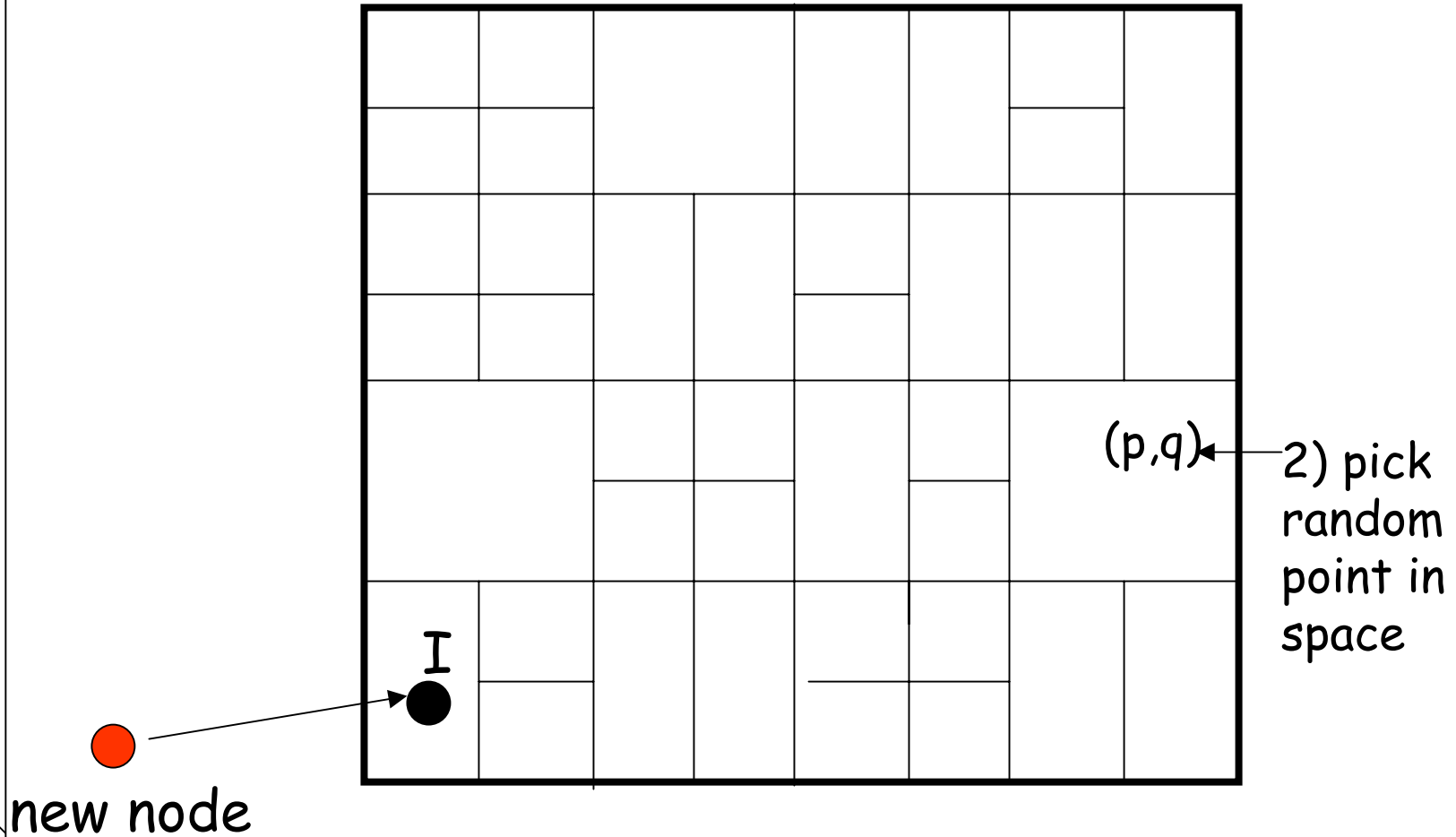
## CAN: node insertion

Bootstrap  
node

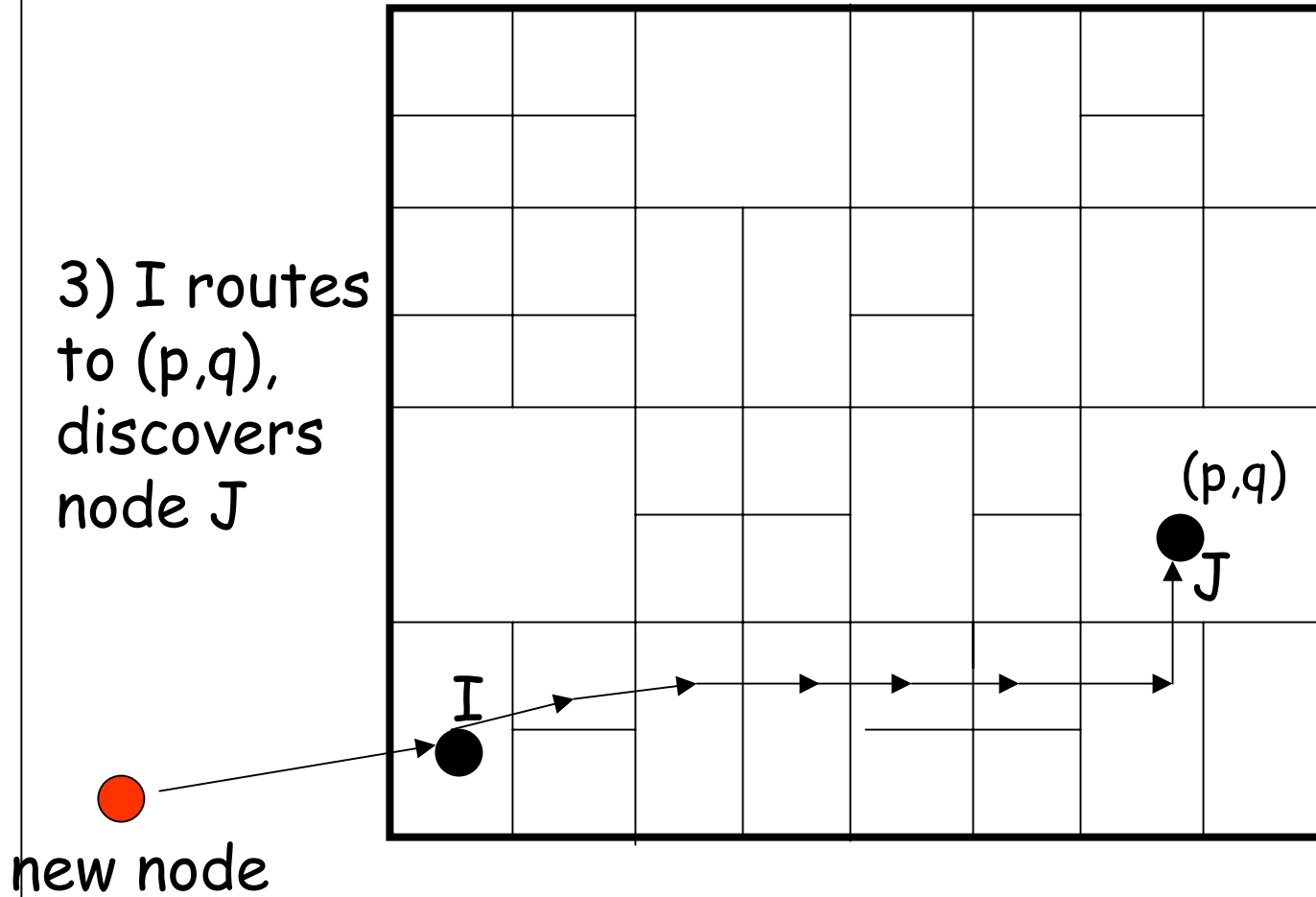


  
new node

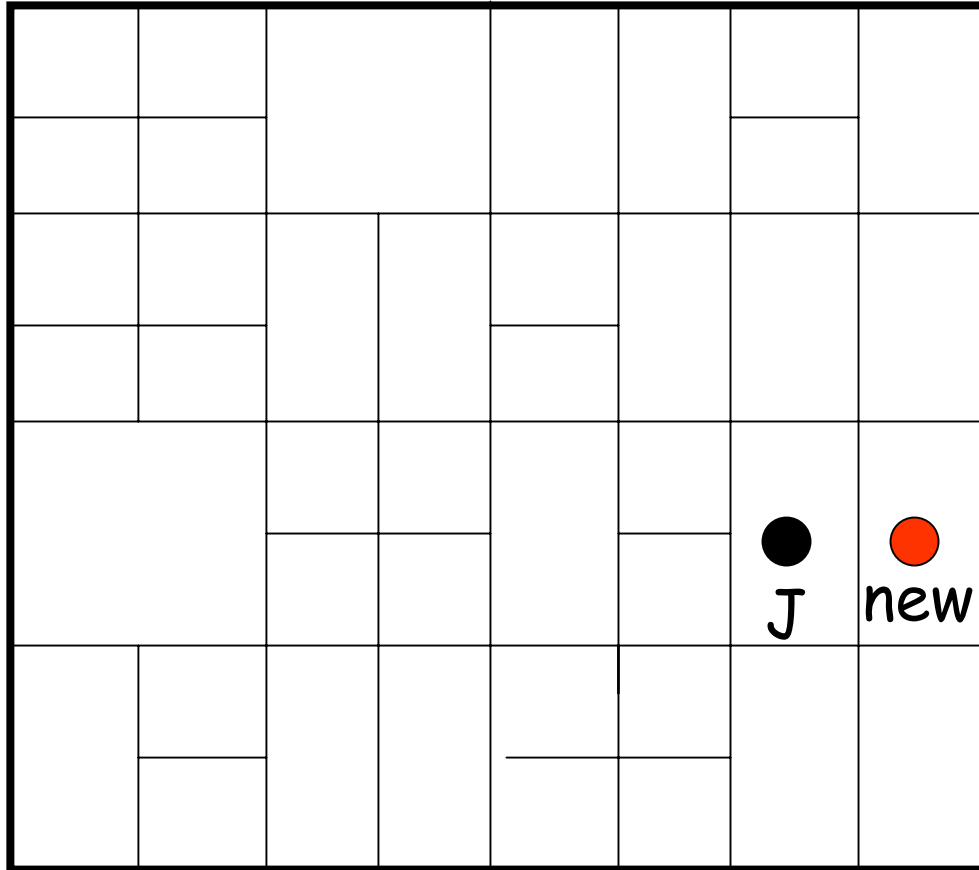
## CAN: node insertion



## CAN: node insertion



## CAN: node insertion



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

## CAN: node failures

- Simple failures
  - know your neighbor's neighbors
  - when a node fails, one of its neighbors takes over its zone
- More complex failure modes
  - simultaneous failure of multiple adjacent nodes
  - scoped flooding to discover neighbors
  - hopefully, a rare event
- Background zone reassignment algorithm

## CAN: scalability

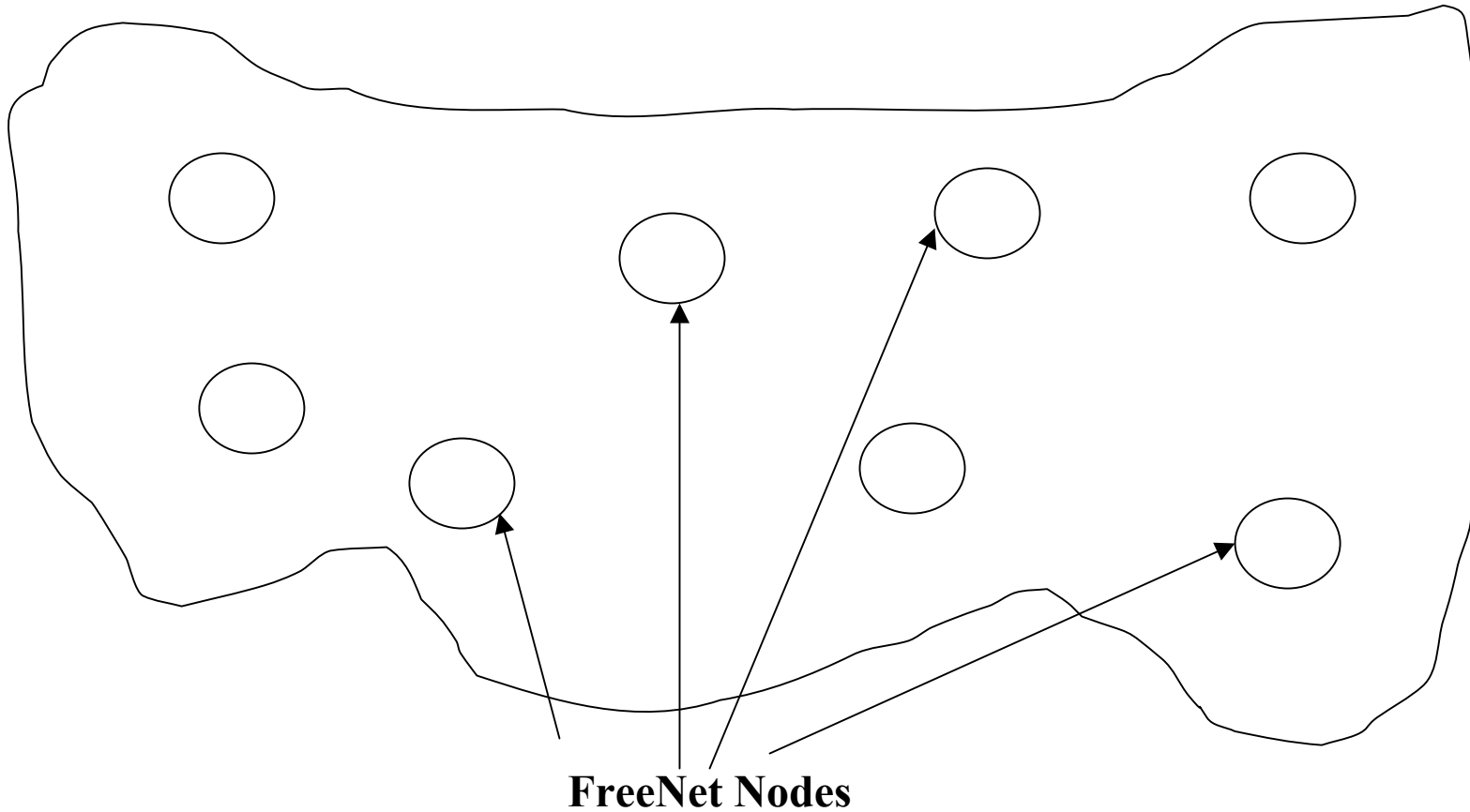
- For a uniformly partitioned space with  $n$  nodes and  $d$  dimensions
  - per node, number of neighbors is  $2d$
  - average routing path is  $(dn^{1/d})/4$  hops
    - A hop here is an application-level hop
    - 1 app-level hop = (possibly) multiple IP-level hops
  - simulations show that the above bounds hold for imperfectly partitioned spaces

Can scale the network without increasing per-node state

## Summary

- Two similar approaches to locating objects by “computed routing”
  - Similar to Manhattan Street Networks
- Both are scalable, reasonably robust
- All these P2P networks ignore underlying topology!

## FreeNet - Serverless, Symmetric, Secure, Parallel Internet File System





## **Content Summary**

**Conceptual Elements**

**File Insertion**

**File Retrieval**

**File Update**

**Large File Management**

**Node Join**

**Communication Protocols**

## Conceptual Elements of FreeNet

### Network Nodes

**Storage**

**Files**

**Routing Tables**

### Files

**Byte String**

**Path Name**

**Search Key**

**Keyword Signed Key**

**Content Key**

**File Signature**

**Directories - Personal Name**

**Spaces**

**Signed SubSpace Key**

### Interactions/Transactions

**Join Network**

**Insert File**

**Request File**

**Storage Management Algorithm**

**File Storage, Retrieval and Retention**

**Routing Tables**

## Files

**File Content = string of bytes**

**File Name/ File Description - Unix text string**

**/text/philosophy/sun-tzu/art-of-war**

**/peertopeer/books/oram**

**Keys associated with each file.**

**{Public Key, Private Key} = PPKP-Generator(File Descriptor)**

**File Keys**

**Keyword Signed Key = Secure Hash Algorithm(Public Key)**

**Content Hash Key = SHA(File Content)**

**File Encryption Key - Use file descriptor as an encryption key.**

**File Signature = Private Key**

## Definitions:

### 1. Node Locality Set

- \* Each node has a locality set of neighbors.

### 2. Node Routing Table

- \* A Node Routing Table is a set of Search Key - Node pairs.

## Node Resource Management

1. Each node allocates the amount of storage to be assigned to FreeNet Functions.
2. Storage is partitioned for file data storage and routing table storage.
3. Node data storage is managed as an LRU Cache. Files which are not accessed are eventually deleted from a node.
4. Routing table entries are also managed as an LRU cache but entries of the routing tables persist after eviction of file data.

## Personal Name Spaces/Directories

**File keys may be coupled to a unique identifier for a personal name space. The personal name space is created as follows.**

- 1. Randomly generate a public/private key pair (The signed subspace key.)**

$$\{\text{Public Key, Private Key}\} = \text{PPKP}(\text{Random}())$$

- 2. Hash the public key**

$$\text{HPK} = \text{SHA}(\text{Public Key})$$

- 3. Hash the file path name/file descriptor**

$$\text{HFPN} = \text{SHA}(\text{Path Name})$$

- 4. File Key = XOR(HFPN, HPK)**

- 5. File Signature = Private Key**

- 6. Publish the file descriptor, the public key and the File Encryption Key.**

- 7. This file is a “directory file” to which only the owner of the private key can add files.**

## **Name Space “Directory Structure”**

**Create a file using the signed-subspace key process with hypertext links to other files, perhaps signed-subspace key files.**

**The linked files may themselves contain links, etc.**

**Used to implement file update.**

## 3. File Insertion Algorithm

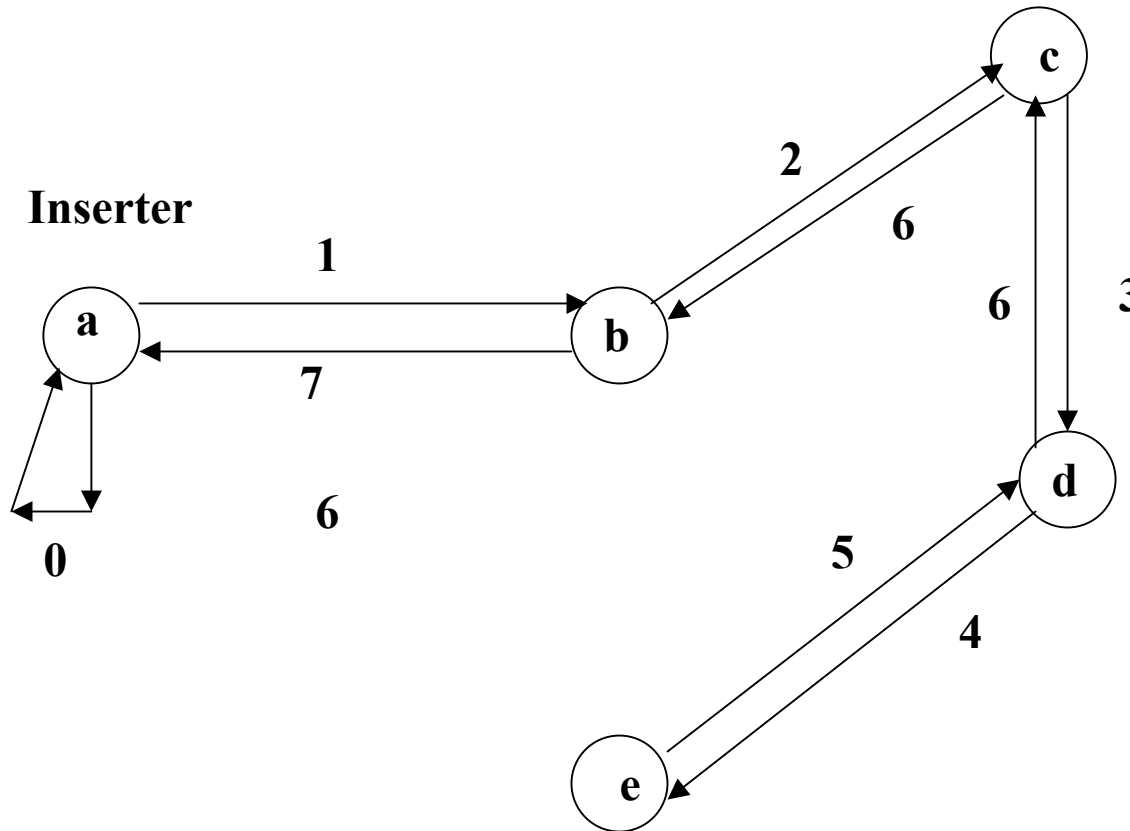
- a) Create a path name and a binary file key for the file.
- b) The creating node sends a message to itself including a “hops to live” counter for this insertion.
- c) The originating node then checks to see if that key is already in use.
- d) If a match is found then the file associated with the previously defined key is returned as though a request had been made and the node knows this key has already been used. The node generates a new key for this file and again attempts the insertion process.
- e) If no match is found then the originating node finds the node with the nearest key in its routing table and sends it an “insert” message with the key and the “hops to live” counter.
- f) If this insert message causes a collision then the file is returned to the upstream inserter and again behave as though a request had been made. (Cache the file locally and create a routing table entry for the data source.

## 3. File Insertion Algorithm - Continued

- g) If the “hops to live” limit is reached without a collision being detected then the insertion is successful and the “all clear” message is sent back to the originator of the file insertion.**
- h) The originator of the file insertion then sends the file itself which is propagated along the path of the key collision search. Each node will create a local copy of the file and establish a routing table entry matching the inserting node and the key. (A forwarding node may choose to arbitrarily change the supposed source when it forwards the insert message.)**
- i) The descriptive string (path name) is published by some out of band mechanism or in the case of name space encoded files, the descriptive string name and the subspace public key.**



# Search/Discovery (and Insertion) in Distributed Systems



**Successful  
insertion of  
file key in five  
nodes. File  
data will  
follow the  
same path.**

## Properties of Insertion Algorithm

### 1. Key-space locality

**Files are cached on nodes with similar keys.**

**2. New nodes can use inserts to extend their network locality space.**

**3. Discourages fake file propagation. (A fake file is a file with “junk” or malicious content with a key identical to some file with actual data.)**

**The original files are propagated upon collision.**

## File Search and Retrieval Algorithm

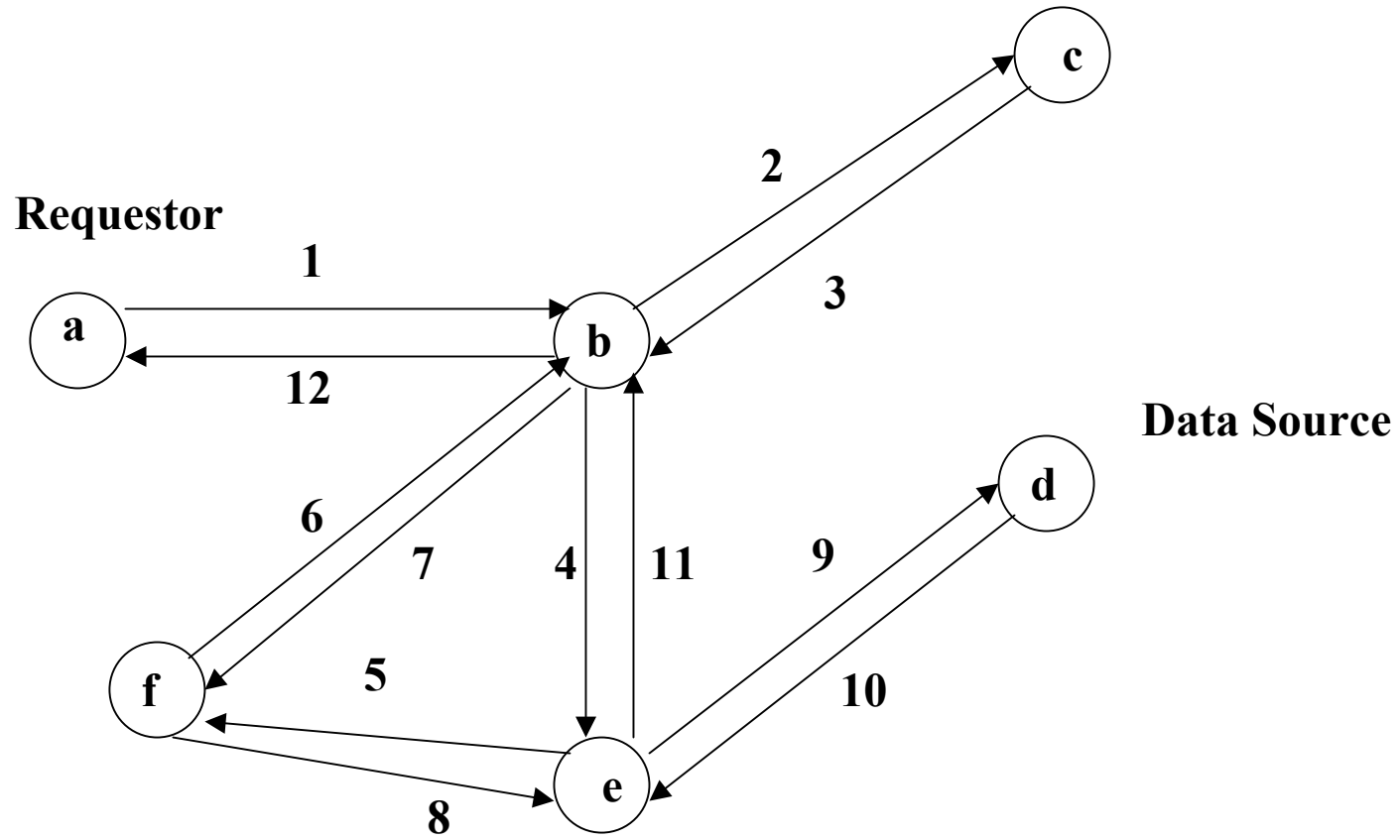
- a) Obtain or calculate the binary file search key.
- b) Send a message to yourself with the key and a “hops to live” counter.
- c) If the data is stored locally the request is satisfied.
- d) If the data is not found then the search is continued at the node in the local routing table associated with the key “nearest” to the search key.
- e) If the search is successful then the data is returned to the upstream requestor which will cache the data and create a new entry in its routing table associating the actual data source with the key.

## File Search and Retrieval Algorithm- Continued

**f) If a node in the search path cannot forward a request to its preferred downstream node then the node with the second nearest key is selected as the downstream target, etc. If a node runs out of possible downstream paths without finding the file then it reports failure to the immediately upstream requestor. This upstream requestor then chooses its second choice target. This process is recursively followed until the originator runs out of downstream targets at which time failure is reported.**

**g) If the “hops to live” limit is exceeded at any time in this search process then the failure result is propagated back to the original requestor.**

# Search/Discovery (and Insertion) in Distributed Systems



## Properties of Search Algorithm

- 1. Nodes will tend to accumulate similar keys and similar files.**
- 2. This should lead to better routing as well as better caching performance.**
- 3. The greater the number of requests the more copies of the file which will exist.**
- 4. The more localized the requests the more localized the number of copies of the file will become.**
- 5. Network connectivity is increased as routing tables are built by requests and inserts.**
- 6. Nodes with popular files will preferentially appear in routing tables of other nodes.**
- 7. While files are clustered by key they are dispersed with respect to subject.**

## File Search and Retrieval for Up-datable Files

- a) An up-datable file should be stored under a “content hash” key.
- b) A file with the content hash key is inserted in a personal name space (signed-subspace key)
- c) The file is encrypted with a random key which is published with the file key.
- d) Retrieval is by first retrieving the file containing the content hash key and using this key to search for the file.

## File Update Algorithm

- a) A new version of an up-datable file is inserted under its content hash key.**
- B) The insert algorithm is executed for the new indirect file under the original signed subspace key.**
- C) When the insert reaches a node with the old version of the file a key collision will occur.**
- D) The node will check the signature on the new version, verify that it is correct and replace the old version of the file with the new version. Then the signed subspace key will always lead to the new version while a content hash search will still lead to the old version.**



## Management of Large Files

- 1. Partition a large file.**
- 2. Create a content hash key for each partition.**
- 3. Create a file to serve as the indirect access file for the partitions of the large file.**
- 3. Insert the content hash key for each partition separately as an entry in the indirect file.**
- 5. Insert the indirect file.**

## Publication of Names and Search Mechanisms

- 1. Create a search engine specific to FreeNet.**
- 2. Create for each actual file a family of “lightweight indirect files” each named by a search keyword relating to the actual file and containing a pointer to the actual file.**
- 3. Encourage users to create directories of “favorites”**

## Node Join Protocol

- 1. The joining node must obtain the address of at least one node which is already a member of FreeNet by some out of band means.**
- 2. The joining node chooses a random seed and sends an announcement message containing the hash of that seed, its address and a “hops to live” counter to the existing nodes for which it knows addresses.**
- 3. When a node receives an announcement message it generates a random seed, XORs that with the hash it received and hashes the result again to generate a “commitment.”**
- 4. The node which received the announcement message forwards the new hash to some node chosen randomly from its routing table and decrements the hops to live counter.**
- 5. The last node to receive the announcement message just generates a seed.**
- 6. There all the nodes in the announce chain share their seeds and the key for the new node is assigned as the XOR of all of the seeds.**
- 7. Then each of the nodes in the announce chain add the new node to their routing table under that key.**

## FreeNet Communication Protocols

- 1. Each request, insert or join action is a transaction which has a (probably) unique ID associated with it.**
- 2. Each message originated as a result of a transaction has a transaction ID which is carried in each message resulting from a transaction.**
- 3. Node addresses are {transport identifier, transport specific identifier}, for example {tcp/192.168.1.1:19114}**
- 4. Nodes which change addresses frequently may wish to use address resolution keys which are signed subspace keys updated to contain the current actual address of the node.**
- 5. All transactions begin with a Request.Handshake message specifying the return address of the sending node.**
- 6. The receiver of a Request.Handshake message may (or may not) respond with a Reply.Handshake message specifying the protocol it understands.**

## FreeNet Communication Protocols - Continued

- 7. All messages have a 64 bit randomly generated transaction ID, a hops to live counter, and a depth counter.**
- 8. Hops to live and depth are set by the originator of a transaction is decremented by each receiver.**
- 9. The propagation chain is continued with hops to live = 1 with some finite probability.**
- 10. Depth is incremented at each hop and is used by a replying node to set hops to live high enough to reach a requestor. A depth of 1 is not automatically incremented but may be passed unchanged with some finite probability.**
- 11. A time-out is superimposed on transactions.**
- 12. Nodes in a chain may send back Reply.Restart messages based on knowledge of network delays not accessible to the originator.**

## Summary and Conclusions

- 1. Design achieves goals of symmetry, fully distributed control, parallelism, high security and use of “unused resources.”**
- 2. Performance properties are generally unknown except for simulations.**
- 3. Algorithm domain for fully distributed control is largely unexplored.**
- 4. State of the art is about where client-server systems were 10 years ago.**
- 5. Next step - Integration of distributed computation and distributed file system.**  
**Potential of system with 100,000,000 computers linked by Internet.**
- 6. Problem - Application design**