

CSCI-1680

Some Alternatives

Rodrigo Fonseca



Alternatives

- **P2P**
 - Focus on scalable routing on flat names
- **Erasur Coding**
 - Alternative to ACK-based reliability
- **Information-Centric Networking**
 - Alternative to pair-based communication



Peer-to-Peer Systems

- **How did it start?**
 - A killer application: file distribution
 - Free music over the Internet! (*not exactly legal...*)
- **Key idea: share storage, content, and bandwidth of individual users**
 - Lots of them
- **Big challenge: coordinate all of these users**
 - In a scalable way (not $N \times N$!)
 - With changing population (aka *churn*)
 - With no central administration
 - With no trust
 - With large heterogeneity (content, storage, bandwidth,...)



3 Key Requirements

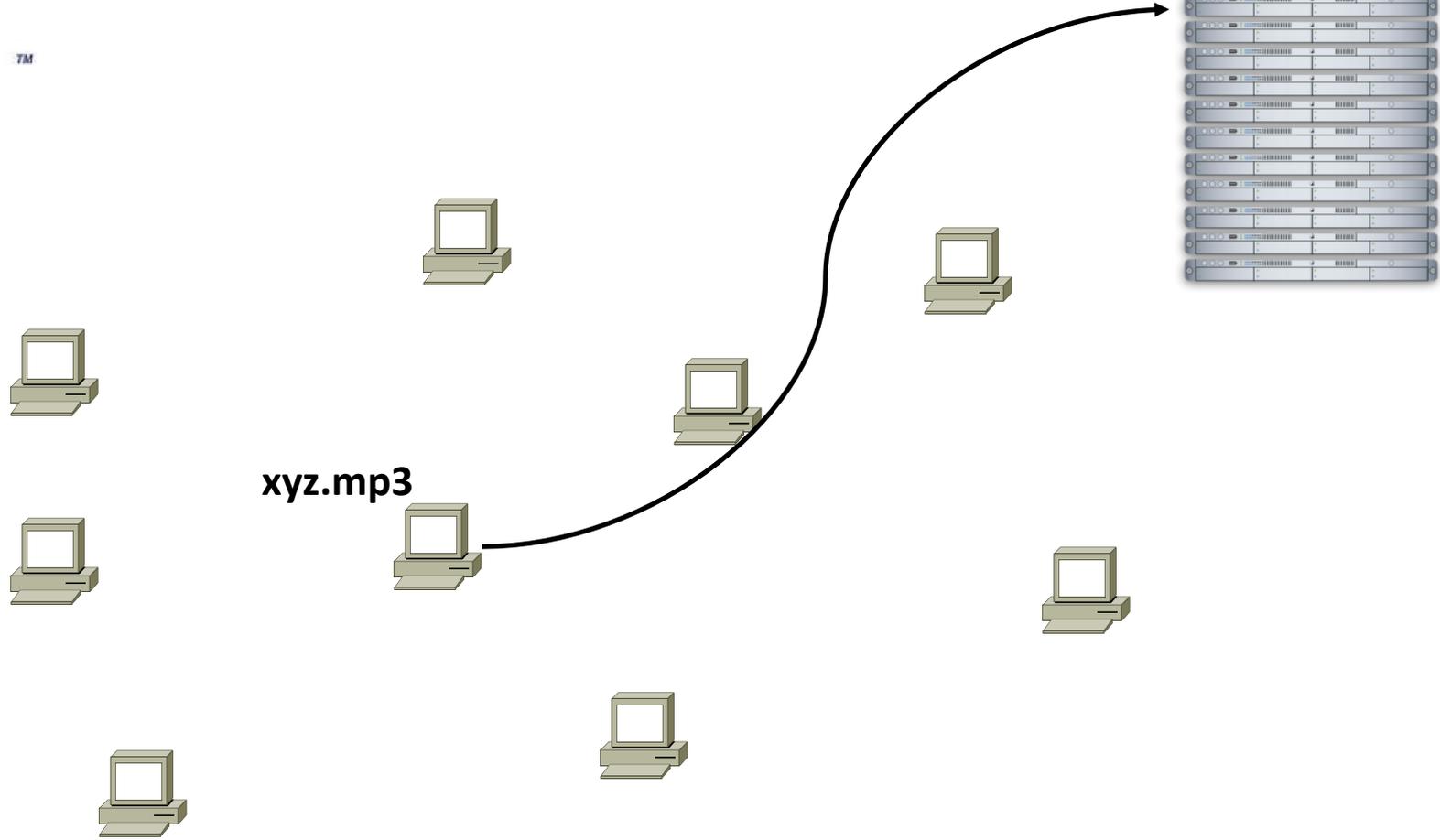
- **P2P Systems do three things:**
- **Help users determine what they want**
 - Some form of search
 - P2P version of Google
- **Locate that content**
 - Which node(s) hold the content?
 - P2P version of DNS (map name to location)
- **Download the content**
 - Should be efficient
 - P2P form of Akamai





TM

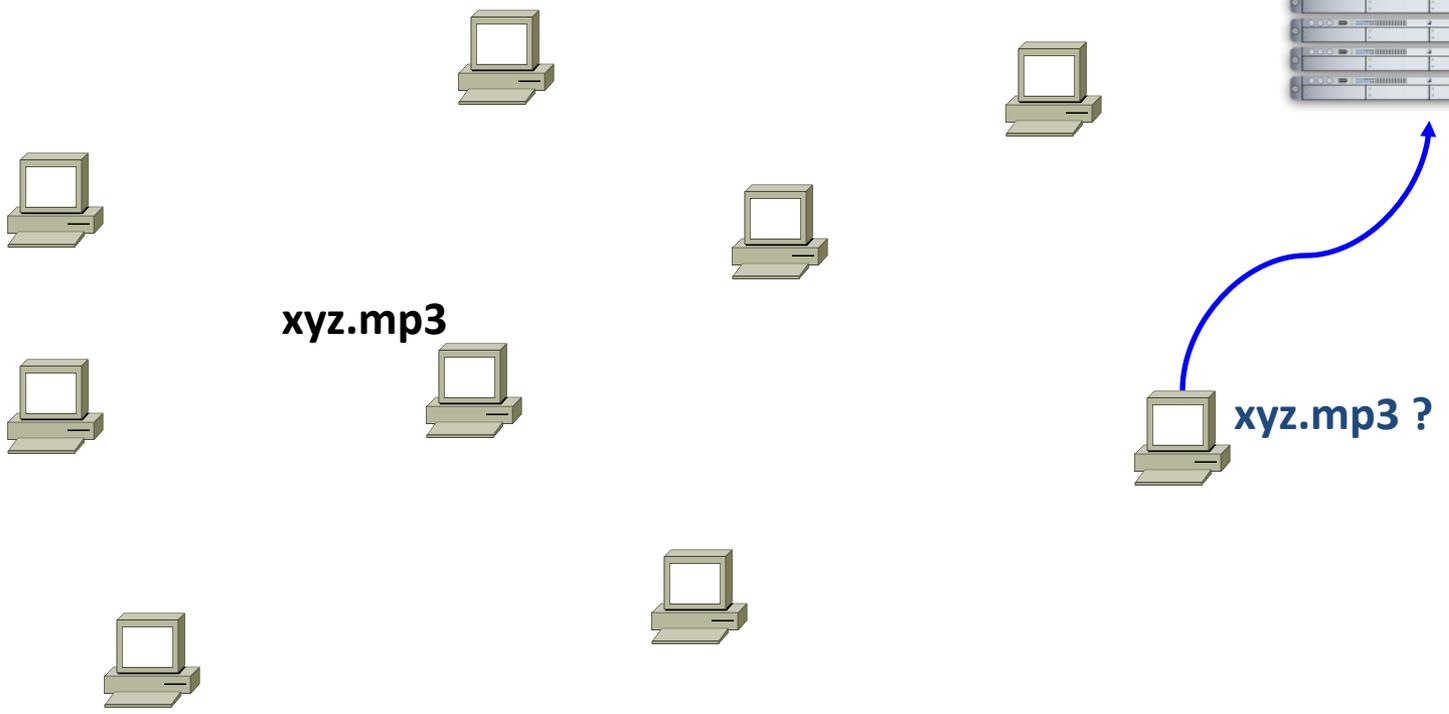
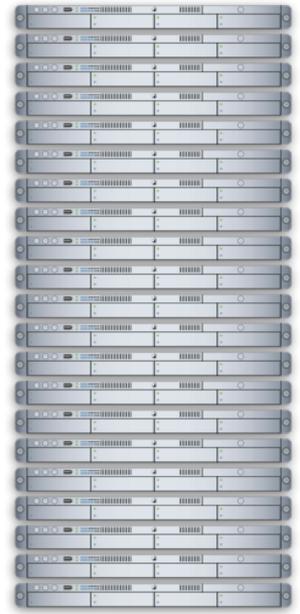
Napster (1999)





TM

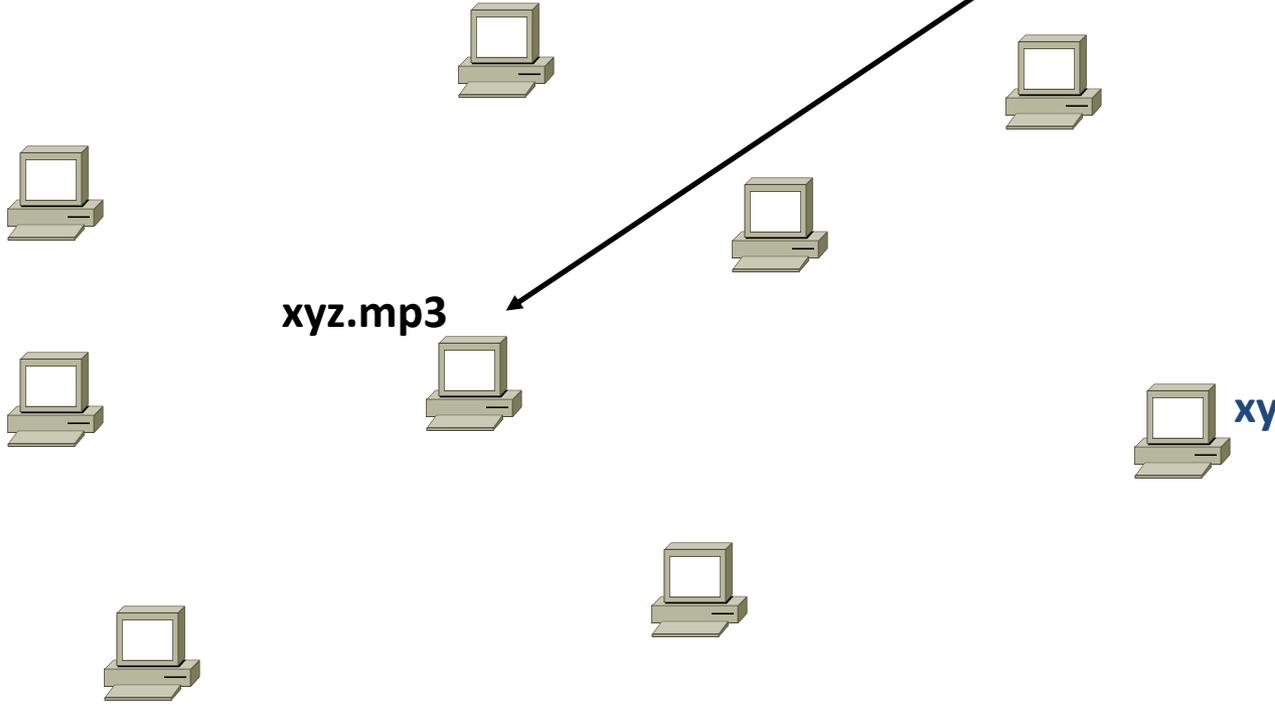
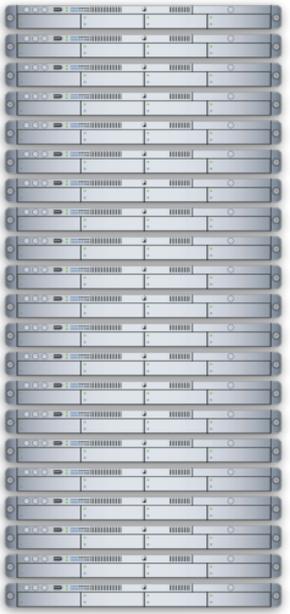
Napster





TM

Napster



xyz.mp3

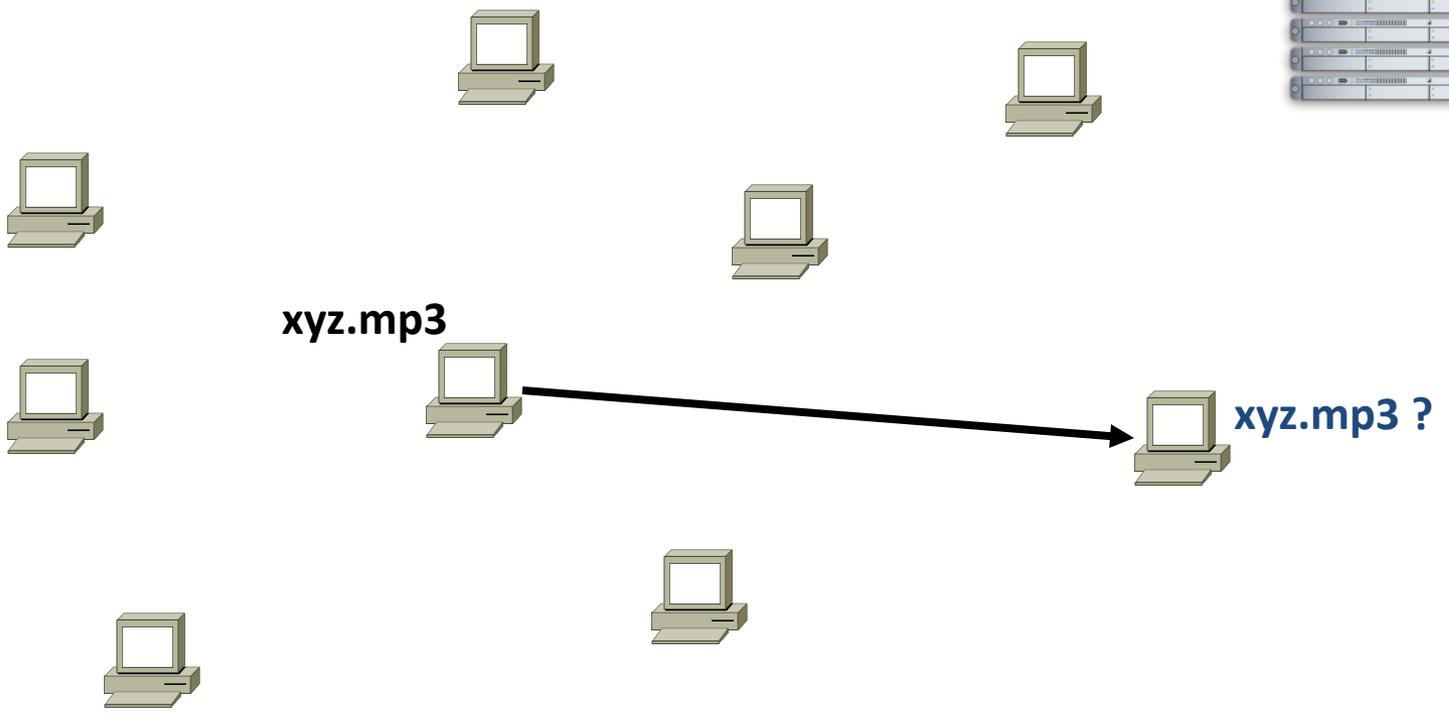
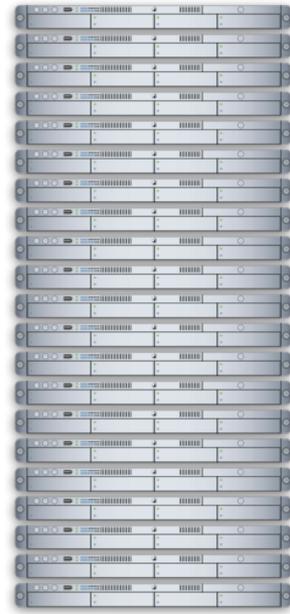
xyz.mp3 ?





TM

Napster



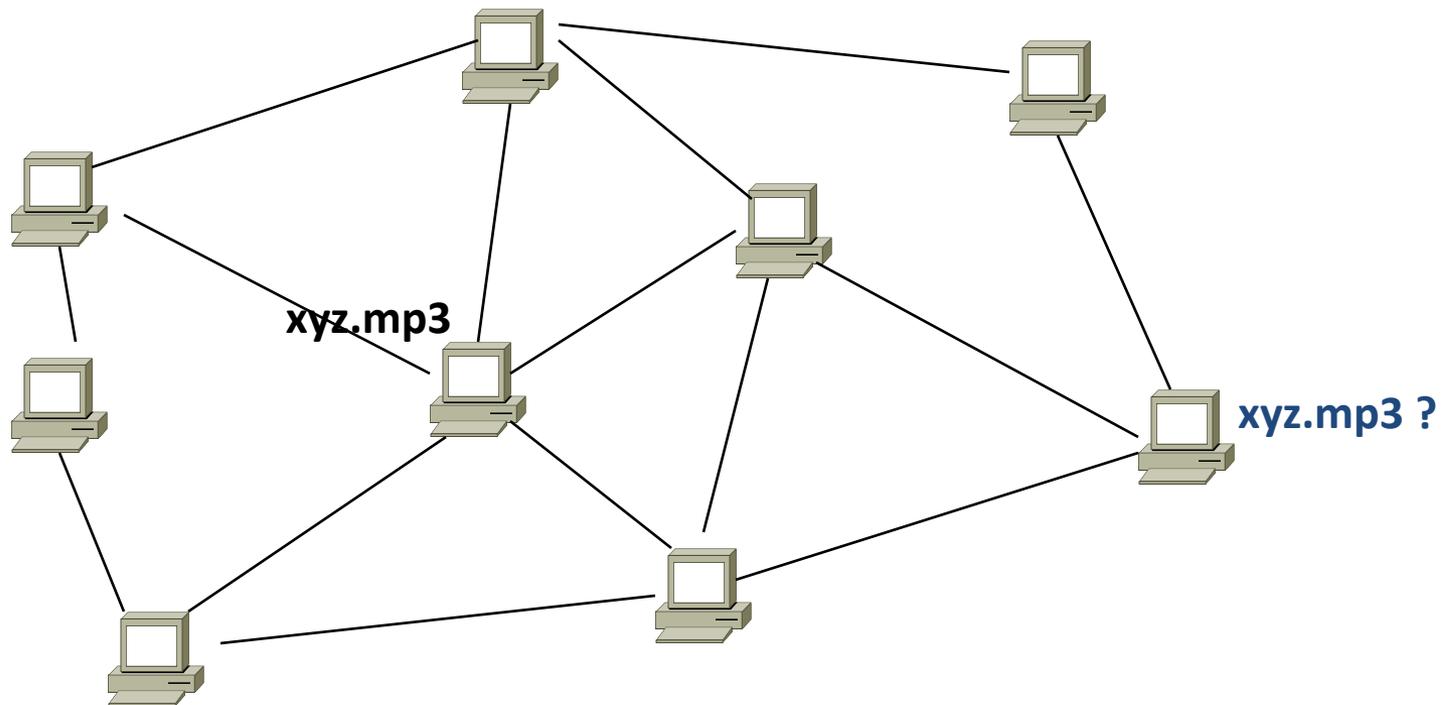
Napster

- **Search & Location: central server**
- **Download: contact a peer, transfer directly**
- **Advantages:**
 - Simple, advanced search possible
- **Disadvantages:**
 - Single point of failure (technical and ... legal!)
 - The latter is what got Napster killed



Gnutella: Flooding on Overlays (2000)

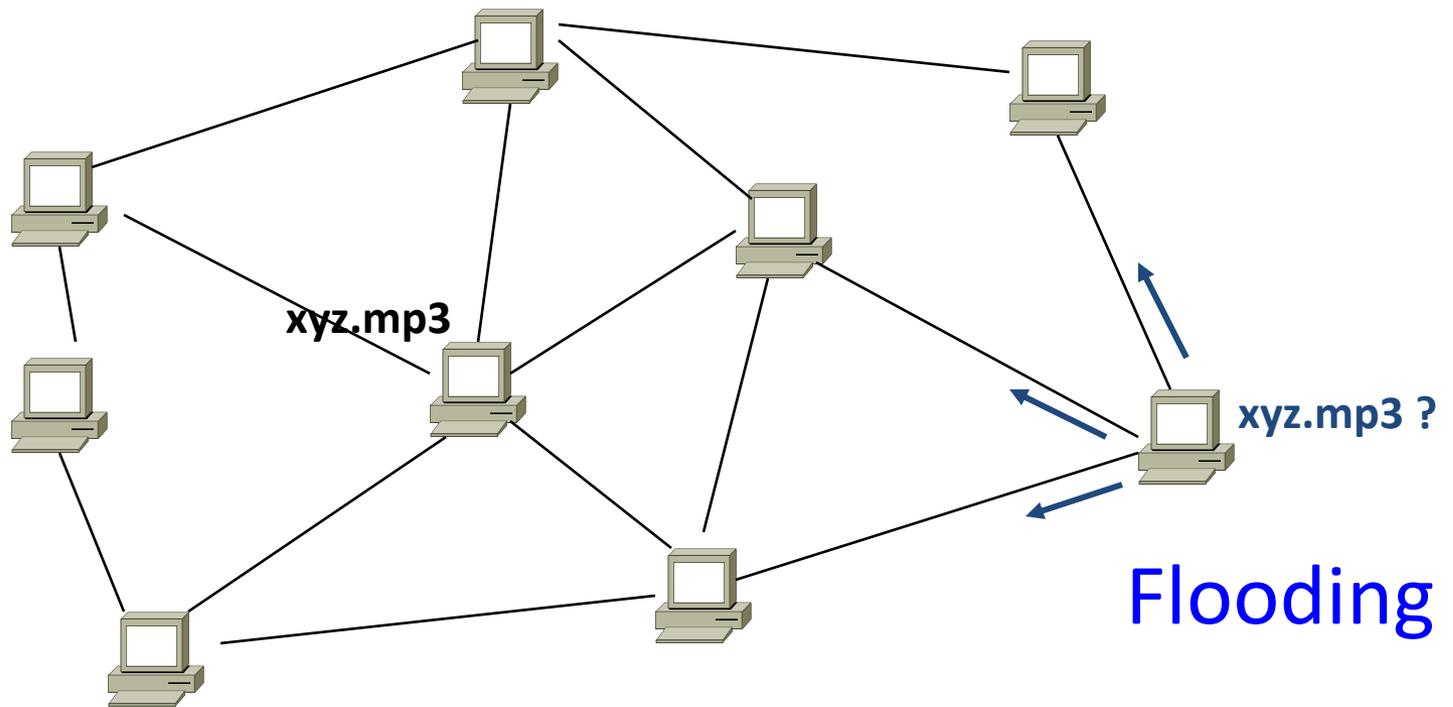
- **Search & Location: flooding (with TTL)**
- **Download: direct**



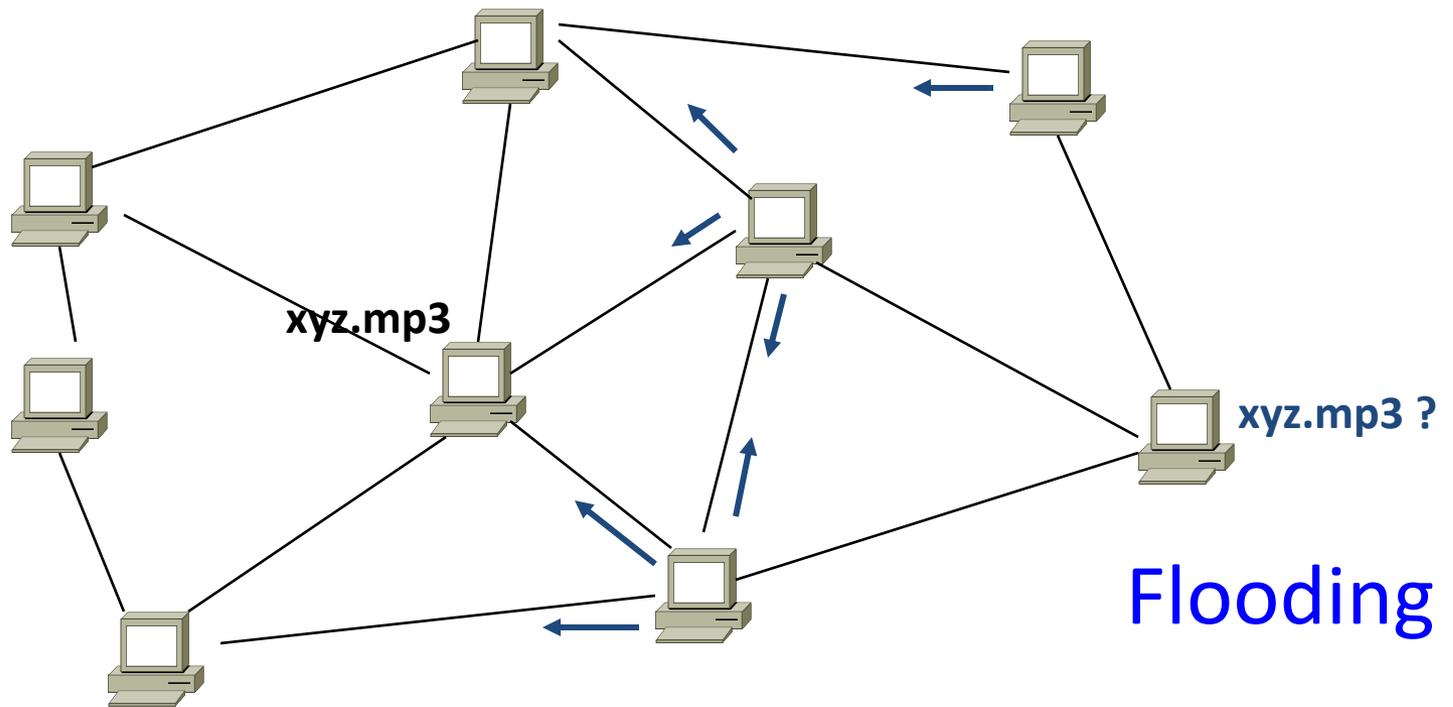
An “unstructured” *overlay network*



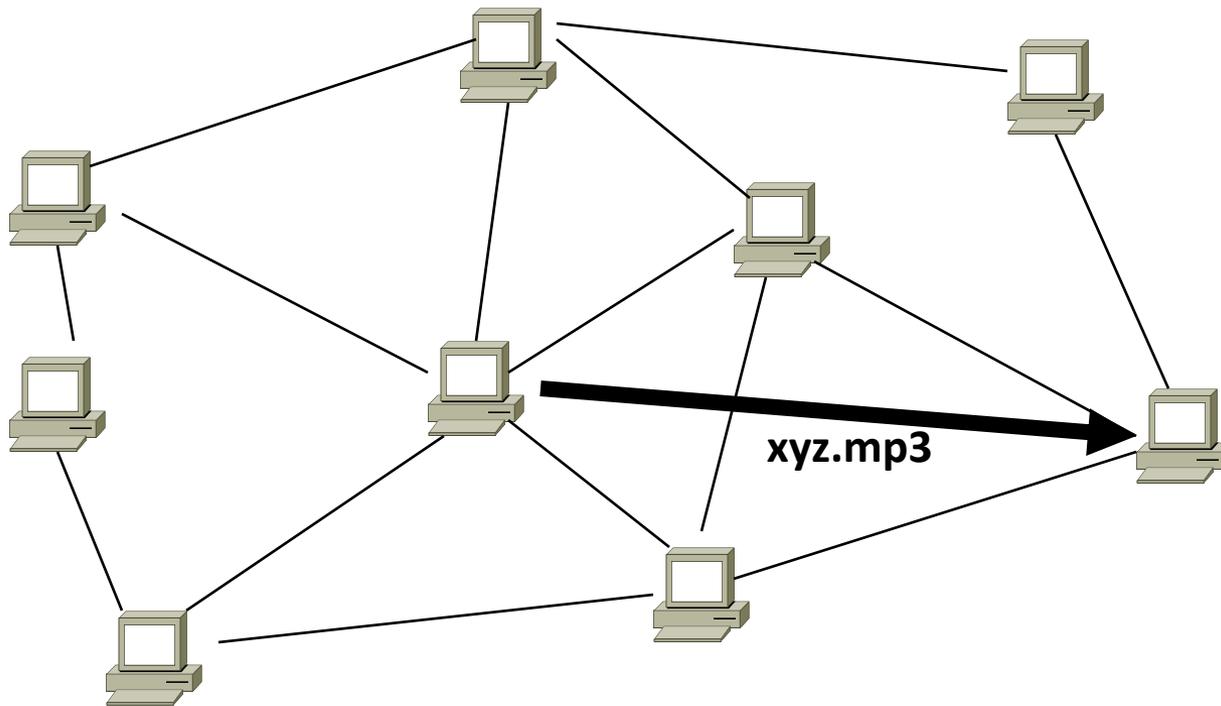
Gnutella: Flooding on Overlays



Gnutella: Flooding on Overlays

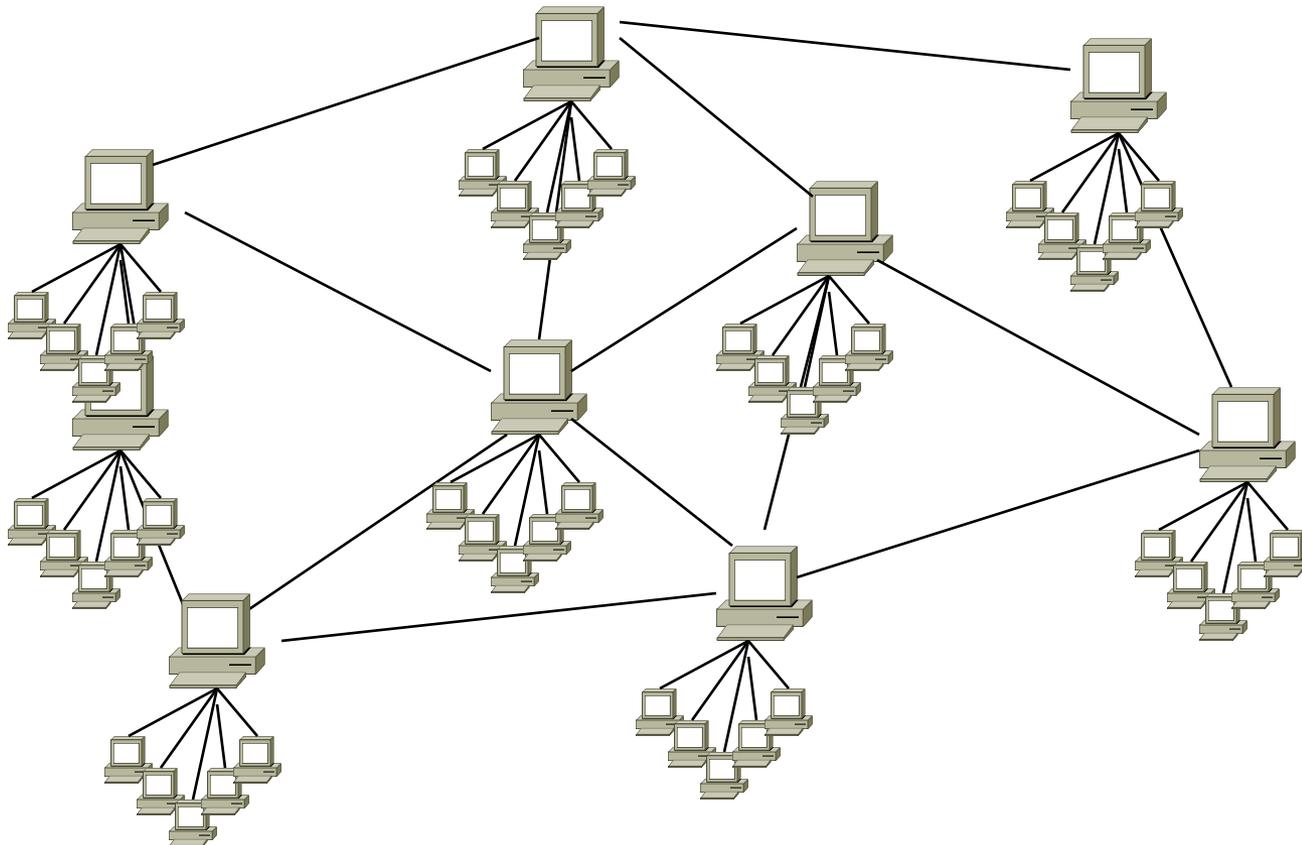


Gnutella: Flooding on Overlays



KaZaA: Flooding w/ Super Peers (2001)

- Well connected nodes can be installed (KaZaA) or self-promoted (Gnutella)



Say you want to make calls among peers

- **You need to find who to call**
 - Centralized server for authentication, billing
- **You need to find where they are**
 - Can use central server, or a decentralized search, such as in KaZaA
- **You need to call them**
 - What if both of you are behind NATs? (only allow outgoing connections)
 - You could use another peer as a relay...



Skype



- **Built by the founders of KaZaA!**
- **Uses Superpeers for registering presence, searching for where you are**
- **Uses regular nodes, outside of NATs, as decentralized relays**
 - This is their killer feature
- **This morning, from my computer:**
 - 29,565,560 people online



Lessons and Limitations

- **Client-server performs well**
 - But not always feasible
- **Things that flood-based systems do well**
 - Organic scaling
 - Decentralization of visibility and liability
 - Finding popular stuff
 - Fancy *local* queries
- **Things that flood-based systems do poorly**
 - Finding unpopular stuff
 - Fancy *distributed* queries
 - Vulnerabilities: data poisoning, tracking, etc.
 - Guarantees about anything (answer quality, privacy, etc.)





BitTorrent (2001)

- **One big problem with the previous approaches**
 - Asymmetric bandwidth
- **BitTorrent (original design)**
 - Search: independent search engines (e.g. PirateBay, isoHunt)
 - Maps keywords -> .torrent file
 - Location: centralized *tracker* node per file
 - Download: chunked
 - File split into many pieces
 - Can download from many peers





BitTorrent

- **How does it work?**
 - Split files into large pieces (256KB ~ 1MB)
 - Split pieces into subpieces
 - Get peers from tracker, exchange info on pieces
- **Three-phases in download**
 - Start: get a piece as soon as possible (random)
 - Middle: spread pieces fast (rarest piece)
 - End: don't get stuck (parallel downloads of last pieces)





BitTorrent

- **Self-scaling: incentivize sharing**
 - If people upload as much as they download, system scales with number of users (no free-loading)
- **Uses *tit-for-tat*: only upload to those who give you data**
 - *Choke* most of your peers (don't upload to them)
 - Order peers by download rate, choke all but P best
 - Occasionally unchoke a random peer (might become a nice uploader)
- **Optional reading:**
[\[*Do Incentives Build Robustness in BitTorrent?* Piatek et al, NSDI'07\]](#)



Structured Overlays: DHTs

- **Academia came (a little later)...**
- **Goal: Solve efficient decentralized location**
 - Remember the second key challenge?
 - Given ID, map to host
- **Remember the challenges?**
 - Scale to millions of nodes
 - Churn
 - Heterogeneity
 - Trust (or lack thereof)
 - Selfish and malicious users



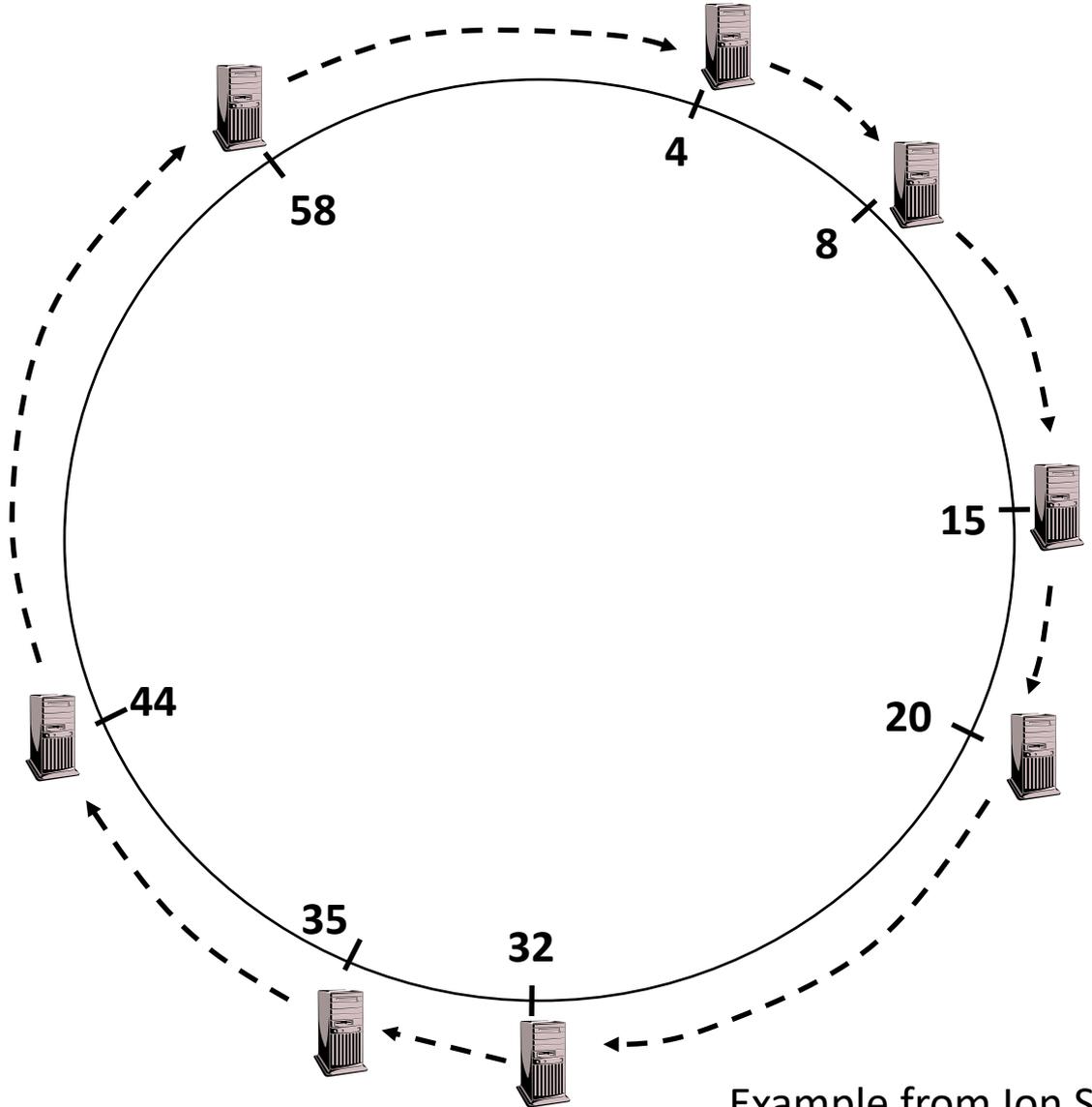
DHTs

- **IDs from a *flat* namespace**
 - Contrast with hierarchical IP, DNS
- **Metaphor: hash table, but distributed**
- **Interface**
 - Get(key)
 - Put(key, value)
- **How?**
 - Every node supports a single operation:
Given a *key*, route messages to node holding *key*



Identifier to Node Mapping Example

- Node 8 maps [5,8]
 - Node 15 maps [9,15]
 - Node 20 maps [16, 20]
 - ...
 - Node 4 maps [59, 4]
-
- Each node maintains a pointer to its successor

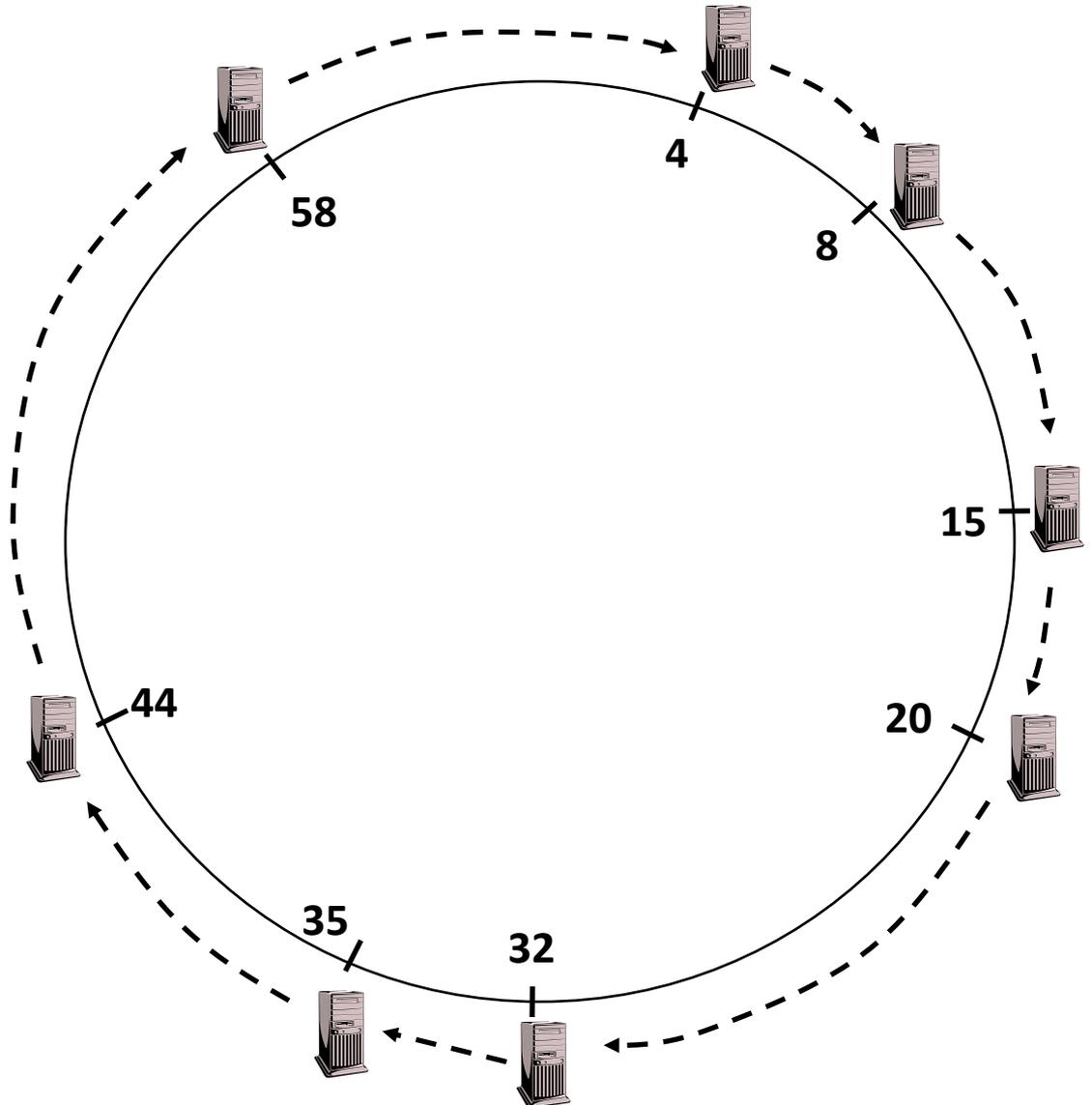


Example from Ion Stoica



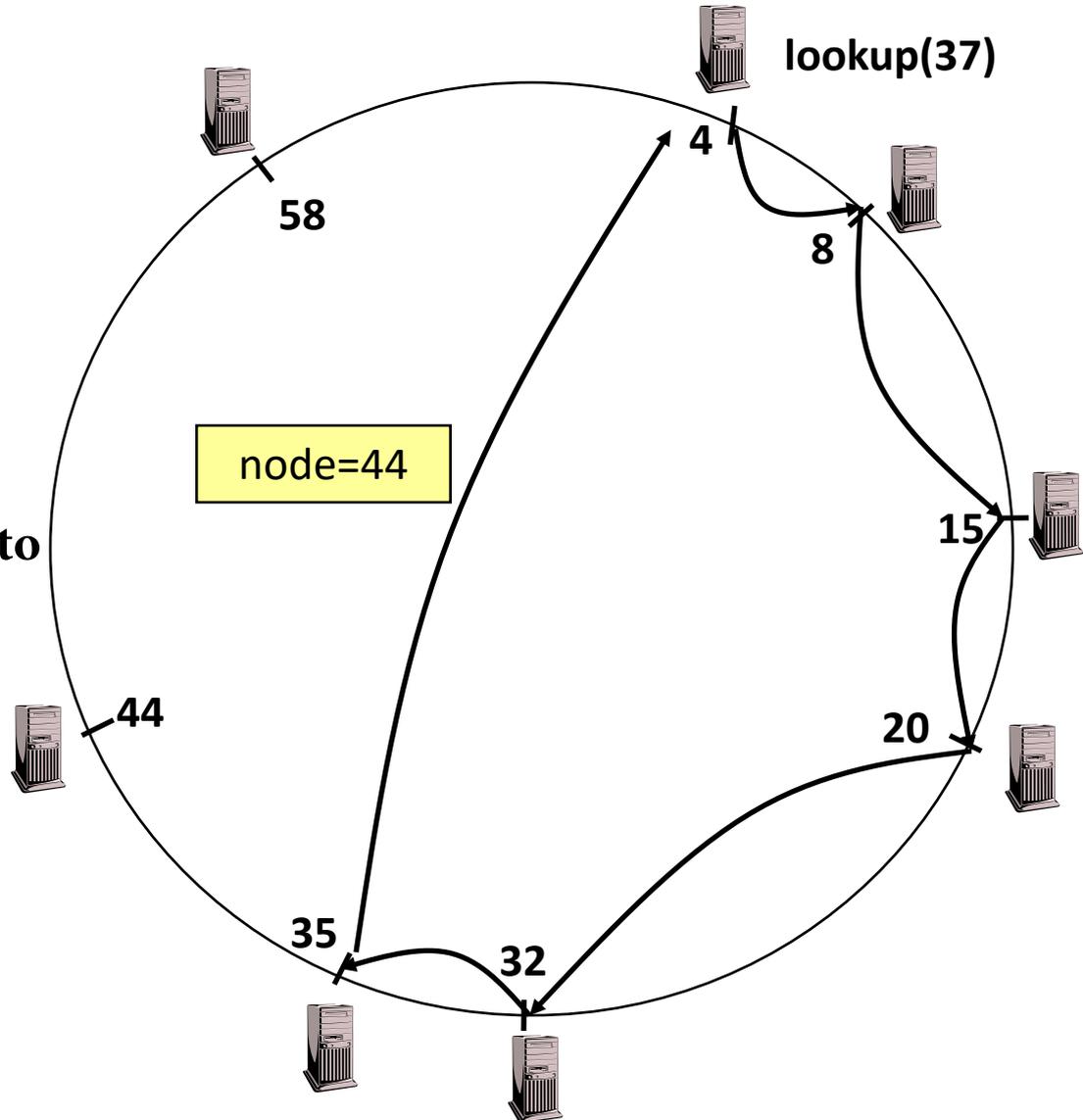
Consistent Hashing-like

- But each node only knows about a small number of other nodes (so far only their successors)



Lookup

- Each node maintains its successor
- Route packet (ID, data) to the node responsible for ID using successor pointers



Optional: DHT Maintenance



Stabilization Procedure

- **Periodic operations performed by each node N to maintain the ring:**

STABILIZE() [N.successor = M]

N->M: *"What is your predecessor?"*

M->N: *"x is my predecessor"*

if x between (N,M), N.successor = x

N->N.successor: NOTIFY()

NOTIFY()

N->N.successor: *"I think you are my successor"*

M: upon receiving NOTIFY from N:

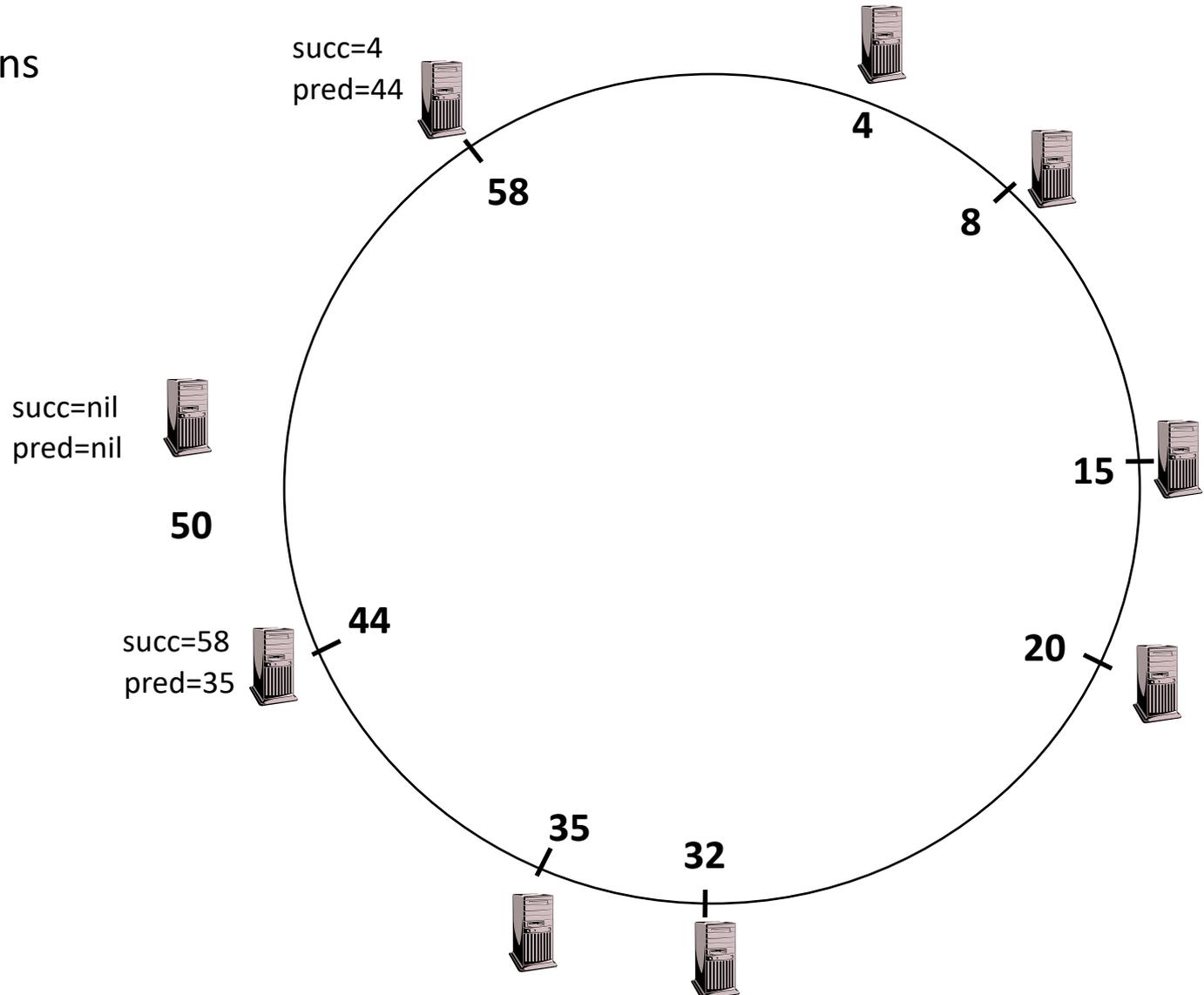
If (N between (M.predecessor, M))

M.predecessor = N



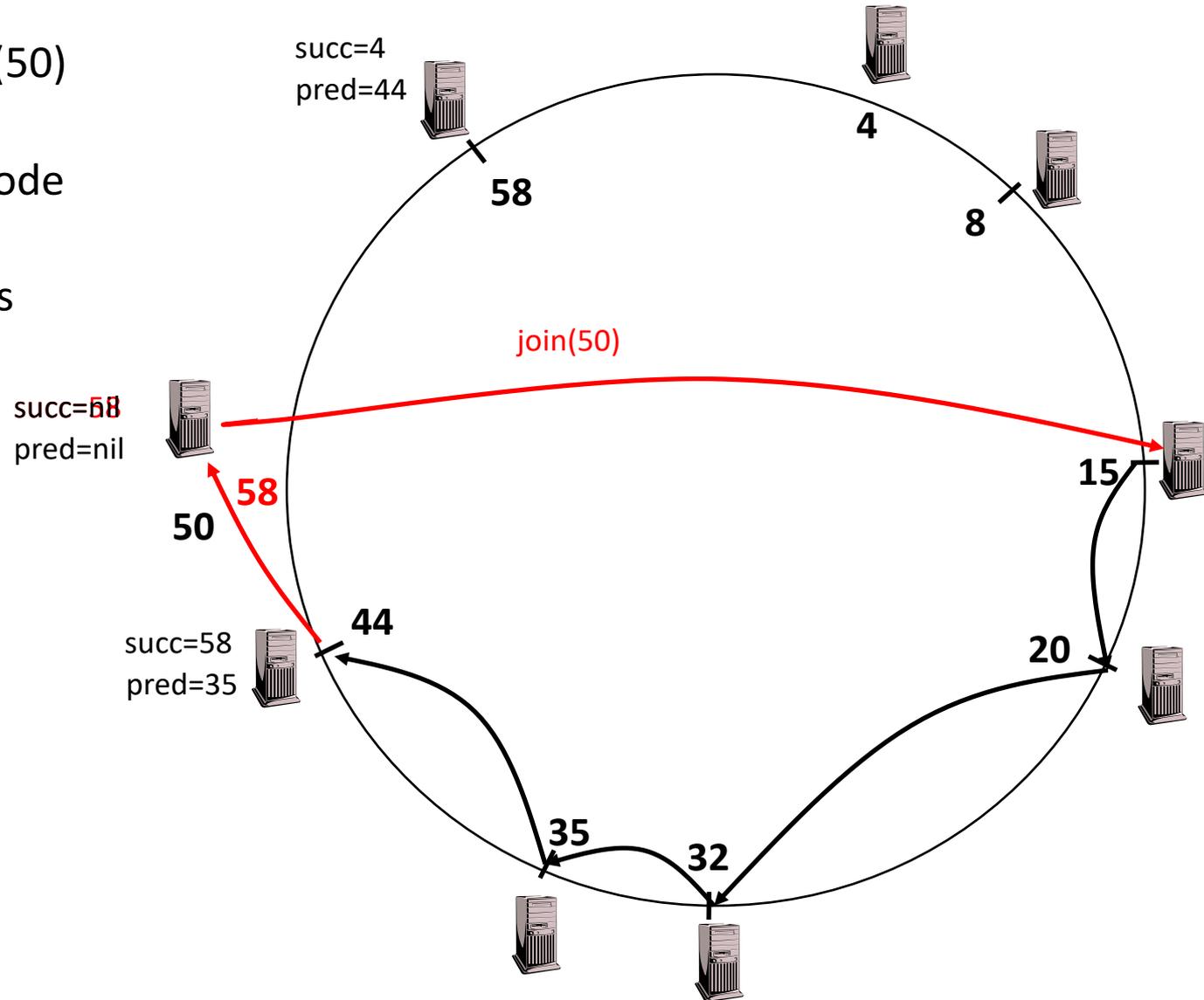
Joining Operation

- Node with id=50 joins the ring
- Node 50 needs to know at least one node already in the system
- Assume known node is 15



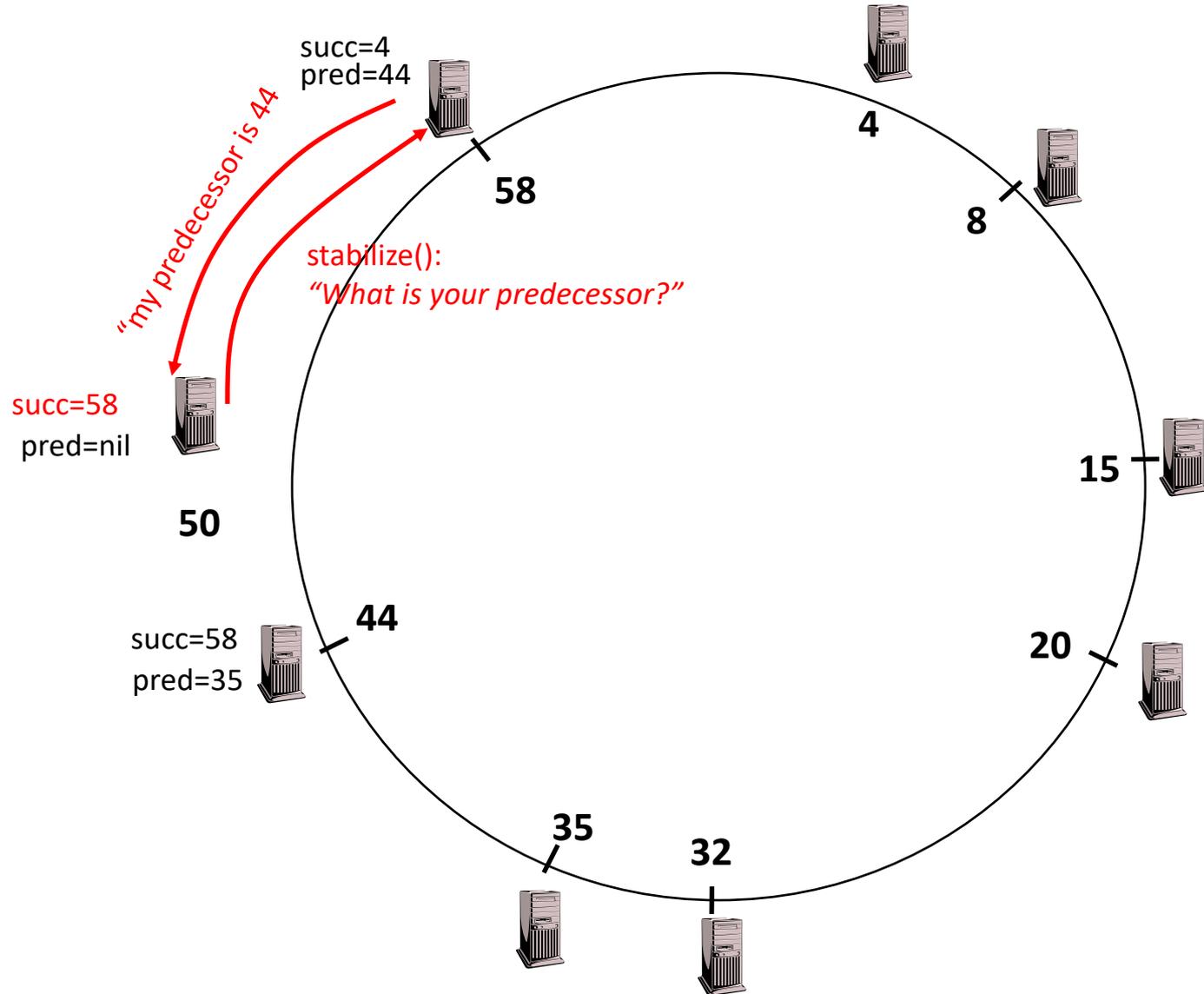
Joining Operation

- Node 50: send join(50) to node 15
- Node 44: returns node 58
- Node 50 updates its successor to 58



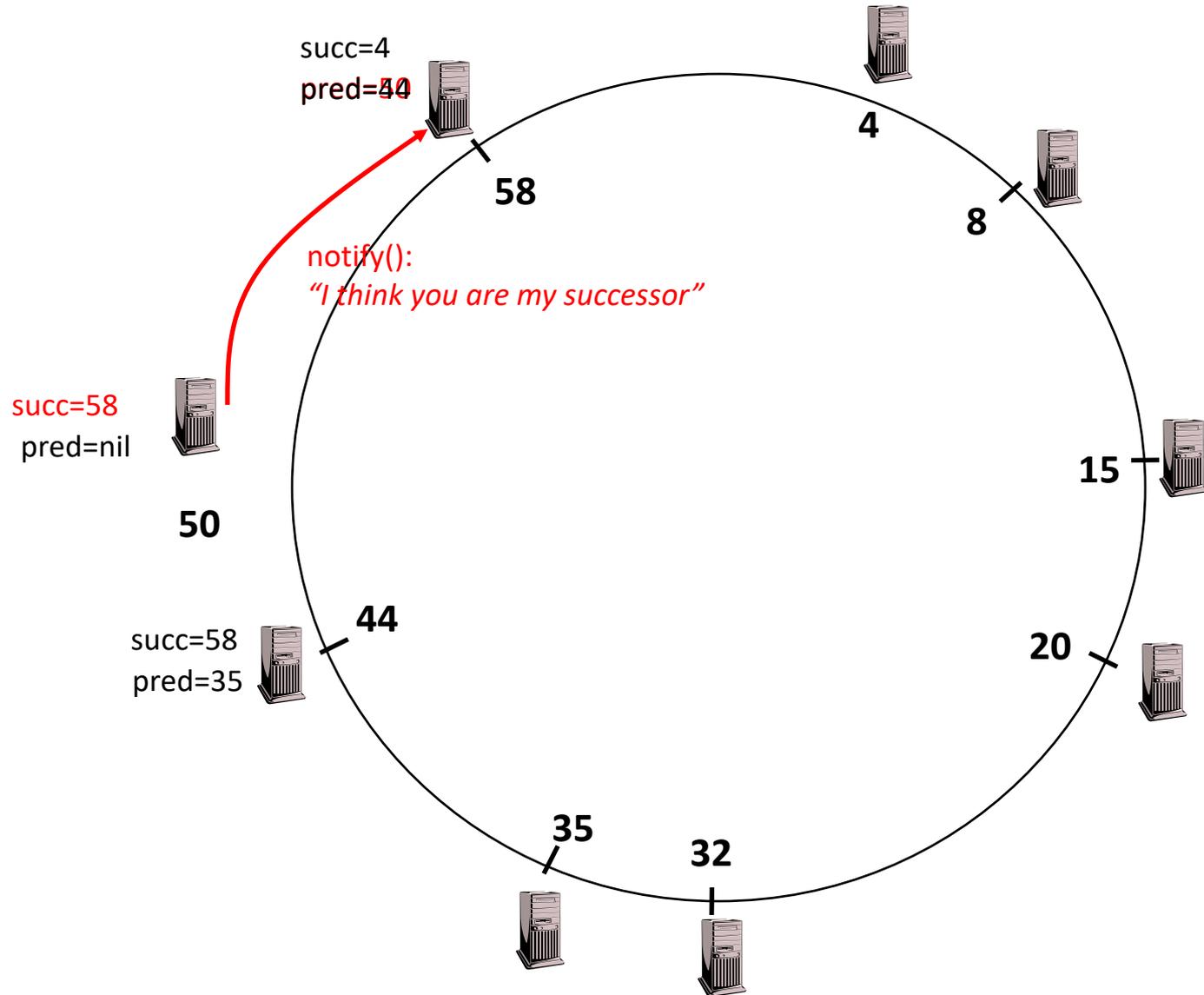
Joining Operation

- Node 50: send stabilize() to node 58
- Node 58:
 - Replies with 44
 - 50 determines it is the right predecessor



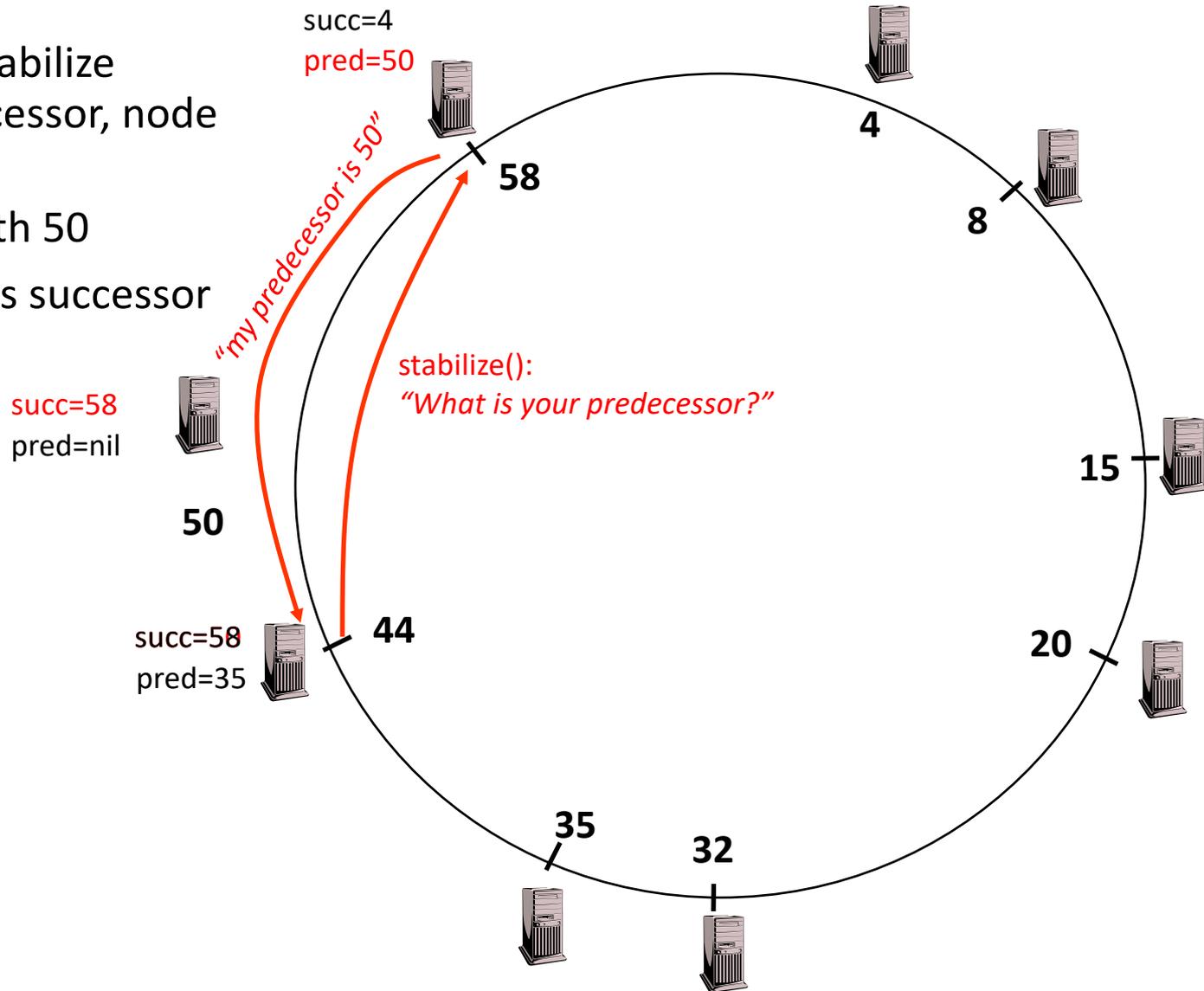
Joining Operation

- Node 50: send notify() to node 58
- Node 58:
 - update predecessor to 50



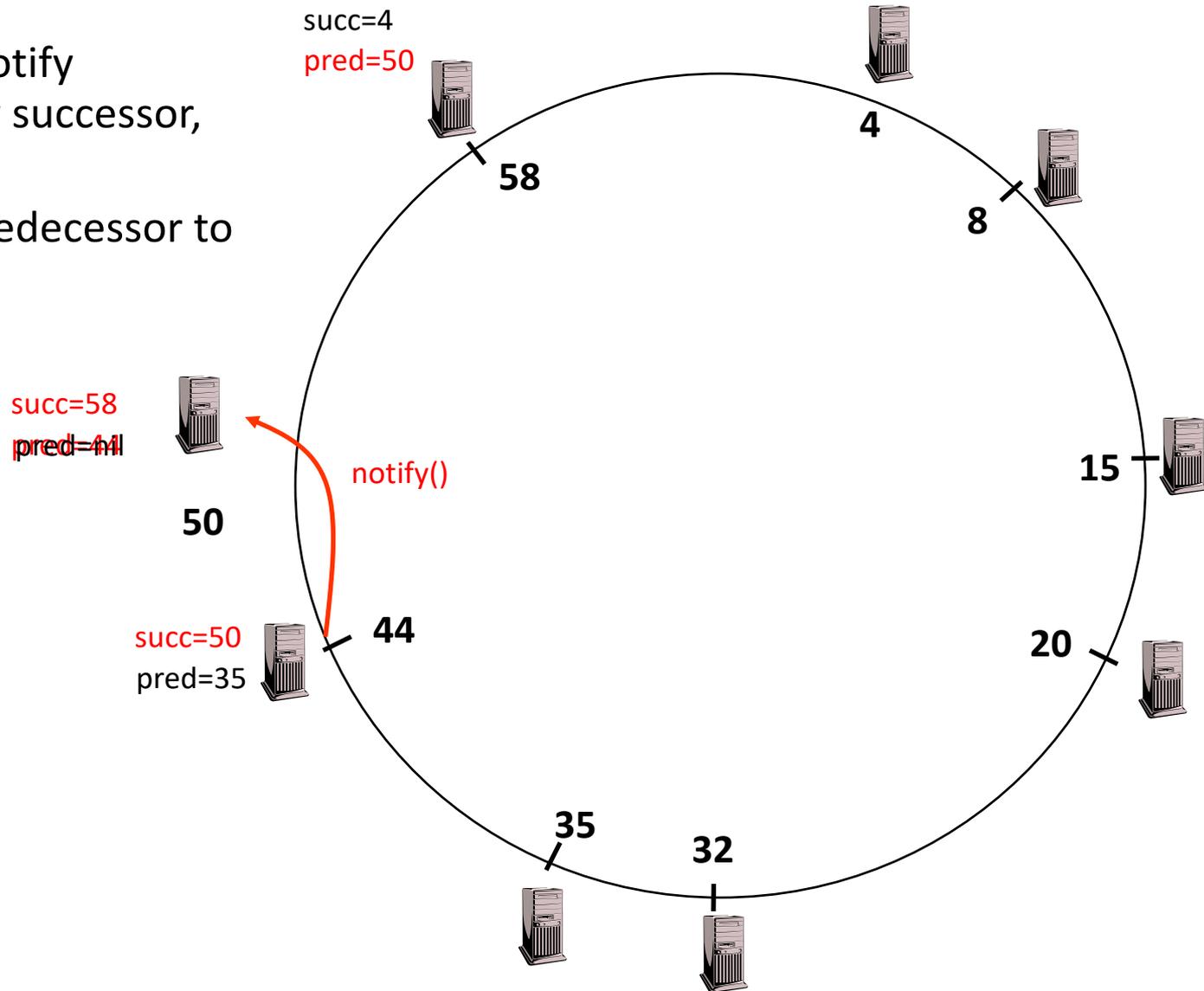
Joining Operation

- Node 44 sends a stabilize message to its successor, node 58
- Node 58 replies with 50
- Node 44 updates its successor to 50



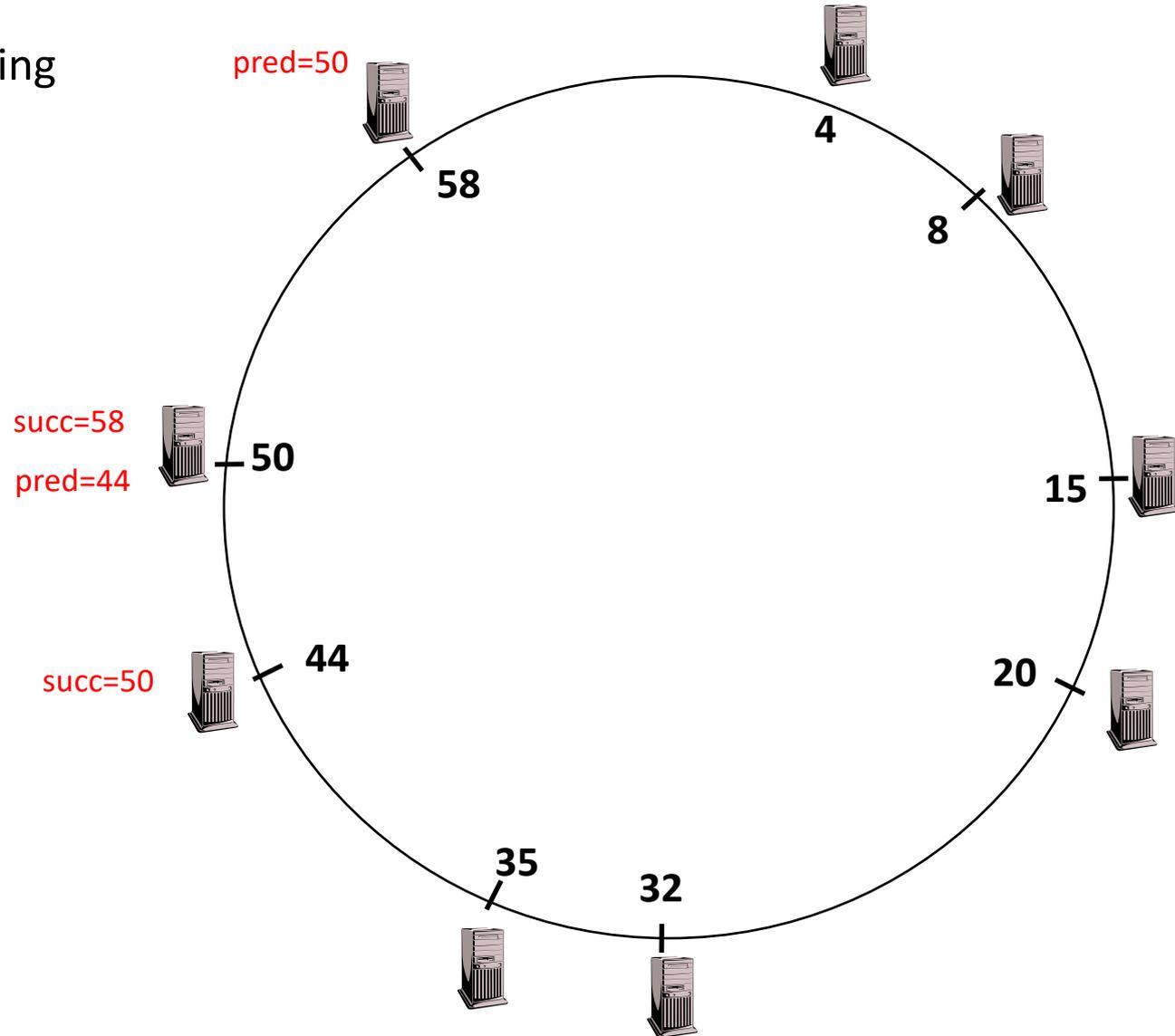
Joining Operation

- Node 44 sends a notify message to its new successor, node 50
- Node 50 sets its predecessor to node 44



Joining Operation (cont'd)

- This completes the joining operation!

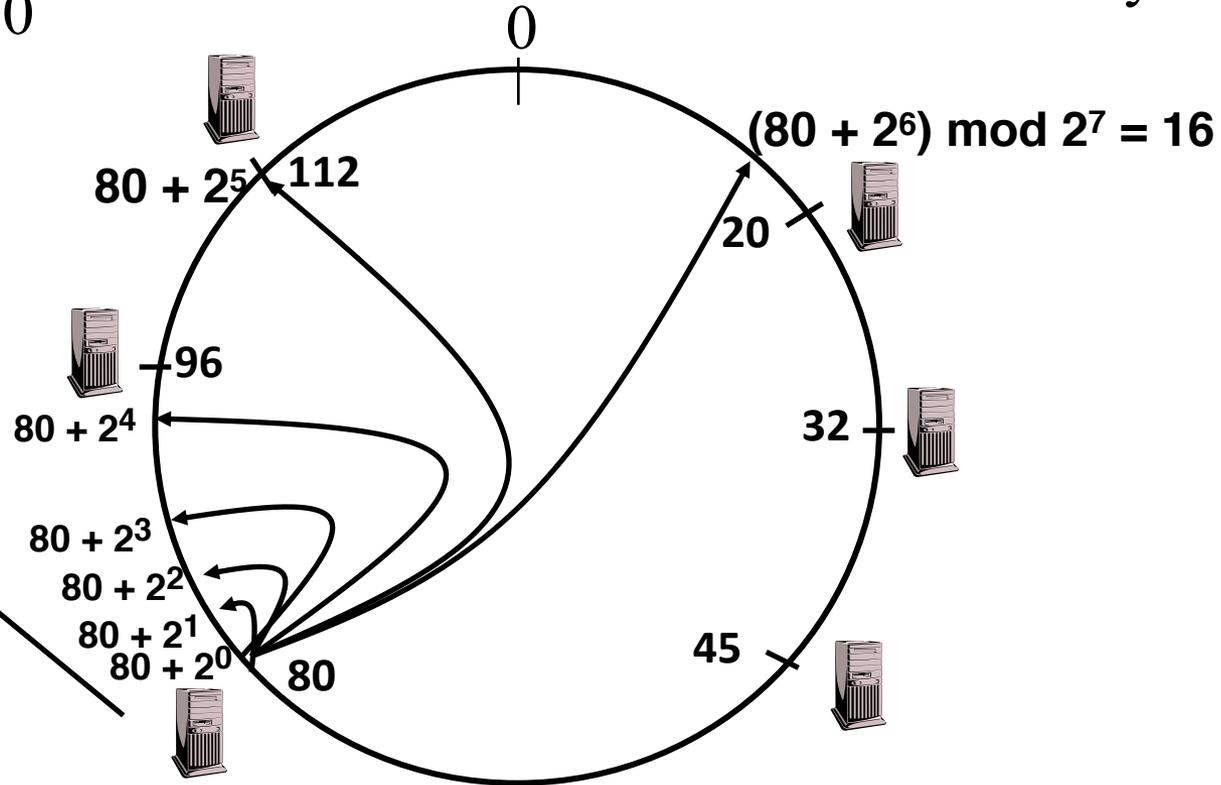


Achieving Efficiency: *finger tables*

Say $m=7$

Finger Table at 80

| i | $ft[i]$ |
|-----|---------|
| 0 | 96 |
| 1 | 96 |
| 2 | 96 |
| 3 | 96 |
| 4 | 96 |
| 5 | 112 |
| 6 | 20 |



i th entry at peer with id n is first peer with id $\geq n + 2^i \pmod{2^m}$



Chord

- **There is a tradeoff between routing table size and diameter of the network**
 - You can achieve diameter $O(1)$ with $O(n)$ -entry routing tables
 - Max diameter with $O(1)$ routing tables (random walks)
- **Chord achieves diameter $O(\log n)$ with $O(\log n)$ -entry routing tables**



Many other DHTs

- **CAN**
 - Routing in n-dimensional space
- **Pastry/Tapestry/Bamboo**
 - (Book describes Pastry)
 - Names are fixed bit strings
 - Topology: hypercube (plus a ring for fallback)
- **Kademlia**
 - Similar to Pastry/Tapestry
 - But the ring is ordered by the XOR metric
 - Used by BitTorrent for distributed tracker
- **Viceroy**
 - Emulated butterfly network
- **Koorde**
 - DeBruijn Graph
 - Each node connects to $2n, 2n+1$
 - Degree 2, diameter $\log(n)$
- ...



Discussion

- **Query can be implemented**
 - Iteratively: easier to debug
 - Recursively: easier to maintain timeout values
- **Robustness**
 - Nodes can maintain ($k > 1$) successors
 - Change notify() messages to take that into account
- **Performance**
 - Routing in overlay can be worse than in the underlay
 - Solution: flexibility in neighbor selection
 - Tapestry handles this implicitly (multiple possible next hops)
 - Chord can select any peer between $[2^n, 2^{n+1})$ for finger, choose the closest in latency to route through



Where are they now?

- **Many P2P networks shut down**
 - Not for technical reasons!
 - Centralized systems work well (or better) sometimes
- **But...**
 - Vuze network: Kademlia DHT, millions of users
 - Concepts incorporated into many systems (e.g., Amazon's DynamoDB)
 - Skype used to use a P2P network similar to KaZaA
- **Shown that you can have scalable routing *without* hierarchy**



Where are they now?

- **DHTs allow coordination of MANY nodes**
 - Efficient *flat* namespace for routing and lookup
 - Robust, scalable, fault-tolerant
- **If you can do that**
 - You can also coordinate co-located peers
 - Now dominant design style in datacenters
 - E.g., Amazon's Dynamo storage system
 - DHT-style systems everywhere
- **Similar to Google's philosophy**
 - Design with failure as the common case
 - Recover from failure only at the highest layer
 - Use low cost components
 - Scale out, not up



An alternative for reliability

- **Erasur coding**
 - Assume you can detect errors
 - Code is designed to tolerate entire missing packets
 - Collisions, noise, drops because of bit errors
 - Forward error correction
- **Examples: Reed-Solomon codes, LT Codes, Raptor Codes**
- **Property:**
 - From K source frames, produce $B > K$ encoded frames
 - Receiver can reconstruct source with *any* K' frames, with K' *slightly* larger than K
 - Some codes can make B as large as needed, on the fly



Erasure Codes

- **Motivation: scalability of reliable multicast**
 - Problem: in large multicast groups, where each receiver misses specific packets, how to coordinate retransmissions?
- **Erasure codes:**
 - Any K out of N messages reconstruct original content
- **Initially:**
 - Fixed-rate codes (e.g. Reed-Solomon ~ 1960)
 - Solve for polynomial of degree K with N linearly independent equations



LT Codes

- **Luby Transform Codes**
 - Michael Luby, circa 1998
- **Encoder: repeat B times**
 1. Pick a degree d (*)
 2. Randomly select d source blocks. Encoded block $t_n =$ XOR of selected blocks

* The degree is picked from a distribution, *robust soliton distribution*, that guarantees that the decoding process will succeed with high probability



More on encoding

- **Picking the degree d of encoded blocks**
 - Robust Soliton Distribution
 - Balances the probability that there is at least one block of degree 1 in each decoding iteration
 - While trying to minimize the probability of decoding failing
- **In practice, you don't encode the list of source blocks on each block, but the state of a pseudo-random number generator**
 - From this you can generate the next numbers in the sequence: degree d , and the next ids of the d source blocks in the encoded block

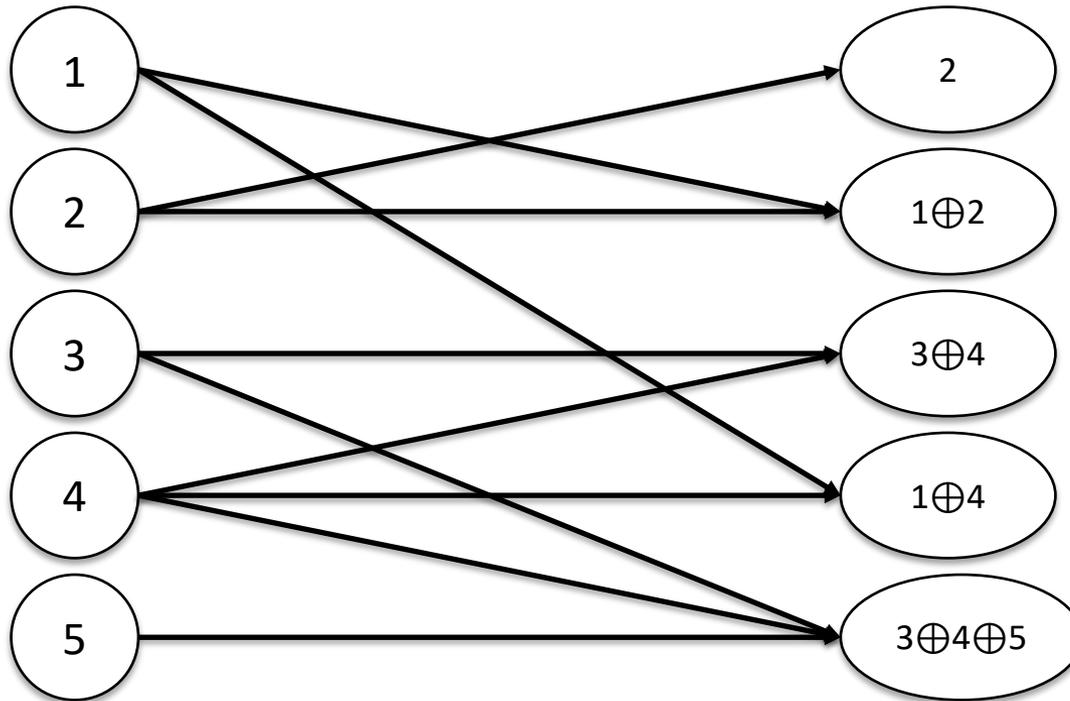


LT Decoder

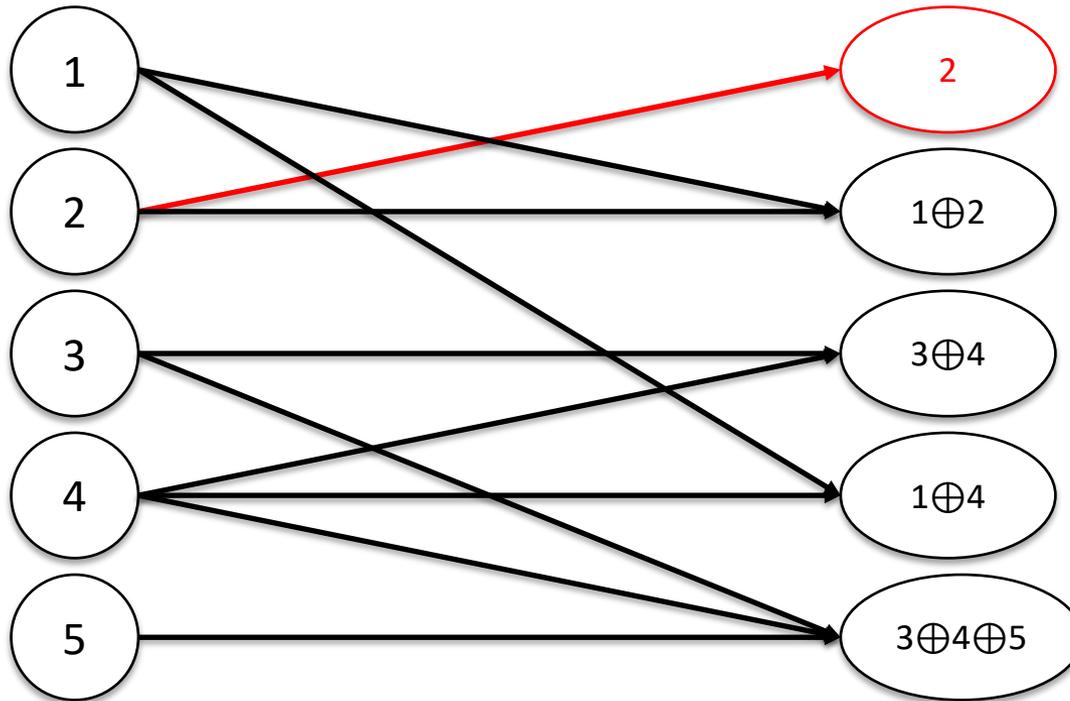
- Find an encoded block t_n with $d=1$
- Set $s_n = t_n$
- For all other blocks t_n , that include s_n ,
set $t_n' = t_n \text{ XOR } s_n$
- Delete s_n from all encoding lists
- Finish if
 1. You decode all source blocks, or
 2. You run out of blocks of degree 1



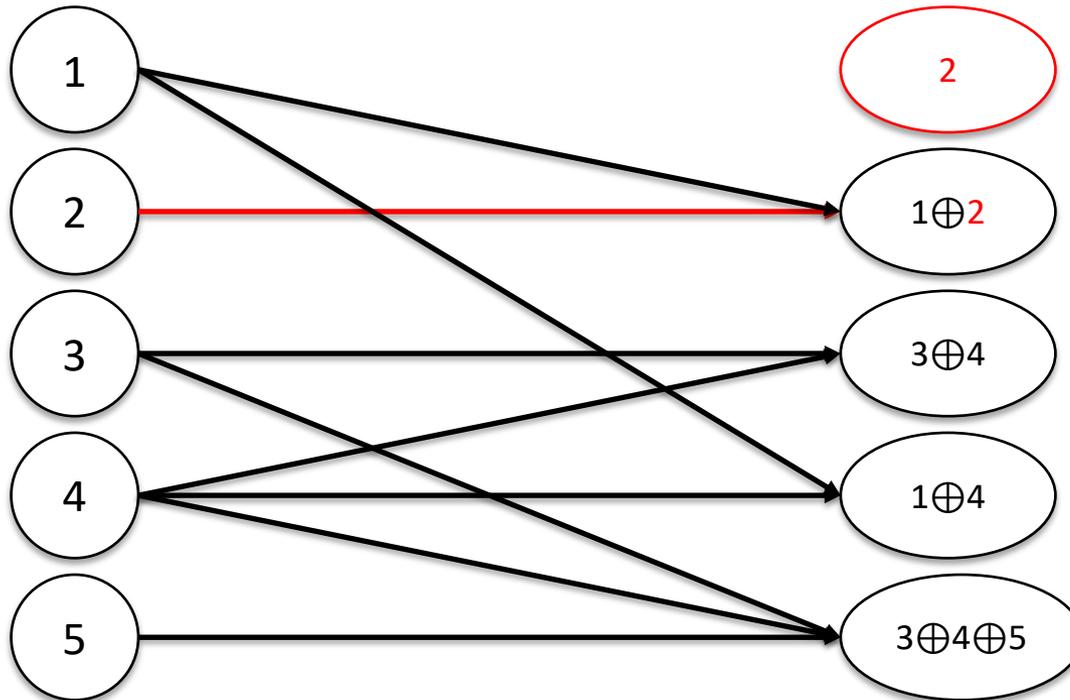
Decoding Example



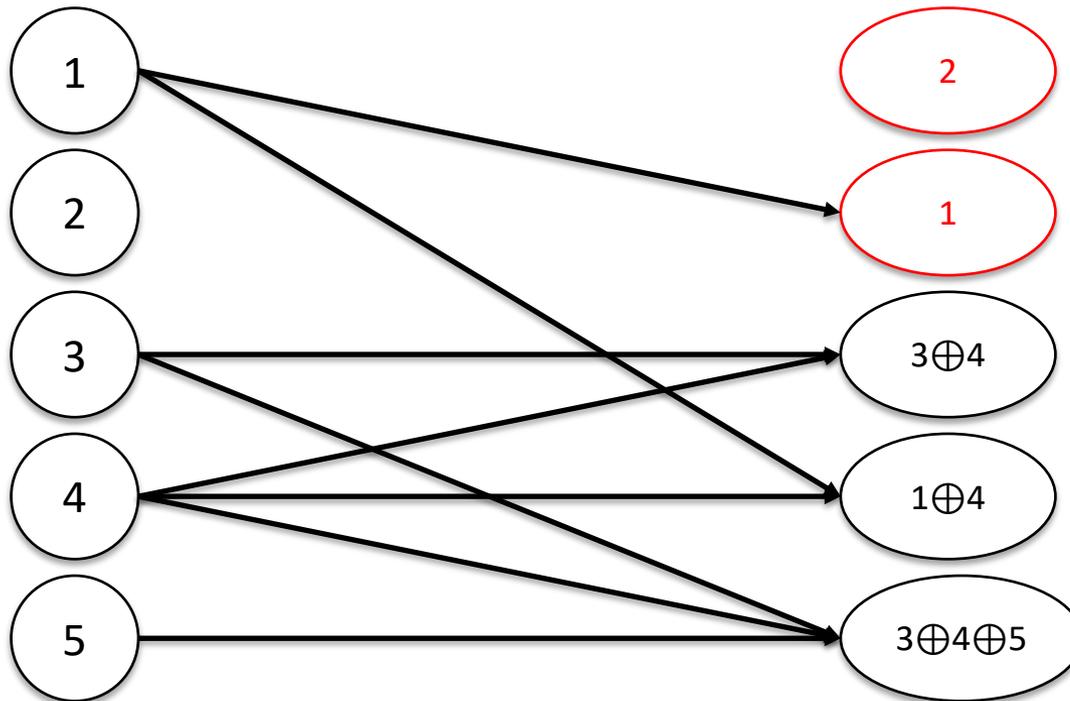
Decoding Example



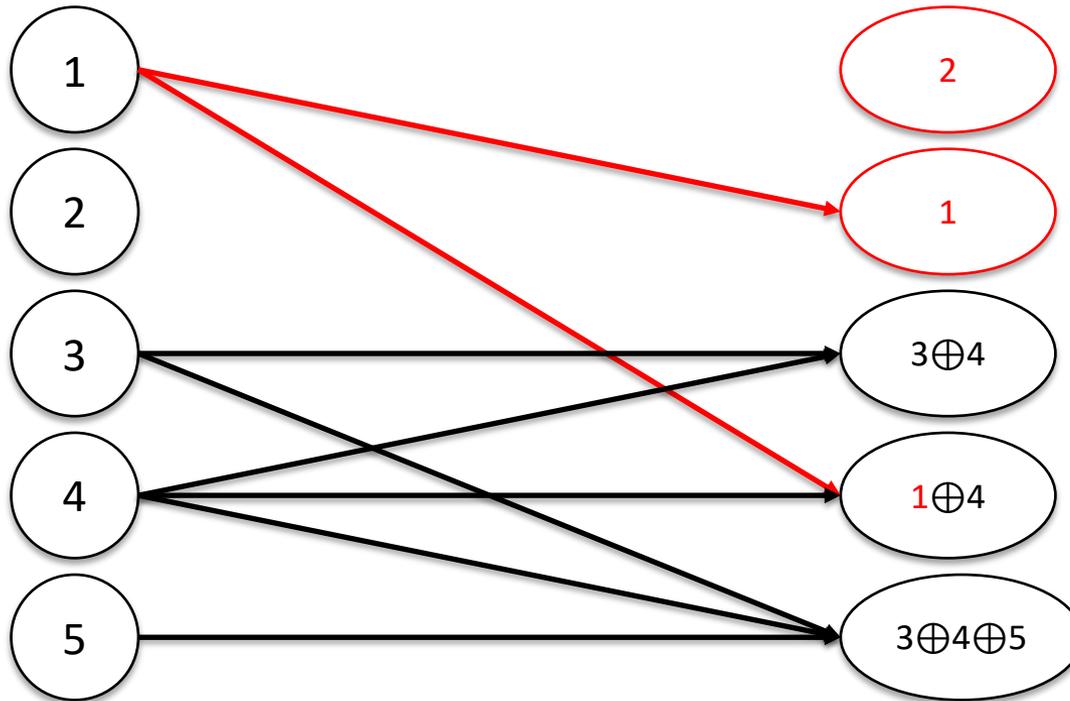
Decoding Example



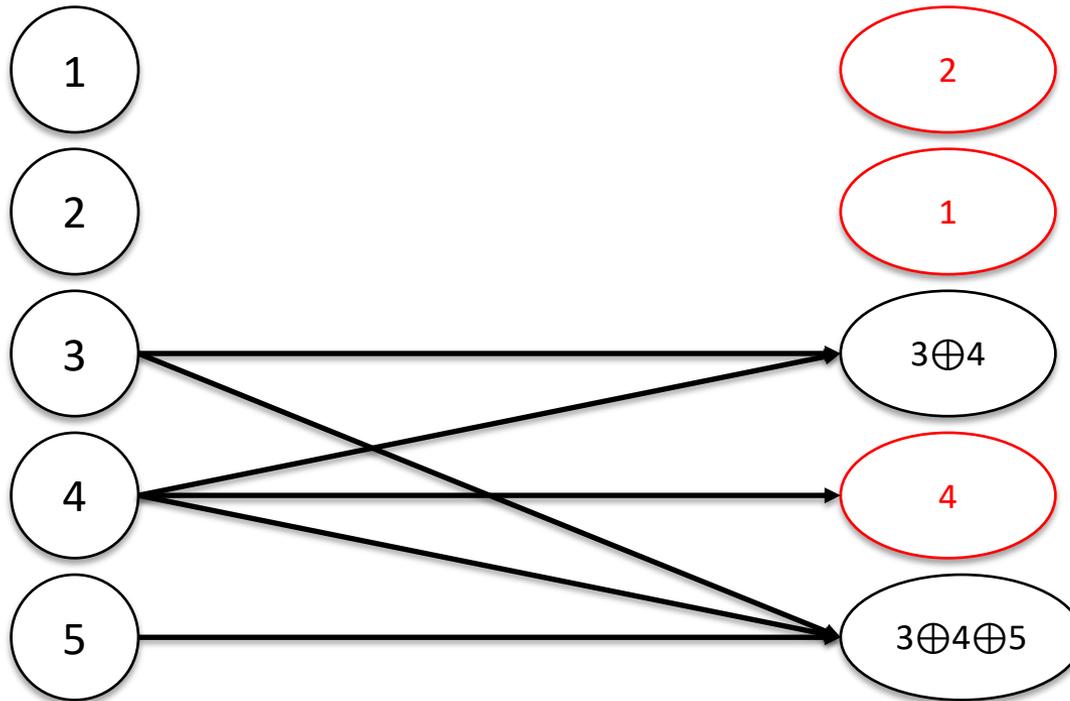
Decoding Example



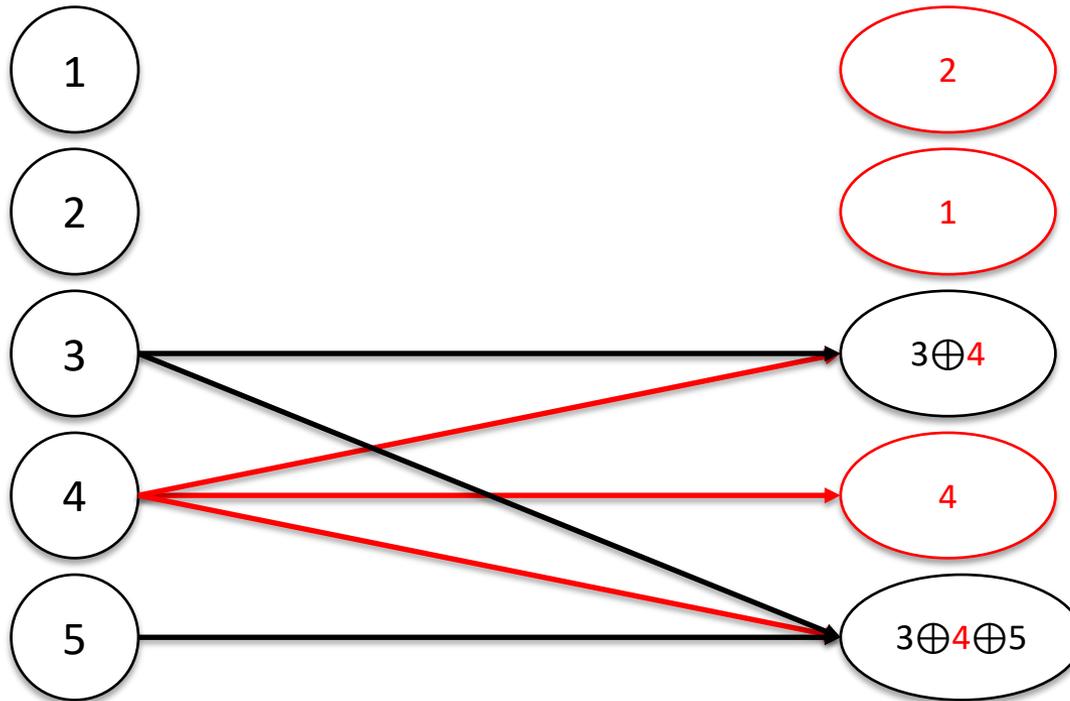
Decoding Example



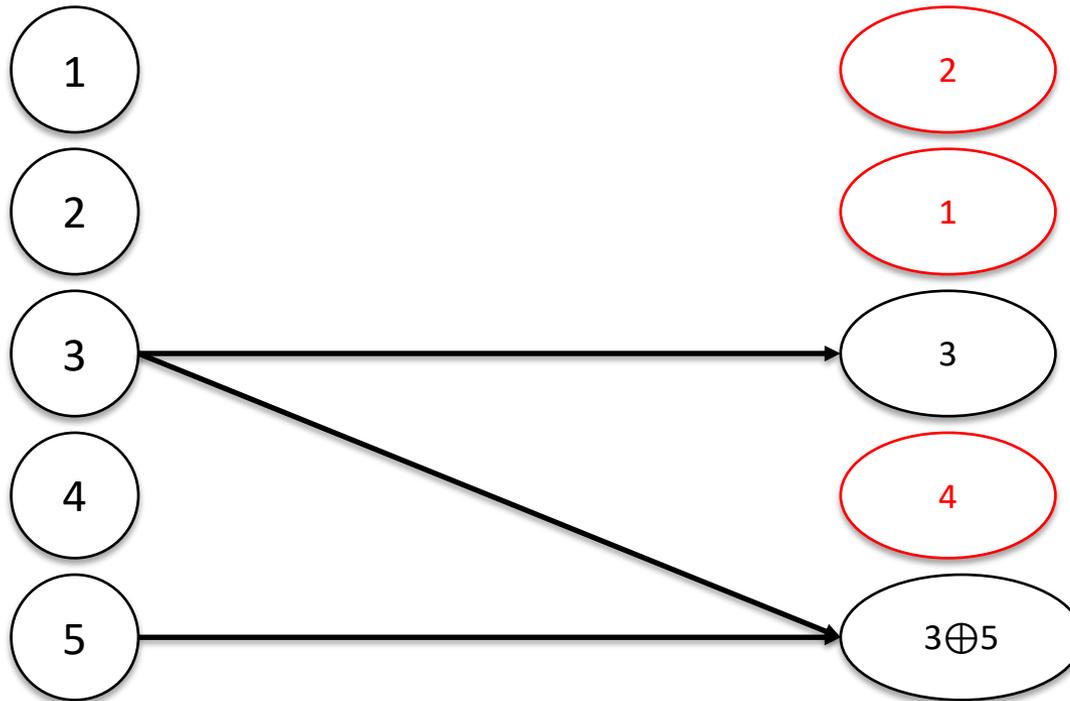
Decoding Example



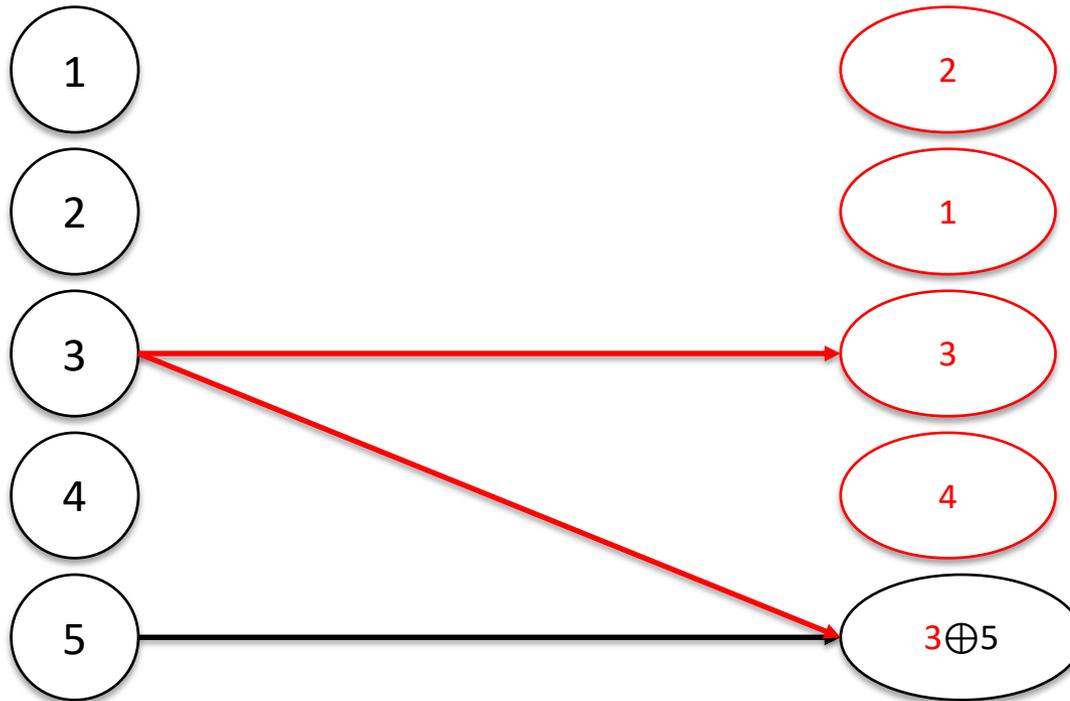
Decoding Example



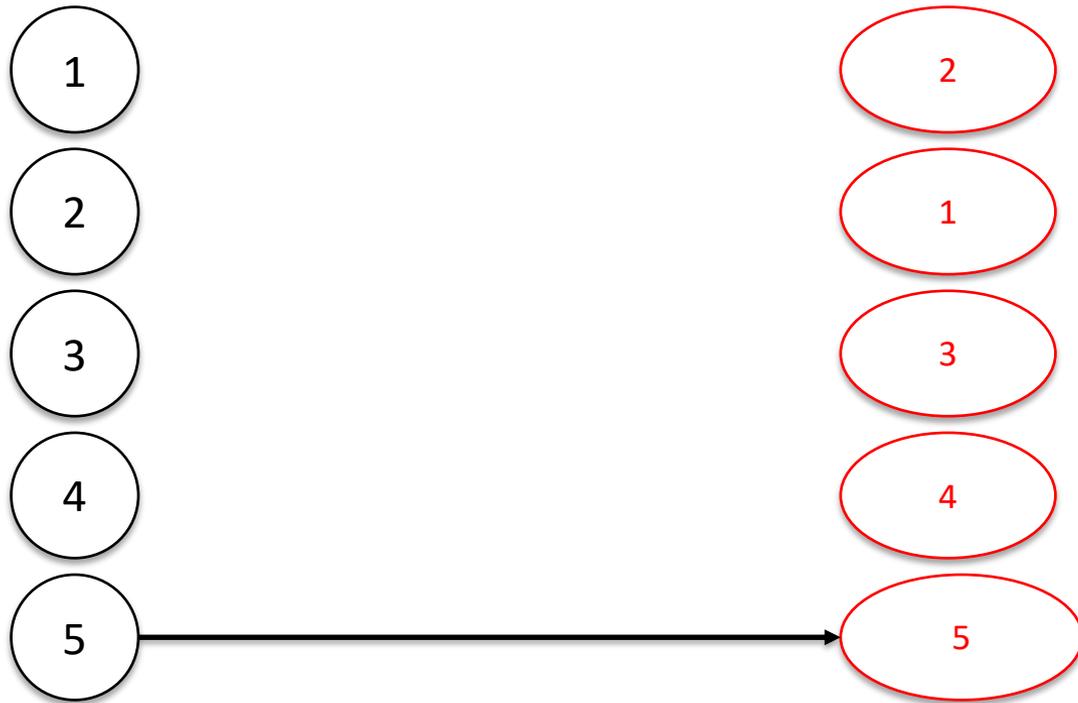
Decoding Example



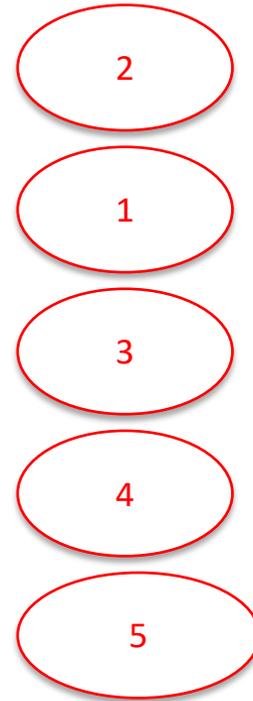
Decoding Example



Decoding Example



Decoding Example



Uses

- **IPTV, defense, postal service, satellite systems**
- **3G/4G/5G Multicast service (Qualcomm)**
- **Storage systems**
- **...**

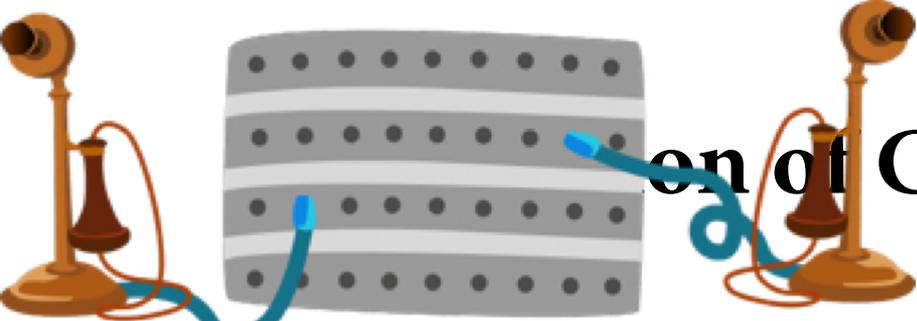


Information-Centric Networking (ICN) **(a vision)**

Named-Data Networking (NDN)
(a specific architecture proposal)

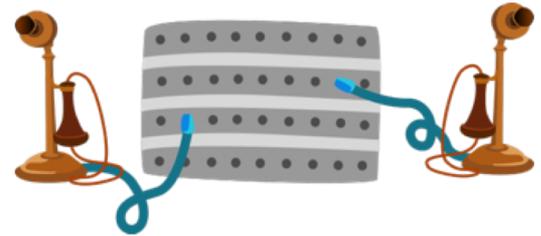
Content-Centric Networking (CCN)
(an earlier project)



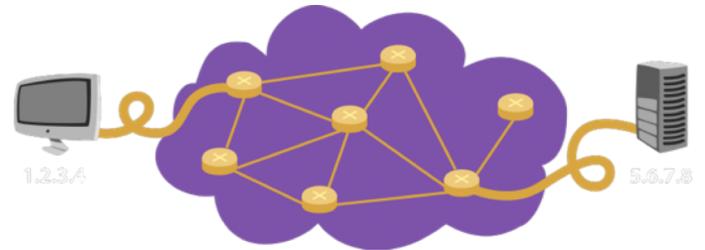


Evolution of Communications

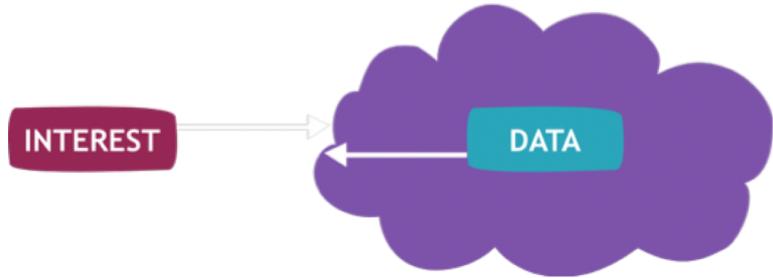
- **Telegraph, Telephony: pairwise communication**

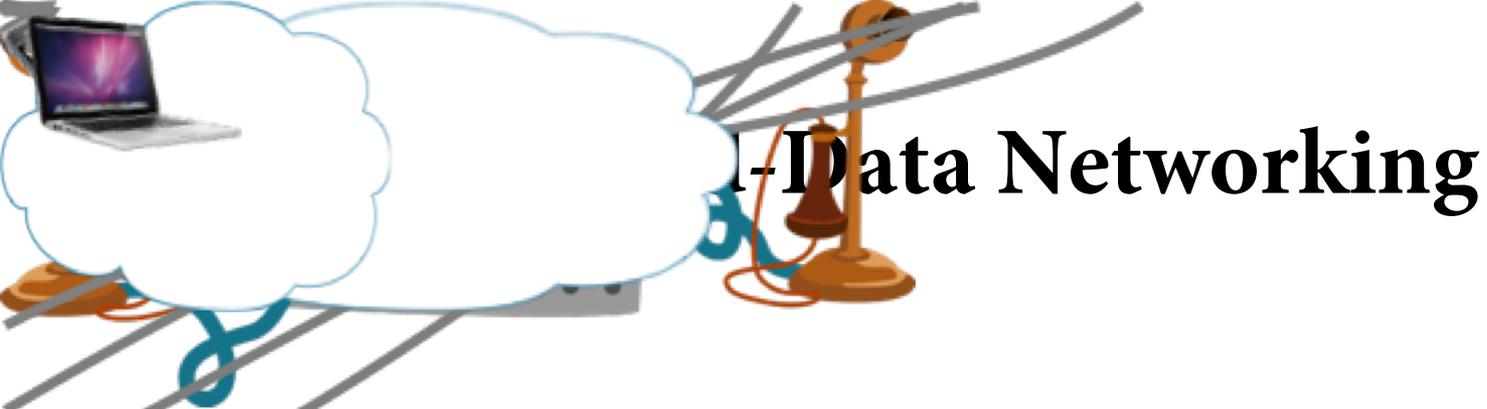


- **IP Networking: node-to-node communication**



- **Today: most users retrieve objects, don't care which server they come from**





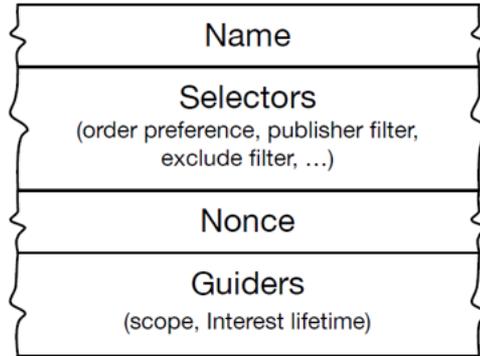
Data Networking

- Focus on *data* rather than on *endpoints*
- All content is named
- In-network storage and multicast arise naturally
- Secure the data rather than the process
 - Each data packet is immutable and signed

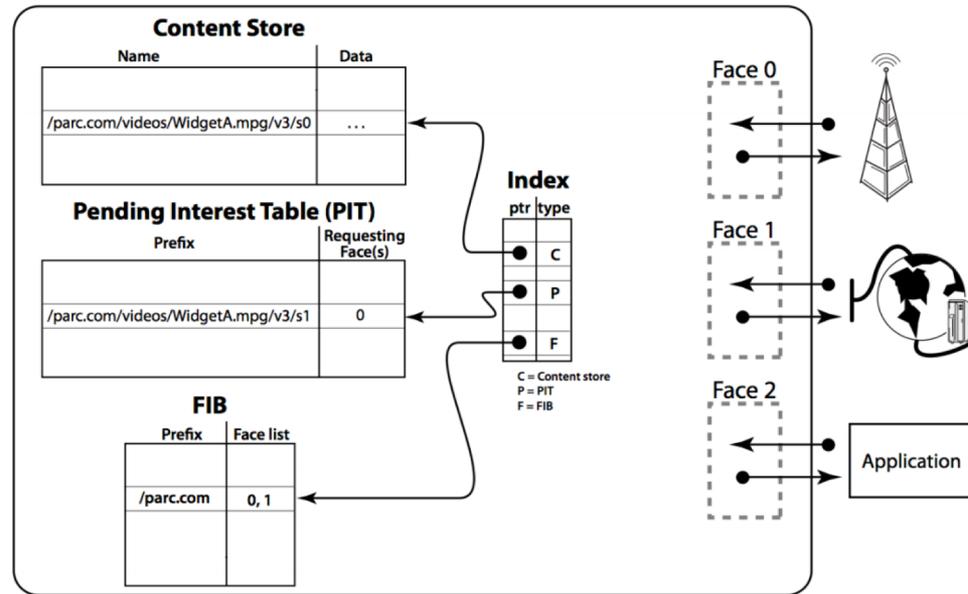
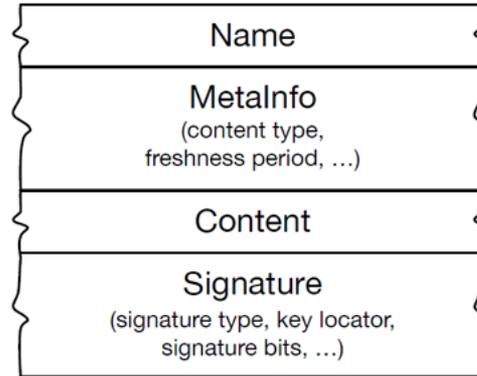


How does this work?

Interest Packet

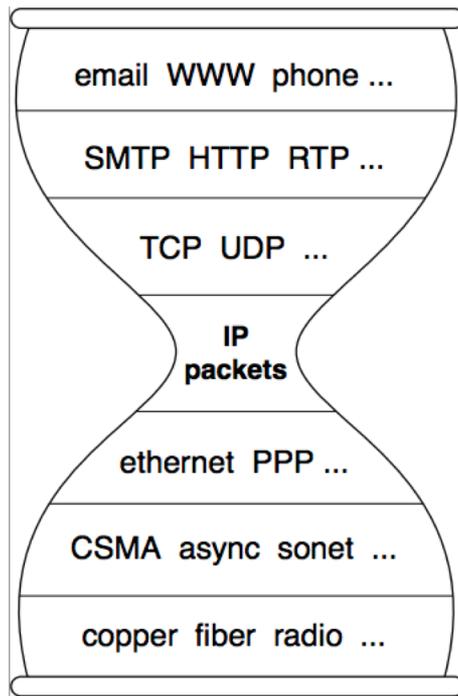


Data Packet

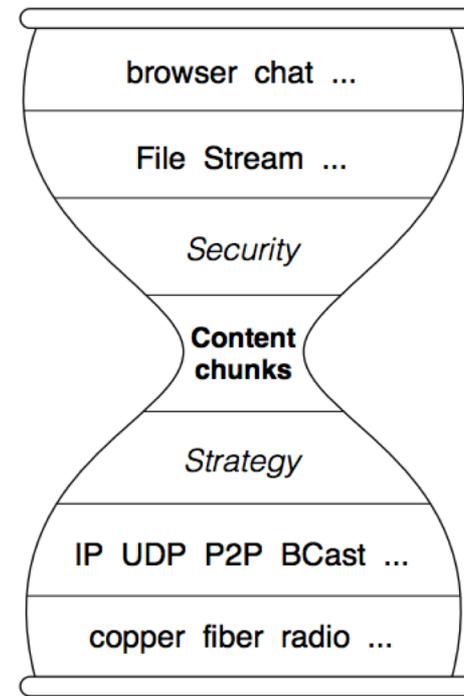


A new architecture

Today's architecture



NDN architecture

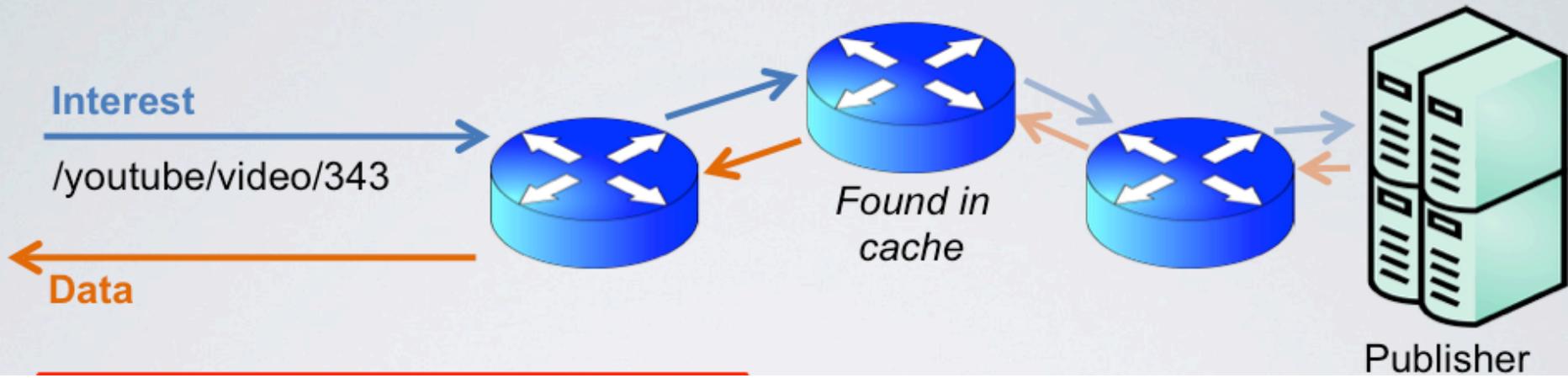


Individual apps

Every node

Individual links





Some details (and questions)

- **Names are hierarchical**
 - e.g. /brown.edu/courses/cs168/f19/videos/l25.mpg/5
 - Can name anything (including endpoints)
- **Routing can work similarly to IP prefix-based routing**
 - Aggregation on prefixes, longest-prefix matching
- **Signatures enable caching anywhere**
 - Hierarchical names provide context for trust management
- **Pull-based model**
 - Not “always on”, no unsolicited packets
 - Eliminates some types of DDoS attacks



Will this work?

- **Many challenges**
 - Does the current architecture work “well enough”?
 - Can we route efficiently on names of unbounded length?
 - How does trust management work? Yet another PKI?
 - What is the role of CDNs?
- **Proponents view this as the underlying architecture in 20 years, with IP-like communications as a special case**
 - Much like telephony today is a special case of IP communications

