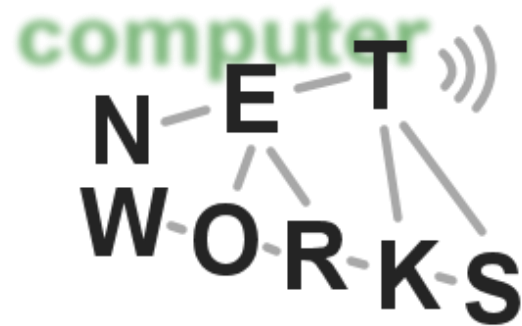# P2P-Networks

David Koll

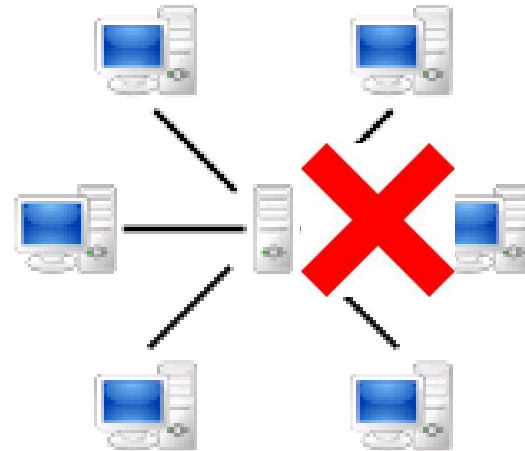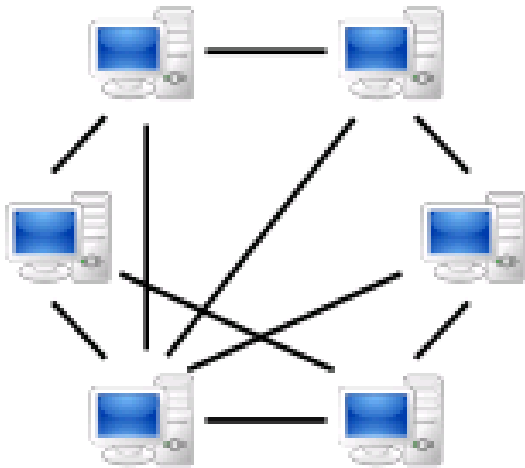Advanced Computer Networks

Summer Semester 2013

# Introduction to P2P Systems

- **What?**
  - What is a P2P system?
- **Why?**
  - Why do we need P2P systems?
- **How?**
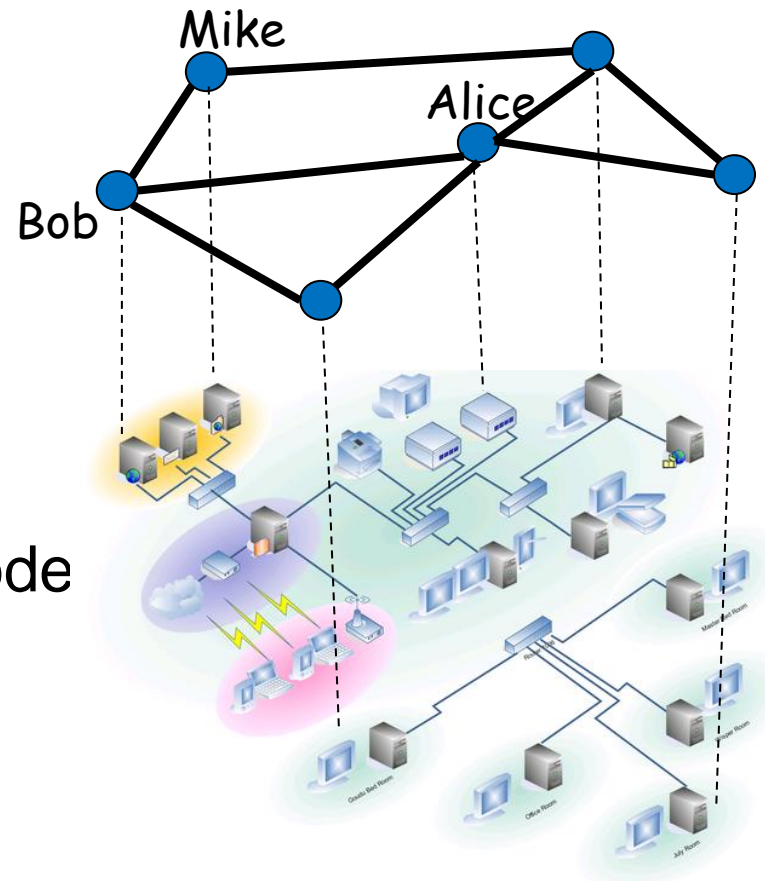  - How are P2P systems working?

# What is a P2P System?

o Peer-to-peer (abbreviated to P2P) refers to a computer network in which each computer in the network can act as a client or server for the other computers in the network, allowing shared access to files and peripherals without the need for a central server [Wiki]



o The sharing of computer resources by direct exchange, rather than requiring the intermediation of a centralized server.

# Features



o Decentralized:

  o No central component

o Role: "all peers are equal"

o Self-organized

  o Highly dynamic behavior of node

    • Free to come, free to go

o Overlay Network

  • A network built on the top of physical network

  • Nodes are connected by logical links

  • Flat system architecture

# Features (Cont.)

- Large-scale resources
  - Heterogeneous
  - Millions of nodes
- Collaboration
  - Based on voluntary participation
  - Global reach
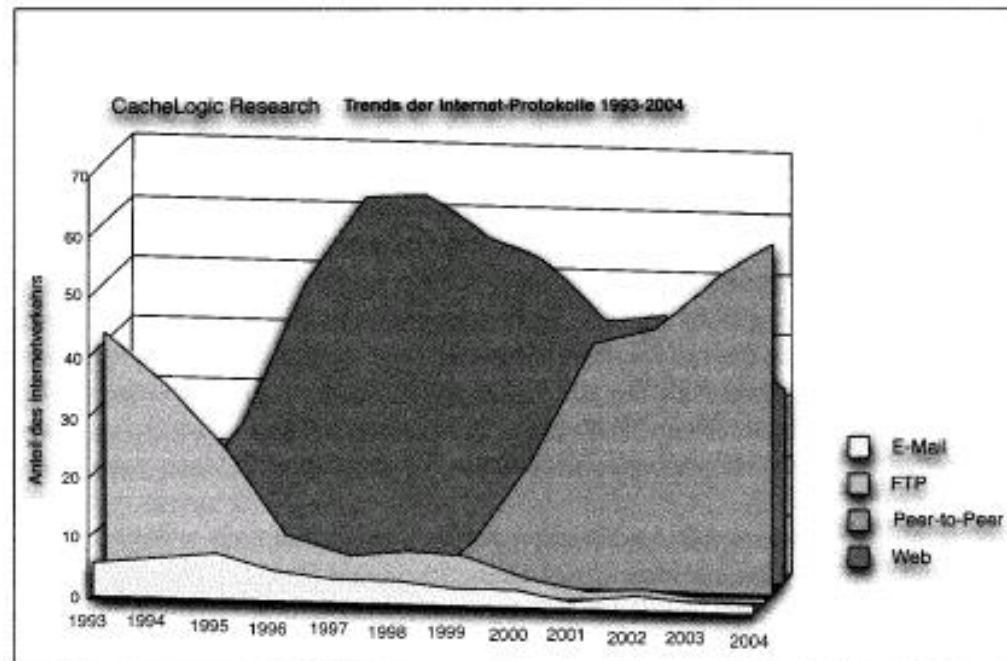- Flexible, resilient to attacks, anonymous
- …

# Why P2P?

# User Requirement of Resource Sharing

○ Text, email: 100 KB

○ Picture: 1MB

○ Music: 10MB

○ Movie: 1GB

○ HD Video: 10GB+

# User Requirement of Resource Sharing

o 2004: More than half of Internet traffic P2P

o 2010: Still 39%, 2014: ~20%

o Absolute P2P traffic increases (by 100% 2010-2014)



CacheLogic Research — Trends der Internet-Protokolle 1993-2004

Legend: E-Mail, FTP, Peer-to-Peer, Web

# Developing of Network Nandwidth

o Maximum bandwidth of common Internet access technologies [Wiki]


o Current network bandwidth cannot satisfy user demand


o Need a more efficient way to share resources

| 56 kbit/s | Modem / Dialup |
|---|---|
| 1.5 Mbit/s | ADSL Lite |
| 1.544 Mbit/s | T1/DS1 |
| 10 Mbit/s | Ethernet |
| 11 Mbit/s | Wireless 802.11b |
| 44.736 Mbit/s | T3/DS3 |
| 54 Mbit/s | Wireless 802.11g |
| 100 Mbit/s | Fast Ethernet |
| 155 Mbit/s | OC3 |
| 600 Mbit/s | Wireless 802.11n |
| 622 Mbit/s | OC12 |
| 1 Gbit/s | Gigabit Ethernet |
| 2.5 Gbit/s | OC48 |
| 9.6 Gbit/s | OC192 |
| 10 Gbit/s | 10 Gigabit Ethernet |
| 100 Gbit/s | 100 Gigabit Ethernet |

# Client/Server

o The client arrives and requests a service at any given point in time

o The server is dedicated to the service and responds to the client

Apple.com
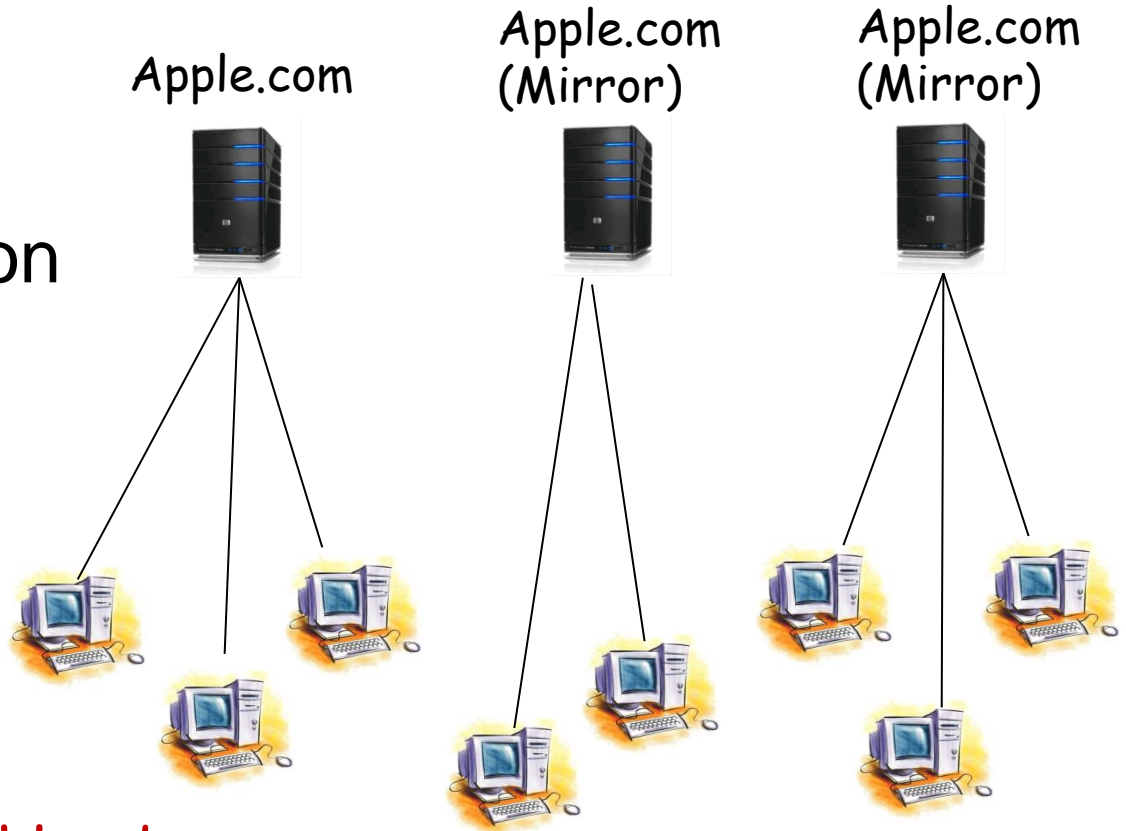
## Problems
- Hot spot-uneven workload
- Bottleneck: bandwidth, CPU, …
- Single point of failure
- Scalability
- ~~Maintenance~~

# Replication
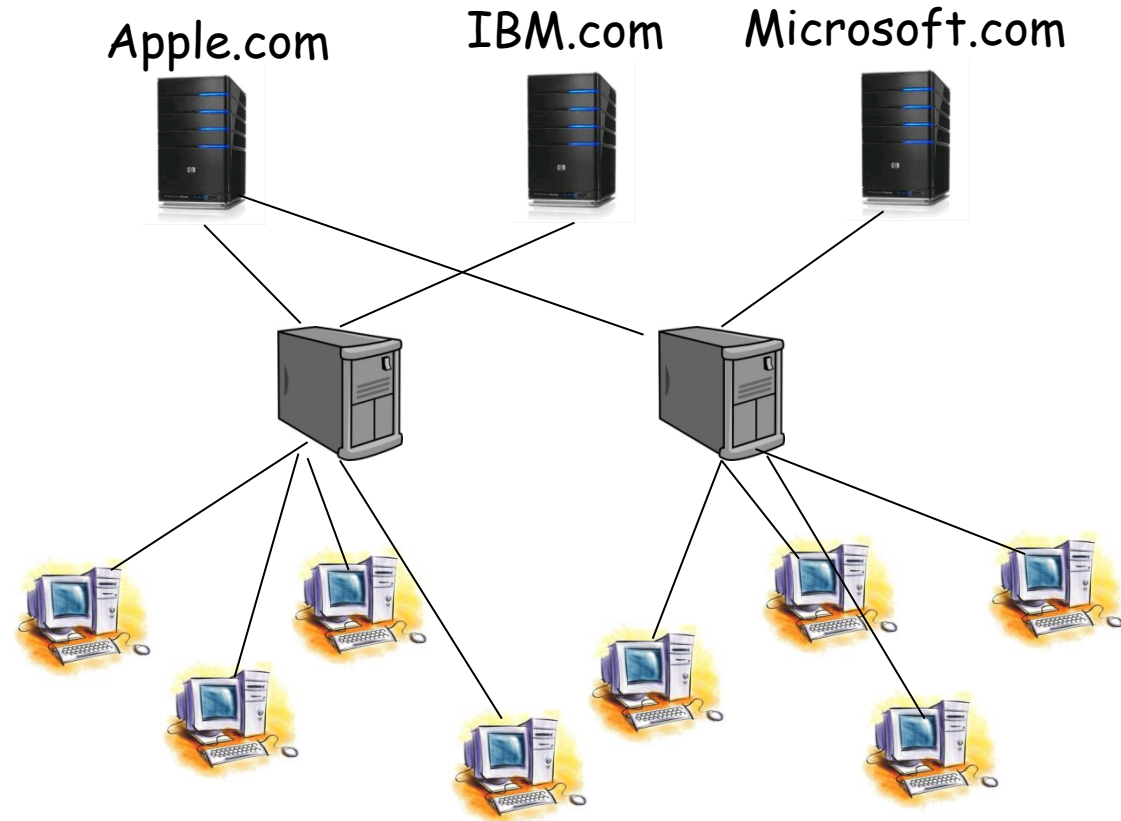
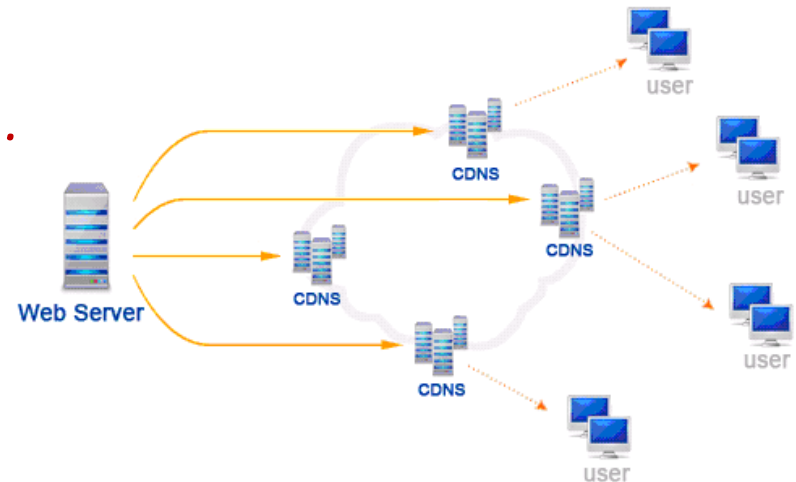Apple.com    Apple.com (Mirror)    Apple.com (Mirror)

o Server Replication

## Problems
- Hot spot-uneven workload
- Bottleneck: bandwidth, CPU, …
- ~~Single point of failure~~
- ~~Scalability~~
- Maintenance

# Proxy, CDN

Apple.com    IBM.com    Microsoft.com

## Problems

- ~~Hot spot uneven workload~~
- Bottleneck: bandwidth, CPU, ...
- ~~Single point of failure~~
- ~~Scalability~~
- Maintenance

Web Server    CDNS    CDNS    CDNS    CDNS    user

# P2P: Advantages

o Changes the way of network bandwidth use

o Easy to deploy, easy to use

o Dynamic for joining and leaving

o Distributed resource sharing

    o Files, data, storage, computation, …

    o Provide something useful and free

    o Anyone can contribute

o Fault tolerant

o Service ability: large scale

o Service of quality: the more users, the better

# P2P: How?

# P2P Applications and Systems

o File sharing

   o Napster, Gnutella, BitTorrent

o Multimedia streaming

   o P2P TV: PeerCast, PPlive, PPStream, TVUZattoo, ….

   o P2P based VOD systems

o Communication

   o Skype, MSN, …

o Computation

   o SETI@home: Search for Extra-Terrestrial Intelligence
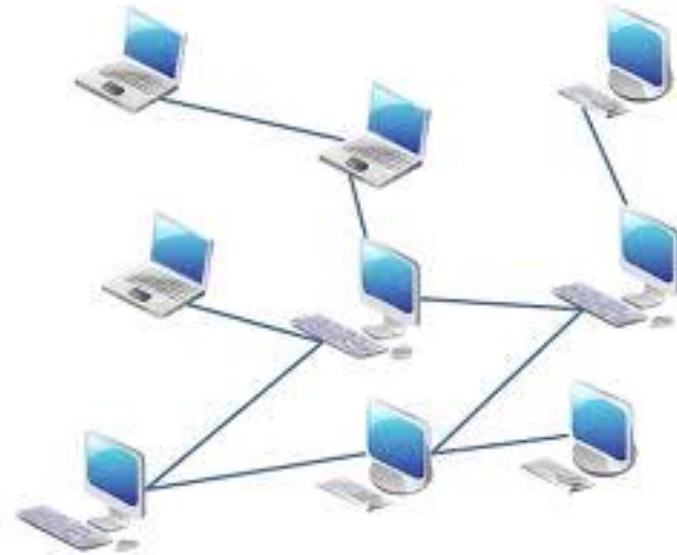
# Current State of P2P

o P2P applications are popular over the world

o P2P networks are mainly used for resource sharing

  o Music, videos, software, …

  o Some are illegal copyrighted materials

o New emerging applications

  o Online media streaming, P2P TV

  o P2P telephone system

  o Software installation and update

  o Decentralized social network applications

# Typical Research Topics

- Structure
  - How to search for information
    - Unstructured P2P
    - Structured P2P

- Security and privacy
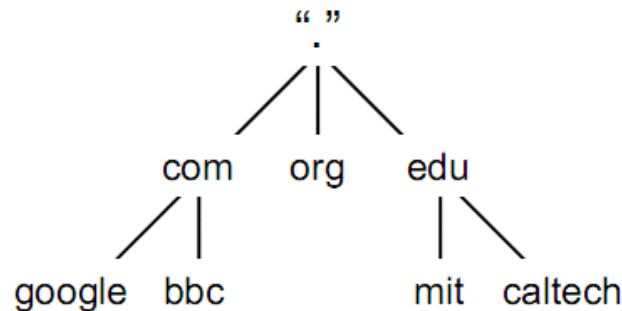  - How to protect system security and user privacy?

- Legal issues

# Search in P2P Networks

o How to locate resources in P2P networks?

# DNS- Domain Name System

○ Translates queries for domain names into IP addresses for the purpose of locating computer services and devices worldwide

○ Distributed database organized in hierarchical structure
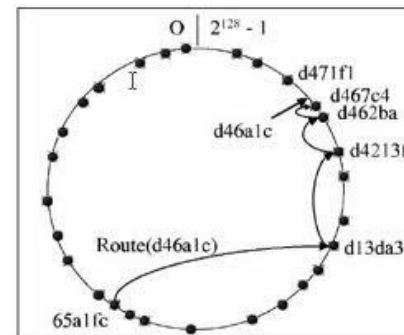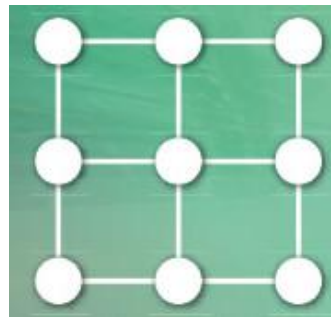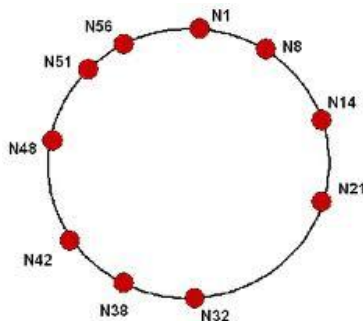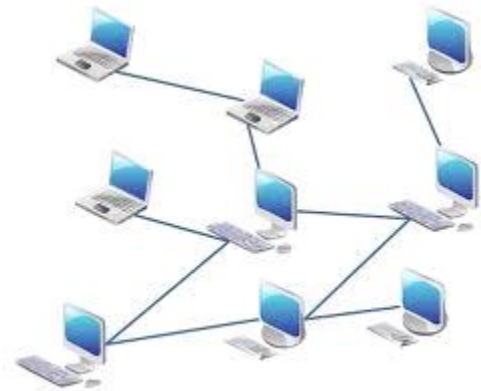
# Information Retrieval System

o Keyword-based

# Search in P2P Networks?

o Unstructured P2P

  o Highly flexible, dynamic, easy to maintain

  o Hard to find information

o Structured P2P

  o Hard to maintain its structure

  o Easy to find information

# Unstructured P2P Networks

# Unstructured P2P Networks

o Example: BitTorrent


o Successor of Napster, Gnutella, …

  o Napster: Pioneer P2P, shut down in 2001

   • easy to manage and search, but relied on central lookup server (drawbacks?)

   • data transfer directly between peers

  o Gnutella: Answer to Napster weaknesses

   • fully distributed P2P network based on overlay network, no central server

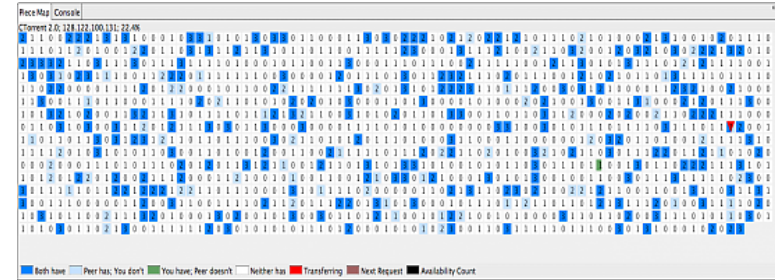   • Search: flooding the network with request (drawbacks?)

# BitTorrent

o A new popular approach to sharing large files

- o Origin of 37% of German internet traffic in 2009

o Originally used for distributing legal content

- o Linux distributions, software updates
- o Official movies

o Goal:

- o Quickly and reliably replicate one file to a large number of clients
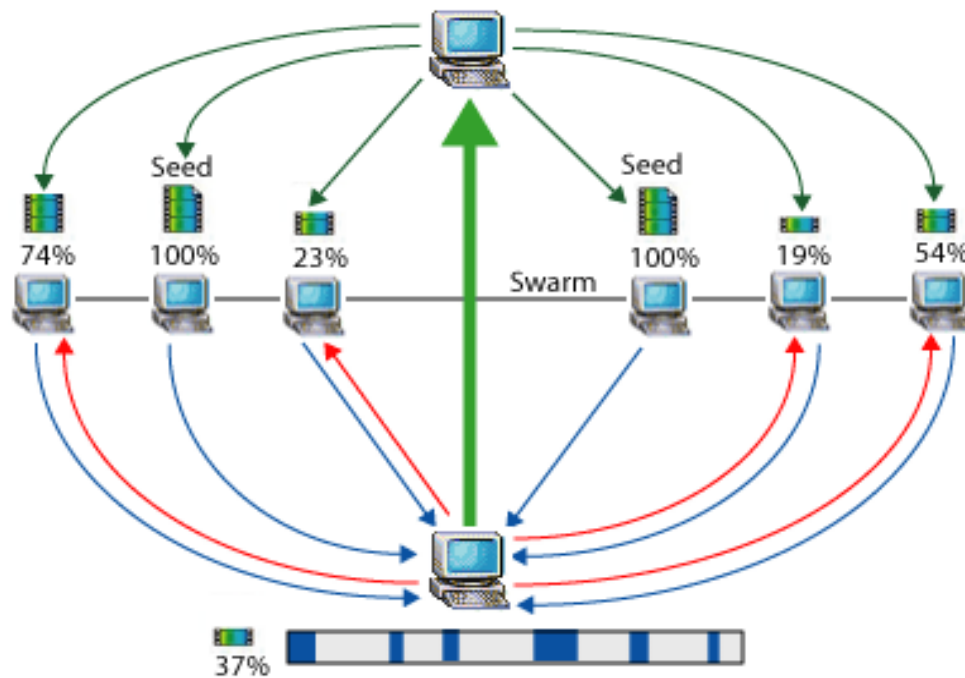
o Call it "P2P content distribution"

# Basic Idea



o Chunking:

- o Files split into smaller pieces or chunks
- o Chunks can be downloaded in parallel
- o Downloading order does not matter

o Swarming

- o Clients join a crowd of peers uploading and downloading the same content
- o Nodes request chunks from neighbors and download content in parallel

o Use a web server to publish content

o Use a central unit to locate resource

# Basic Idea



BitTorrent tracker identifies the swarm and helps the client software trade pieces of the file you want with other computers.

Seed

Seed

74%  100%  23%  Swarm  100%  19%  54%

37%

Computer with BitTorrent client software receives and sends multiple pieces of the file simultaneously.

©2005 HowStuffWorks

# Basic Components

o Web server: for content publication

o Tracker: a special central server for running the content distribution system

    o Tracking active peers

    o Mapping from file name to peers

o Peer

    o Seed: a peer with a complete copy of the file

    o Leecher: peer still downloading the file

o ".torrent" file: metadata and description of the file

    o The number of chunks

    o The tracker's IP

Web server

Tracker

Seed

Torrent-file

Tracker: 127.0.0.1
Chunks: 42
Chunk 1: 12345678
Chunk 2: 90ABCDEF

# **Operation**

Web server

Tracker

Torrent-file

Seed creates torrent-file
and hosts it somewhere

Seed
starts
tracker

Tracker: 127.0.0.1
Chunks: 42
Chunk 1: 12345678
Chunk 2: 90ABCDEF
...

Seed

o Sharing a file:
  - o (1) Seed generates a ".torrent" file from the file
  - o (2) Upload the ".torrent" file to some public web server or sending it to friends by email

o Searching a file:
  - o No dedicate search component
  - o User can search ".torrent" file from web server

o Downloading a file:
  - o (1) Download the ".torrent" file
  - o (2) Connect to the tracker to locate the file
  - o (3) Choose some fast peers to download chunks in parallel

# Tit-for-Tat Policy and Chunk Selection

o Tit-for-Tat policy

  o The more you give, the more you get

  o A peer serves peers that serve it

  o Encourages cooperation, discourage free-riding

o Chunk selection

  o Peers uses <span style="color:red">rarest first</span> policy when downloading chunks

  o Having a rare chunk makes peer attractive to others

  o The goal is to maximize availability of each chunk
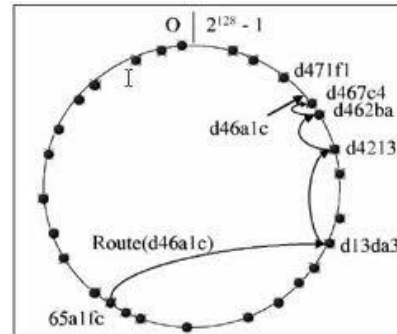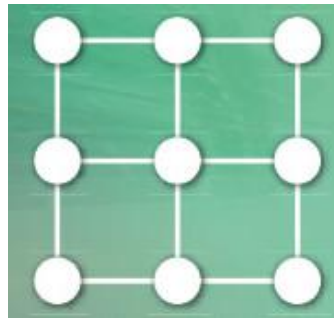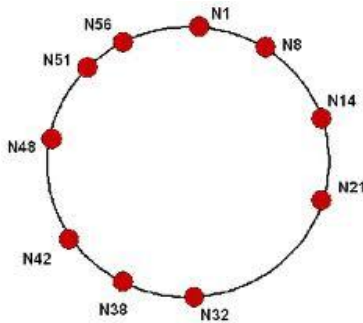
# BitTorrent : Pros and Cons

o Strengths

  o Works well for "hot content", very fast and resilient

  o Proficient in utilizing partially downloaded files

  o Discourages "free-riding"

  o Efficient for distributing large files to a large number of clients

o Weaknesses

  o Assumes all interested peers active at same time

  o Tracker could be single point of failure

  o Lack of search feature

# Structured P2P Networks

# Structured P2P Networks

o Routing & Lookup

o DHTs



o The following slides are based on a lecture by Prof. Roscoe, ETH Zürich, and provided with his kind permission

# Problem Space

o Challenge: spread lookup database among P2P participants

o Goals:

  o Scalable – operates with millions of nodes

  o Self-organized – no central, external control

  o Load-distributing – every member should contribute (at least ideally)

  o Fault tolerant – robust against node leaves or failures

  o Robustness – resiliance against malicious activity

# Idea

- o Distributed Hash Tables
    - o Hash content identifiers to machines
    - o Hash IP addresses
    - o Store content (or content locator) at machine with closest hash value

- o Originally 4 papers submitted to SIGCOMM 2001:
    - o CAN, Chord, Pastry, Tapestry

- o Widely used in practice (e.g., BitTorrent uses Kademlia DHT)

# Background: Hash Functions

o Hash function maps arbitrary input sequence to fixed length output:

   o H(m) = x, x of fixed length

o Crypto-Hashes:

   o Small input changes result in large output changes (Avalanche criterium)

   o If H(m1) = x is known, it is hard to find another m2 giving H(m2) (collision resistant)

o Inheritly hash functions span whole $2^k$ space (k bits hash length)

# MD5 / SHA-1

- Message Digest Algorithm 5
    - 128 bit hash values
    - Weak collisions found

- SHA-1 (similar to MD4)
    - 160 bit hash values
    - Stronger than MD5, but „under researcher's attack": find collisions in $2^{69}$

- But: Both algorithms efficiently map input homogeniously to $2^k$ space
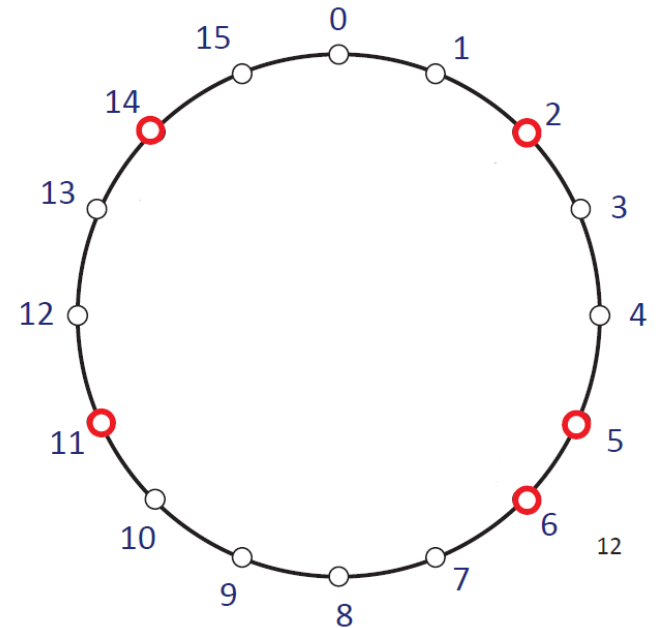
# DHTs

o Index data by hash value

o Assign each node in the network a portion of the hash address space

o DHT provides the lookup function

# Example: Chord

o Published 2001 at SIGCOMM by Stoica et al. „Chord: A Scalable Peer-to-peer Lookup Service for Internet Applications"

o Keys are SHA-1 hashes – 160 bit identifiers

o Key: Identifier of a data item

o Value: Identifier of a node

o Host (key,value) pair at node with ID larger or equal to key – successor(key)

# Identifier Space

o Identifier in $2^4$ space

- o Space from 0..15
- o Nodes pick IDs:
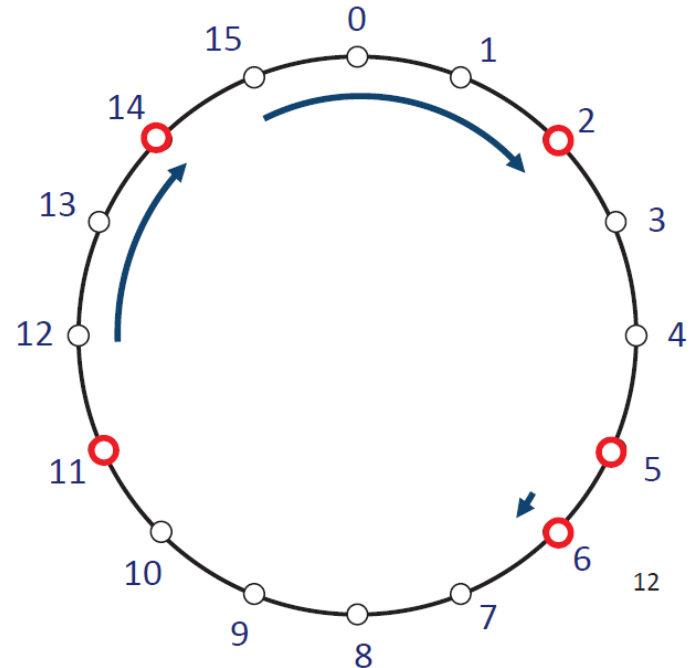  - 2,5,6,11,14 covered by nodes
  - Remaining values are not directly covered by a node

# Successor

o First node in clockwise direction with ID larger or equal the key
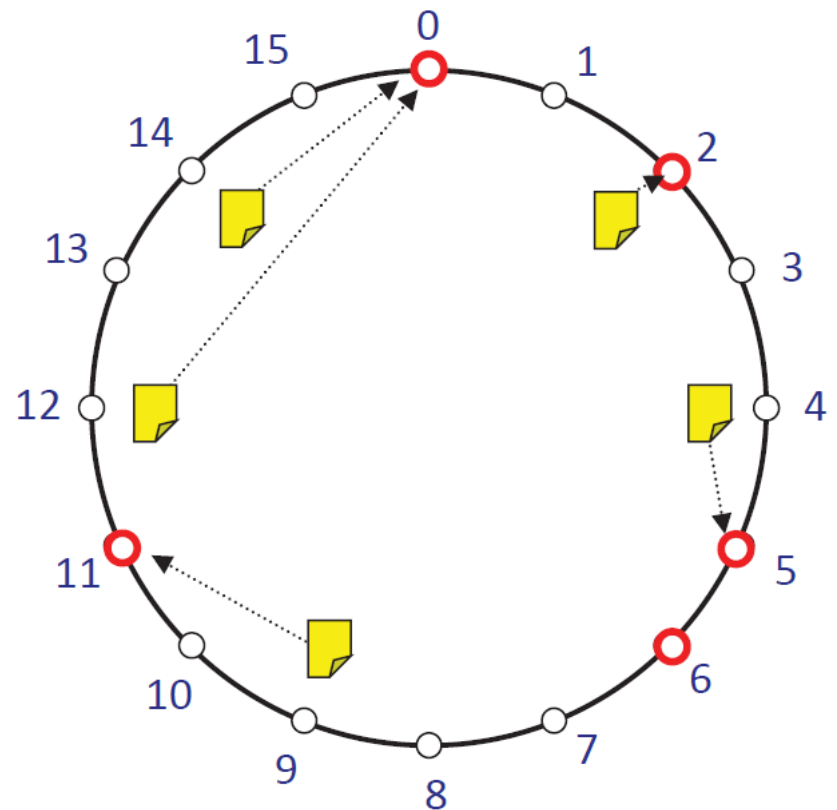
o Examples:
  - o succ(6) = 6
  - o succ(12) = 14
  - o succ(15) = 2

# How to store and locate data?

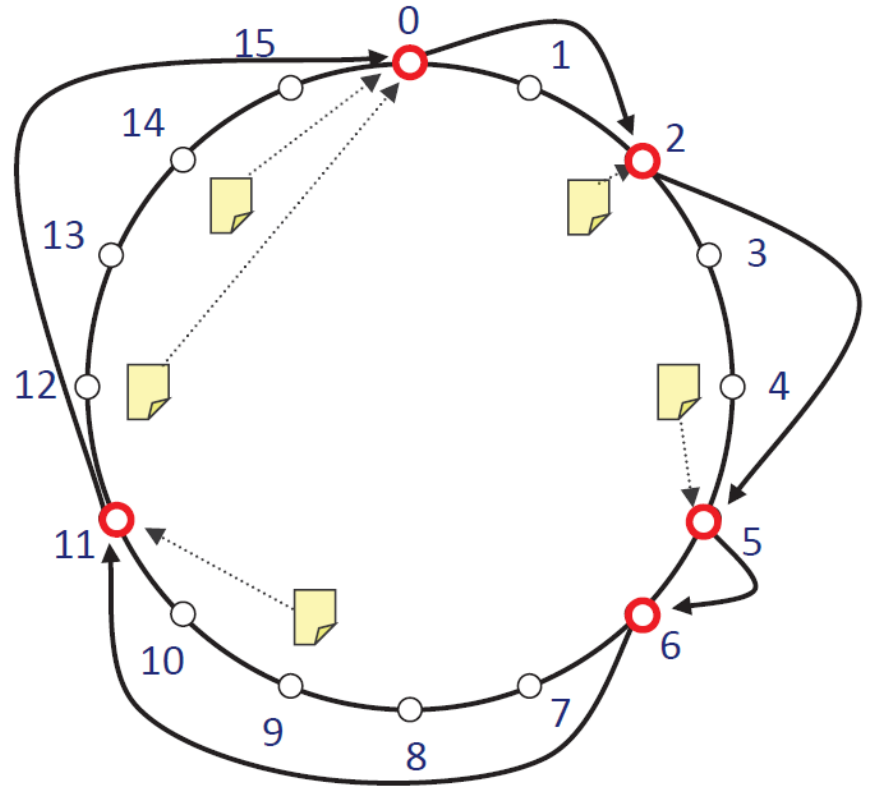o Each (key,value) pair is assigned the identifier H(key)

o Each item is stored at its succ(H(key))

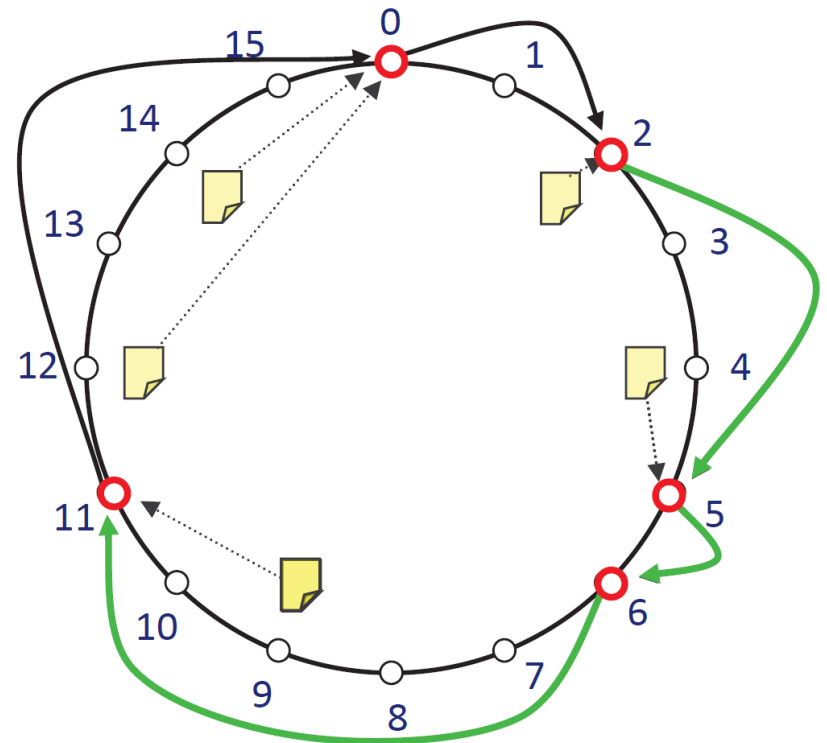| Drink | Location | H(Drink) |
|-------|----------|----------|
| Beer | Göttingen | 12 |
| Wine | France | 2 |
| Whisky | Scotland | 9 |
| Wodka | Russia | 14 |

# Successor Pointer

o Each node points to its successor

    o Known as node's succ pointer

o Example:

    o 0's succ = 2

    o 2's succ = 5
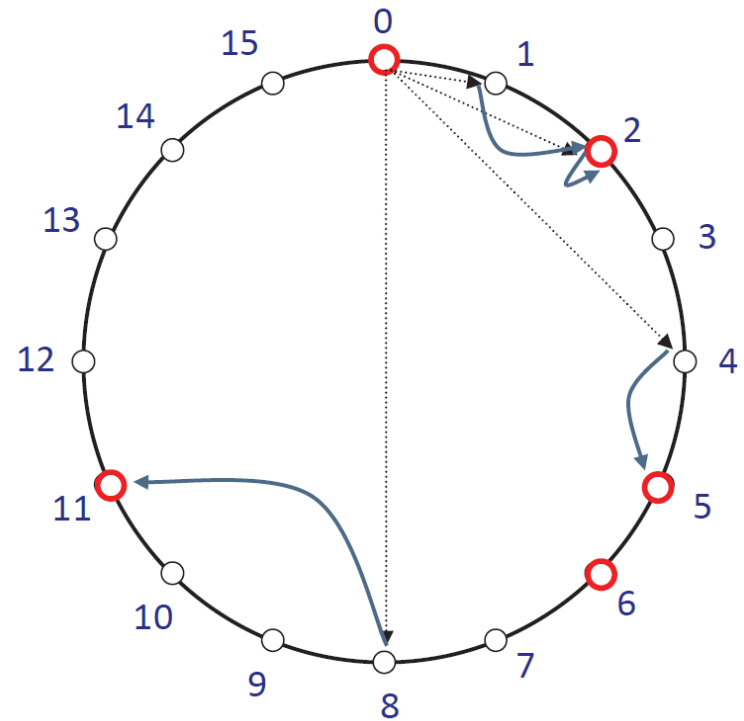
    o ...

# Basic Lookup of Data

o Lookup key:
  - o Calculate H(key)
  - o Follow succ pointers until key is found
  - o Lookup time: O(n)

o Example:
  - o „Where can I drink Whisky?"
  - o Calculate H(Whisky) = 9
  - o Traverse nodes:
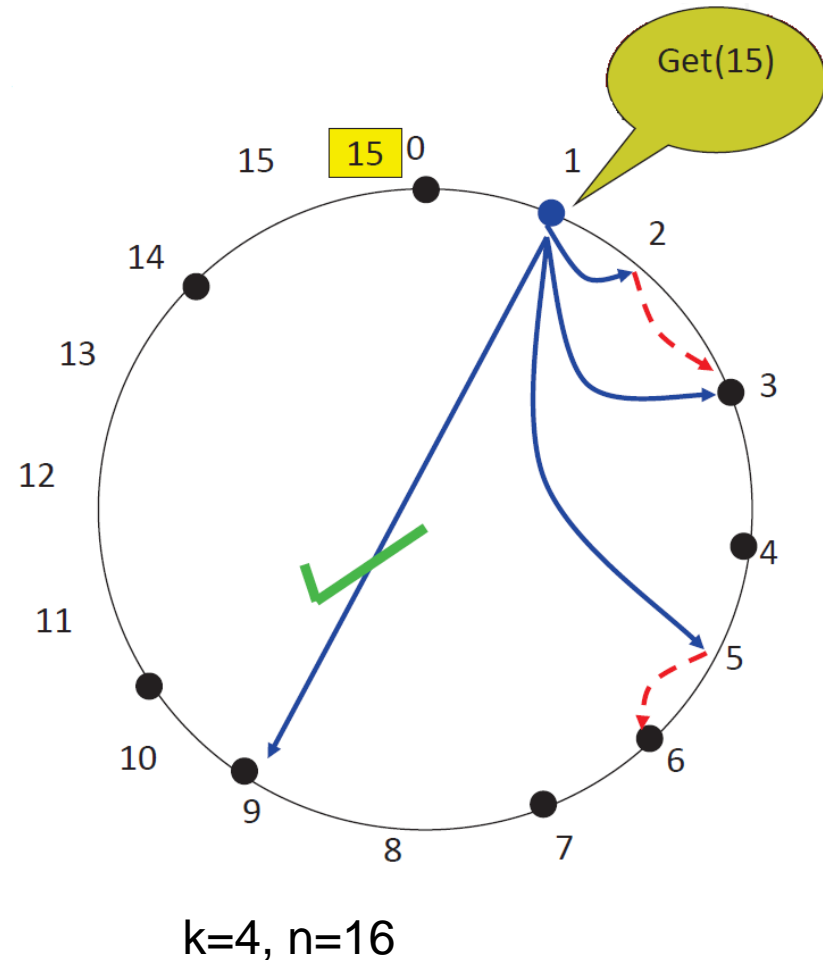    - • 2,5,6,11
  - o Return „Scotland"

# Scalable Lookup

o Each node maintains finger table (max k entries)

o for i in 0..k-1: finger[i] = succ(n+$2^{i-1}$)

  o Point to succ(n+1)
  o Point to succ(n+2)
  o Point to succ(n+4)

  o ...
  o Point to succ(n+$2^{i-1}$)
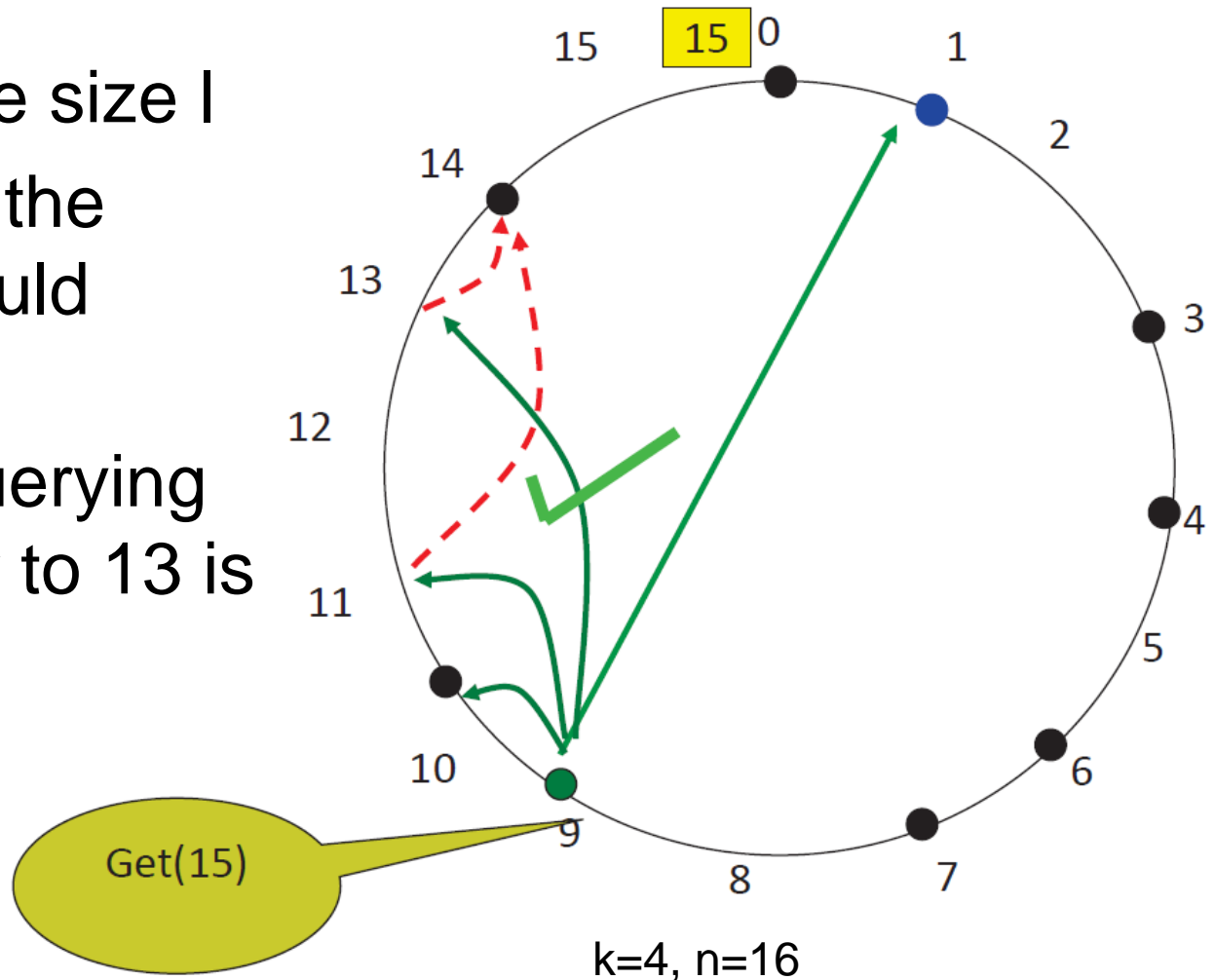
o Makes lookup time logarithmic!

  o O(log n)

# Routing

o Determines the next hop

o Each node n knows
  $succ(n+2^{i-1})$ for all $i=1..k$

o Forward queries for key
  to then highest
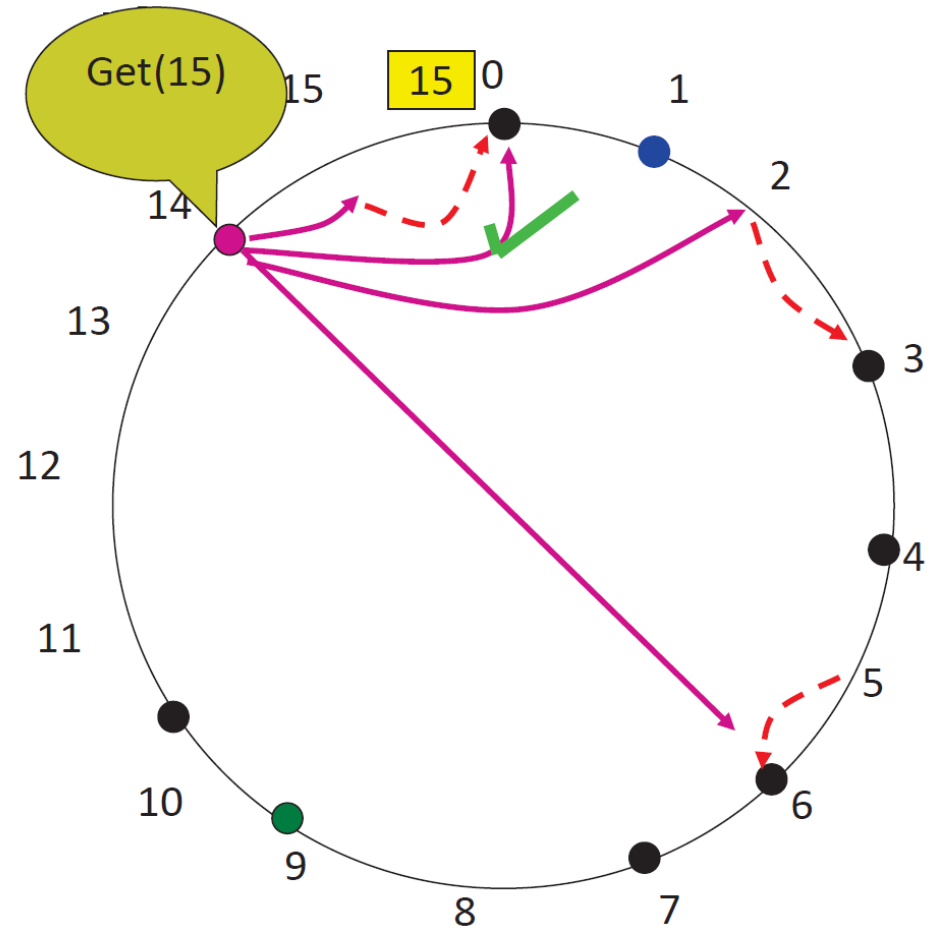  predecessor of key

o Routing entries = $\log_2(n)$



Get(15)

k=4, n=16

# Routing cont'd

o Routing table size l

o Node 9 was the highest 1 could reach

o Node 9 is querying again, finger to 13 is best



k=4, n=16

# Routing cont'd

o 13 is handled by 14

o 14 completes the route:

    o 15 is found at 0

# Routing cont'd

o From node 1, 3 hops to node 0 where item 15 is stored

o k=4 equals an ID space of 16, therefore the maximum number of hops is:

    o $Log_2(16) = 4$

o Average complexity is ½ log(n)

# Routing cont'd

o Such routing algorithms solve the lookup problem

o General concept:

  o Each node has only a limited view on the network

  o A node that receives a message containing a destination ID that is not managed by that node, it just forwards the request to the closest hop

o Here, algorithm is based on numeric closeness


o In Gummadi et al., „*The Impact of DHT Routing Geometry on Resilience and Proximity*", SIGCOMM 2003, implications are discussed

# Recursive vs. Iterative Lookup

o Recursive: Each node forwards the request (as shown) to the next hop

   o Fast, efficient
   o Each node can optimize forwarding


o Iterative: The requesting client queries the next hop iteratively from the nodes

   o Allows the lookup client to keep in control
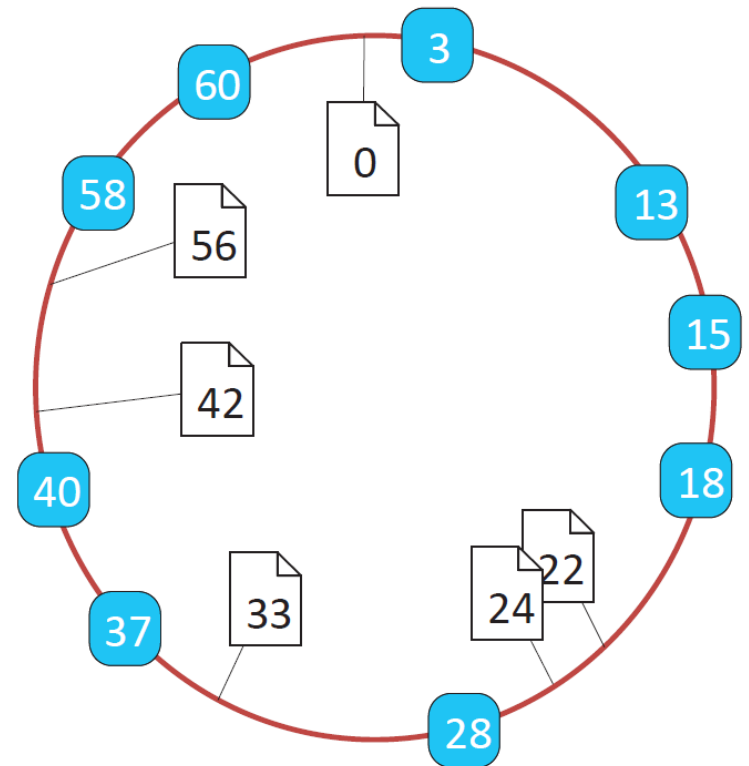   o Lookup client detects and localizes failures

# Achieved goals

o The DHT is scalable, as operations are performed in log(n)

o It is self-organized as each node directly knows its position (thanks to the hash function) and learns about the next hops

o On average load-distributing

o What about joins and especially leaves?

# Node Join and Leave

o **Node join:**

1. Bootstrap: a new node contacts a known node in the DHT
2. The new node gets a partion of the address space
3. Routing information is updated
4. The new node retrieves all tuples for which it is responsible

o **Node departure:**

  o Replication and load balancing

o **Node failure:**

  o Reactive or proactive recovery
  o Maintenance, load balancing, redistribution of tuples
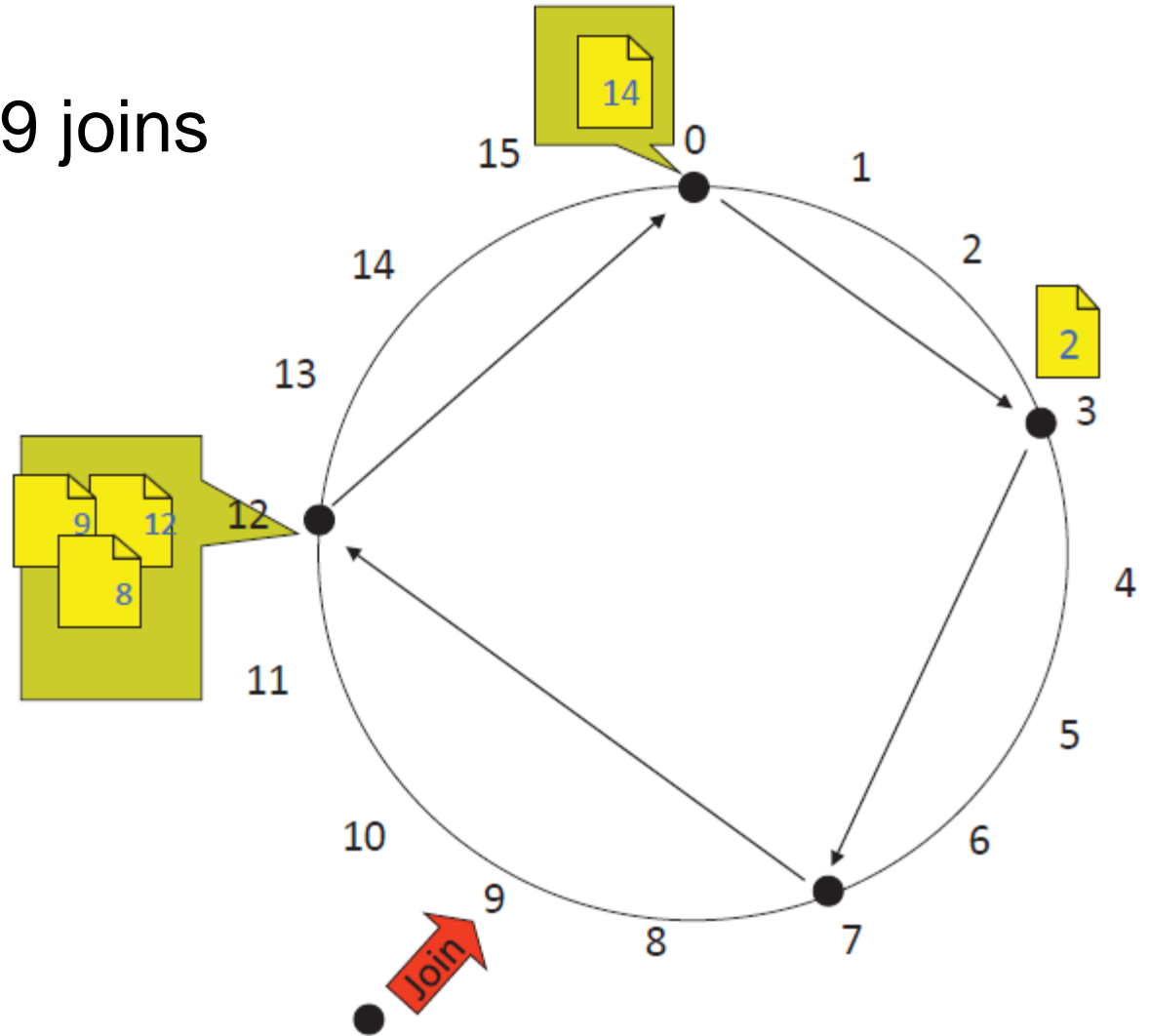  o Data is lost if not replicated!

# Node Join and Leave

- Join:
  - Lookup of own ID's successor
  - Contact that to get successors and predecessor

- Leaves:
  - Ping successors regularly
  - Always ensure x live nodes in successor set
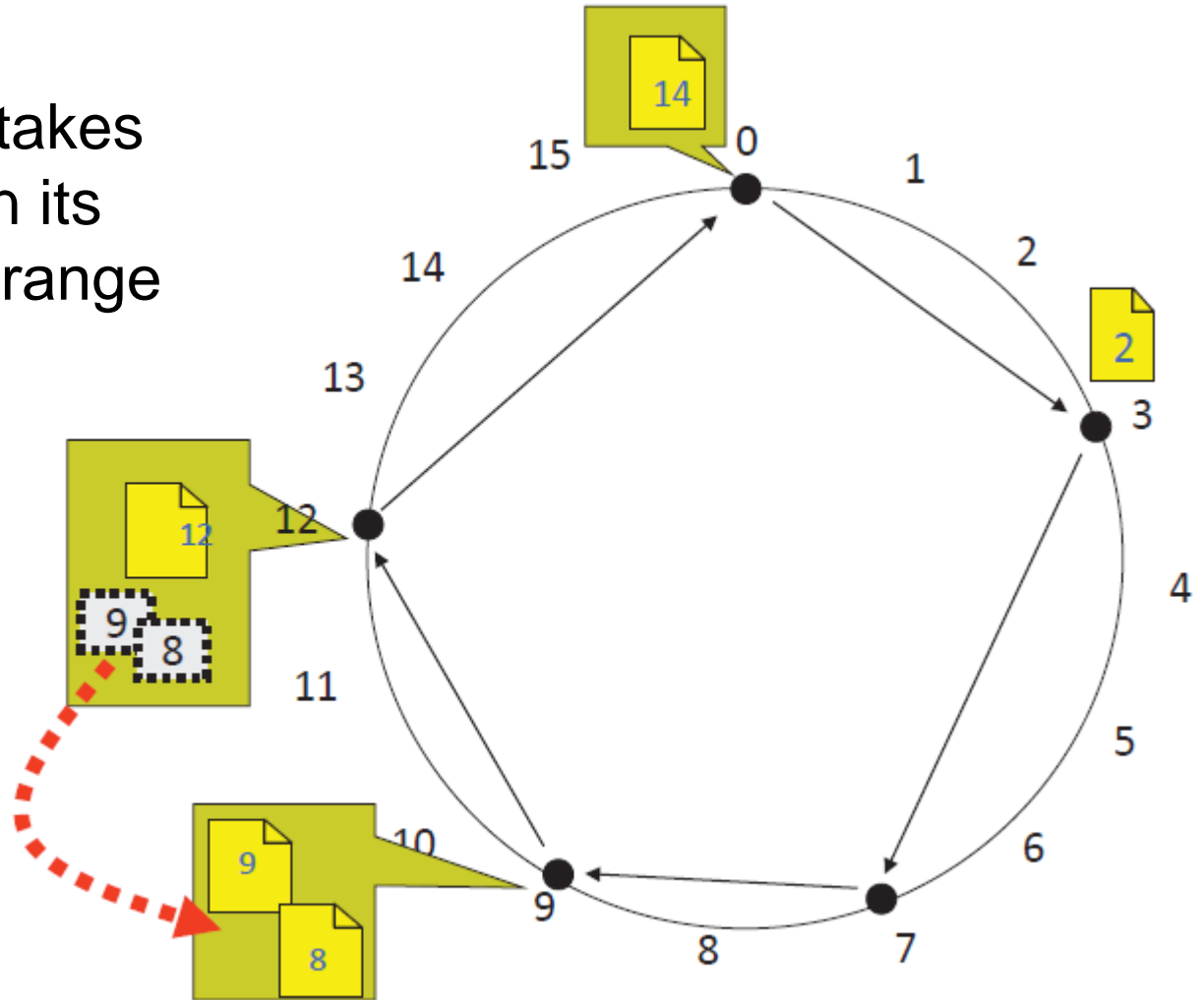
- Thereby, failures are treated as „normal"

# Node Join Example
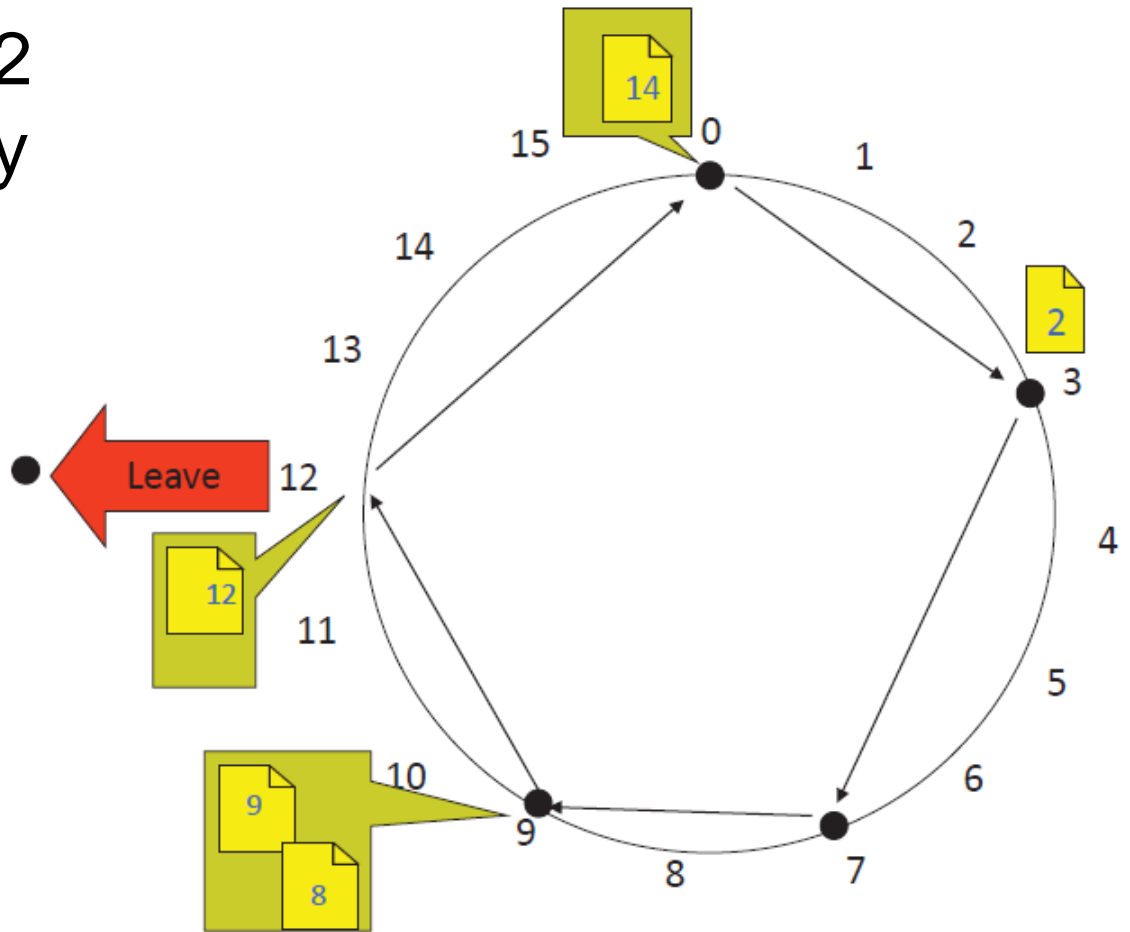
o Assume node 9 joins

# Node Join Example cont'd

o The new node takes over the docs in its "responsibility" range
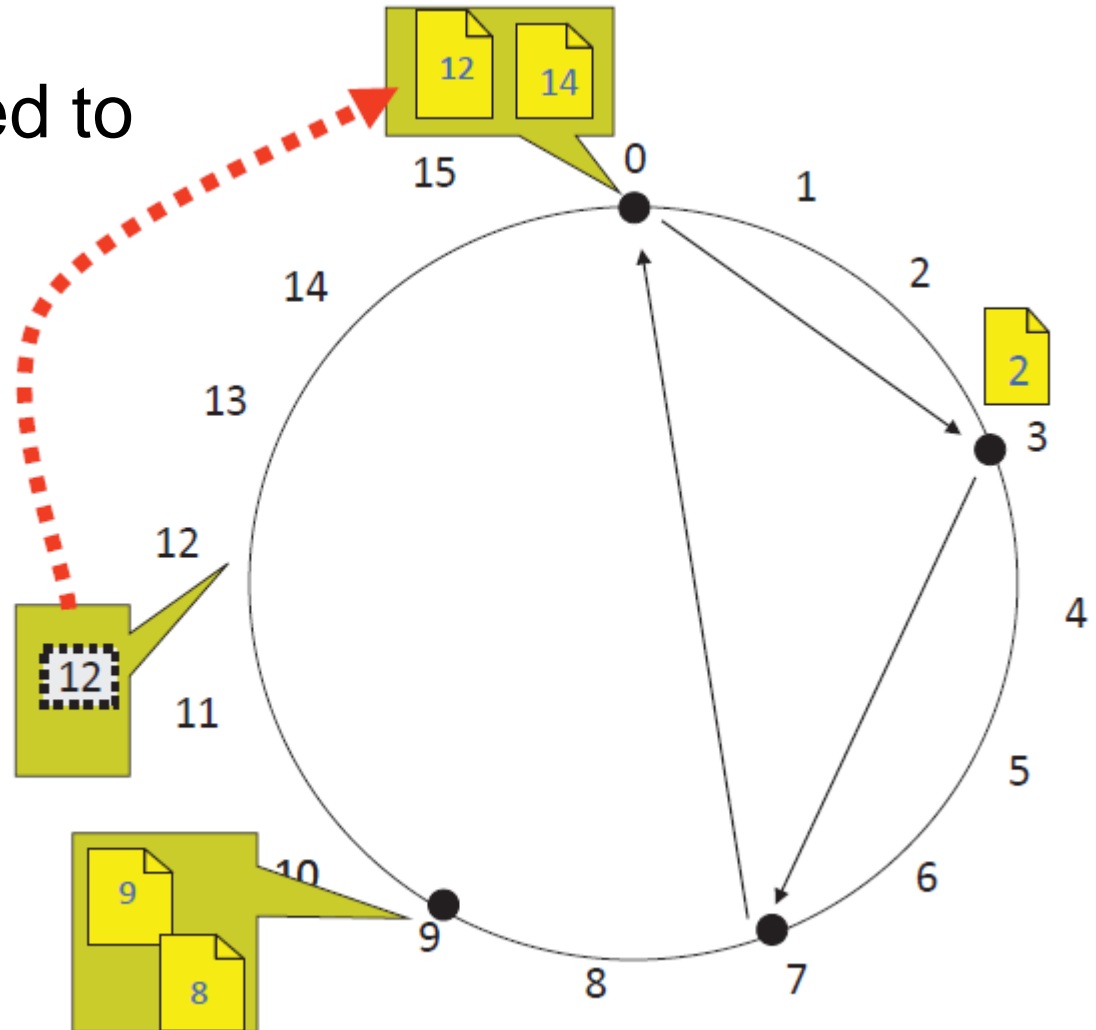
o Docs 9,8 from its successor

# Node Leave

o Assume node 12 leaves gracefully

# Node Leave cont'd

o Data is transferred to succ(12) = 14

o Node 12 informs predecessor and successor, who update their finger tables
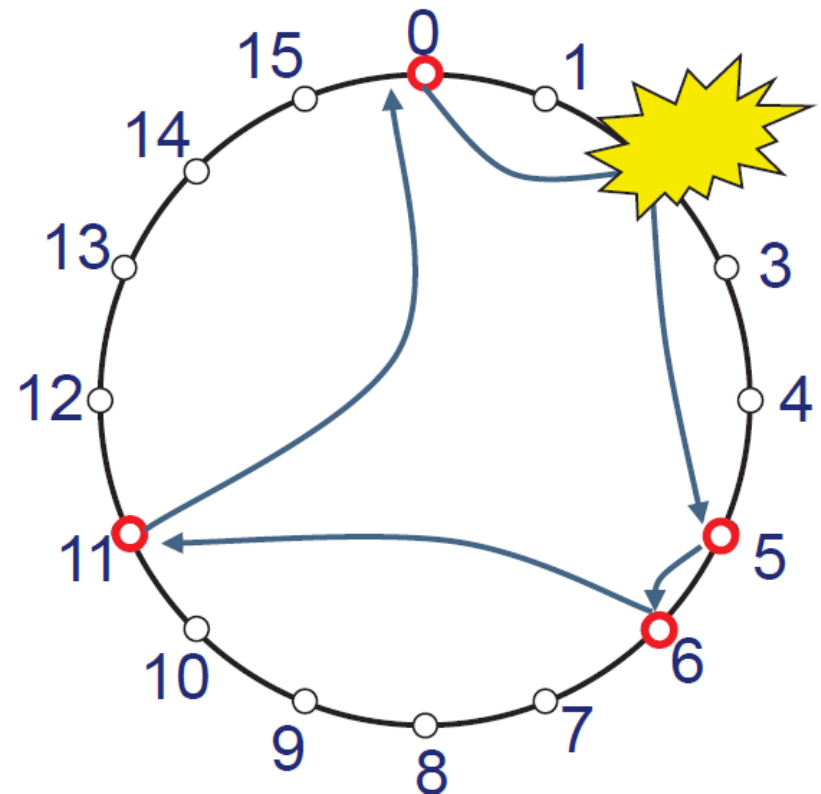
# Direct vs. Indirect Storage

o Direct storage:

- o Actual data is stored at the node responsible for it
- o The data is copied towards the responsible node upon node join
- o The node that contributed the data can leave without loss of its data
- o But: High storage and communication overhead!

o Indirect storage:

- o Instead of data, the references to the data are stored
- o The inserting node keeps the data
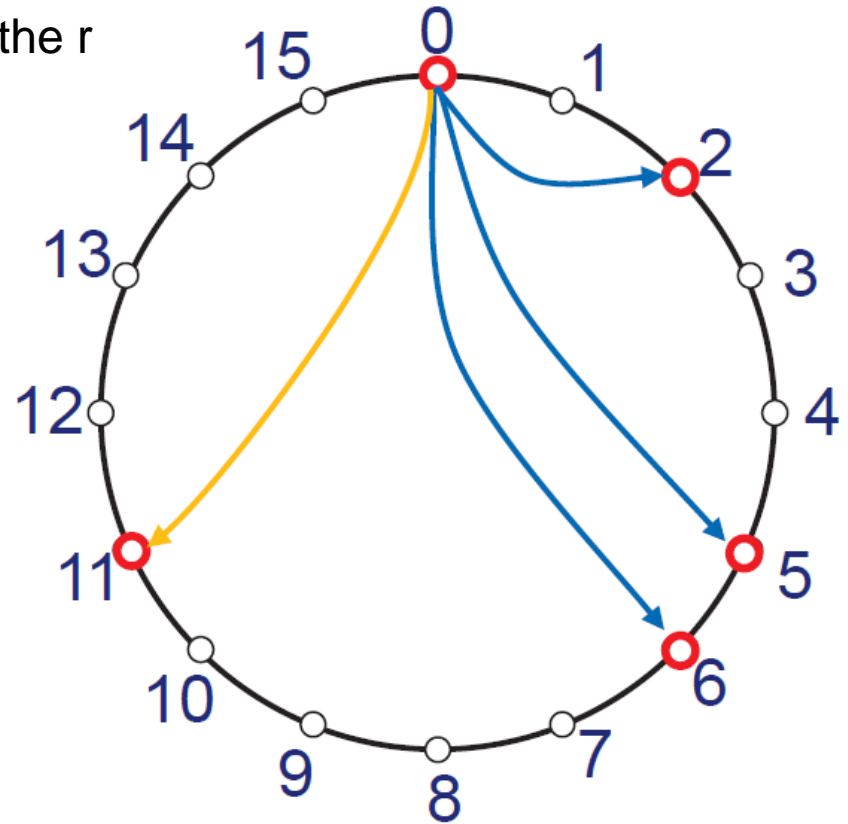- o Lower load on the DHT

# The Fragile Ring

o Problem: Everything is organized in a fragile ring structure

  o Failure of a node breaks the ring and data is lost

  o No way to recover as previous predecessor and successor don't know about each other!

# Successor Sets

○ As a solution, each node keeps:
  ○ A Successor set with pointers to the r closest successors
  ○ Predecessor pointer

○ If successor fails, replace with closest alive successor

○ If predecessor fails, set pointer to nil

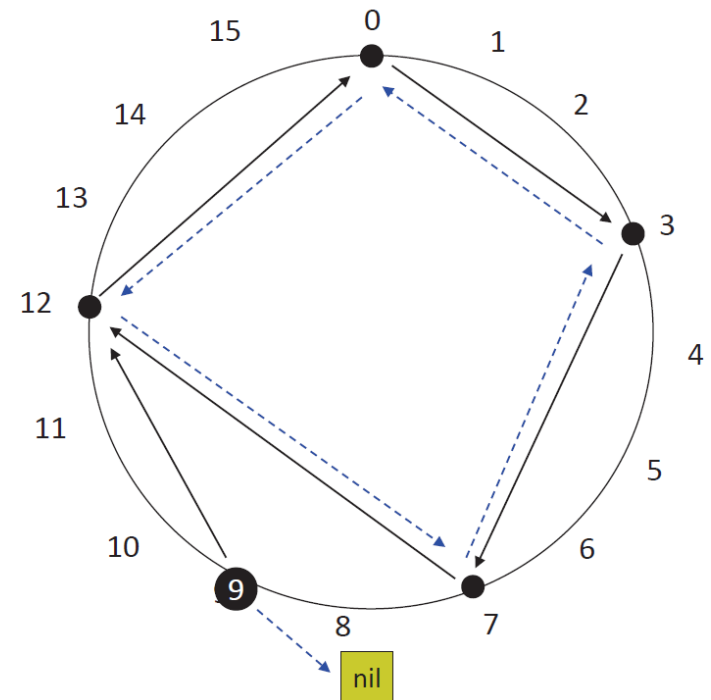○ Replicate objects throughout the successor set

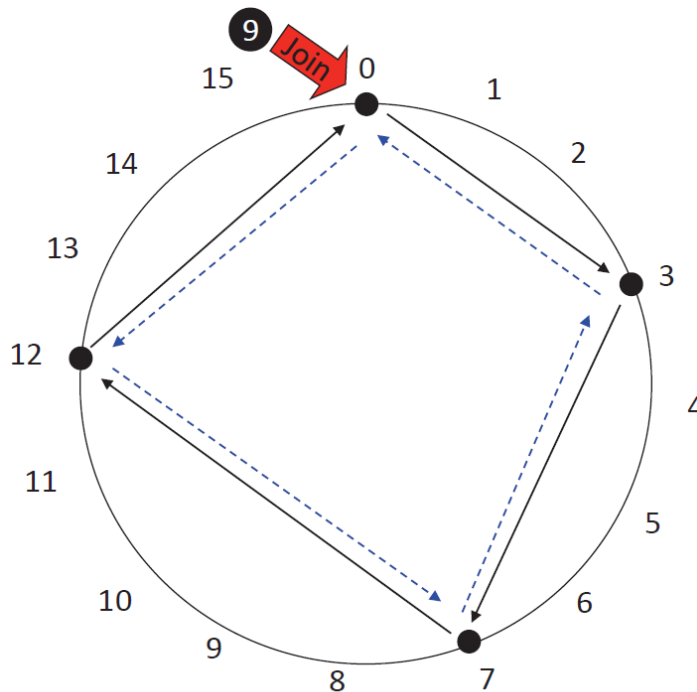# **Further Challenges**

o How does a node learn its:

- o Predecessors?

- o Fingers?

o What if "better" fingers come along later?

- o How would a node find out?

o How does a node react to failing or leaving fingers?

o All basically the same problem

# Periodic Stabilization

- Used to make pointers eventually correct
- Requires an additional predecessor pointer
  - First node met in anti-clockwise direction starting at n-1

- A node n joins the DHT through a node o:
  - Find n's successor by lookup(n)
  - n sets its successor to the found successor
  - Stabilization fixes the rest
    - stabilize() function is run peridically by each node
  - The new node does not determine its predecessor: its predecessor detects and fixes inconsistencies
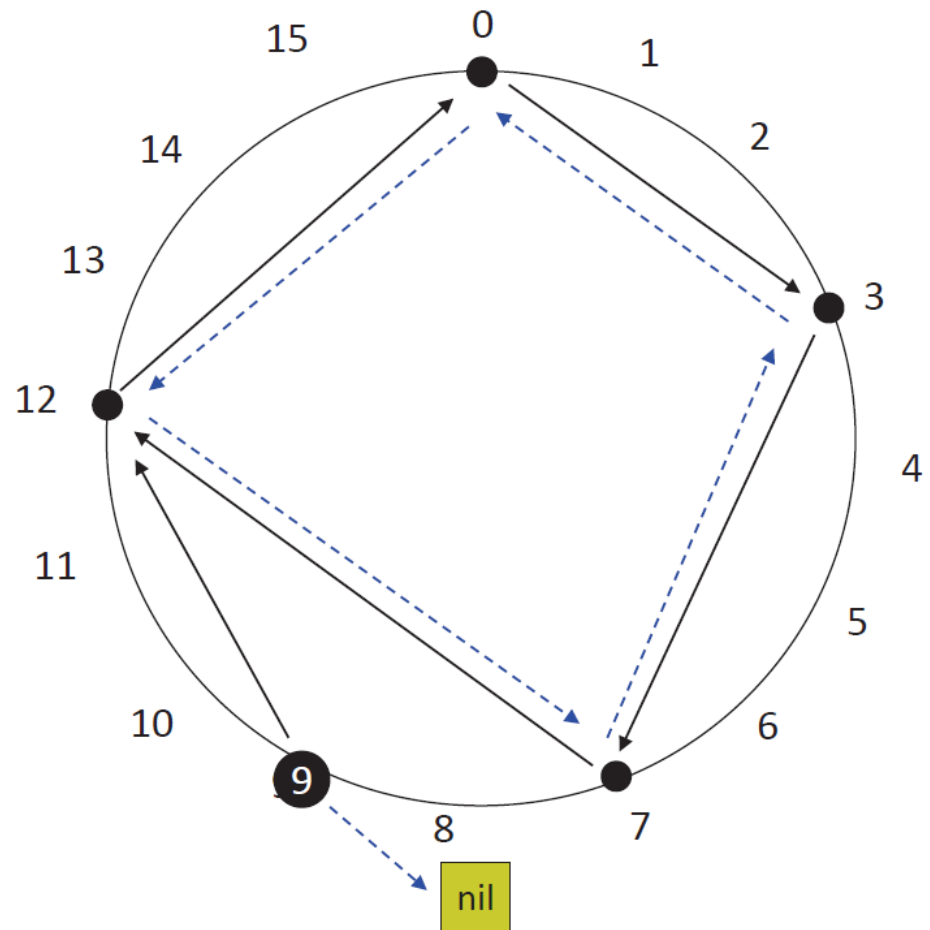
# Periodic Stabilization Example

1. 9 joins through node 0
2. 9 sets its predecessor to nil
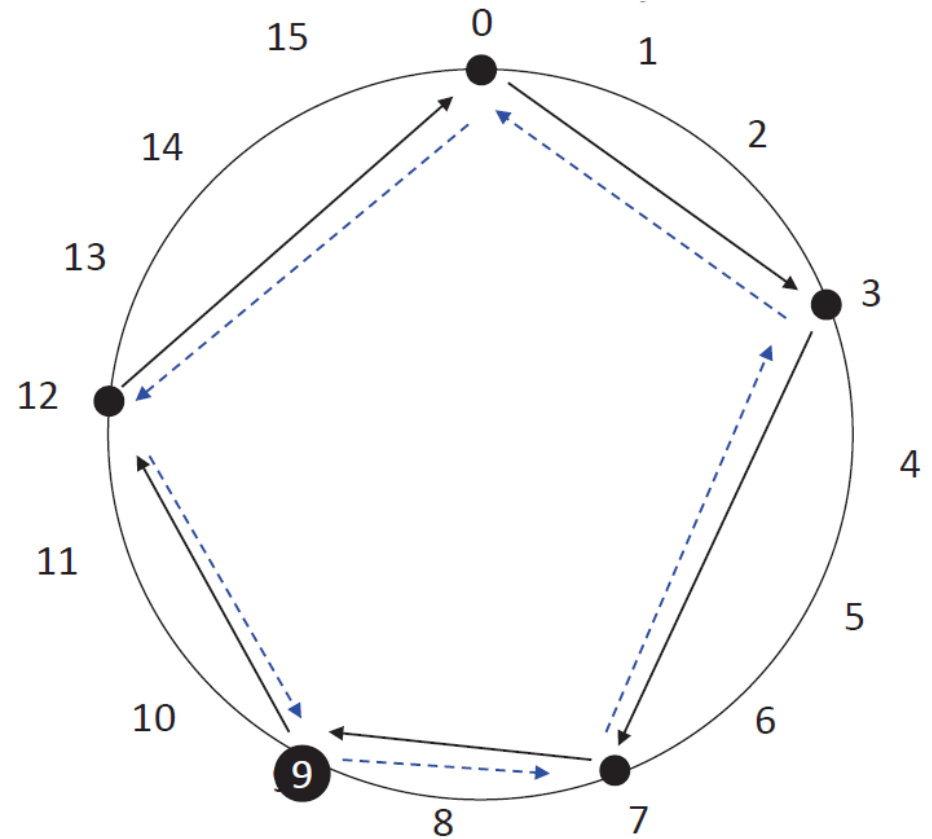3. 9 asks 0 for succ(9). Receives "12"
4. 9 sets its succ to 12

# Periodic Stabilization Example

o 9 runs stabilize()

1. 9 asks 12 for its predecessor
2. 12 replies with "7"
3. 9 notifies 12 that 9 is now its predecessor

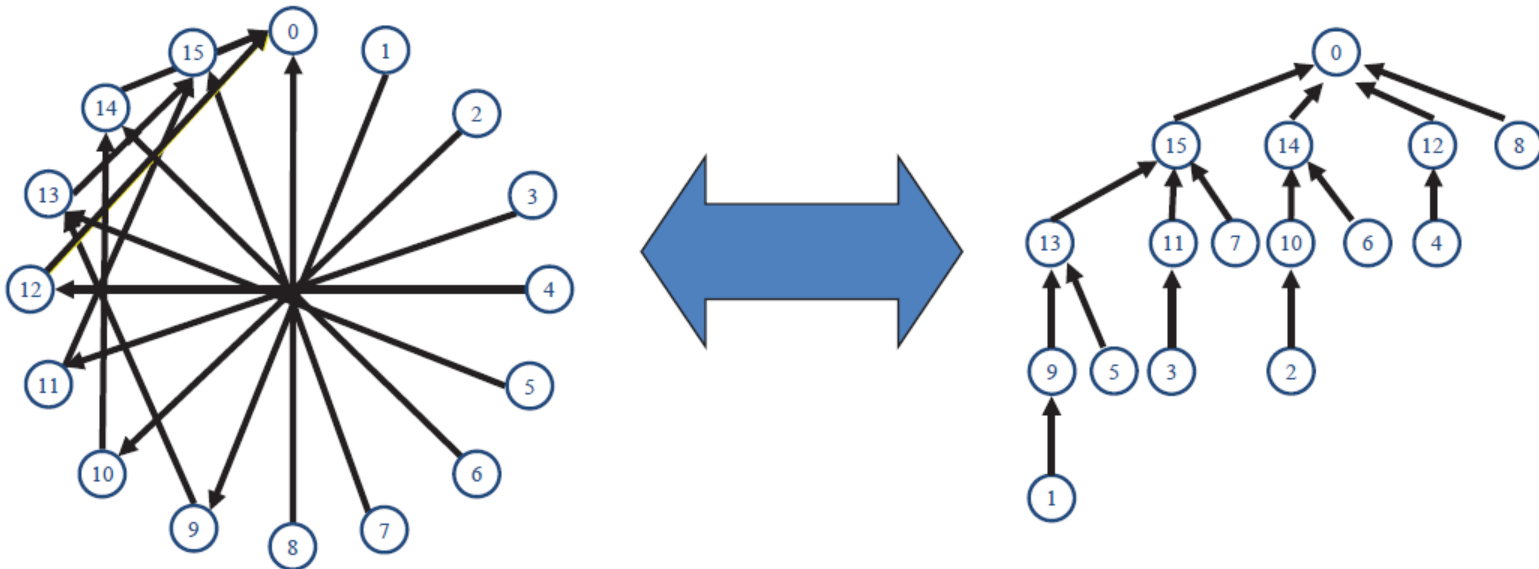# Periodic Stabilization Example

o 7 runs stabilize()

1. 7 discovers from 12 that pred(12) is now 9
2. 7 sets successor to 9
3. 7 notifies 9
4. 9 sets pred(9) to 7

# Chord in a "Tree View"

o Finger tables are Chord's core
  o Providing O(log n) hop routing by at least halving the distance to the target by each hop
  o Forest of binomial trees rooted at each key

# **Chord - Conclusion**

o Lookup time: O(log n)

o Drawbacks:

   o Rigidity

   - Complicates recovery from failed nodes and routing table
   - Precludes proximity-based routing

   o Unidirectional routing

   o Incoming traffic is not used to re-enforce routing
     tables

o Fault-tolerant, but not very robust.

# Other DHTs

o Kademlia (used in BitTorrent)

    o Lookup also done in O(log n) – as with most DHTs

    o Uses distance between two nodes: XOR of both nodes' IDs

    o Nodes still responsible for a part of ID space

    o Location of content basically the same as in Chord

        • Node closest to searched ID

        • O(log n) since XOR can halve distance at each hop

        • Note: This distance is *not* geographical

    o For details: Maymounkov and Mazières, Kademlia: "A Peer-to-peer Information System Based on the XOR metric", 2002

# References

- [1] Stephanos Androutsellis-Theotokis and Diomidis Spinellis. A survey of peer-to-peer content distribution technologies. ACM Comput. Surv. 36(4), 335-371. 2004.

- [2] Hari Balakrishnan, M. Frans Kaashoek, David Karger, Robert Morris, and Ion Stoica. Looking up data in p2p systems. Comm. ACM 46,2(Feb.), 43–48. 2003.

- [4] Pouwelse, Johan; et al. "The Bittorrent P2P File-Sharing System: Measurements and Analysis". Peer-to-Peer Systems IV. Berlin: Springer. pp. 205–216. 2005.

- [5] Erik Nygren, Ramesh K. Sitaraman, and Jennifer Sun, The Akamai Network: A Platform for High-Performance Internet Applications, ACM SIGOPS Operating Systems Review, Vol. 44, No.3, July 2010.