## Peer-to-Peer Systems and Distributed Hash Tables
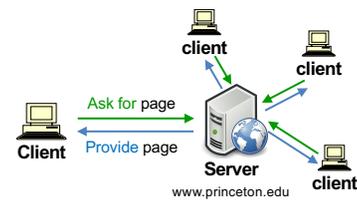
COS 418: Distributed Systems
Lecture 9

Mike Freedman
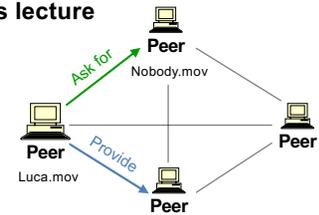
---

## Distributed Application Architecture



**This lecture**

Ask for page / Provide page

Client-Server

Ask for / Provide
Nobody.mov / Luca.mov

Peer-to-Peer

---

## Today

1. **Peer-to-Peer Systems**
   – **What, why, and the core challenge**

2. Distributed Hash Tables (DHT)

3. The Chord Lookup Service

4. Concluding thoughts on DHTs, P2P

---

## What is a Peer-to-Peer (P2P) system?



Node / Internet / Node

- A **distributed** system architecture:
  – **No centralized control**
  – Nodes are **roughly symmetric** in function

- **Large** number of **unreliable** nodes

---

## P2P adoption

Successful adoption in **some niche areas**

1. Client-to-client (legal, illegal) **file sharing**
   1. Napster (1990s), Gnutella, BitTorrent, etc.

2. **Digital currency:** no natural single owner (Bitcoin)

3. **Voice/video telephony:** user to user anyway (Skype in old days)
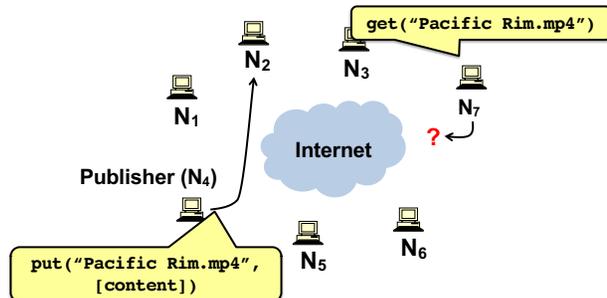   – Issues: Privacy and control

## Why might P2P be a win?

- **High capacity for services** through parallelism and scalability:
  – More disks, network connections, CPUs, etc. as peers join
  – Data are divided and duplicated, accessible from multiple peers concurrently

- **Absence of a centralized server** may mean:
  – **Less chance** of service overload as load increases
  – Easier **deployment**
  – A single failure **won't wreck** the whole system (no single point of failure)
  – System as a whole is **harder to attack**
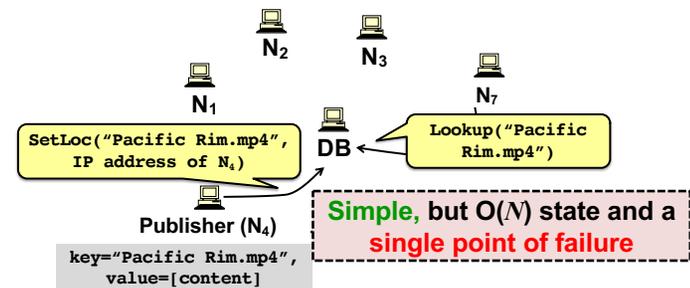    • Both technically *and* legally
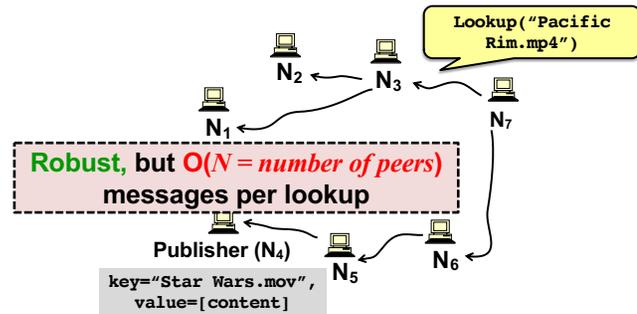
## The lookup problem: locate the data



## Centralized lookup (Napster)



**Simple**, but $O(N)$ state and a **single point of failure**

## Flooded queries (original Gnutella)

Lookup("Pacific Rim.mp4")

$N_2$   $N_3$

$N_1$

$N_7$

**Robust,** but **O**($N = number\ of\ peers$) messages per lookup

Publisher ($N_4$)

key="Star Wars.mov", value=[content]

$N_5$   $N_6$

---

## Tradeoffs in distributed systems

Gnutella — Nearly no state / Many msgs

\# msgs

Napster — High state / Good performance / Single PoF

\# state

---

## Tradeoffs in distributed systems

Gnutella — Nearly no state / Many msgs

\# msgs

Ideal

Napster — High state / Good performance / Single PoF

\# state

---

## Tradeoffs in distributed systems

Gnutella — Nearly no state / Many msgs

\# msgs

DHT (Chord) — msgs < Gnutella / state < Napster

Napster — High state / Good performance / Single PoF

\# state

## Today

1. Peer-to-Peer Systems

2. **Distributed Hash Tables (DHT)**

3. The Chord Lookup Service

4. Concluding thoughts on DHTs, P2P

## What is a DHT (and why)?

- Distributed Hash Table: an abstraction of hash table in a distributed setting

  ```
  key = hash(data)
  lookup(key) → IP addr (Chord lookup service)
  send–RPC(IP address, put, key, data)
  send–RPC(IP address, get, key) → data
  ```

- **Partitioning data** in **large-scale distributed systems**
  - Tuples in a global database engine
  - Data blocks in a global file system
  - Files in a P2P file-sharing system

## Cooperative storage with a DHT

## DHT is expected to be

- Decentralized: no central authority

- Scalable: low network traffic overhead

- Efficient: find items quickly (latency)

- Dynamic: nodes fail, new nodes join

## Today

1. Peer-to-Peer Systems

2. Distributed Hash Tables (DHT)

3. **The Chord Lookup Service**

**17**

## Chord identifiers

- **Hashed values (integers) using the same hash function**
  - **Key identifier** = SHA-1(key)
  - **Node identifier** = SHA-1(IP address)

- **How does Chord partition data?**
  - i.e., map key IDs to node IDs

- **Why hash key and address?**
  - Uniformly distributed in the ID space
  - Hashed key → load balancing;     hashed address → independent failure

**18**

## Alternative:  mod (n) hashing

- **System of n nodes:  1…n**
  - Key identifier = *hash(key) mod n*
  - Good load balancing

- **What if a node fails?**
  - Instead of n nodes, now *n -1* nodes
  - Mapping of all keys change, as now *hash(key) mod (n-1)*

    - **N = 5**
      - 12594   mod 5 = 4
      - 28527   mod 5 = 2
      - 816     mod 5 = 1
      - 716565 mod 5 = 0
    - **N = 4**
      - 12594   mod 4 = **2**
      - 28527   mod 4 = **3**
      - 816     mod 4 = **0**
      - 716565 mod 4 = **1**

**19**

## Consistent hashing [Karger '97] – data partition

**Identifiers have m = 3 bits**
**Key space: [0, $2^3$-1]**

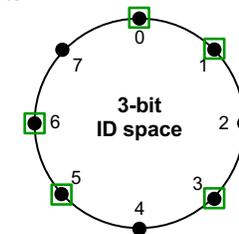● **Identifiers/key space**
☐ **Node**



**3-bit ID space**

**20**

## Consistent hashing [Karger '97] – data partition

Identifiers have m = 3 bits
Key space: [0, $2^3$-1]

● Identifiers/key space
□ Node

3-bit ID space

Stores key 7, 0
Stores key 1
Stores key 6
Stores keys 2, 3
Stores keys 4, 5

Key is stored at its **successor:** node with next-higher ID

21

---

## Consistent hashing [Karger '97] – basic lookup

Identifiers have m = 3 bits
Key space: [0, $2^3$-1]

● Identifiers/key space
□ Node
- ➔ Successor pointer

3-bit ID space

Stores key 7, 0
Stores key 1
Stores key 6
Stores keys 2, 3
Stores keys 4, 5
At Node 1
Key 1 ?

O(N) messages and hops!

Look up key 1

22

---

## Chord – finger tables

Identifiers have **m = 3 bits**
Key space: [0, $2^3$-1]

● Identifiers/key space
□ Node

3-bit ID space

Each node keeps **m** states
Key space ➔ **m** ranges via
**(N+$2^{k-1}$) mod $2^m$, 1<=k<=m**

2, [2,3], node 3
3, [3,5], node 3
5, [5,1], node 5

Separators
Key ranges
Successors of separators

24

---

## Chord – finger tables

Identifiers have **m = 3 bits**
Key space: [0, $2^3$-1]

● Identifiers/key space
□ Node

3-bit ID space
Node 1

Each node keeps **m** states
Key space ➔ **m** ranges via
**(N+$2^{k-1}$) mod $2^m$, 1<=k<=m**

1, [1,2), node 1
2, [2,4), node 3
4, [4,0), node 5

2, [2,3), node 3
3, [3,5), node 3
5, [5,1), node 5

4, [4,5), node 5
5, [5,7), node 5
7, [7,3), node 0

O(logN) messages and hops!

Look up key 1

25

6

## Implication of finger tables

- A **binary lookup tree** rooted at every node
  - Threaded through other nodes' finger tables

- Better than arranging nodes in a single tree
  - Every node acts as a root
    - So there's **no root hotspot**
    - **No single point** of failure
    - But a **lot more state** in total

## Chord lookup algorithm properties

- **Interface:** lookup(key) $\rightarrow$ IP address

- **Efficient:** O(log N) messages per lookup
  - N is the total number of nodes (peers)

- **Scalable:** O(log N) state per node
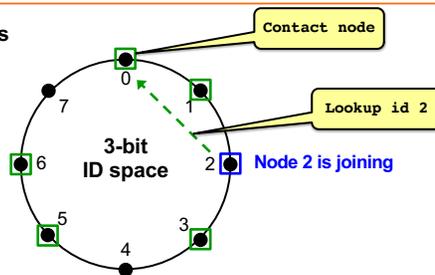
- **Robust:** survives massive failures

## Chord – node joining

**Identifiers have m = 3 bits**
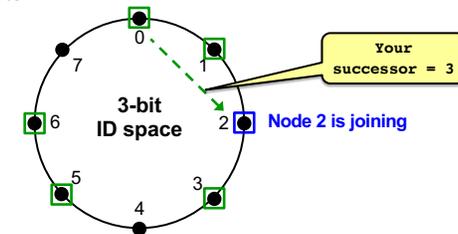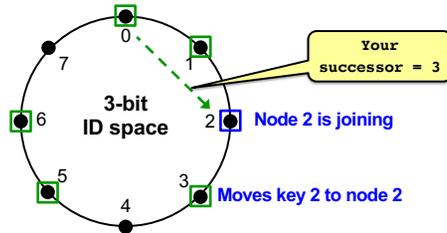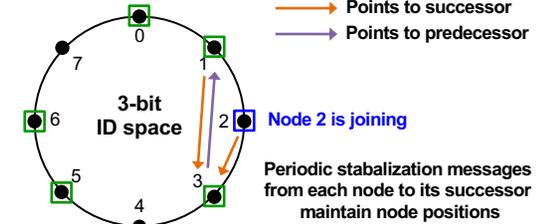**Key space: [0, $2^3$-1]**

- Identifiers/key space
- Node

## Chord – node joining

**Identifiers have m = 3 bits**
**Key space: [0, $2^3$-1]**
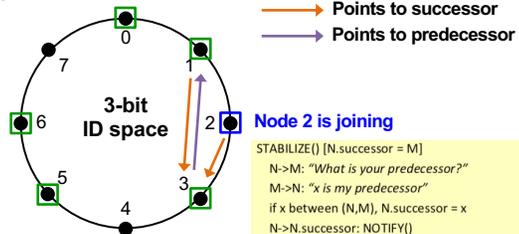
- Identifiers/key space
- Node

## Chord – node joining

**Identifiers have m = 3 bits**
**Key space: [0, $2^3$-1]**

● Identifiers/key space

□ Node

**3-bit ID space**

Your successor = 3

Node 2 is joining

Moves key 2 to node 2

## Chord – node joining

**Identifiers have m = 3 bits**
**Key space: [0, $2^3$-1]**

→ Points to successor
→ Points to predecessor

● Identifiers/key space

□ Node

**3-bit ID space**

Node 2 is joining

**Periodic stabalization messages from each node to its successor maintain node positions**

## Chord – node joining

**Identifiers have m = 3 bits**
**Key space: [0, $2^3$-1]**

→ Points to successor
→ Points to predecessor

● Identifiers/key space

□ Node

**3-bit ID space**

Node 2 is joining

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

## Chord – node joining

**Identifiers have m = 3 bits**
**Key space: [0, $2^3$-1]**

→ Points to successor
→ Points to predecessor

● Identifiers/key space

□ Node

**3-bit ID space**

Node 2 is joining

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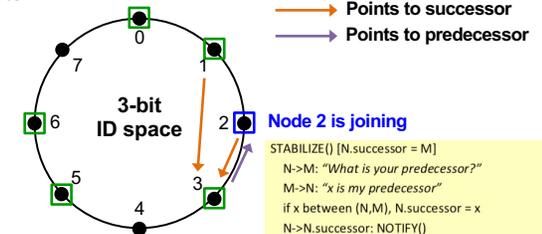  If (N between (M.predecessor, M))
    M.predecessor = N

## Chord – node joining

**Identifiers have m = 3 bits**
**Key space: [0, 2³-1]**

● Identifiers/key space

☐ Node

**3-bit ID space**

0
7
6
5
4
3
2
1

→ Points to successor
→ Points to predecessor

**Node 2 is joining**

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))
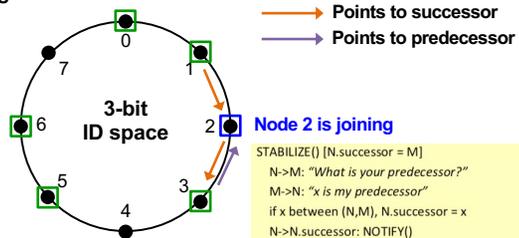    M.predecessor = N

35

---

## Chord – node joining

**Identifiers have m = 3 bits**
**Key space: [0, 2³-1]**

● Identifiers/key space

☐ Node

**3-bit ID space**

0
7
6
5
4
3
2
1

→ Points to successor
→ Points to predecessor

**Node 2 is joining**

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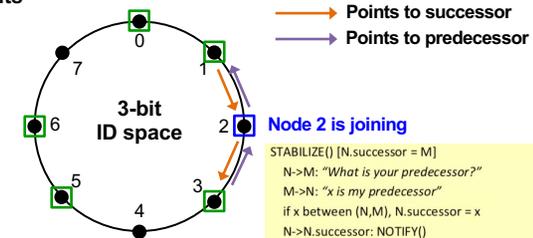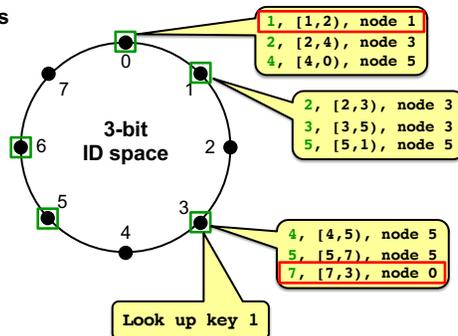  If (N between (M.predecessor, M))
    M.predecessor = N

36

---

## Chord – failures and successor list

**Identifiers have m = 3 bits**
**Key space: [0, 2³-1]**
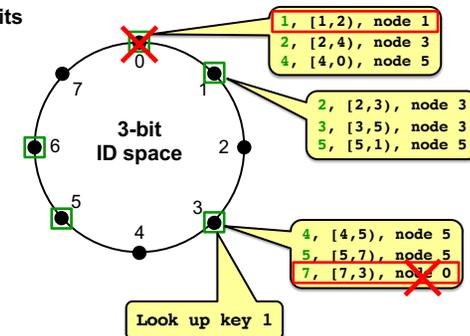
● Identifiers/key space

☐ Node

**3-bit ID space**

0
7
6
5
4
3
2
1

```
1, [1,2), node 1
2, [2,4), node 3
4, [4,0), node 5
```

```
2, [2,3), node 3
3, [3,5), node 3
5, [5,1), node 5
```

```
4, [4,5), node 5
5, [5,7), node 5
7, [7,3), node 0
```

**Look up key 1**

37

---

## Chord – failures and successor list

**Identifiers have m = 3 bits**
**Key space: [0, 2³-1]**

● Identifiers/key space

☐ Node

**3-bit ID space**

0
7
6
5
4
3
2
1

```
1, [1,2), node 1
2, [2,4), node 3
4, [4,0), node 5
```

```
2, [2,3), node 3
3, [3,5), node 3
5, [5,1), node 5
```

```
4, [4,5), node 5
5, [5,7), node 5
7, [7,3), node 0
```

**Look up key 1**

38

9

## Slide 39

# Chord – failures and successor list

**Identifiers have m = 3 bits**
**Key space: [0, 2³-1]**

● Identifiers/key space

□ Node

→ Points to successor

**3-bit ID space**

```
1, [1,2), node 1
2, [2,4), node 3
4, [4,0), node 5
```

```
2, [2,3), node 3
3, [3,5), node 3
5, [5,1), node 5
```

```
4, [4,5), node 5
5, [5,7), node 5
7, [7,3), node 0
```
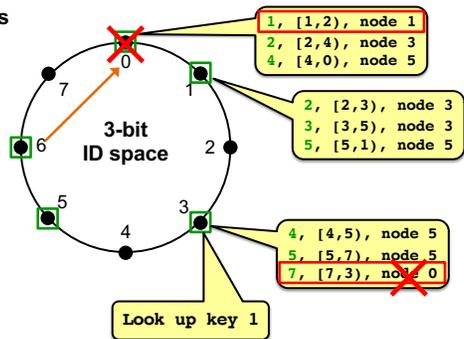
**Look up key 1**

39

## Slide 40

# Chord – failures and successor list

**Identifiers have m = 3 bits**
**Key space: [0, 2³-1]**

● Identifiers/key space

□ Node

→ Points to successor

**3-bit ID space**

```
1, [1,2), node 1
2, [2,4), node 3
4, [4,0), node 5
```

```
2, [2,3), node 3
3, [3,5), node 3
5, [5,1), node 5
```

```
4, [4,5), node 5
5, [5,7), node 5
7, [7,3), node 0
```

**Look up key 1**

**Succ. of id 7 (Succ. Of node 6)**
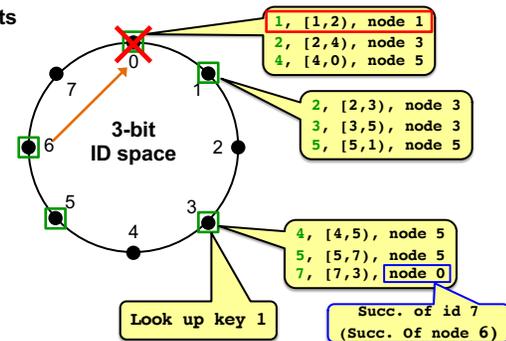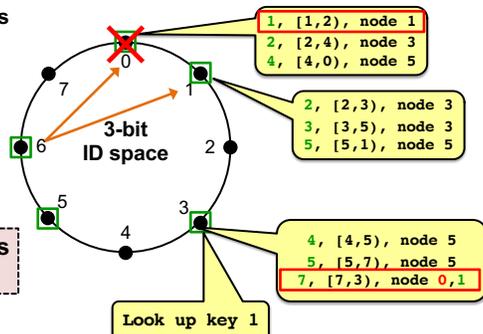
40

## Slide 41

# Chord – failures and successor list

**Identifiers have m = 3 bits**
**Key space: [0, 2³-1]**
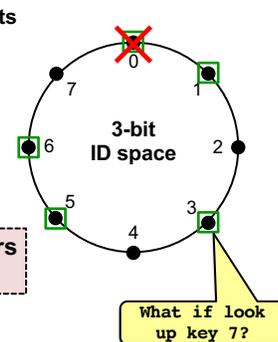
● Identifiers/key space

□ Node

→ Points to successor

**3-bit ID space**

```
1, [1,2), node 1
2, [2,4), node 3
4, [4,0), node 5
```

```
2, [2,3), node 3
3, [3,5), node 3
5, [5,1), node 5
```

```
4, [4,5), node 5
5, [5,7), node 5
7, [7,3), node 0,1
```

**r-nearest successors (r = logN)**

**Look up key 1**

41

## Slide 42

# Chord – failures and successor list

**Identifiers have m = 3 bits**
**Key space: [0, 2³-1]**

● Identifiers/key space

□ Node

**3-bit ID space**

**r-nearest successors (r = logN)**

**What if look up key 7?**

42

**10**

## DHash replicates blocks at *r* successors

**Identifiers have m = 3 bits**
**Key space: [0, $2^3$-1]**

● Identifiers/key space
□ Node

**3-bit ID space**

Key 7

Key 7

**"Adjacent" nodes in the ring may be far away in the network → Independent failures**

**r-nearest successors (r = logN)**

What if look up key 7?

---

## Today

1. Peer-to-Peer Systems

2. Distributed Hash Tables

3. The Chord Lookup Service

4. **Concluding thoughts on DHT, P2P**

---

## Why don't all services use P2P?

- **High latency and limited bandwidth** between peers (vs. intra/inter-datacenter, client-server model*)
  - 1 M nodes = 20 hops; 50 ms / hop gives 1 sec lookup latency (assuming no failures / slow connections…)

- User computers are **less reliable** than managed servers

- **Lack of trust** in peers' correct behavior
  - Securing DHT routing hard, unsolved in practice

---

## DHTs in retrospective

- Seem promising for finding data in large P2P systems
- Decentralization seems good for load, fault tolerance

- **But:** the **security problems** are difficult
- **But:** **churn** is a problem, particularly if log(n) is big

- DHTs have not had the hoped-for impact

## What DHTs got right

- **Consistent hashing**
  - Elegant way to divide a workload across machines
  - Very useful in clusters: actively used today in Amazon Dynamo and other systems
- **Replication** for high availability, efficient recovery
- **Incremental scalability**
  - Peers join with capacity, CPU, network, etc.
- **Self-management:** minimal configuration

47

47