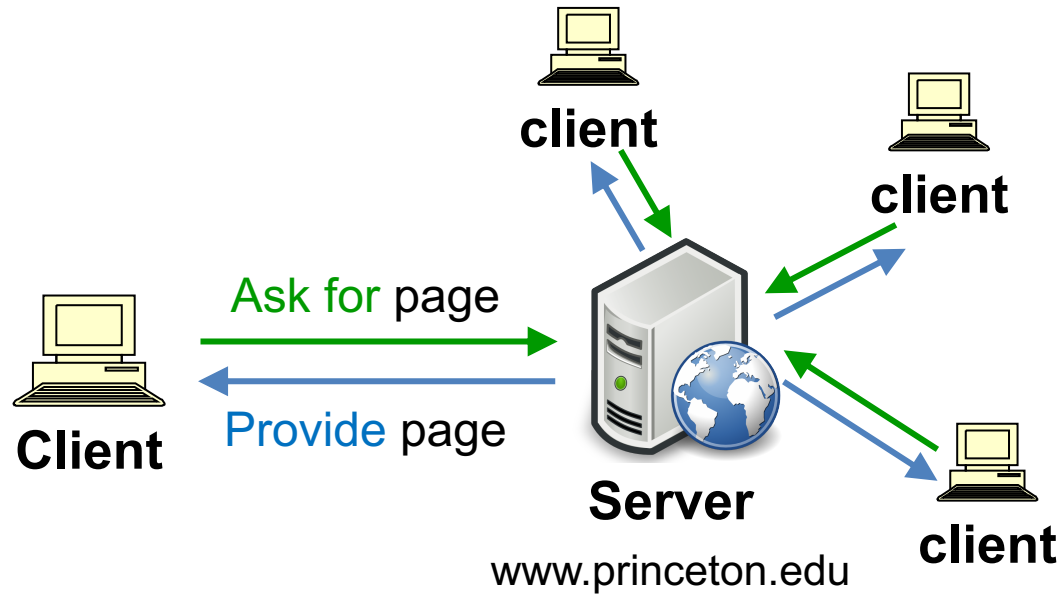# Peer-to-Peer Systems and Distributed Hash Tables

COS 418: Distributed Systems
Lecture 9
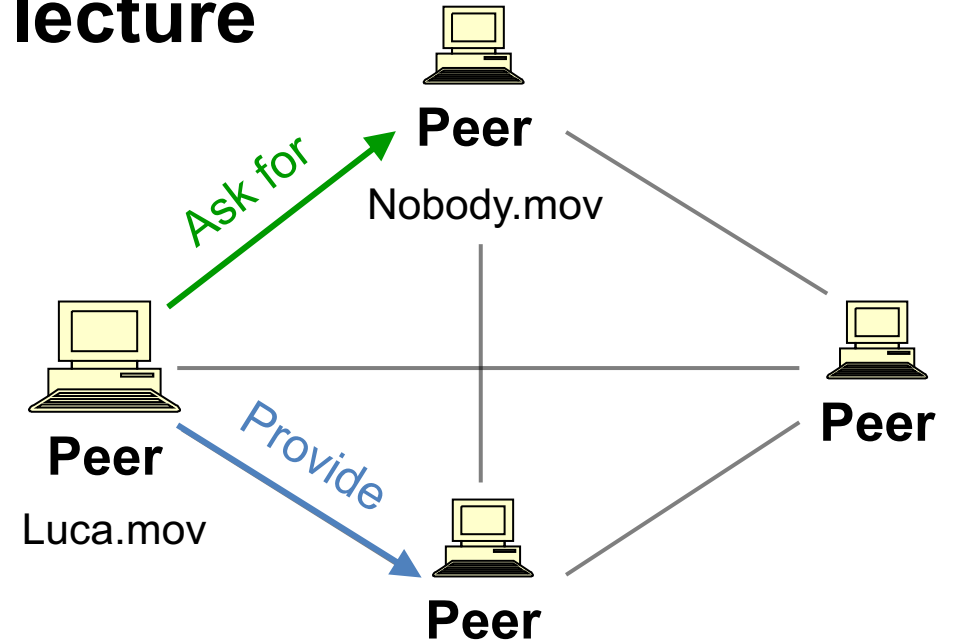
Haonan Lu

# Distributed Application Architecture



**Client-Server**

**Peer-to-Peer**

This lecture

Ask for page

Provide page

Client

Server

www.princeton.edu

client

client

client

Peer

Nobody.mov

Ask for

Provide

Luca.mov

Peer

Peer

Peer

# Today

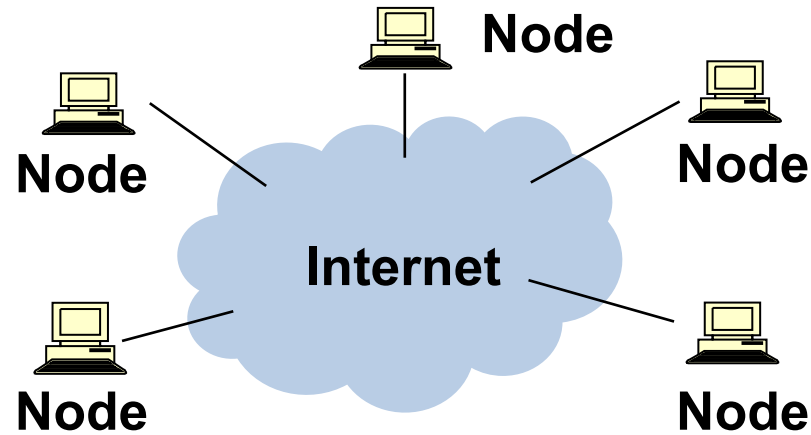1. **Peer-to-Peer Systems**
   - **What, why, and the core challenge**

2. Distributed Hash Tables (DHT)

3. The Chord Lookup Service

4. Concluding thoughts on DHTs, P2P

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



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


- **Large** number of **unreliable** nodes
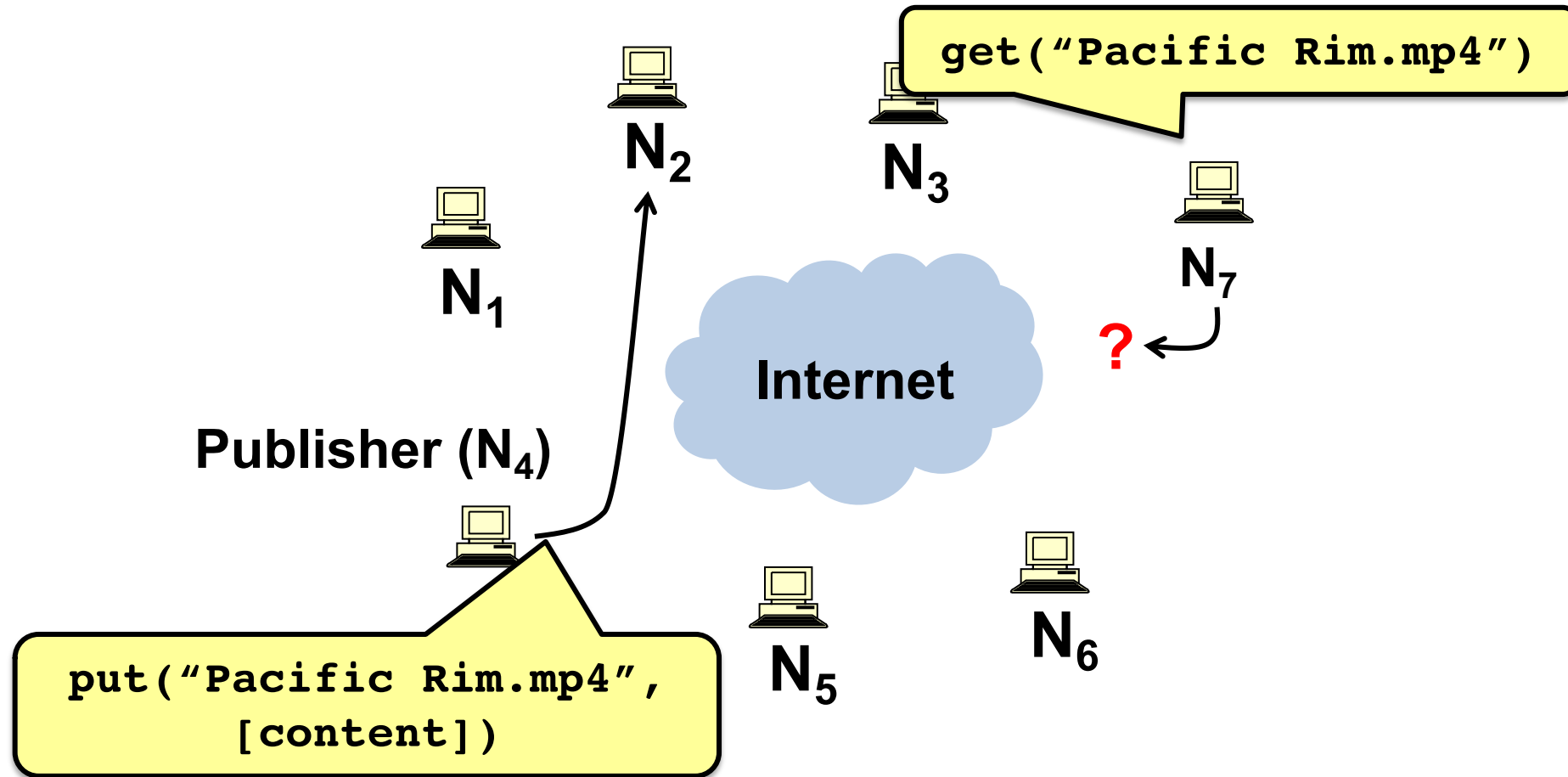
# P2P adoption

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)
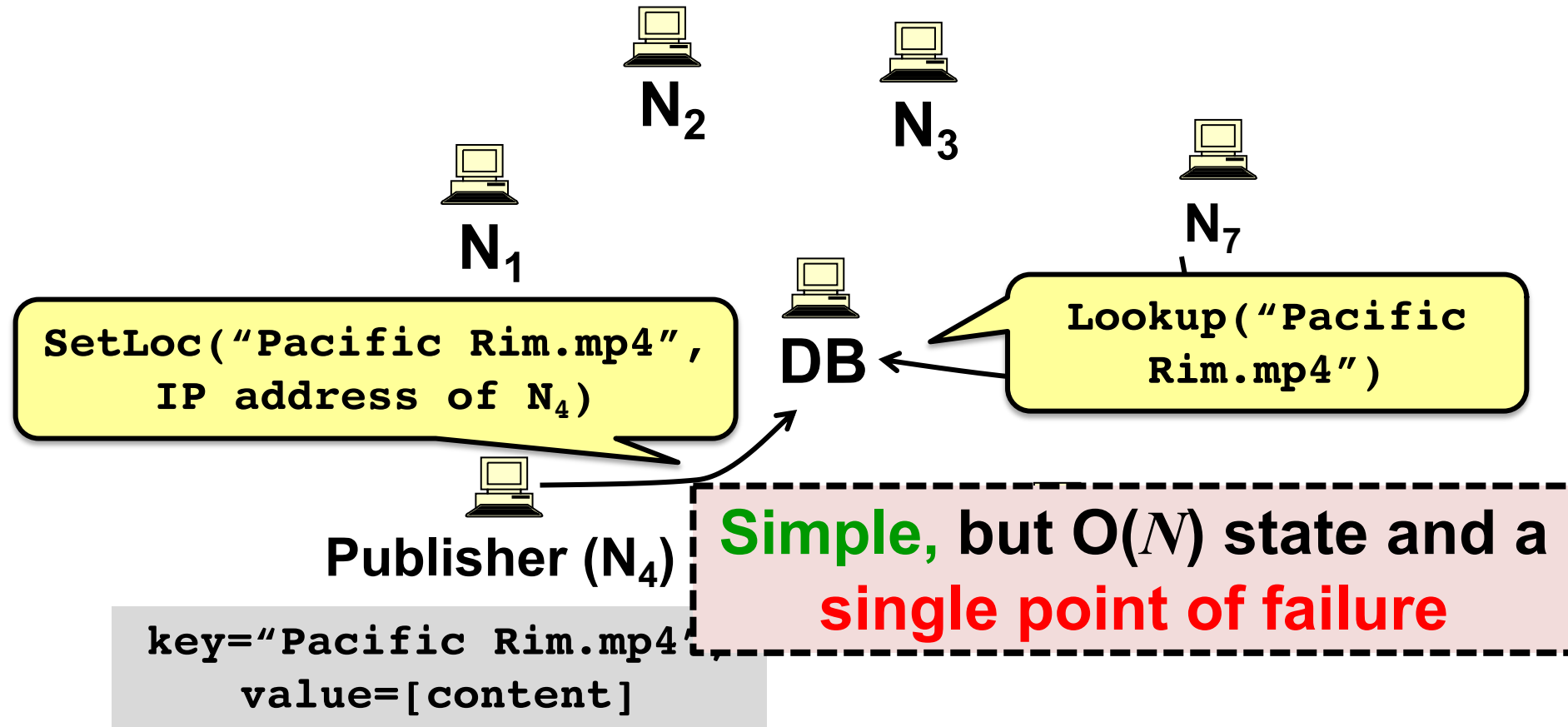   – Issues: Privacy and control

# 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**

# The lookup problem: locate the data



get("Pacific Rim.mp4")

N2

N3

N7

N1

?

Internet

Publisher (N4)

put("Pacific Rim.mp4",
[content])

N5

N6

# Centralized lookup (Napster)



N₂

N₃

N₁

DB

N₇

**SetLoc("Pacific Rim.mp4", IP address of N₄)**

**Lookup("Pacific Rim.mp4")**

**Publisher (N₄)**

**key="Pacific Rim.mp4", value=[content]**

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

# Flooded queries (original Gnutella)



Lookup("Pacific Rim.mp4")

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

Publisher ($N_4$)

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

# Tradeoffs in distributed systems

# Tradeoffs in distributed systems



**# msgs**

**# states**

Gnutella

Ideal

Napster

Nearly no states
Many msgs

Many states
Good performance
Single PoF

# Tradeoffs in distributed systems



# msgs

# states

Gnutella

**Nearly no states
Many msgs**

DHT
(Chord)

**msgs < Gnutella
states < Napster**

Napster

**Many states
Good performance
Single PoF**

# Today

1. Peer-to-Peer Systems

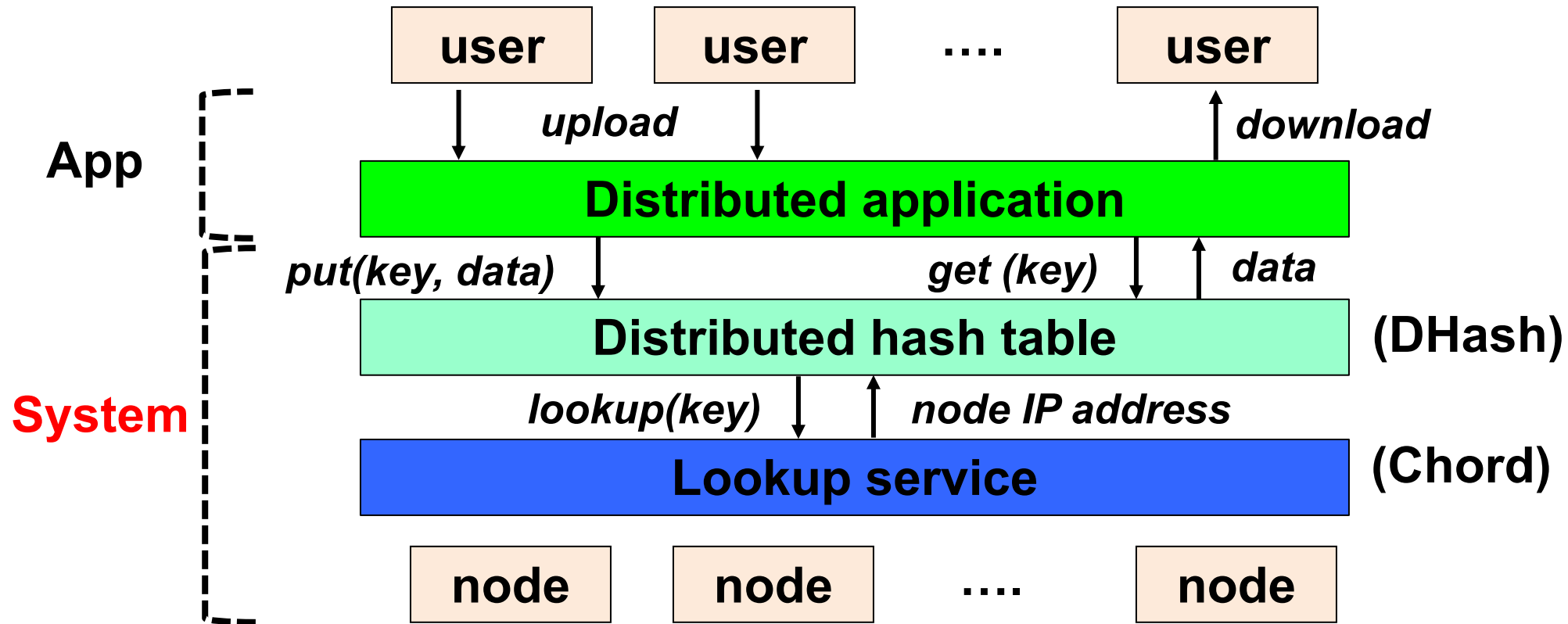2. **Distributed Hash Tables (DHT)**

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



| user | user | …. | user |

App

*upload* *download*

**Distributed application**

*put(key, data)* *get (key)* *data*

System

**Distributed hash table** (DHash)

*lookup(key)* *node IP address*

**Lookup service** (Chord)

| node | node | …. | node |

# 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**

# 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

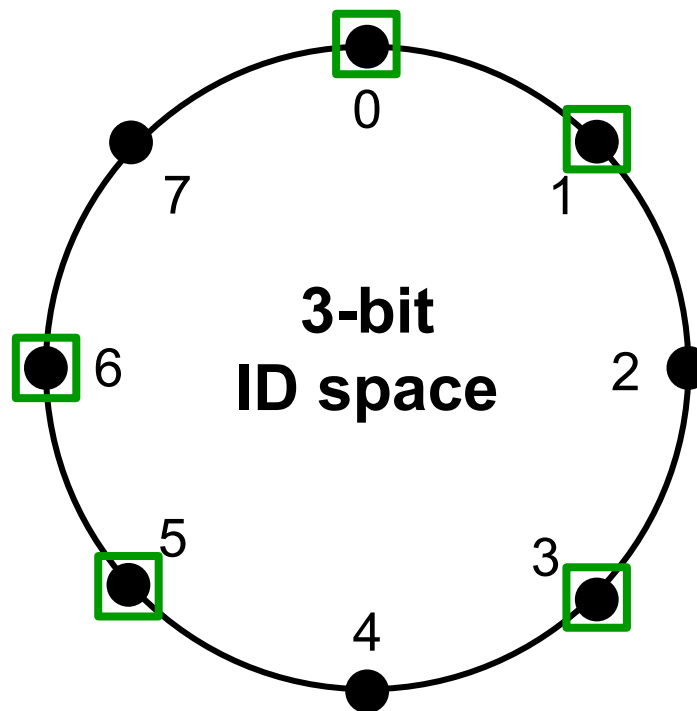# Consistent hashing [Karger '97] – data partition

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

●    **Identifiers/key space**

□    **Node**
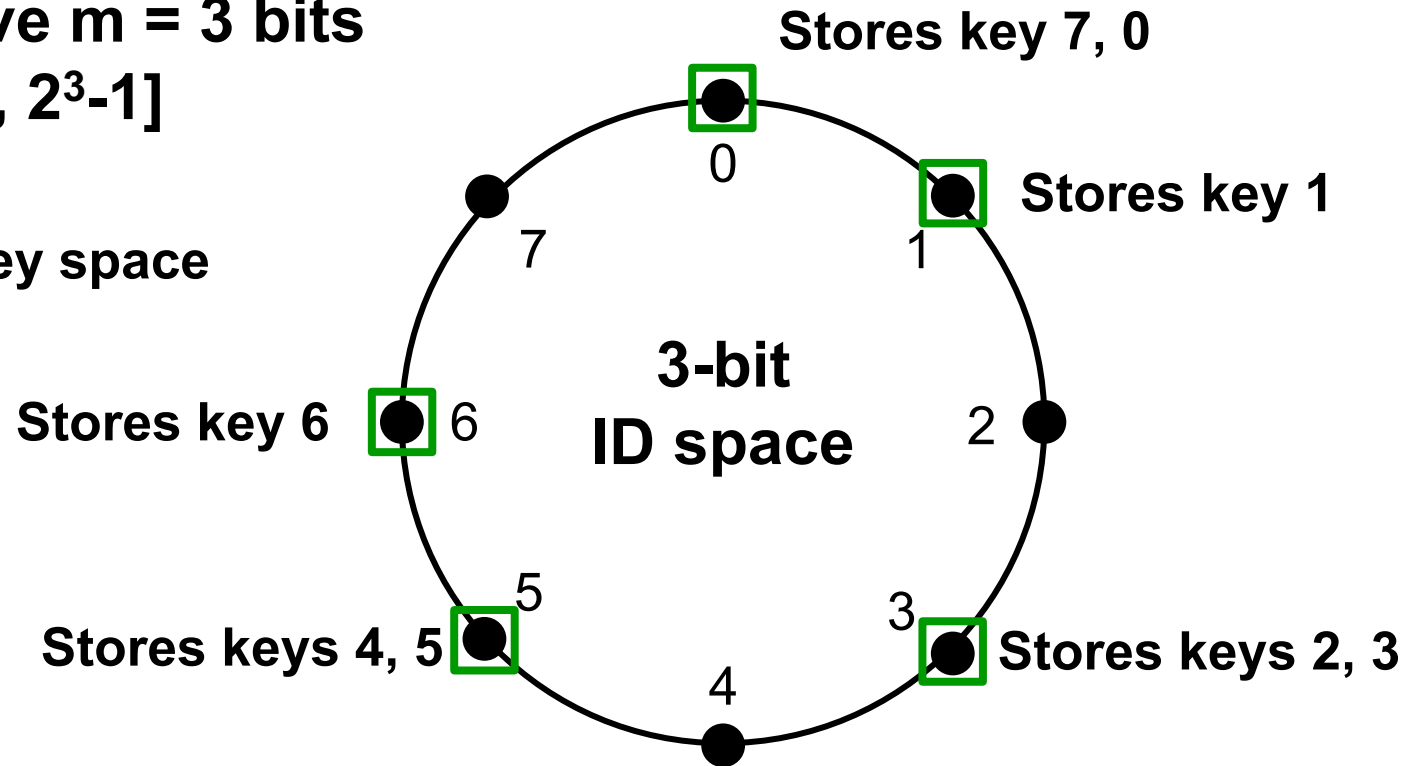
**3-bit ID space**

0
1
2
3
4
5
6
7

# Consistent hashing [Karger '97] – data partition

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

● **Identifiers/key space**

□ **Node**



Stores key 7, 0

Stores key 1

Stores key 6

Stores keys 4, 5

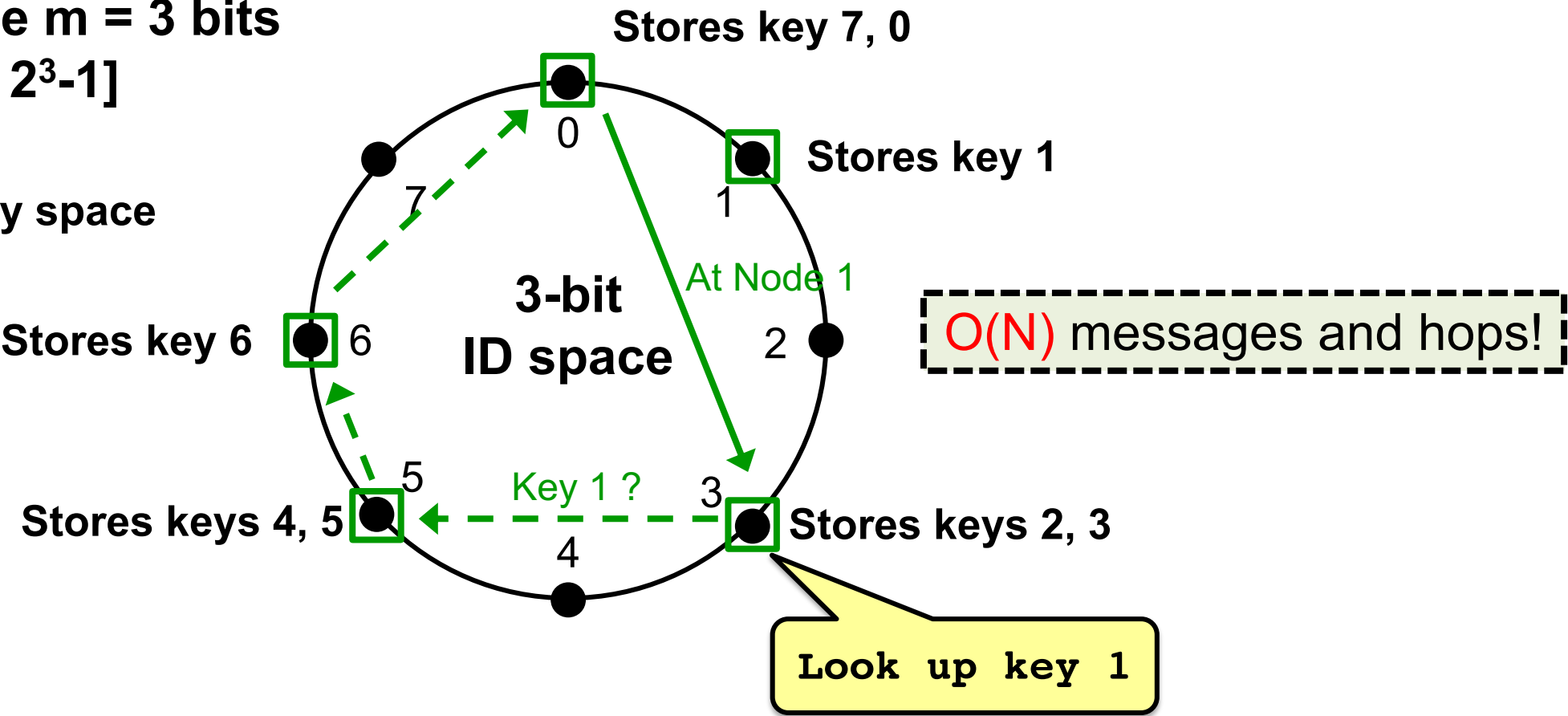Stores keys 2, 3

**3-bit
ID space**

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

# Consistent hashing [Karger '97] – basic lookup

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

Stores key 7, 0

Stores key 1

●    **Identifiers/key space**

☐    **Node**

---▸   **Successor pointer**

**3-bit ID space**

At Node 1

Stores key 6

Key 1 ?

Stores keys 4, 5

Stores keys 2, 3

O(N) messages and hops!

**Look up key 1**

# Chord – finger tables

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

●     **Identifiers/key space**

☐     **Node**

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

**3-bit
ID space**

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

**Separators**

**Key ranges**

**Successors
of separators**

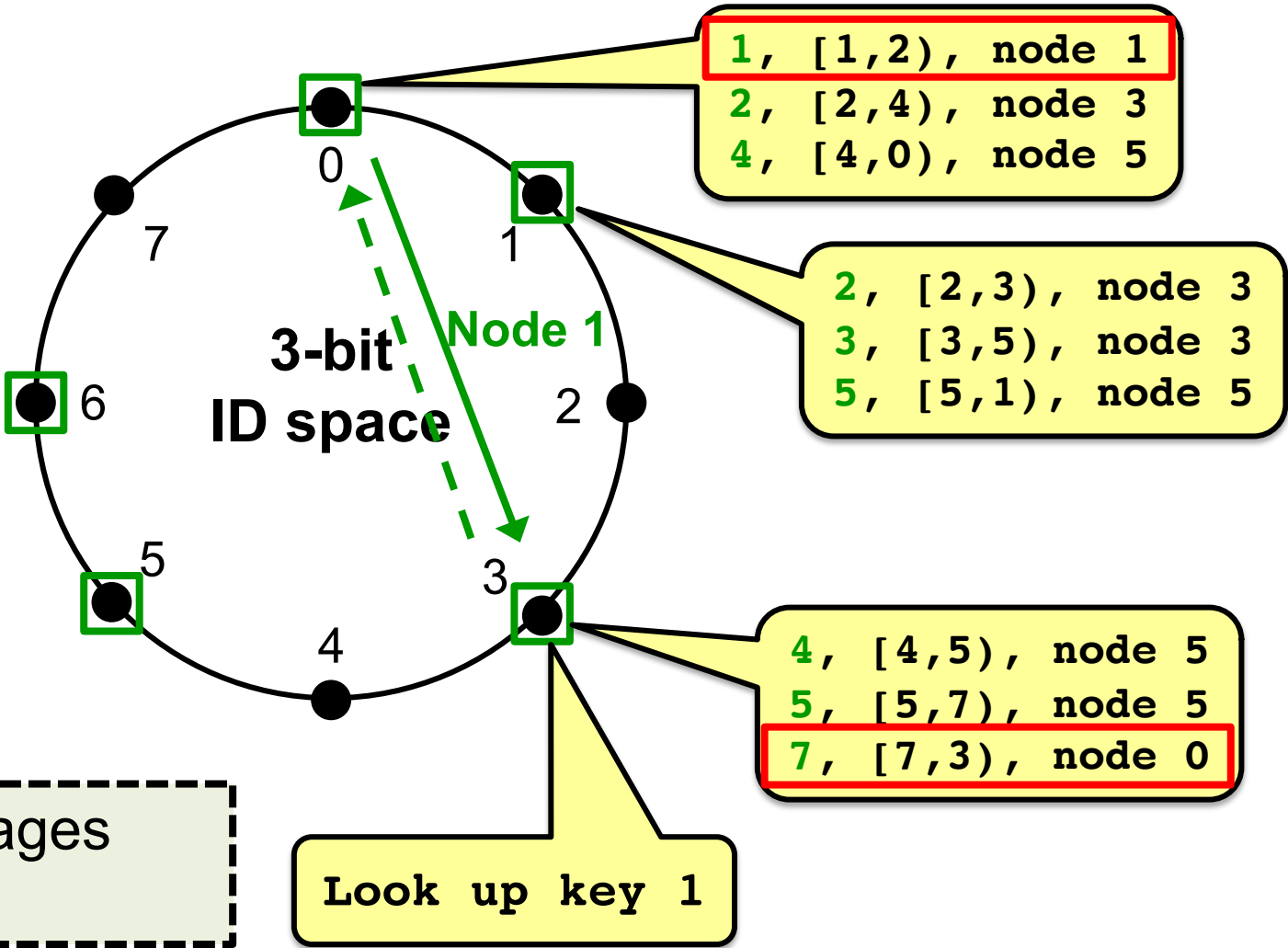# Chord – finger tables

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

●    **Identifiers/key space**

□    **Node**

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

O(logN) messages
and hops!

**3-bit
ID space**

Node 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**

24

# 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
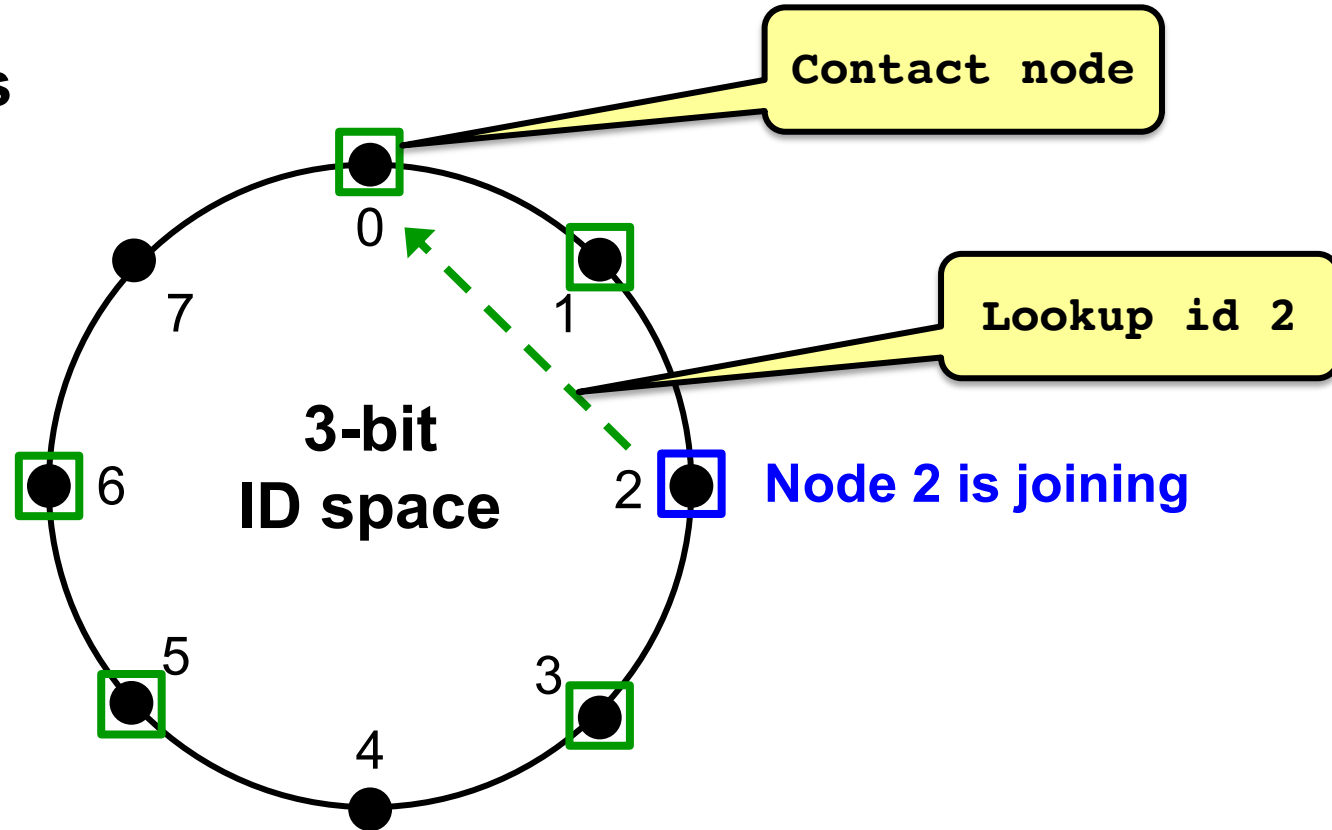
- **Robust:** survives massive failures

# Chord – node joining

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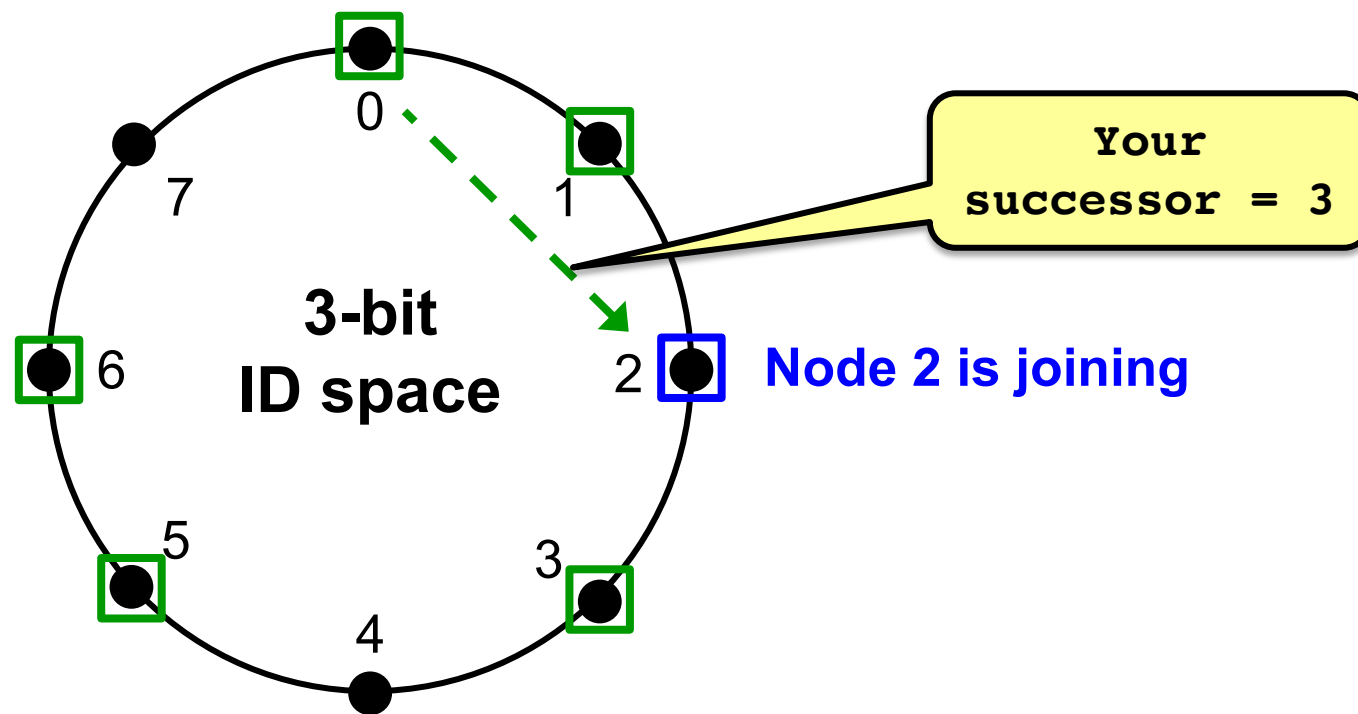
●    **Identifiers/key space**

□    **Node**



Contact node

Lookup id 2

Node 2 is joining

**3-bit ID space**

# 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**

# Chord – node joining

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

# Chord – node joining

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

●    **Identifiers/key space**

☐    **Node**



**3-bit**
**ID space**

→ **Points to successor**
→ **Points to predecessor**

**Node 2 is joining**

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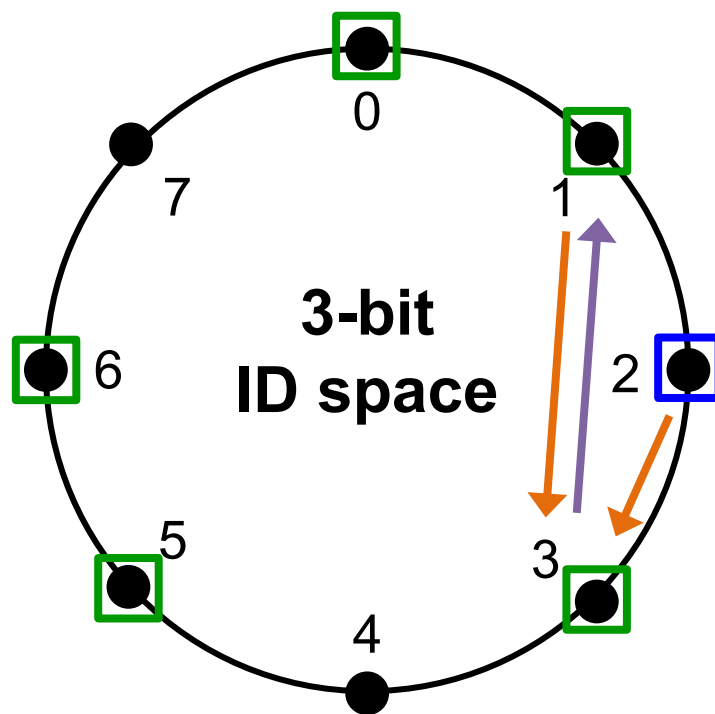**from each node to its successor**
**maintain node positions**

# Chord – node joining

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



● Identifiers/key space

□ Node

**3-bit ID space**

→ **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

*The pseudocode comes from Rodrigo Fonseca's lecture notes

# Chord – node joining

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

●    **Identifiers/key space**

☐    **Node**

**3-bit
ID space**

<span style="color:orange">→</span> **Points to successor**
<span style="color:purple">→</span> **Points to predecessor**

<span style="color:blue">**Node 2 is joining**</span>

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
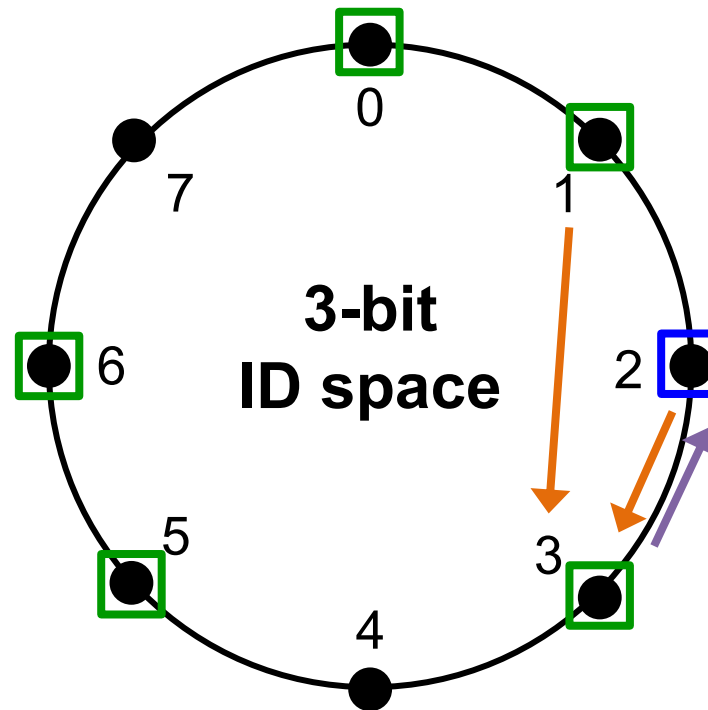
*The pseudocode comes from Rodrigo Fonseca's lecture notes

# Chord – node joining

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

● **Identifiers/key space**

□ **Node**



→ **Points to successor**
→ **Points to predecessor**

**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

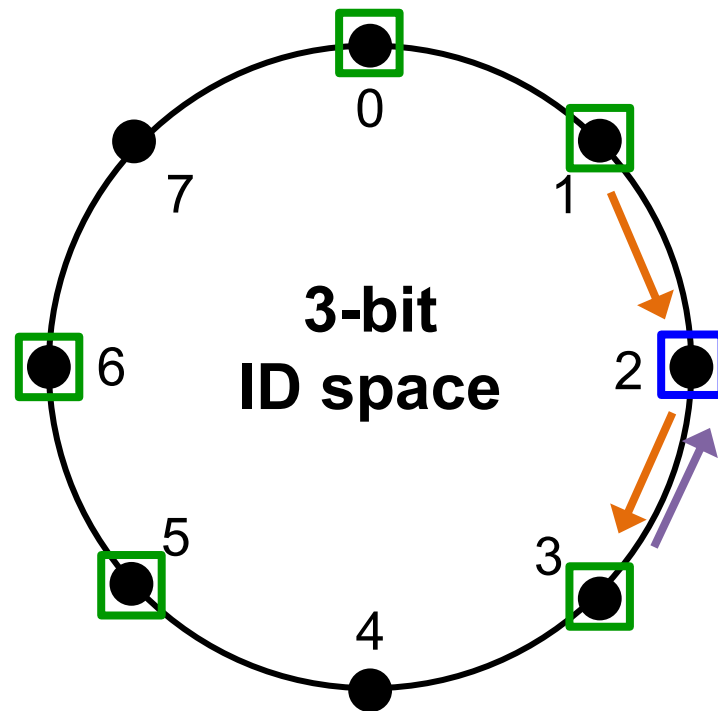*The pseudocode comes from Rodrigo Fonseca's lecture notes

34

# Chord – node joining

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

● **Identifiers/key space**

□ **Node**



→ **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
```

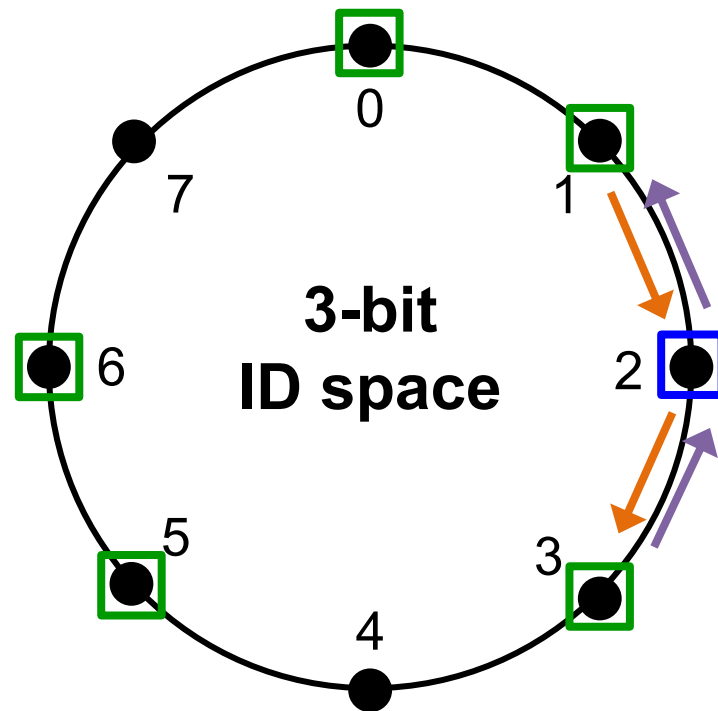*The pseudocode comes from Rodrigo Fonseca's lecture notes

35

# Chord – failures and successor list

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

● Identifiers/key space

□ Node

**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**

# Chord – failures and successor list

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

●    **Identifiers/key space**

□    **Node**

**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**
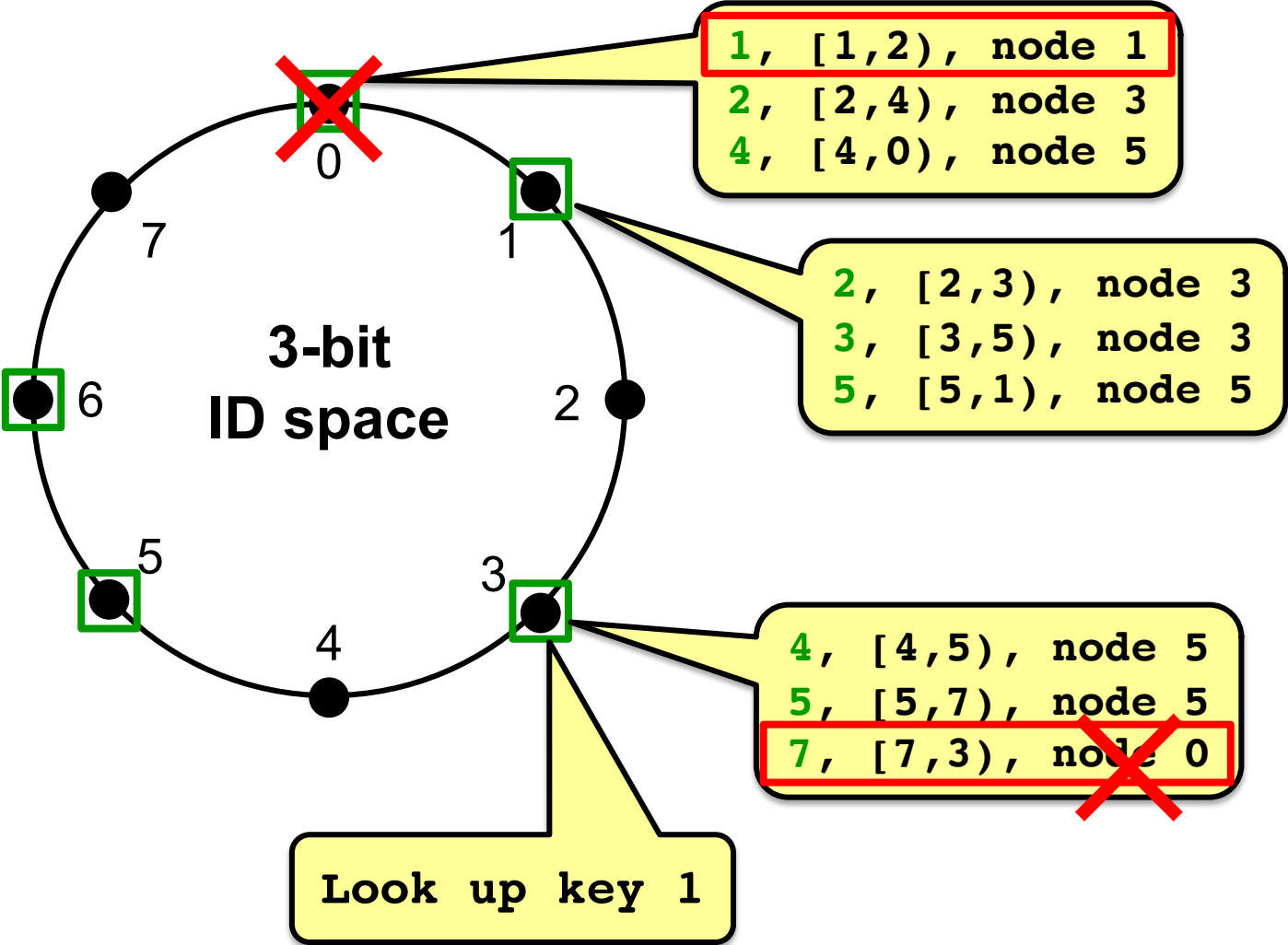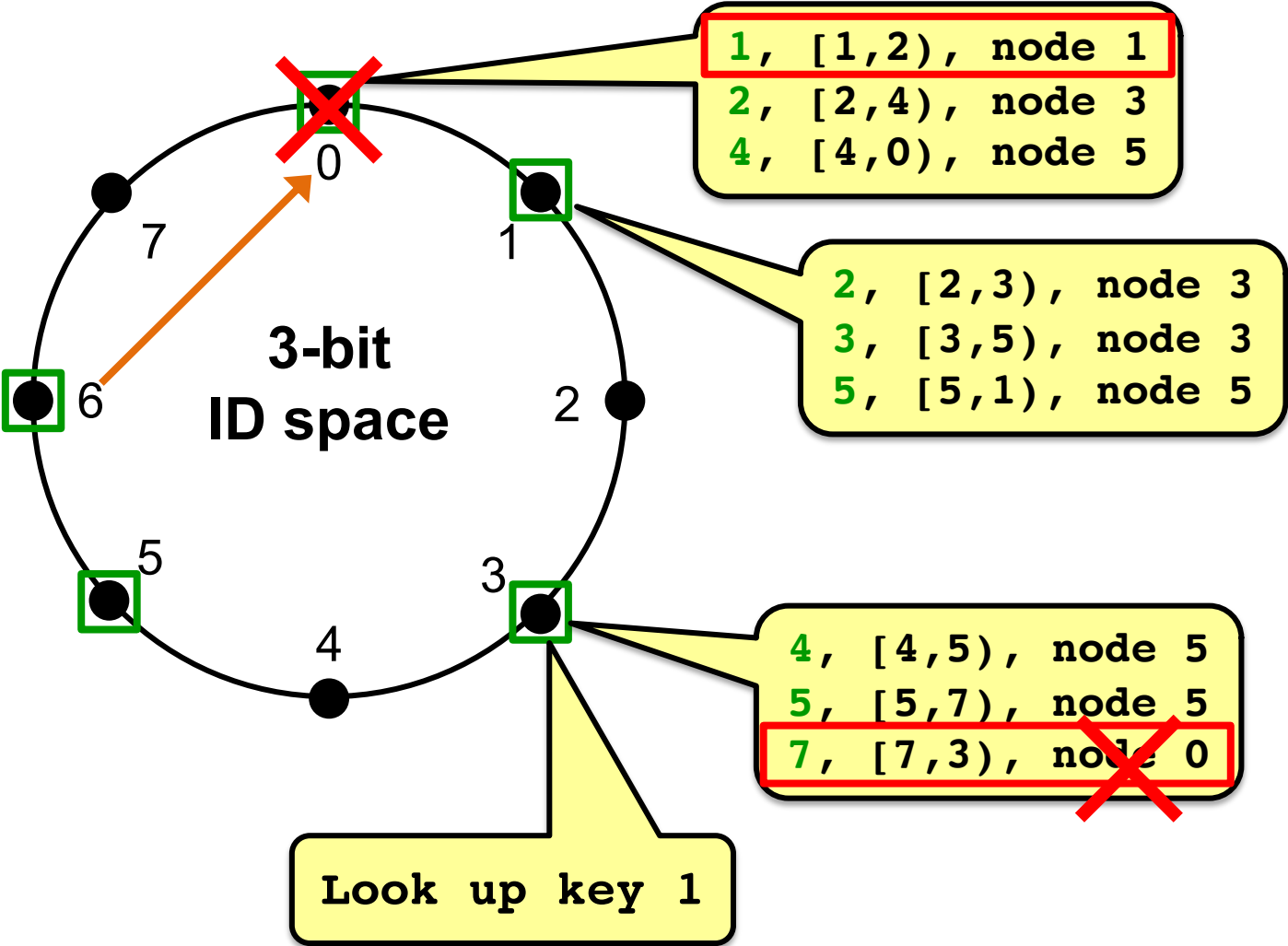
# Chord – failures and successor list

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

# Chord – failures and successor list

**Identifiers have m = 3 bits**

**Key space: [0, $2^3$-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)**

# Chord – failures and successor list

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

●    **Identifiers/key space**

□    **Node**

→    **Points to successor**

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

**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
```

**Look up key 1**

40

# 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?**

# DHash replicates blocks at *r* successors

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

●    **Identifiers/key space**

☐    **Node**



**Key 7**

**Key 7**

**3-bit
ID space**

7
0
1
6
2
5
3
4

**"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?

1. **High latency and limited bandwidth** between peers (vs. intra/inter-datacenter, client-server model)
   1. 1 M nodes = 20 hops; 50 ms / hop gives 1 sec lookup latency

2. User computers are **less reliable** than managed servers

3. **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