

# Peer-to-Peer Systems and Distributed Hash Tables



---

COS 418/518: Distributed Systems  
Lecture 9 & 10

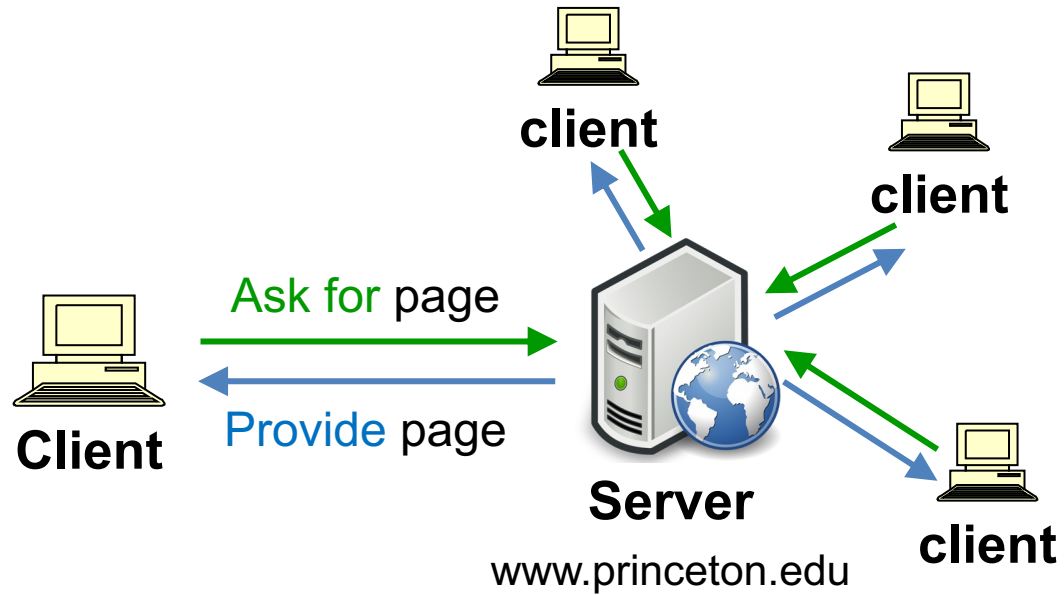
Jialin Ding, Mike Freedman

# Today

---

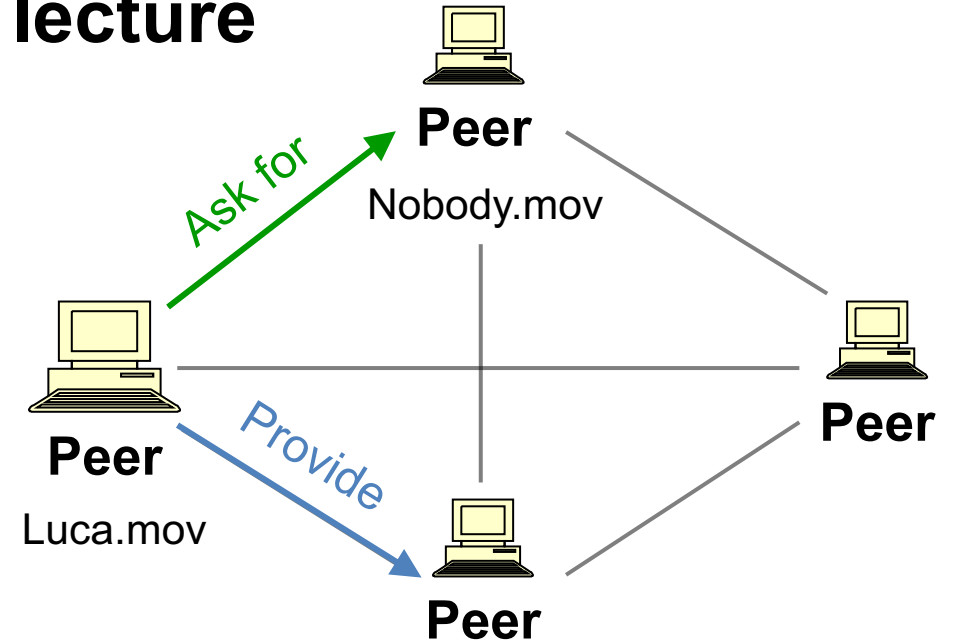
- 1. Peer-to-Peer Systems**
- 2. Distributed Hash Tables (DHT)**
- 3. The Chord Lookup Service**

# Distributed Application Architecture



**Client-Server**

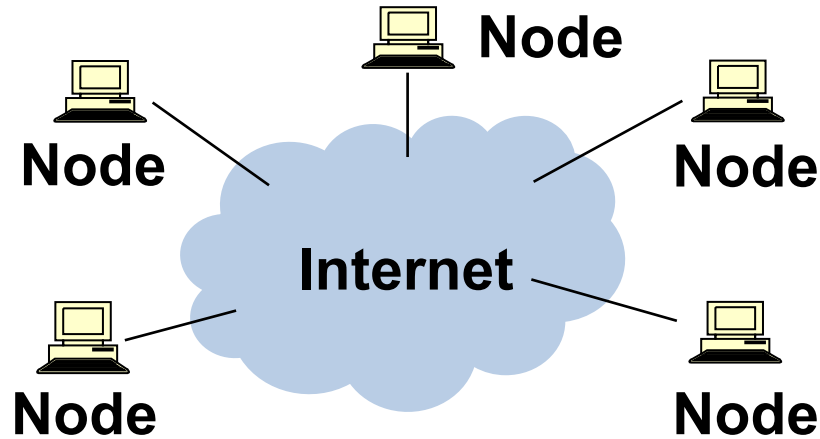
## This lecture



**Peer-to-Peer**

# 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

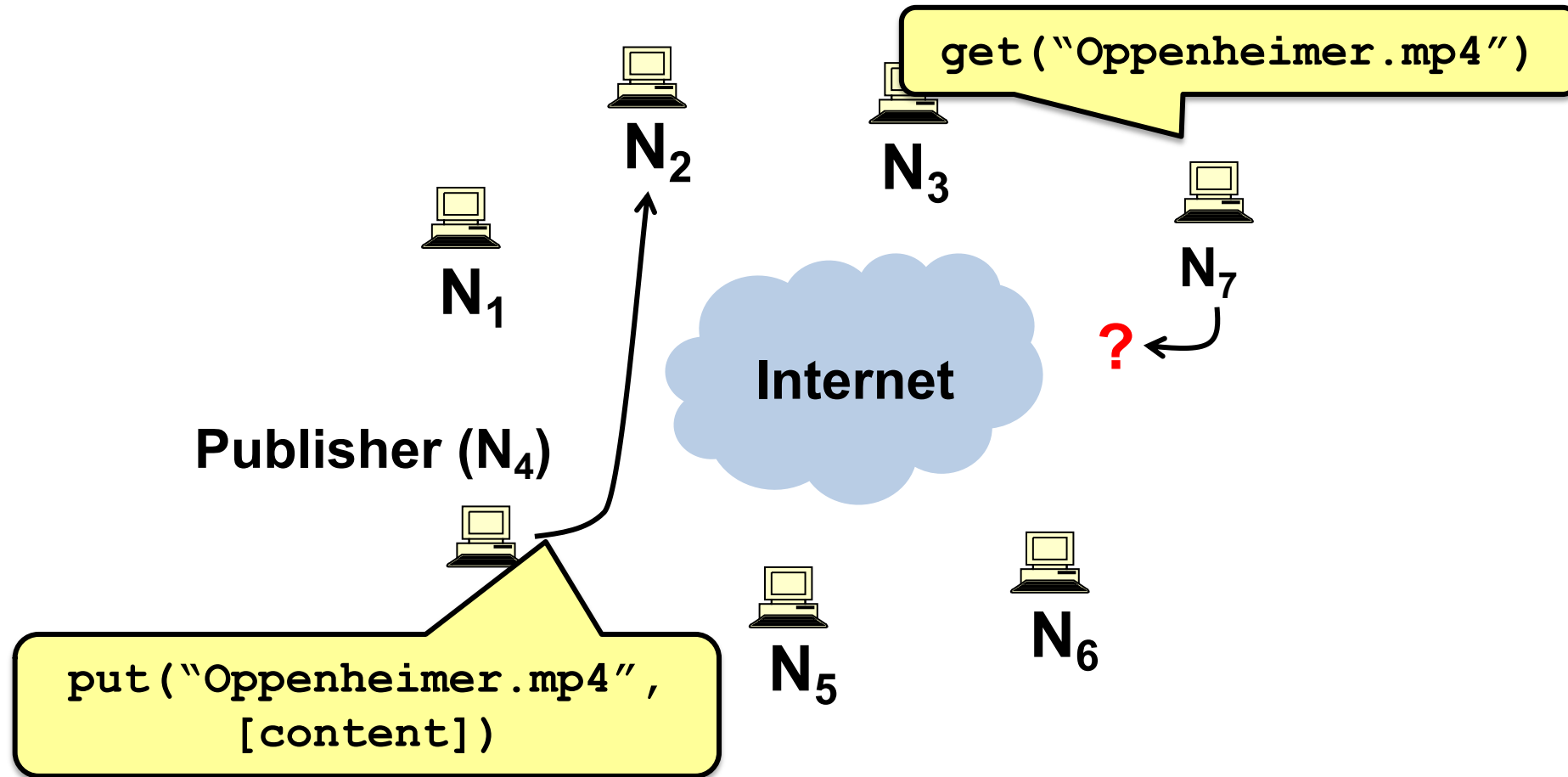
# P2P adoption

---

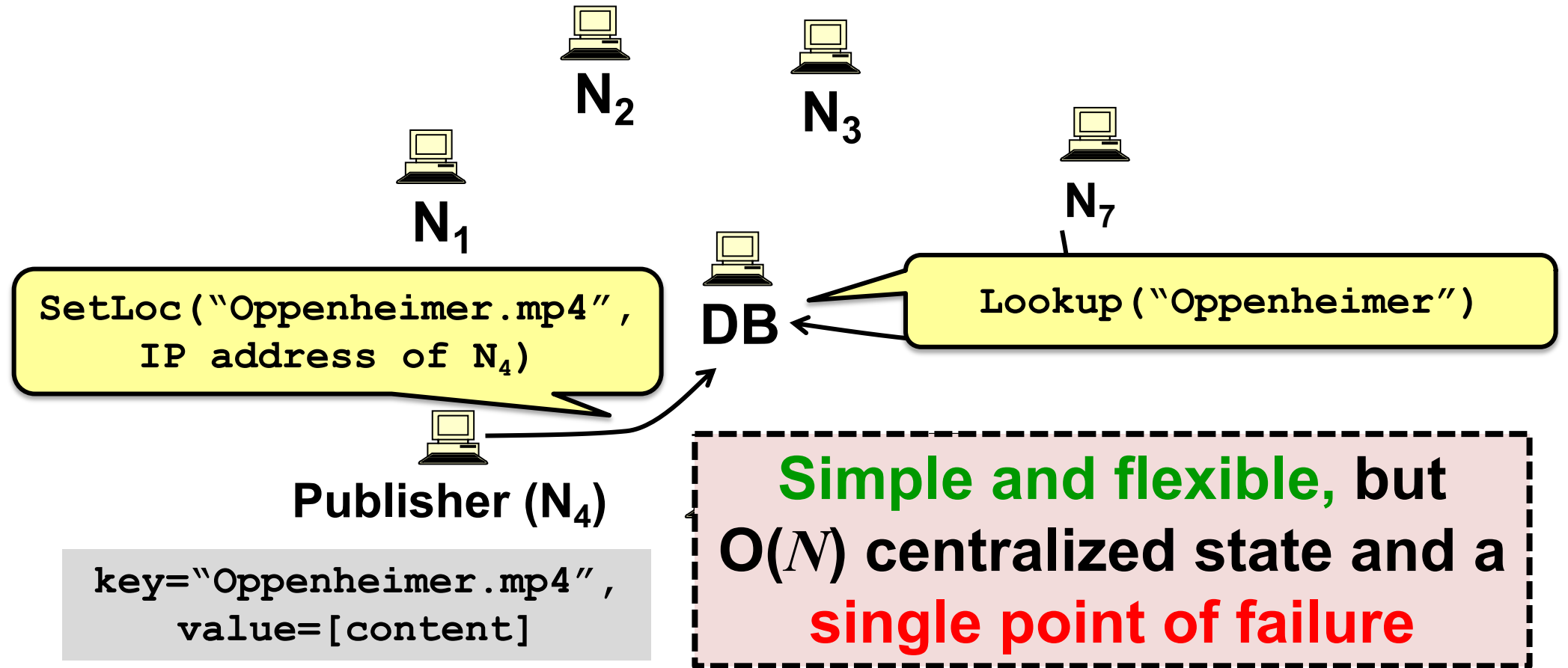
Successful adoption in **some niche areas**

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
  - Issues: Privacy and control

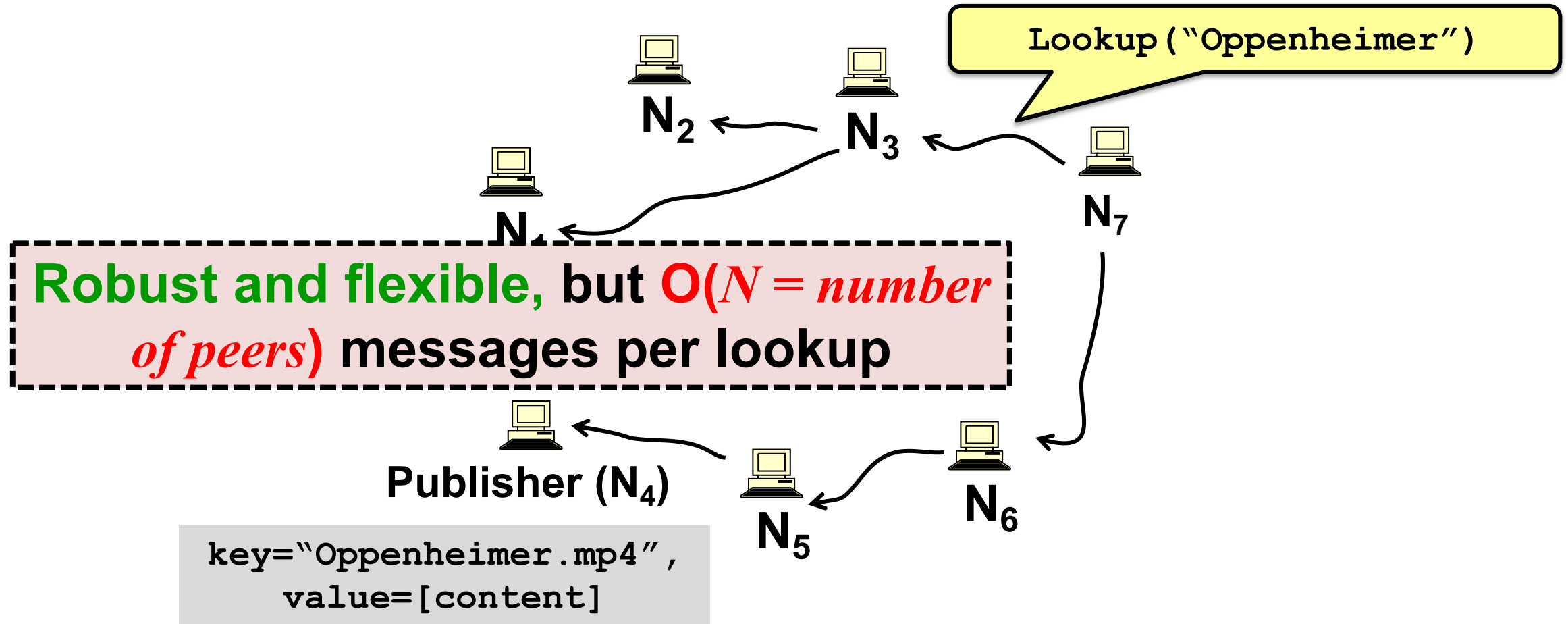
# The lookup problem: locate the data



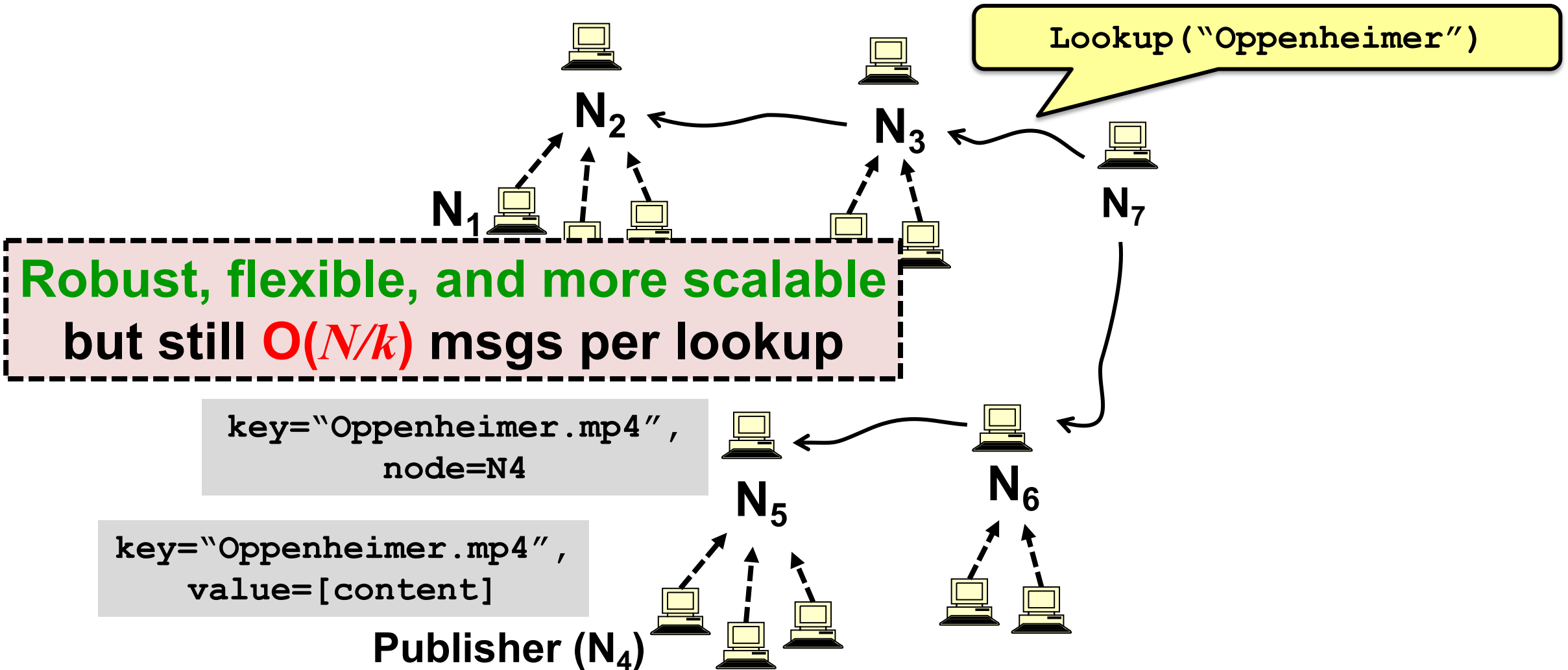
# Centralized lookup (Napster)



# Flooded queries (original Gnutella)

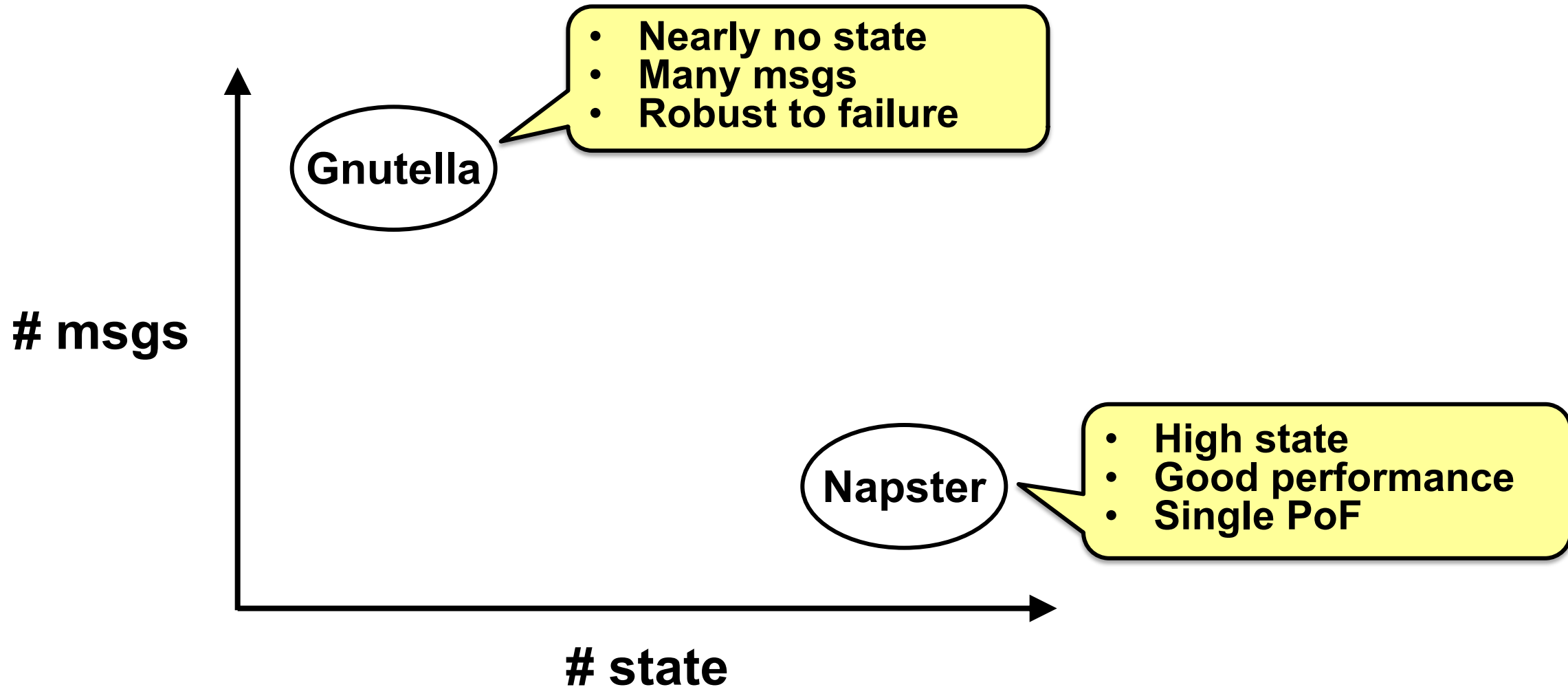


# Flooded queries pt 2 (Gnutella w/ SuperPeers)

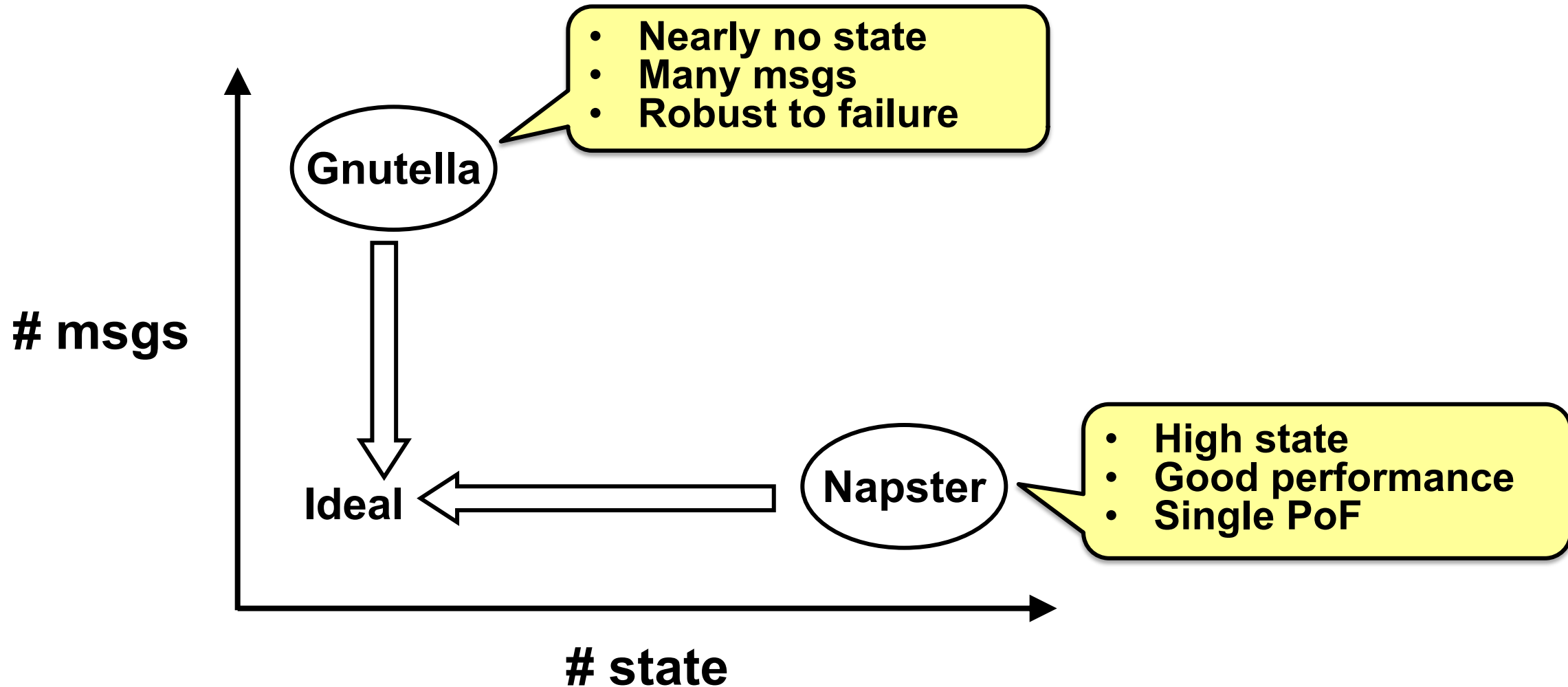


# Tradeoffs in distributed systems

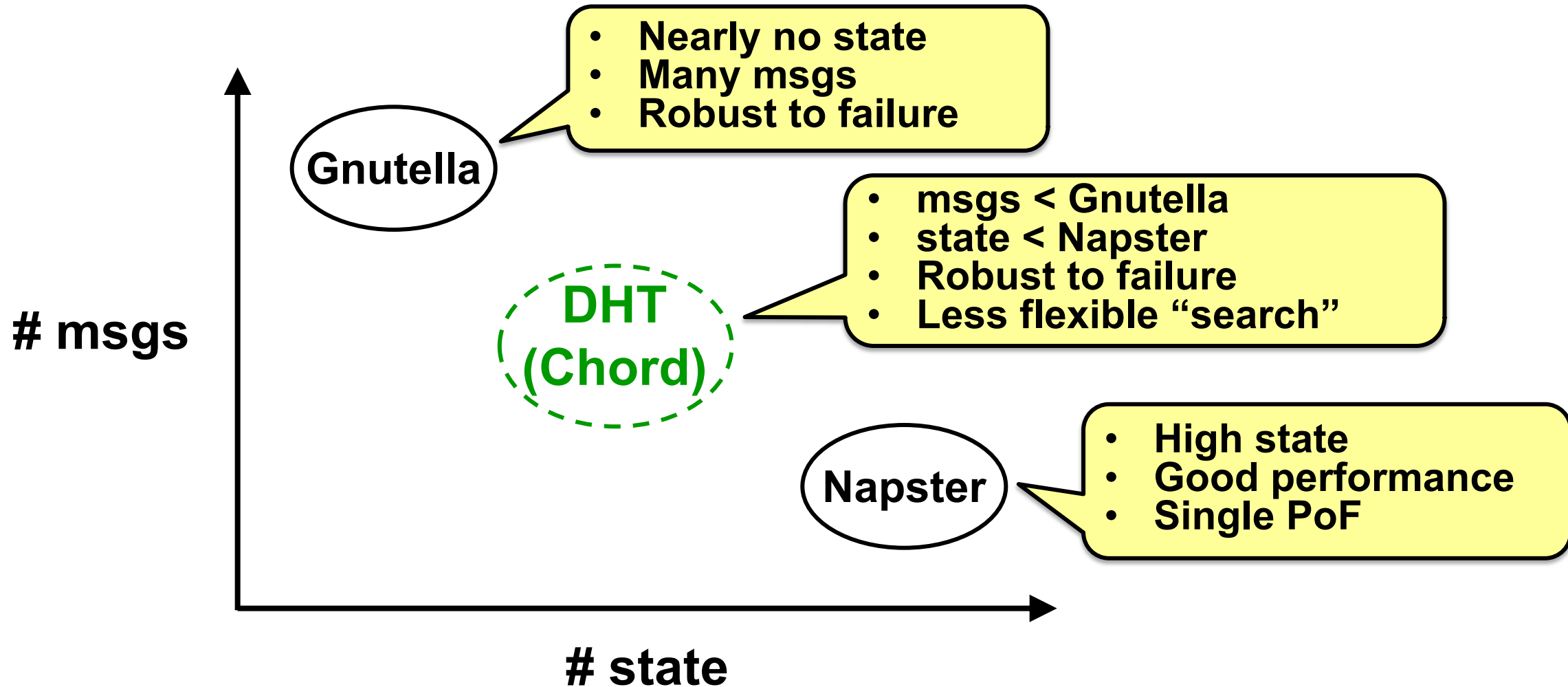
---



# Tradeoffs in distributed systems



# Tradeoffs in distributed systems



# What is a DHT (and why)?

---

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

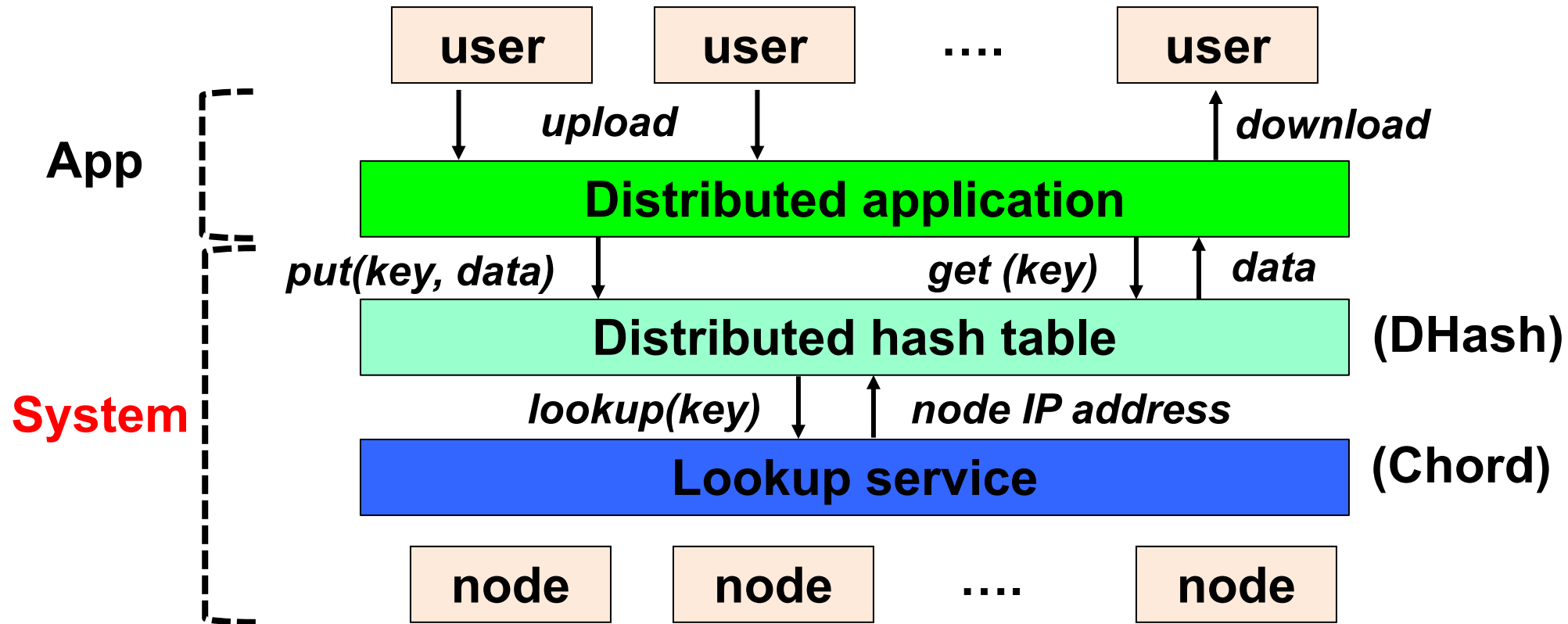
`lookup(key) → IP addr` (Chord lookup service)

`send-RPC(IP address, put, key, data_one)`

`send-RPC(IP address, get, key) → data_one`

- 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**

# Chord identifiers

---

- **Hashed values (integers) using the same hash function**
  - **Key identifier** =  $SHA-1(key) \bmod 2^{160}$
  - **Node identifier** =  $SHA-1(IP\ address) \bmod 2^{160}$
- **What is “SHA-1”?**
  - SHA-1 is a cryptographic hash function that maps input to 160-bit output hash
  - Some properties:
    1. Output hashes looks randomly distributed across output space
    2. Given  $hash1$ , hard to find  $input1$  where  $SHA1(input1) = hash1$
    3. Given  $input1$  and  $hash1$ , hard to find  $input2$  where  $SHA1(input2) = hash1$
    4. Hard to find  $input1$  and  $input2$  where  $SHA1(input1) = SHA1(input2)$

# Chord identifiers

---

- **Hashed values (integers) using the same hash function**
  - Key identifier =  $SHA-1(key) \bmod 2^{160}$
  - Node identifier =  $SHA-1(IP\ address) \bmod 2^{160}$
- **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  $\rightarrow$  load balancing; hashed address  $\rightarrow$  independent failure

# Alternative: mod (n) hashing

---

- **System of n nodes: 1...n**
  - Node that owns key is assigned via  $hash(key) \bmod n$
  - Good load balancing
- **What if a node fails?**
  - Instead of n nodes, now  $n - 1$  nodes
  - Mapping of all keys change, as now  $hash(key) \bmod (n-1)$ 
    - **N = 5**
      - 12594  $\bmod 5 = 4$
      - 28527  $\bmod 5 = 2$
      - 816  $\bmod 5 = 1$
      - 716565  $\bmod 5 = 0$
    - **N = 4**
      - 12594  $\bmod 4 = 2$
      - 28527  $\bmod 4 = 3$
      - 816  $\bmod 4 = 0$
      - 716565  $\bmod 4 = 1$

# Consistent hashing [Karger '97]

## Data partitioning

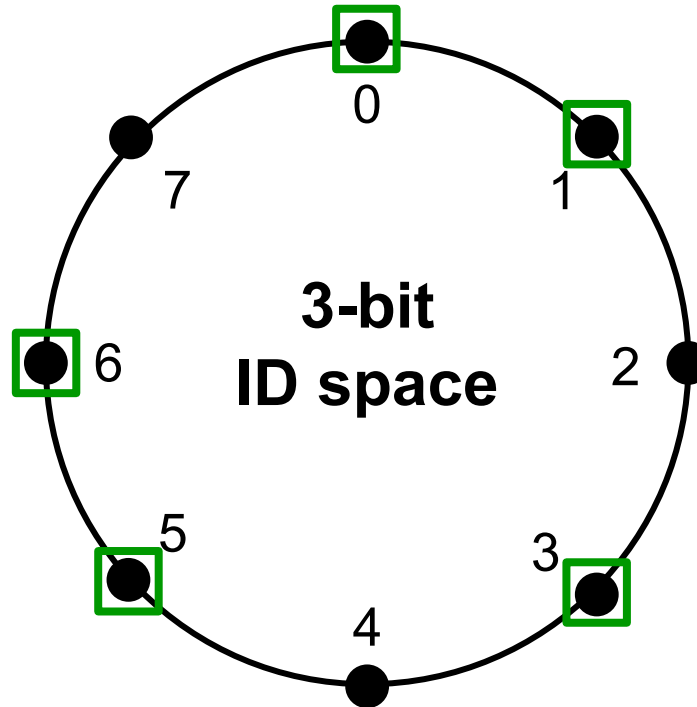
---

Identifiers have  $m = 3$  bits

Key space:  $[0, 2^3-1]$

● Identifiers/key space

□ Node



# Consistent hashing [Karger '97]

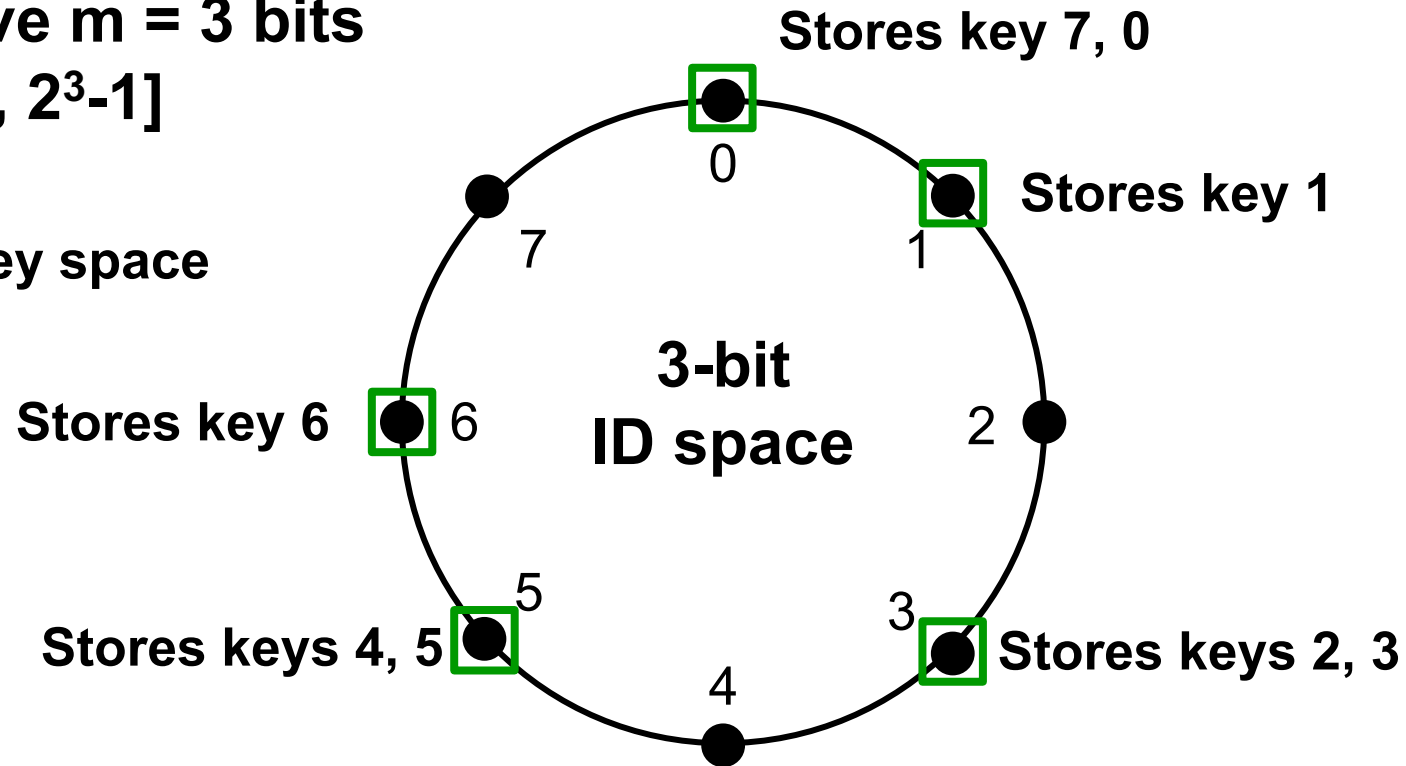
## Data partitioning

Identifiers have  $m = 3$  bits

Key space:  $[0, 2^3-1]$

● Identifiers/key space

□ Node



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

# Consistent hashing [Karger '97]

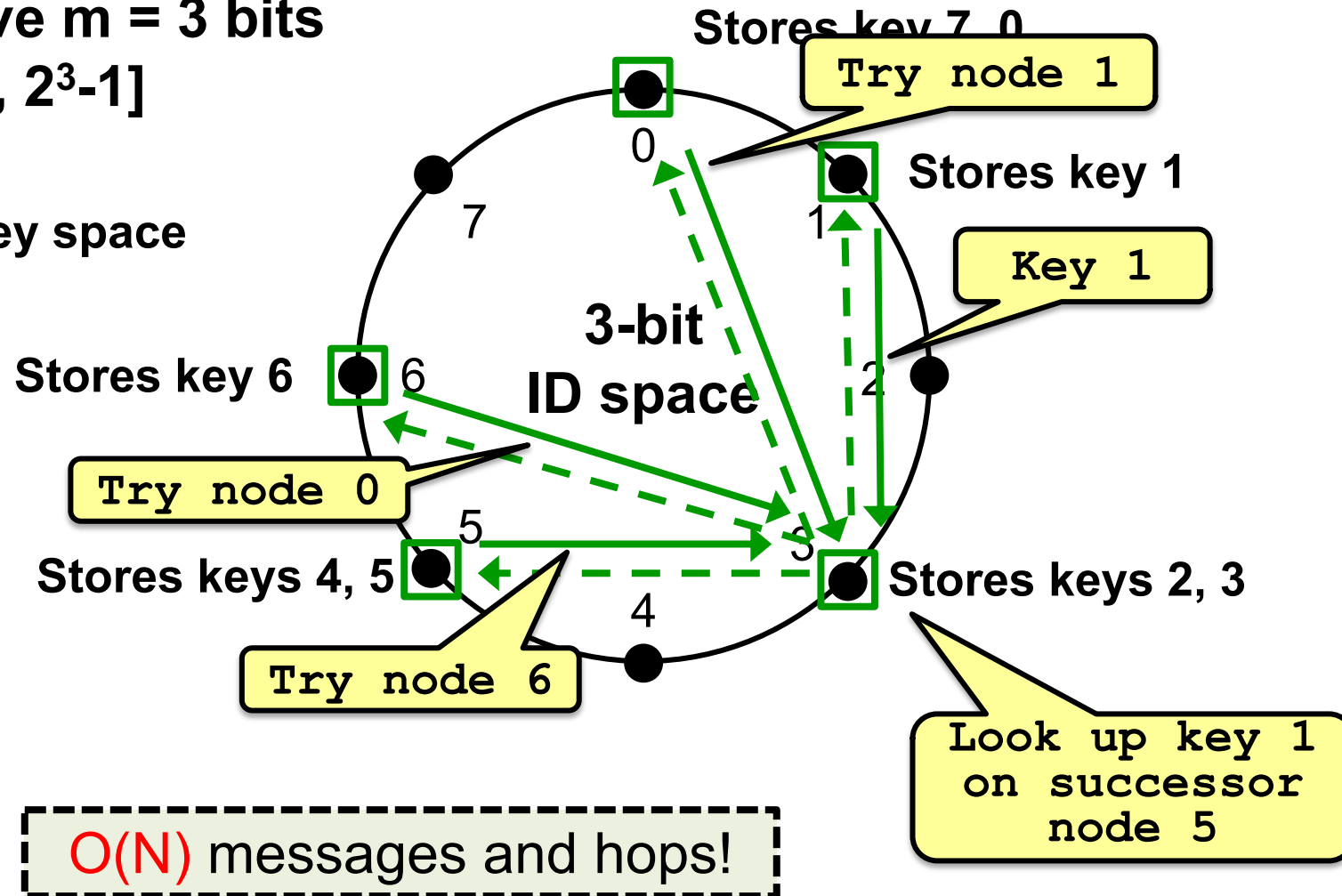
## Strawman lookup vis successors

Identifiers have  $m = 3$  bits

Key space:  $[0, 2^3-1]$

● Identifiers/key space

□ Node



# Consistent hashing [Karger '97]

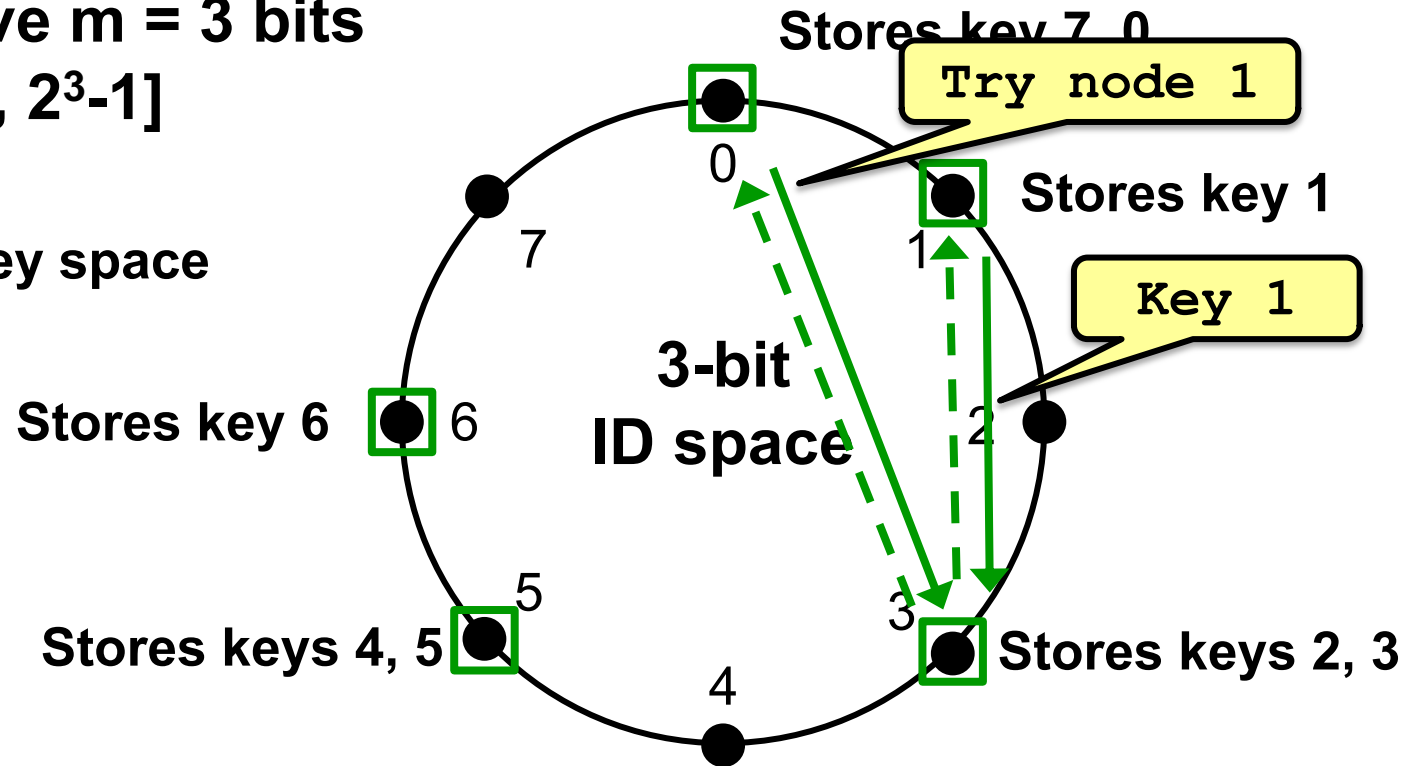
## Observation about last hop

Identifiers have  $m = 3$  bits

Key space:  $[0, 2^3-1]$

● Identifiers/key space

□ Node



Try to find key's **predecessor node** as fast as possible.  
This predecessor will know key's successor (owner).

# Chord – finger tables for *find\_predecessor*

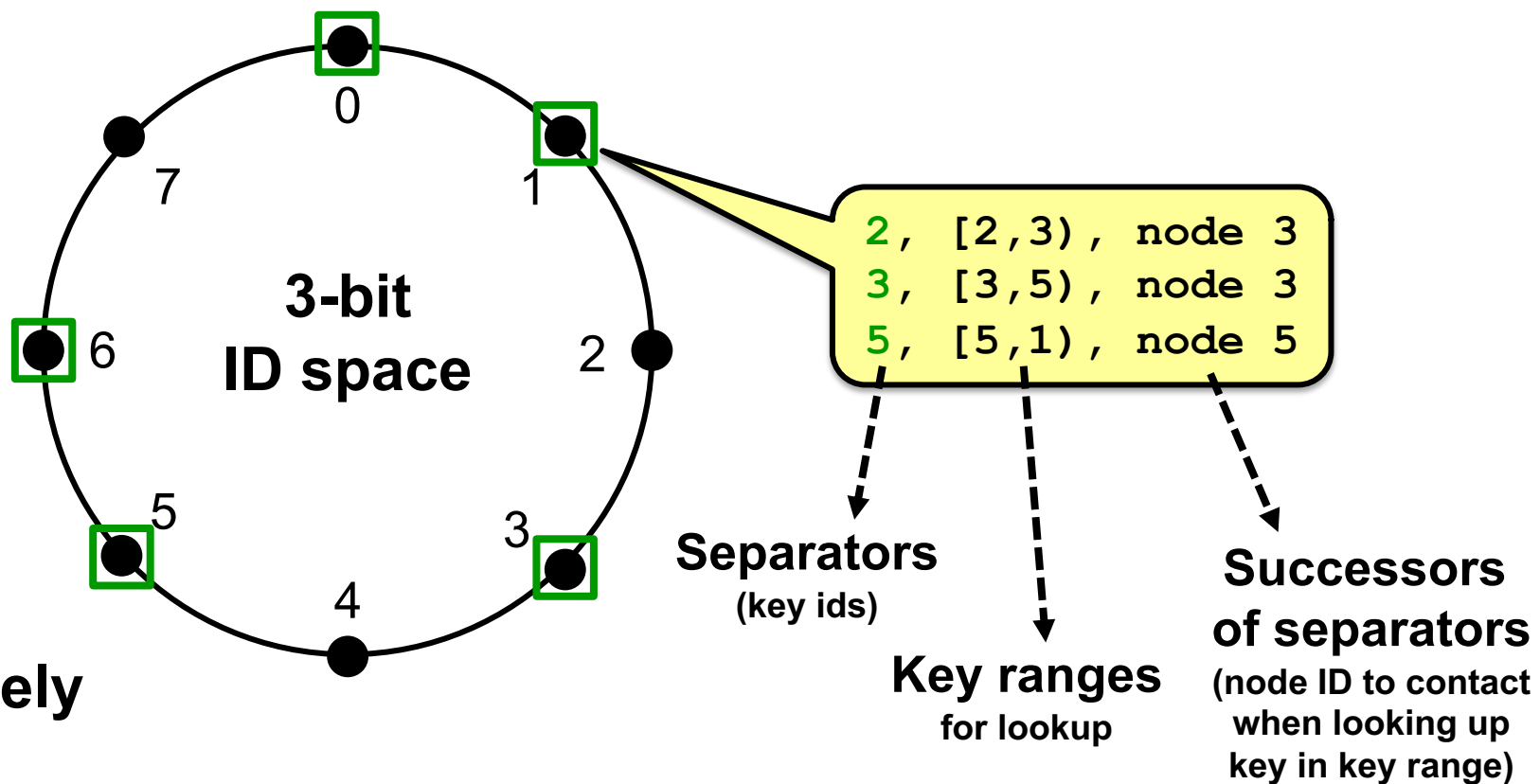
Each node keeps **m** states

Key space  $\rightarrow$  **m** ranges via  
 $(N + 2^{k-1}) \bmod 2^m, 1 \leq k \leq m$

Example for node  $N = 1$ :

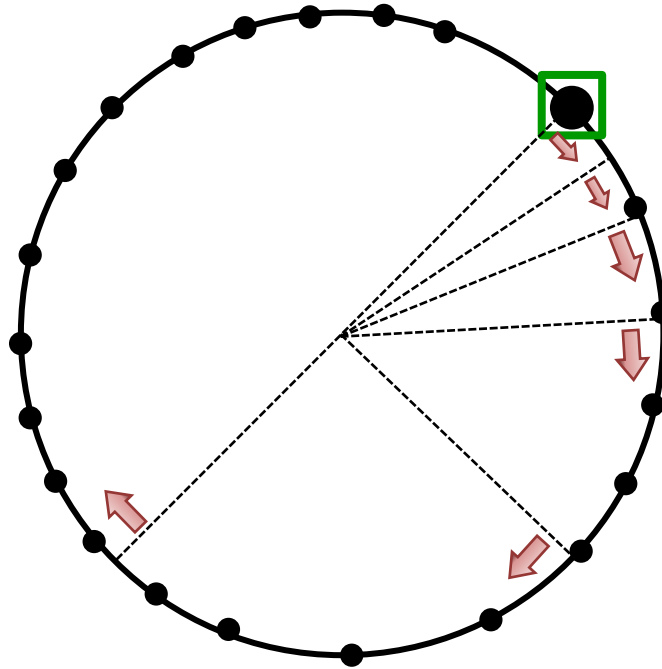
- $1 + 1 \bmod 8 \Rightarrow 2$
  - $1 + 2 \bmod 8 \Rightarrow 3$
  - $1 + 4 \bmod 8 \Rightarrow 5$
- N**       **$2^{k-1}$**

“Finger” is node immediately succeeding separator



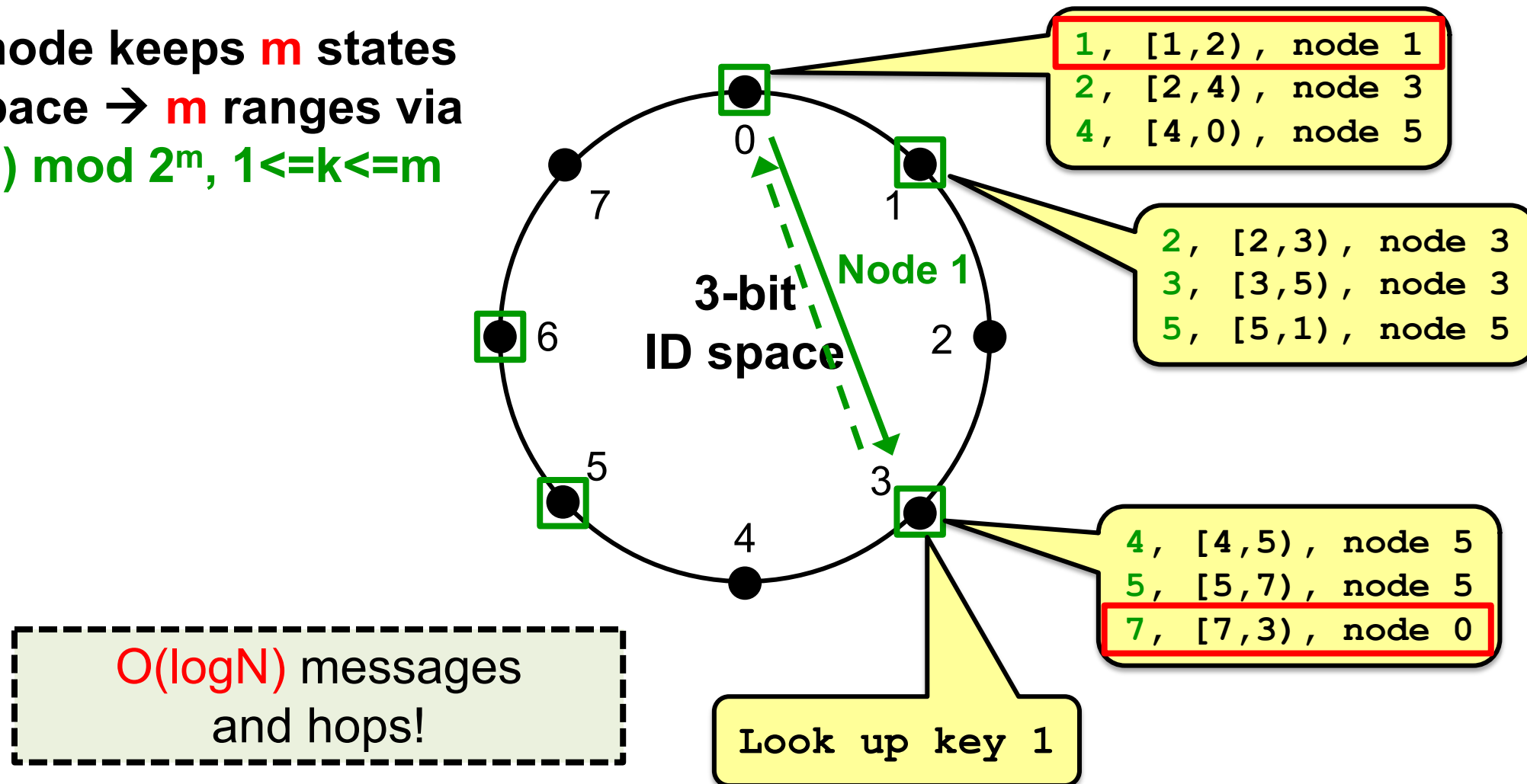
# Chord – finger tables for *find\_predecessor*

---

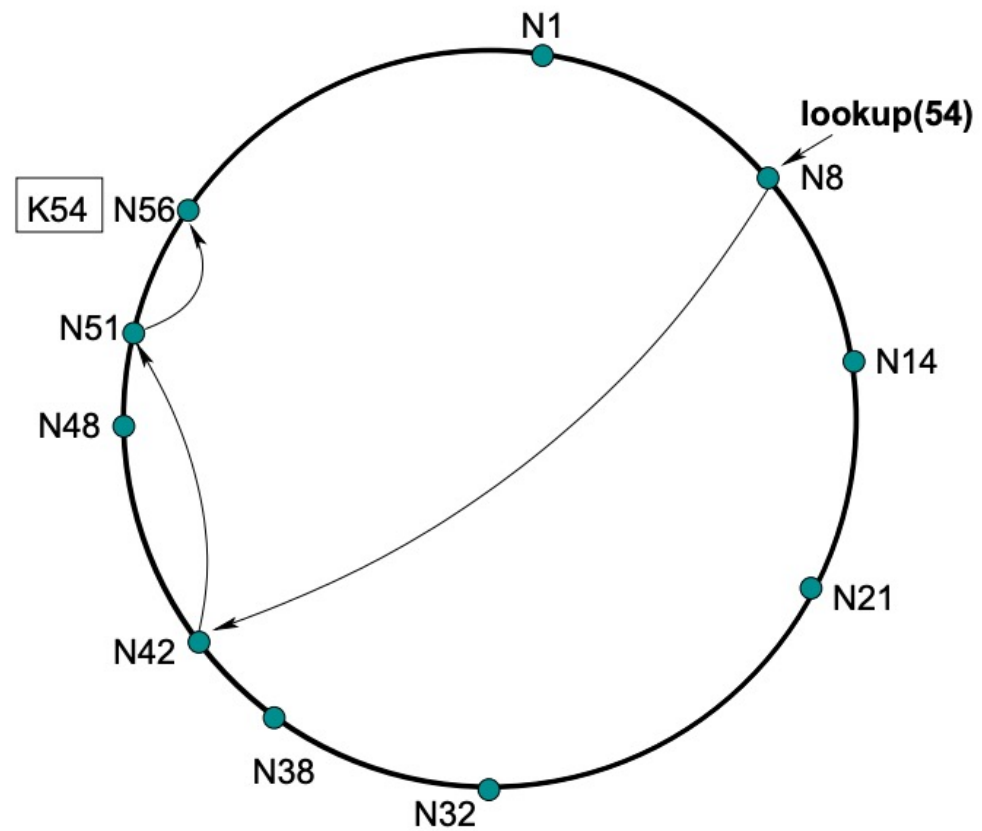
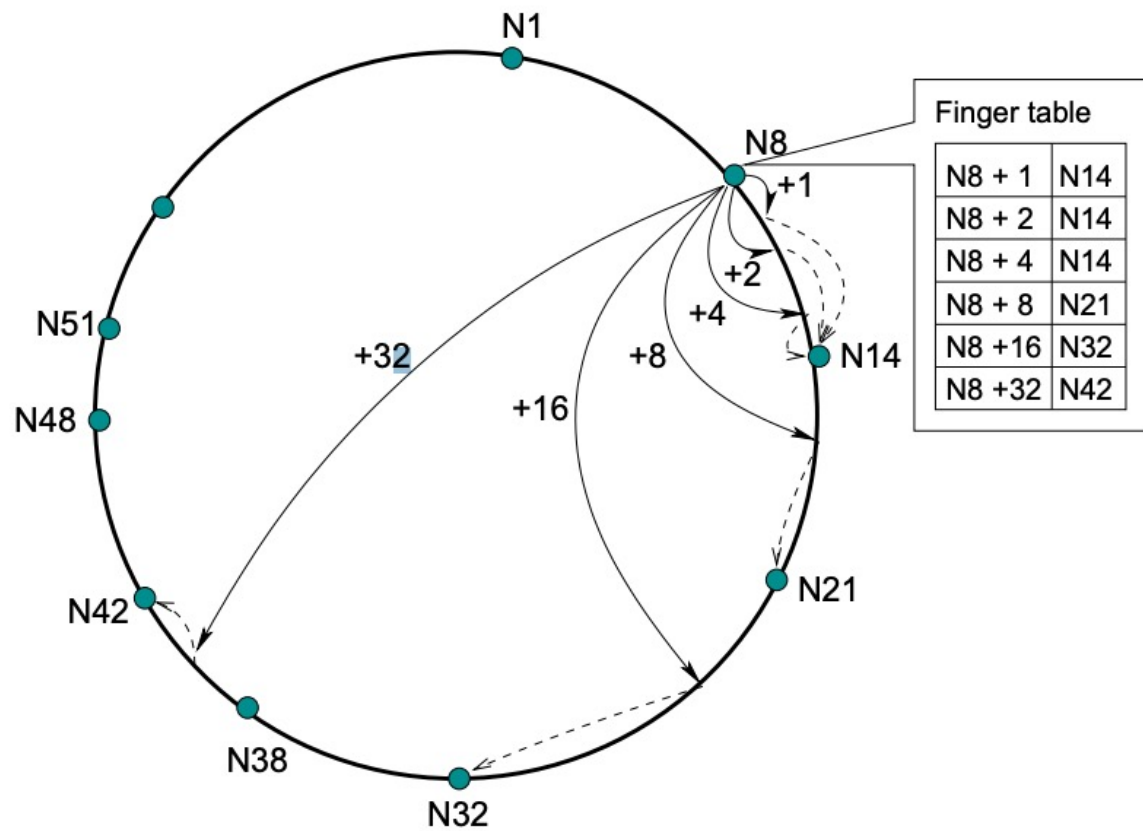


# Chord – finger tables for *find\_predecessor*

Each node keeps **m** states  
Key space  $\rightarrow$  **m** ranges via  
 $(N+2^{k-1}) \bmod 2^m, 1 \leq k \leq m$



# Chord – finger tables



# Implication of finger tables

---

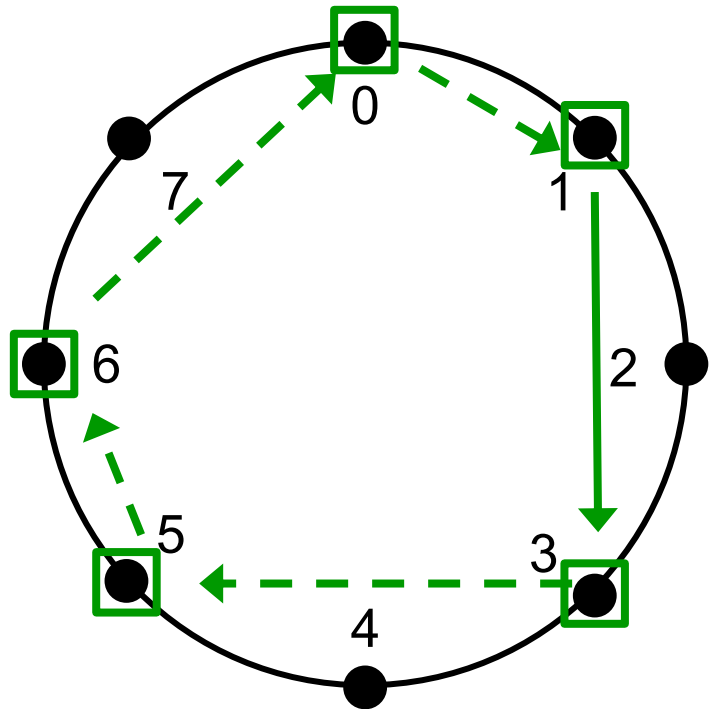
- A **binary lookup tree** rooted at every node
  - Threaded through other nodes' finger tables
- Better than arranging nodes in a single tree
  - Every node acts as a root
    - So there's **no root hotspot**
    - **No single point** of failure
    - But a **lot more state** in total:  $N$  nodes each have  $O(\log N)$

# Chord lookup algorithm properties

---

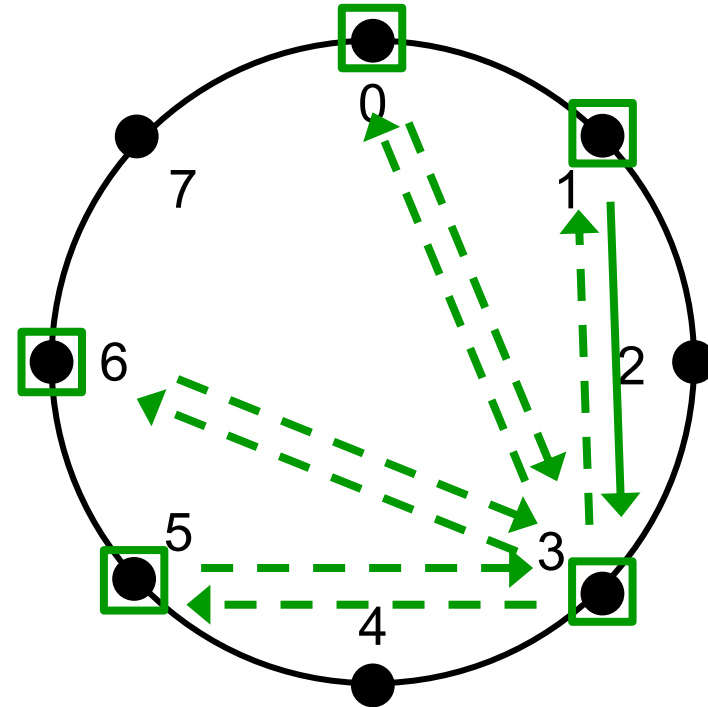
- **Interface:**  $\text{lookup}(\text{key}) \rightarrow \text{IP address}$
- **Efficient:**  $O(\log N)$  messages per lookup
  - $N$  is the total number of nodes (peers)
- **Scalable:**  $O(\log N)$  state per node
- **Robust:** survives massive failures

# Chord – Recursive vs. Iterative Lookup



Recursive Lookup

get(key 1)



Iterative Lookup

# System Dynamics

---

- Handling node joins
- Handling node failures
  - Rebuilding lookup structures
  - Ensure data durability

# Chord – finger tables

Identifiers have  $m = 3$  bits

Key space:  $[0, 2^3-1]$

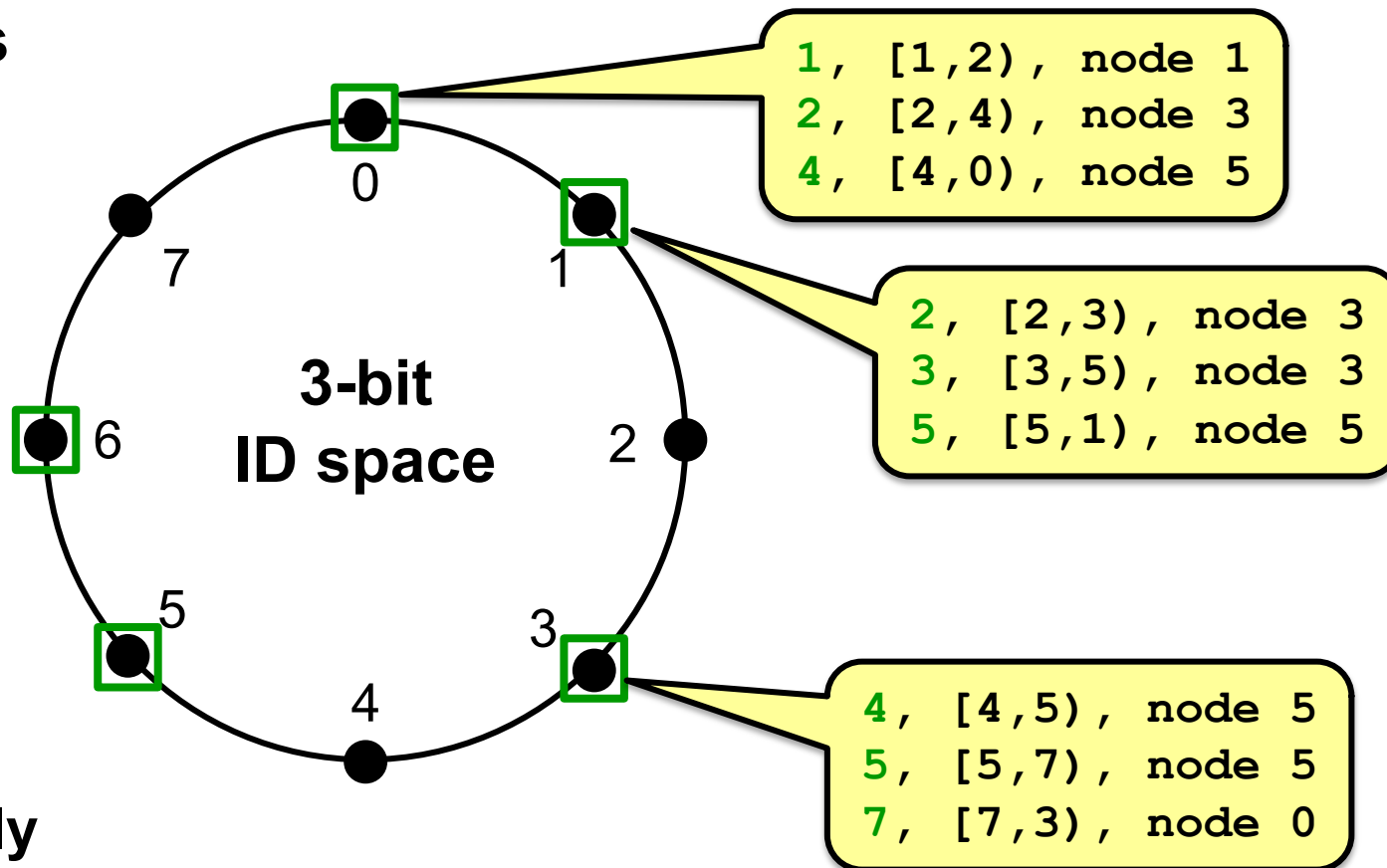
● Identifiers/key space

□ Node

Each node keeps  $m$  states

Key space  $\rightarrow m$  ranges via  
 $(N + 2^{k-1}) \bmod 2^m, 1 \leq k \leq m$

“Finger” is node immediately  
succeeding separator



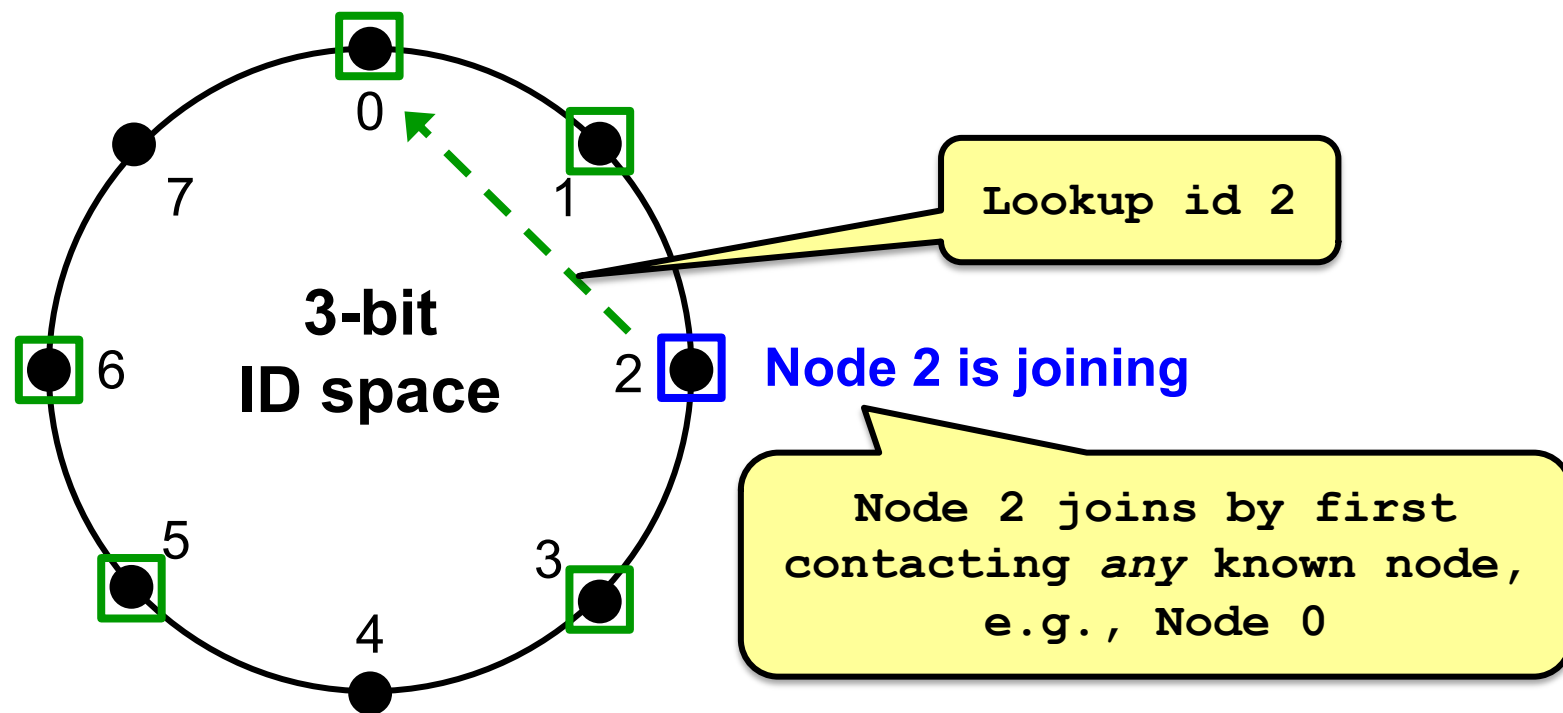
# Chord – node joining

Identifiers have  $m = 3$  bits

Key space:  $[0, 2^3-1]$

● Identifiers/key space

□ Node



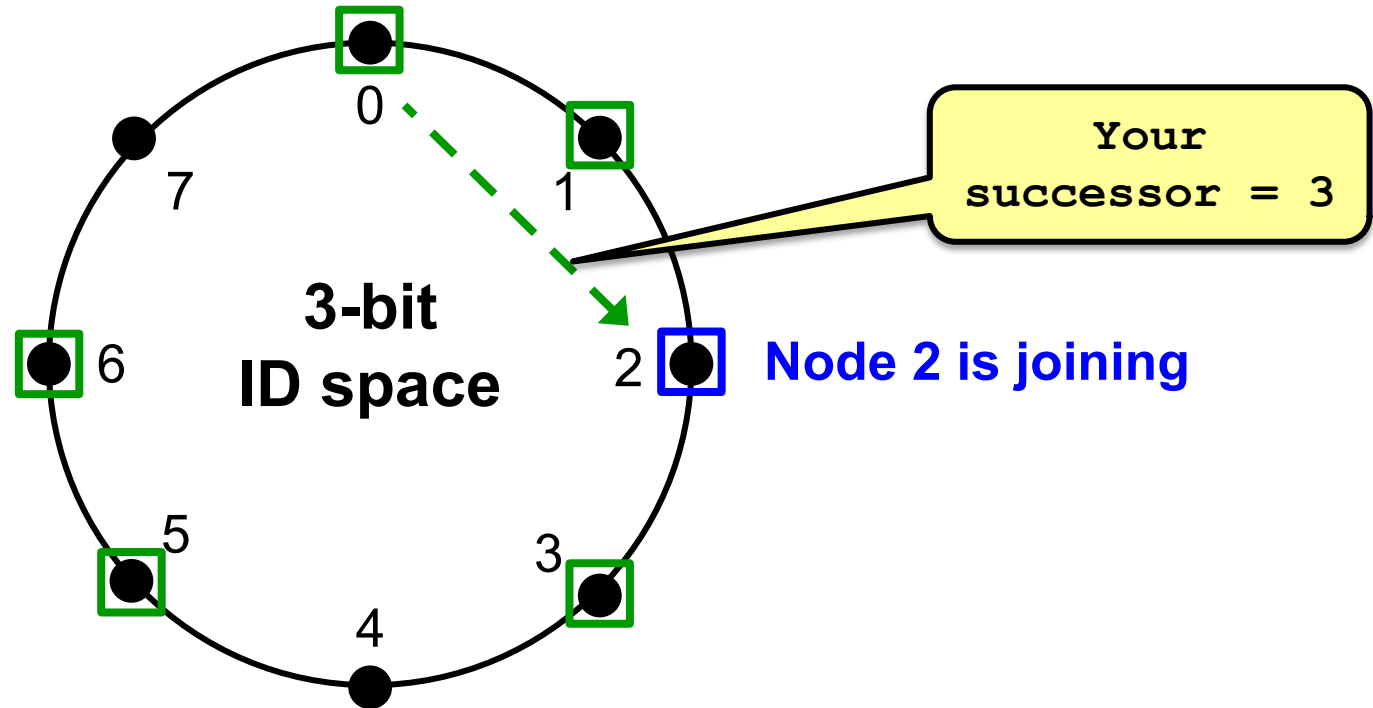
# Chord – node joining

Identifiers have  $m = 3$  bits

Key space:  $[0, 2^3-1]$

● Identifiers/key space

□ Node



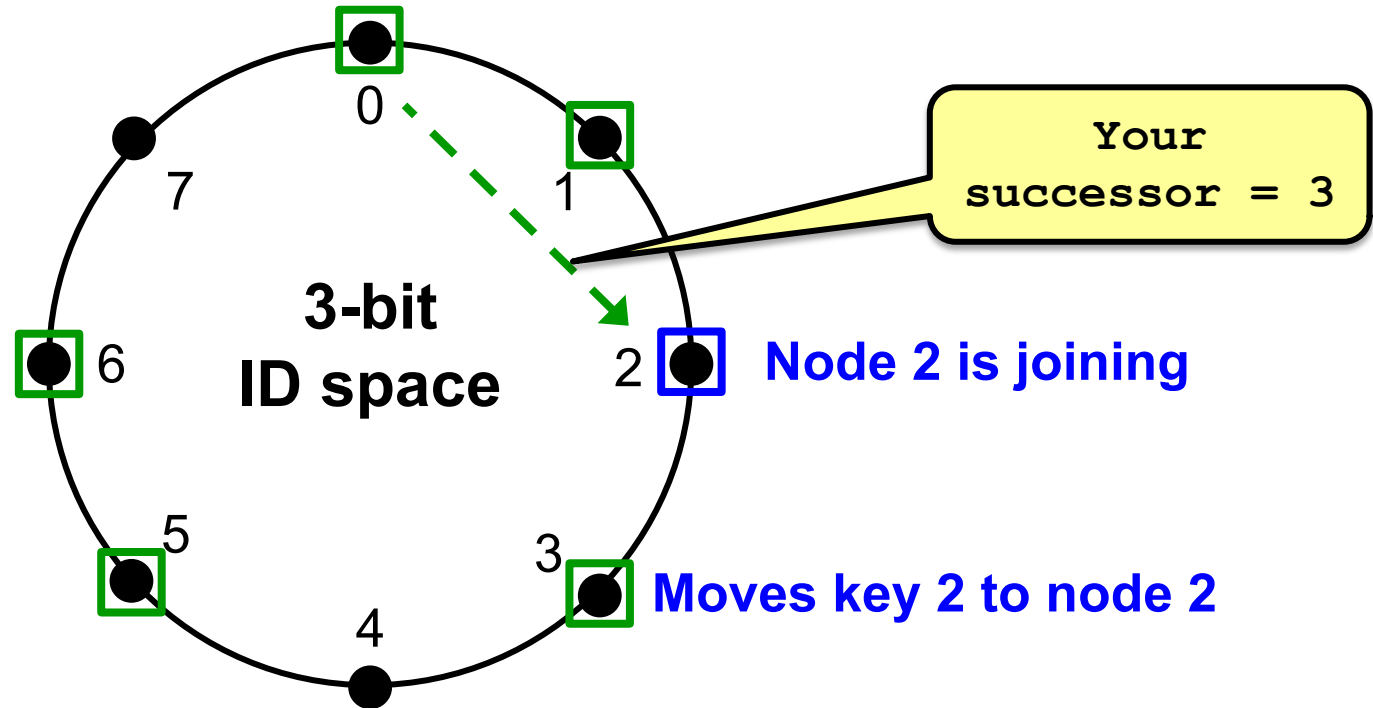
# Chord – node joining

Identifiers have  $m = 3$  bits

Key space:  $[0, 2^3-1]$

● Identifiers/key space

□ Node



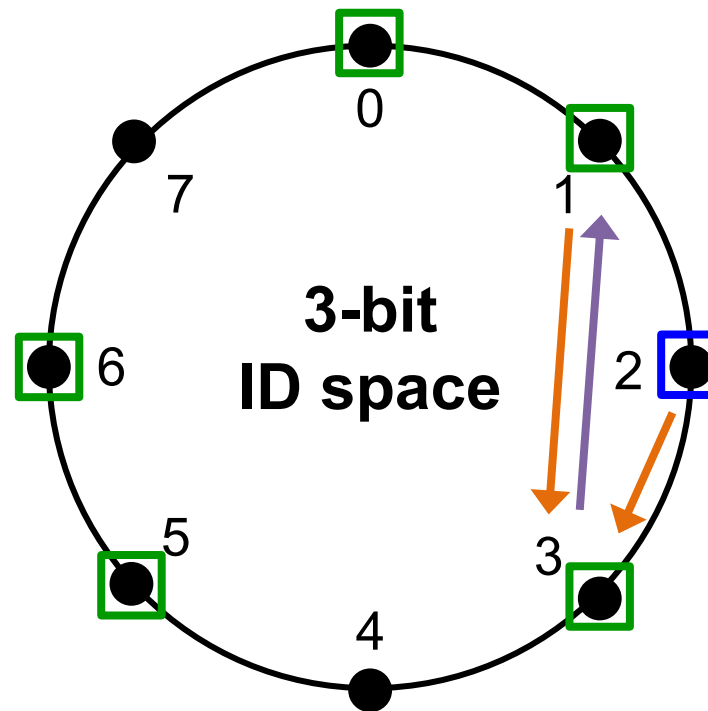
# Chord – node joining

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

Periodic stabilization messages  
from each node to its successor  
maintain node positions

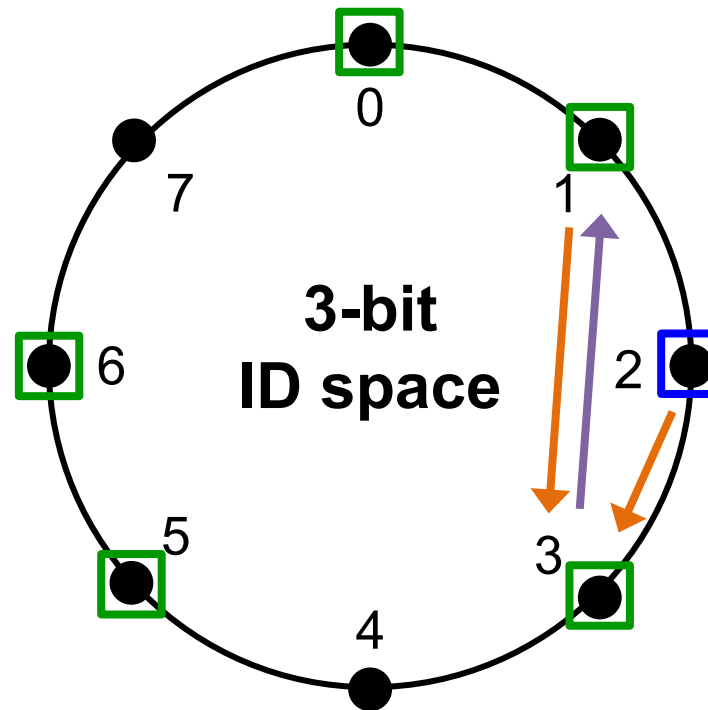
# Chord – node joining

Identifiers have  $m = 3$  bits

Key space:  $[0, 2^3-1]$

● Identifiers/key space

□ Node



→ 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
```

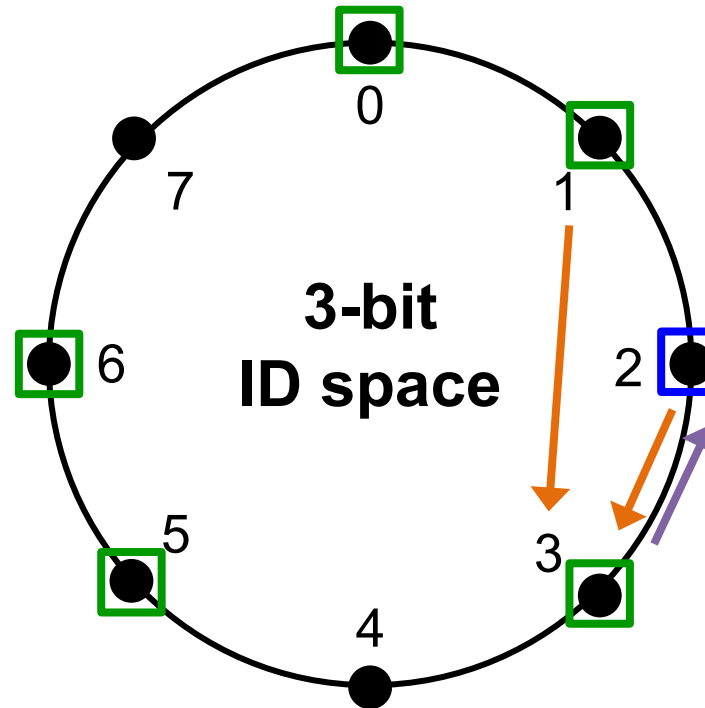
# Chord – node joining

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
```

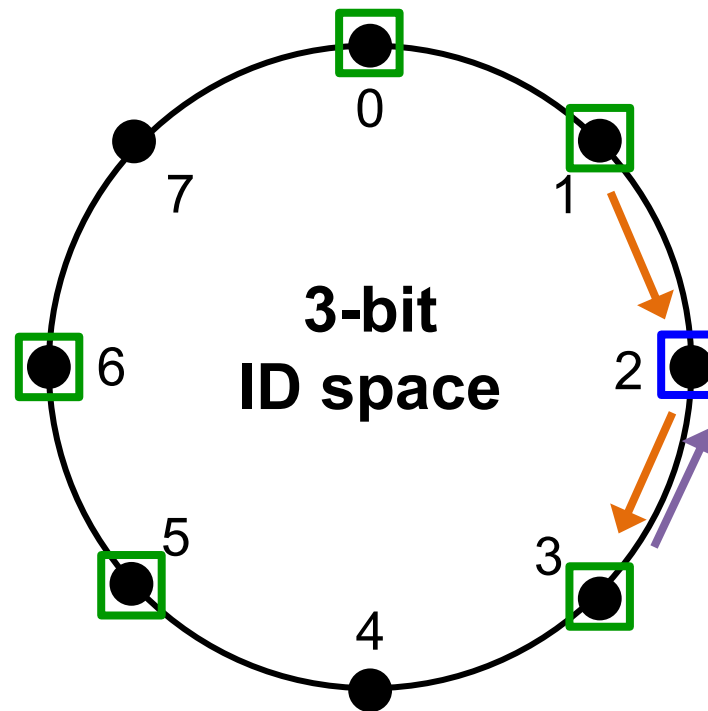
# Chord – node joining

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
```

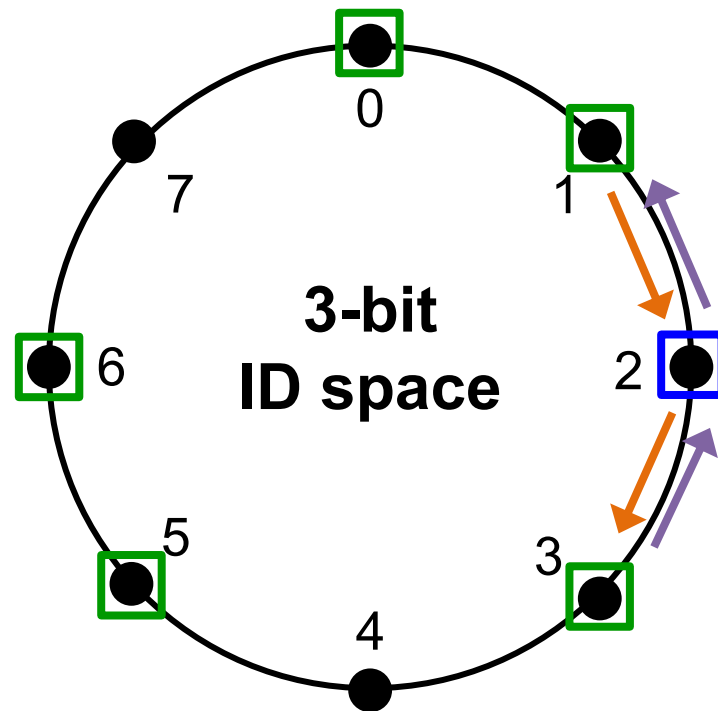
# Chord – node joining

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:

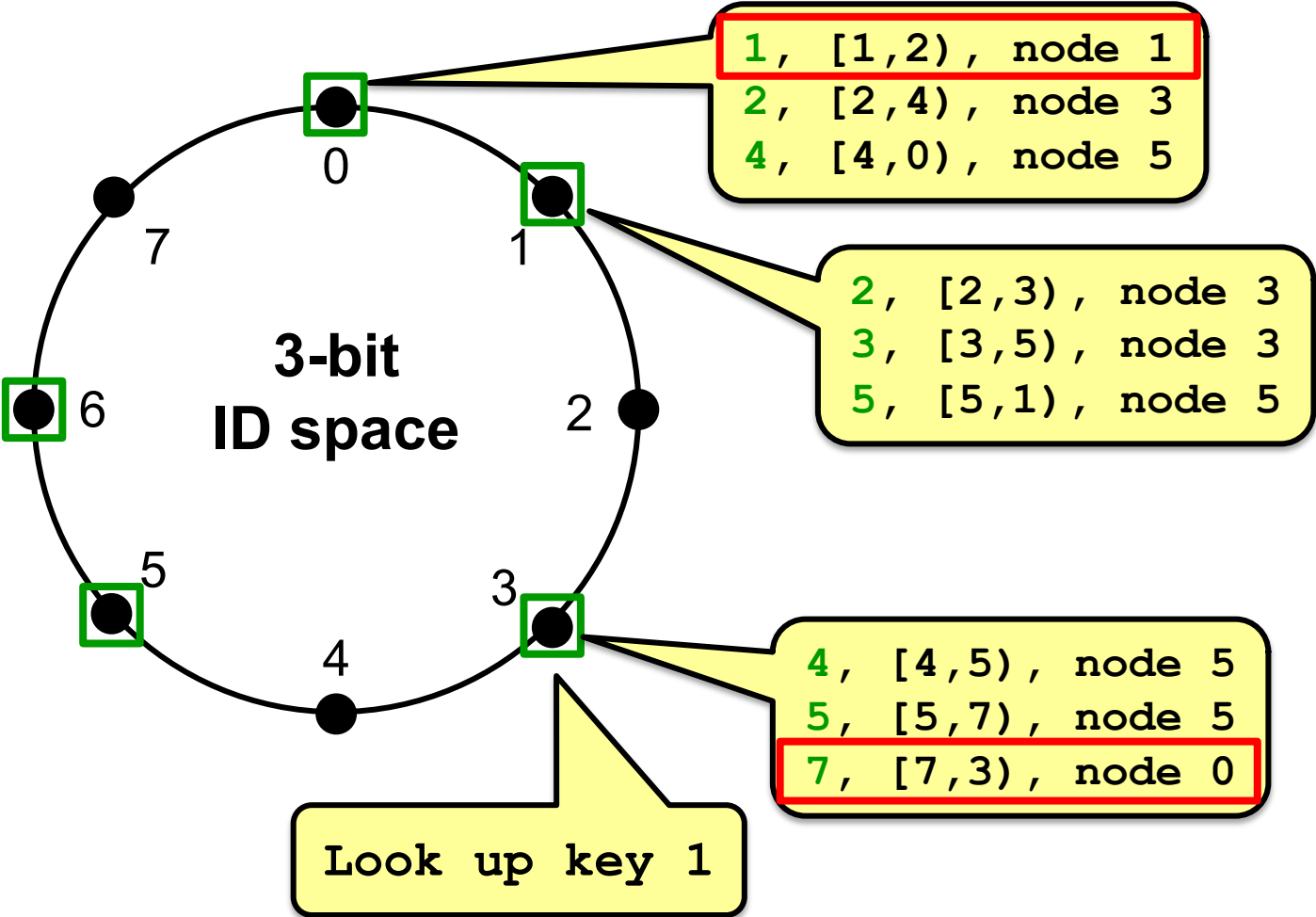
If (N between (M.predecessor, M))

M.predecessor = N

# Chord – failures and successor list

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

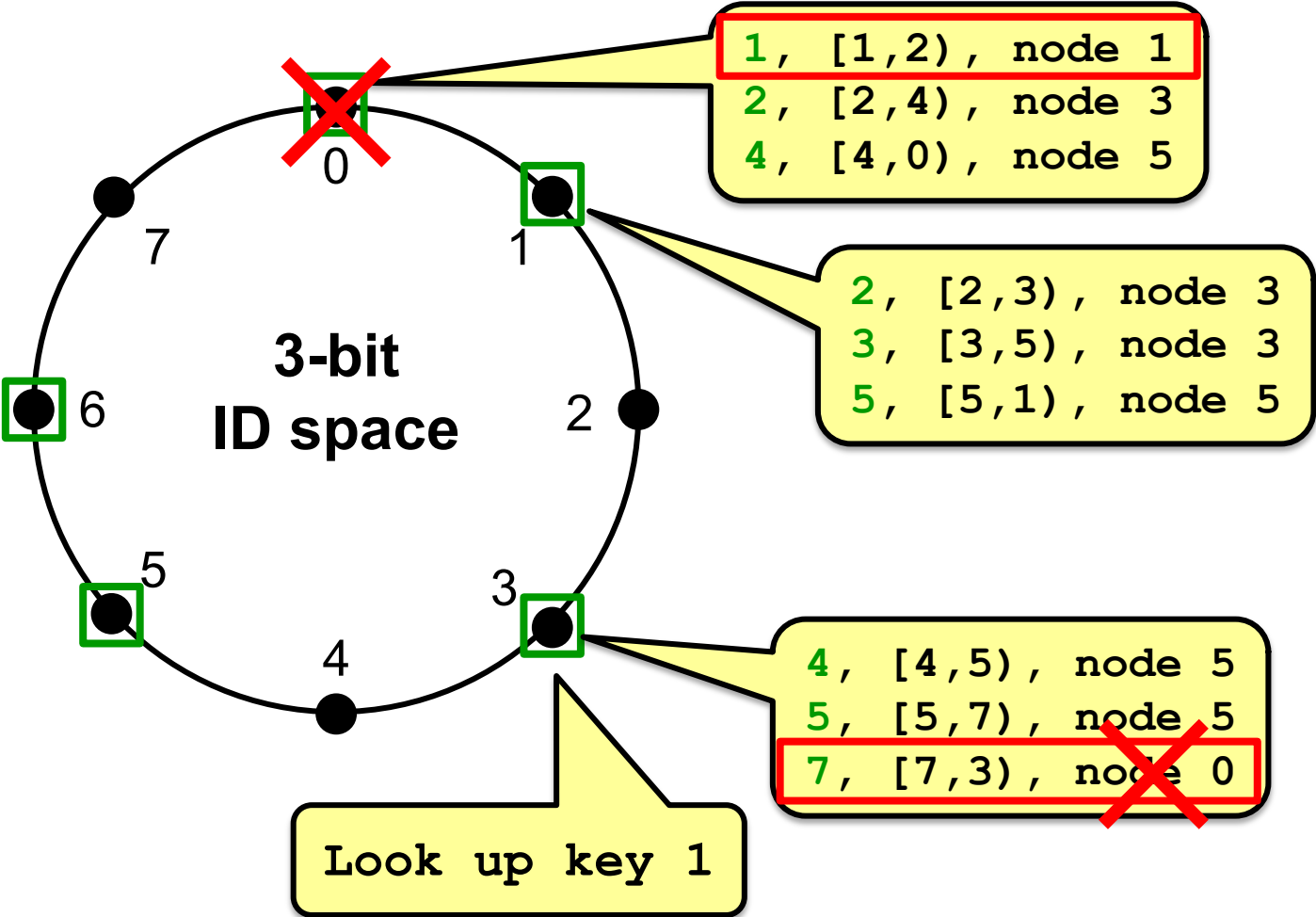
- Identifiers/key space
- Node



# Chord – failures and successor list

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

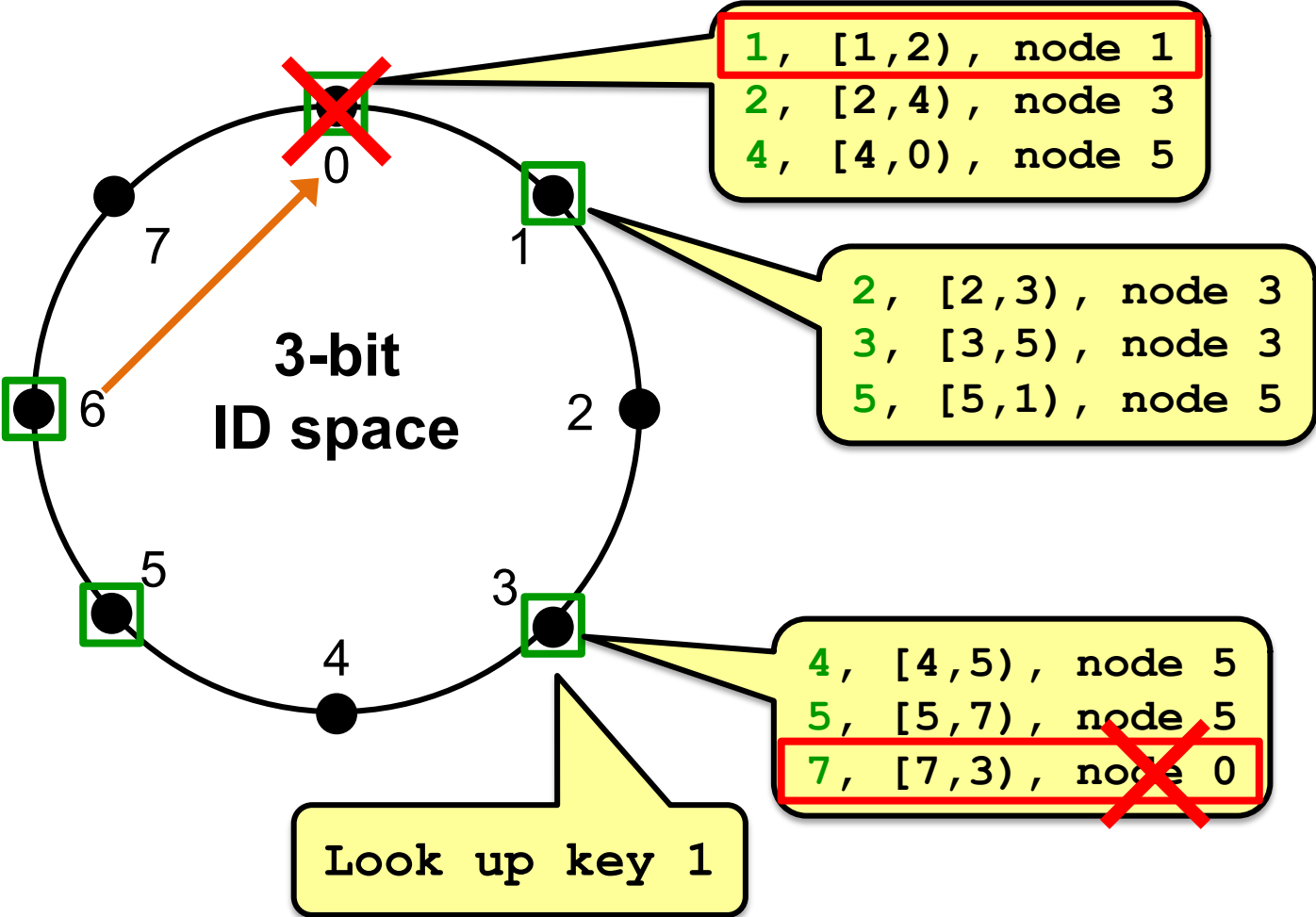
- Identifiers/key space
- Node



# Chord – failures and successor list

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

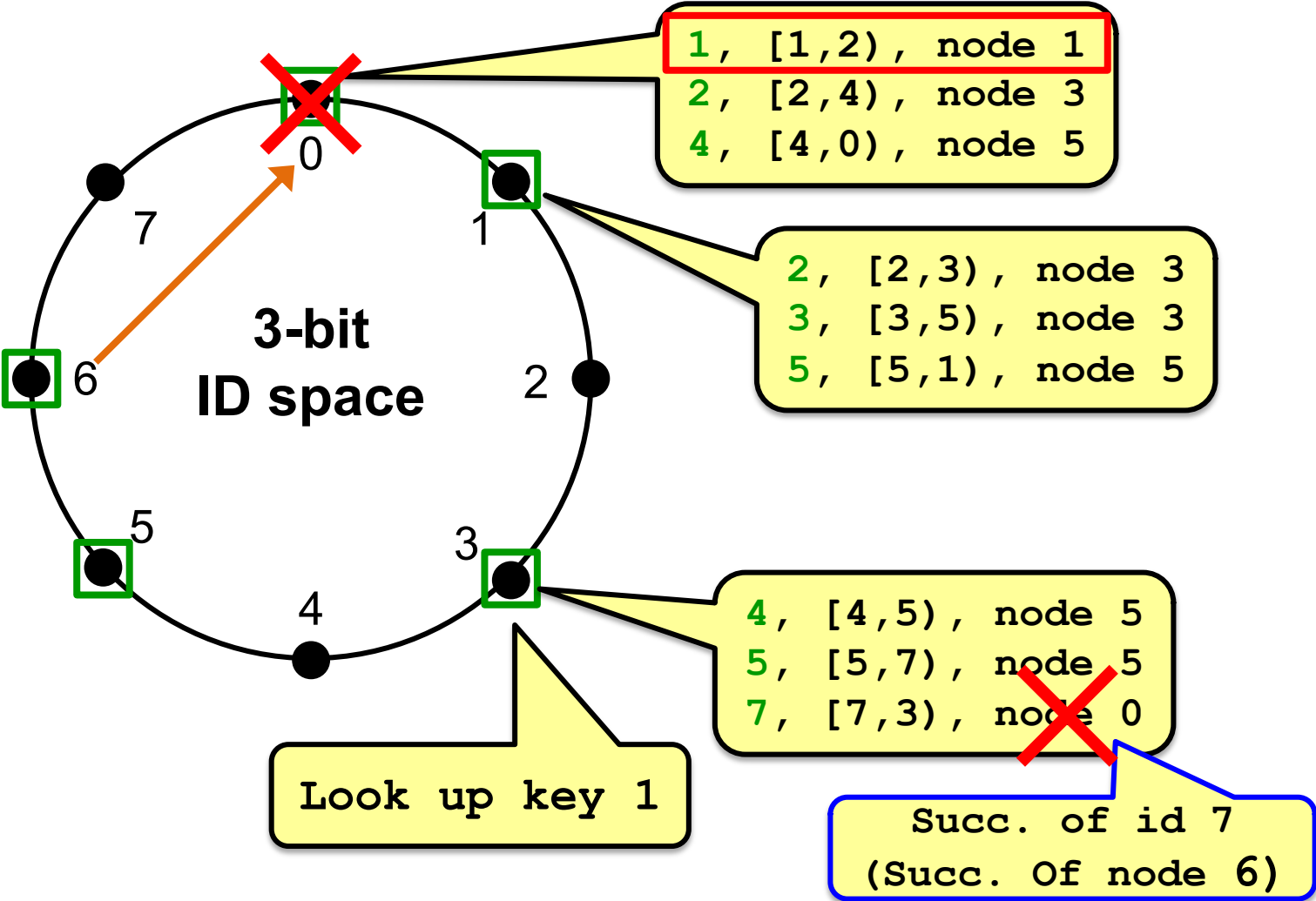
- Identifiers/key space
- Node
- Points to successor



# Chord – failures and successor list

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

- Identifiers/key space
- Node
- Points to successor



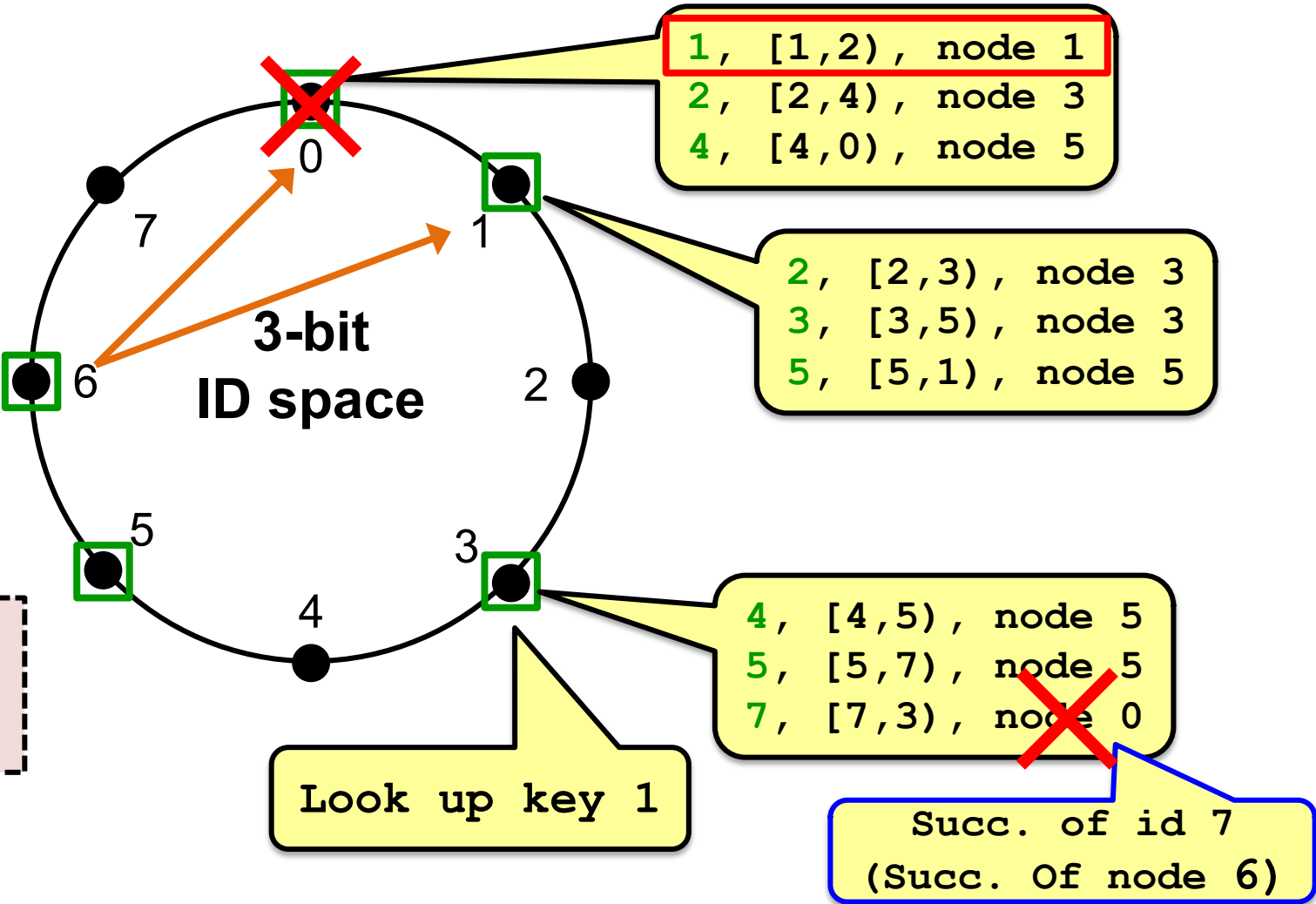
# Chord – failures and successor list

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

- Identifiers/key space
- Node

→ Points to successor

r-nearest successors  
( $r = \log N$ )



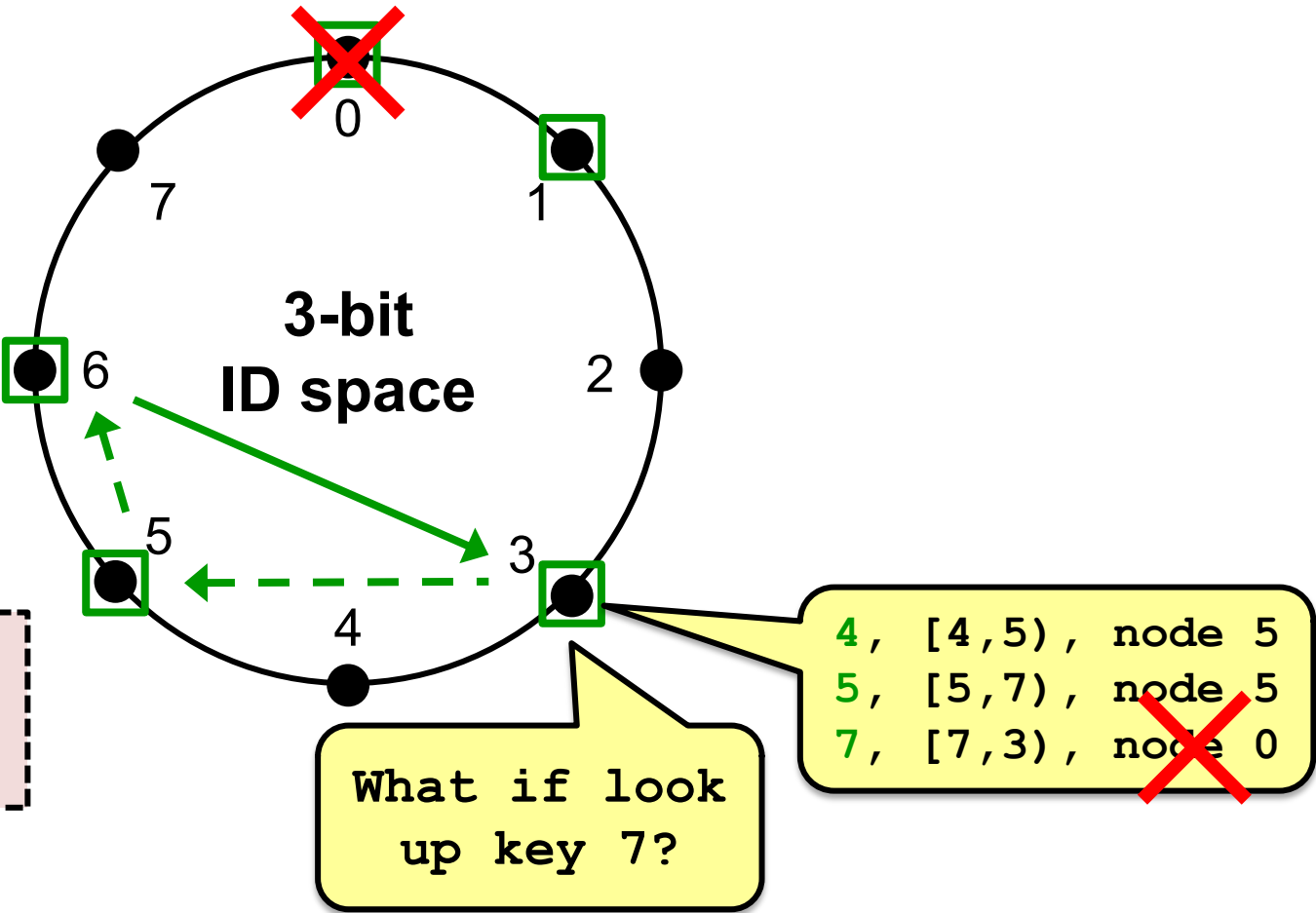
# Chord – failures and successor list

Identifiers have **m** = 3 bits

Key space:  $[0, 2^3-1]$

● Identifiers/key space

□ Node



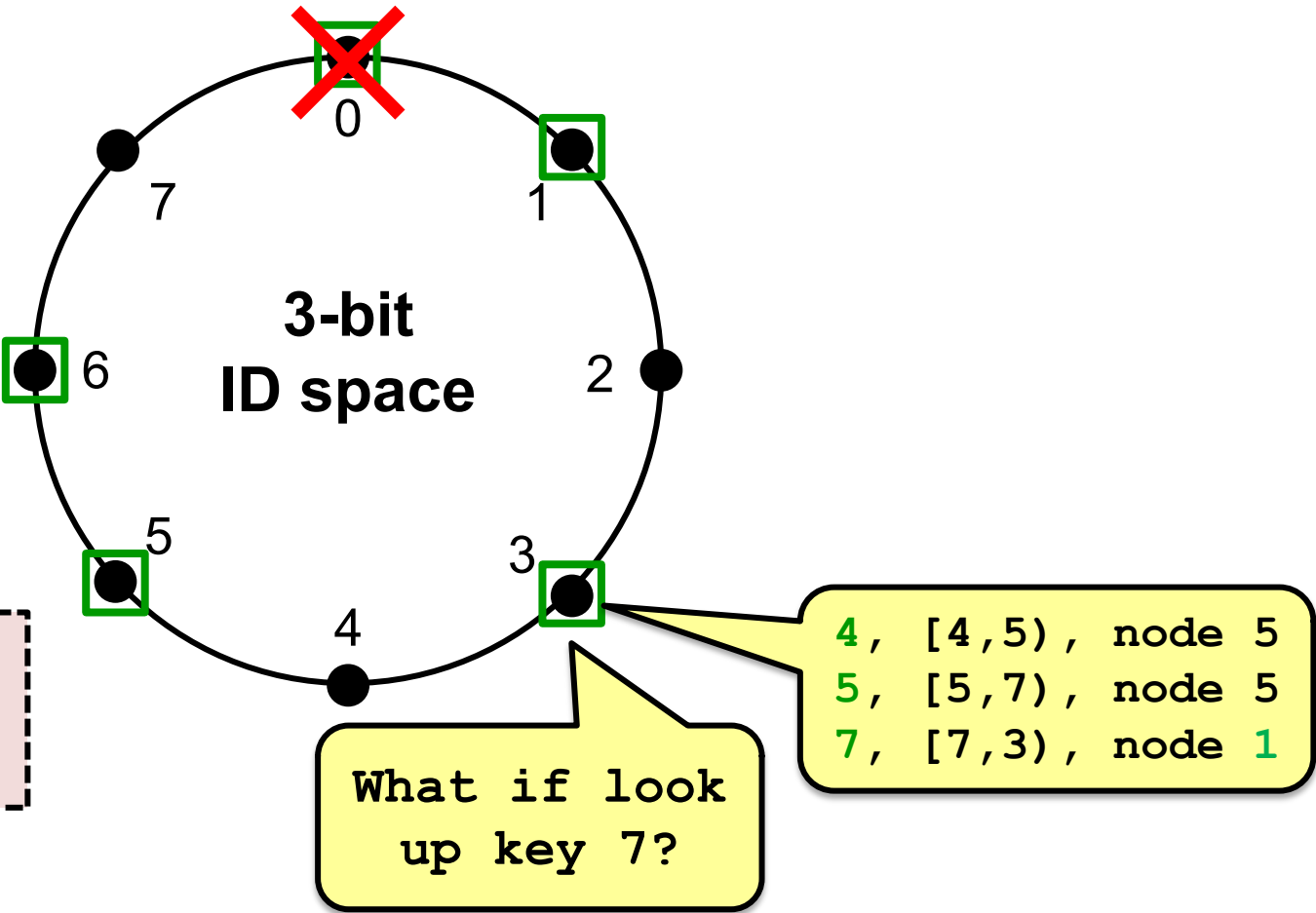
# Chord – failures and successor list

Identifiers have **m** = 3 bits

Key space:  $[0, 2^3-1]$

● Identifiers/key space

□ Node



r-nearest successors  
( $r = \log N$ )

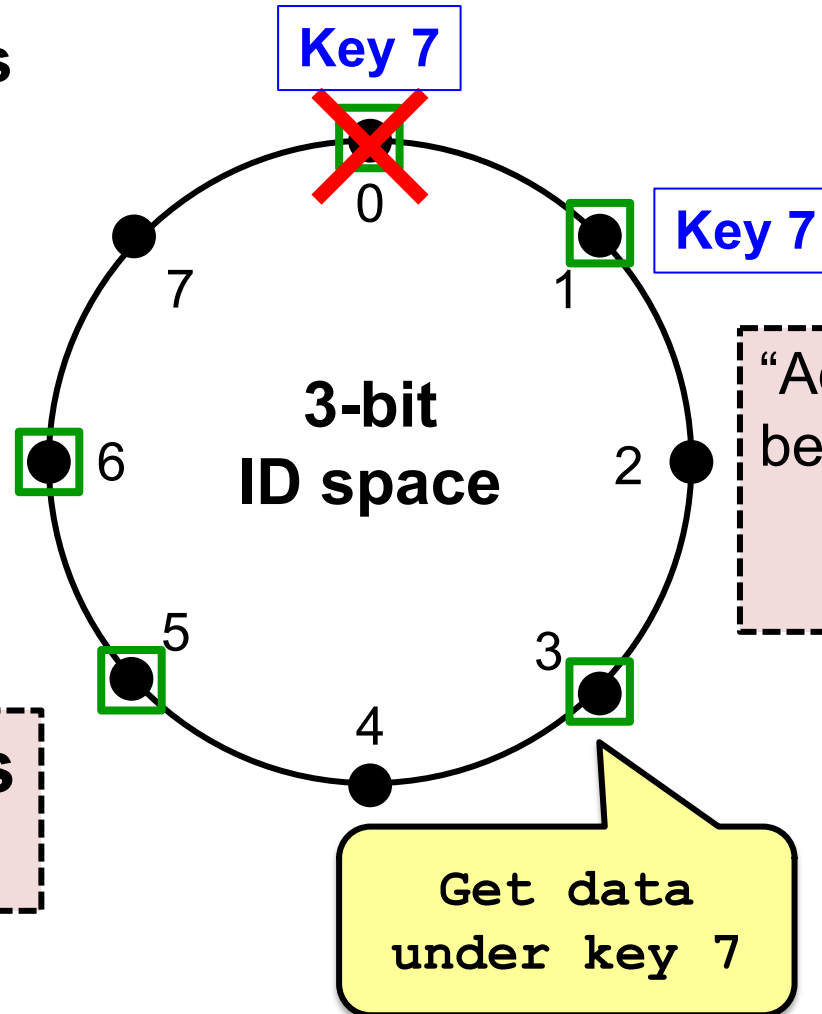
# DHash replicates data blocks at $r$ successors

Identifiers have  $m = 3$  bits

Key space:  $[0, 2^3-1]$

● Identifiers/key space

□ Node



“Adjacent” nodes in the ring may be far away in the network

→ Independent failures

$r$ -nearest successors  
( $r = \log N$ )

# 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