# Distributed Hash Tables & Chord
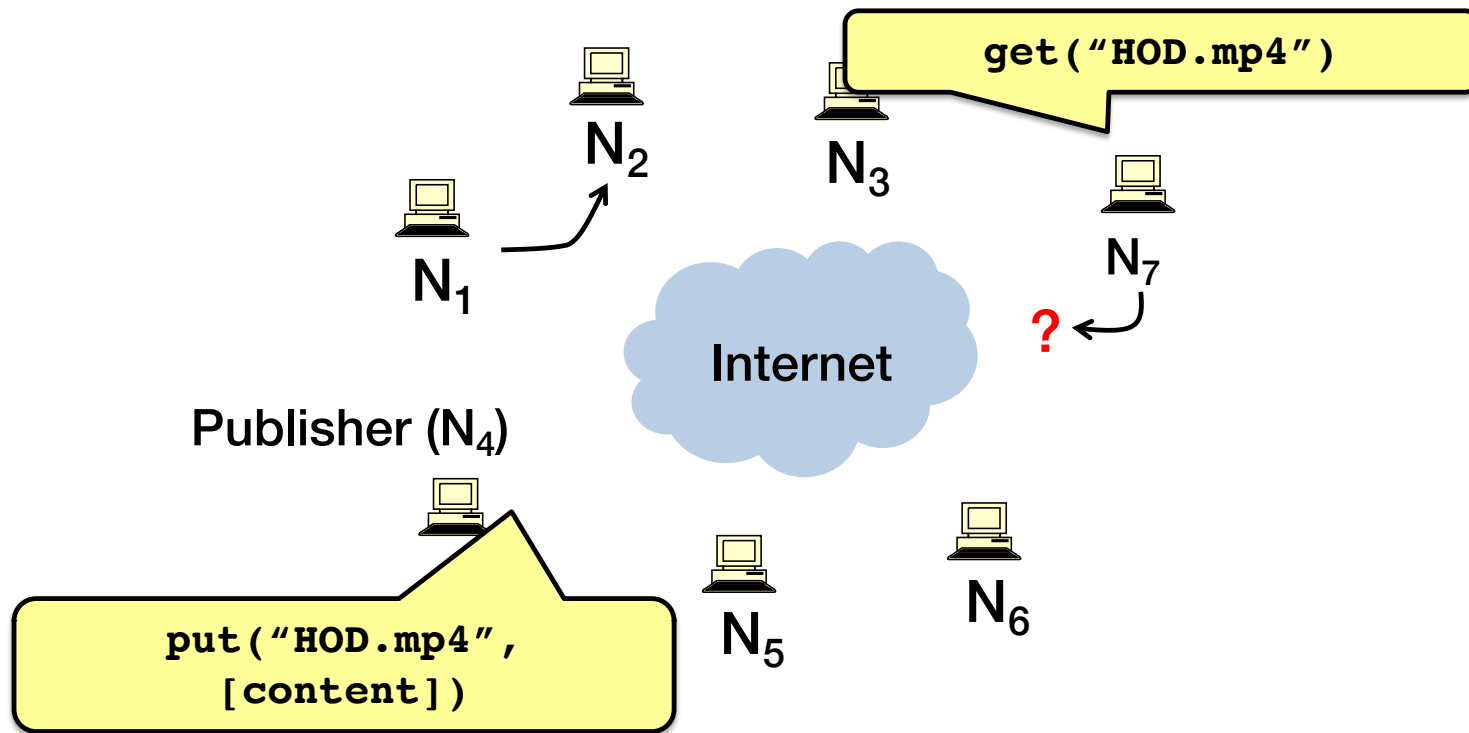
COS 418/518: Distributed Systems

Lecture 8

Wyatt Lloyd

# The lookup problem: locate the data



get("HOD.mp4")

N_2

N_3

N_1

N_7

?

Internet

Publisher (N_4)

put("HOD.mp4",
[content])

N_5

N_6

# What is a DHT (and why)?

- Distributed Hash Table: an abstraction of hash table in a distributed setting
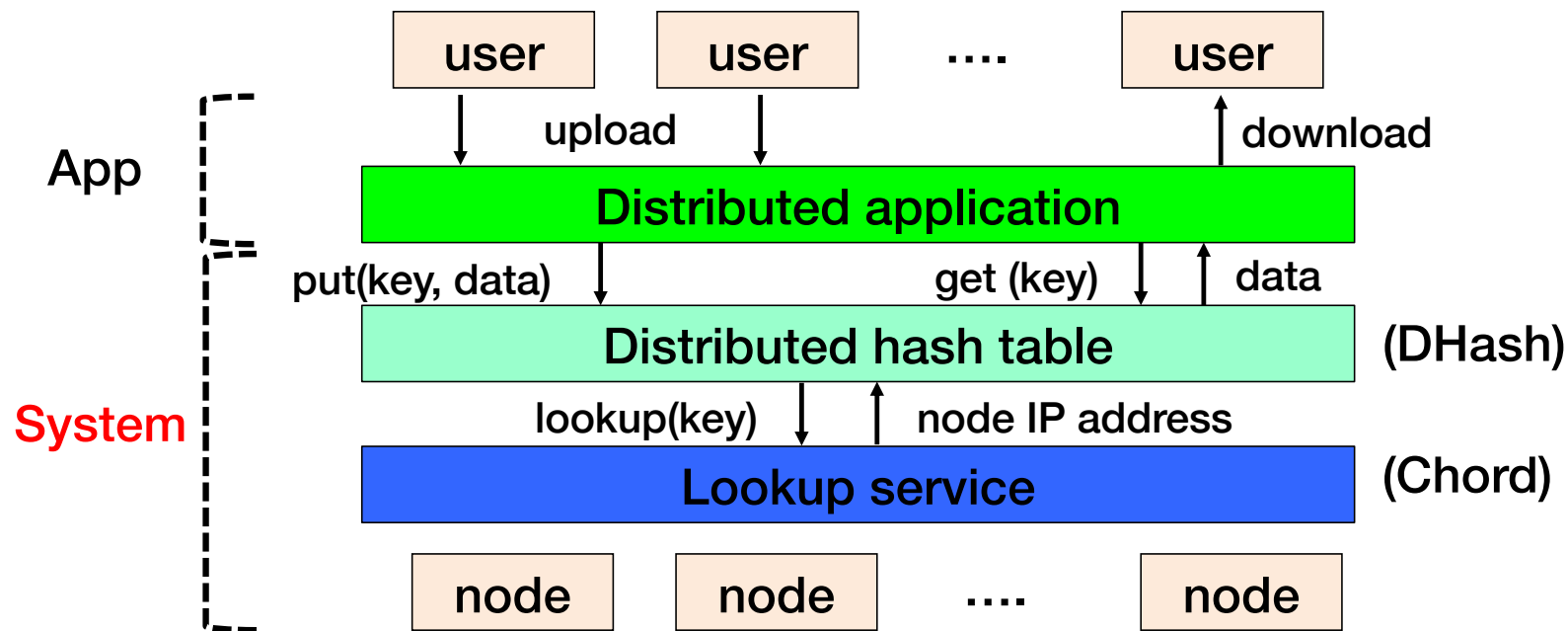
  ```
  key = hash(data)
  lookup(key) → IP addr (Chord lookup service)
  send-RPC(IP address, put, key, data)
  send-RPC(IP address, get, key) → data
  ```

- **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

3

# 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

# 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
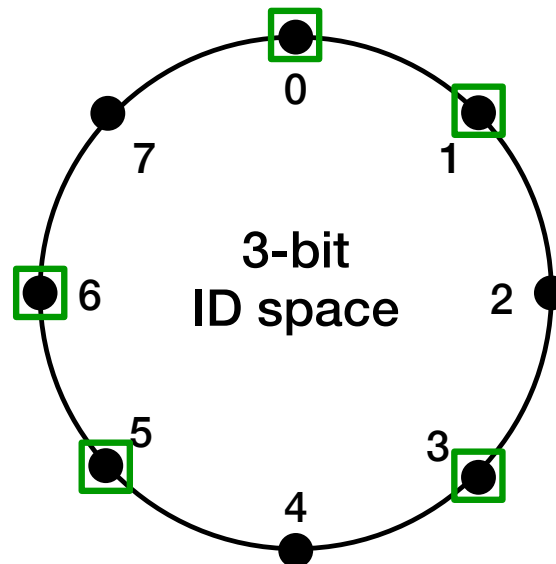  - Hashed key → load balancing; hashed address → independent failure

# Consistent hashing [Karger '97] – data partition

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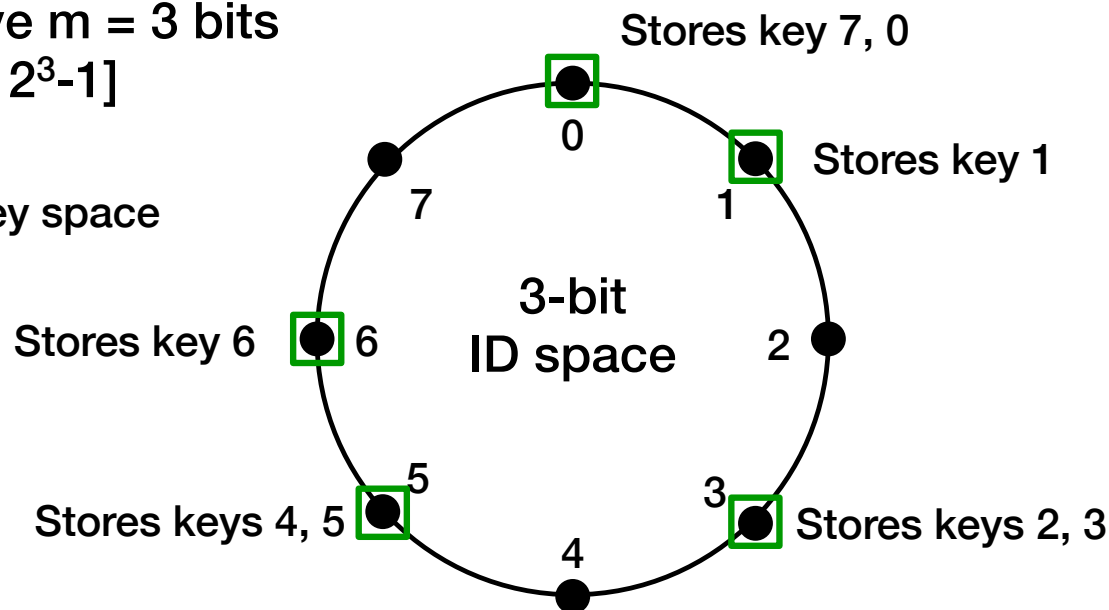
● Identifiers/key space

□ 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

3-bit
ID space

Stores key 7, 0

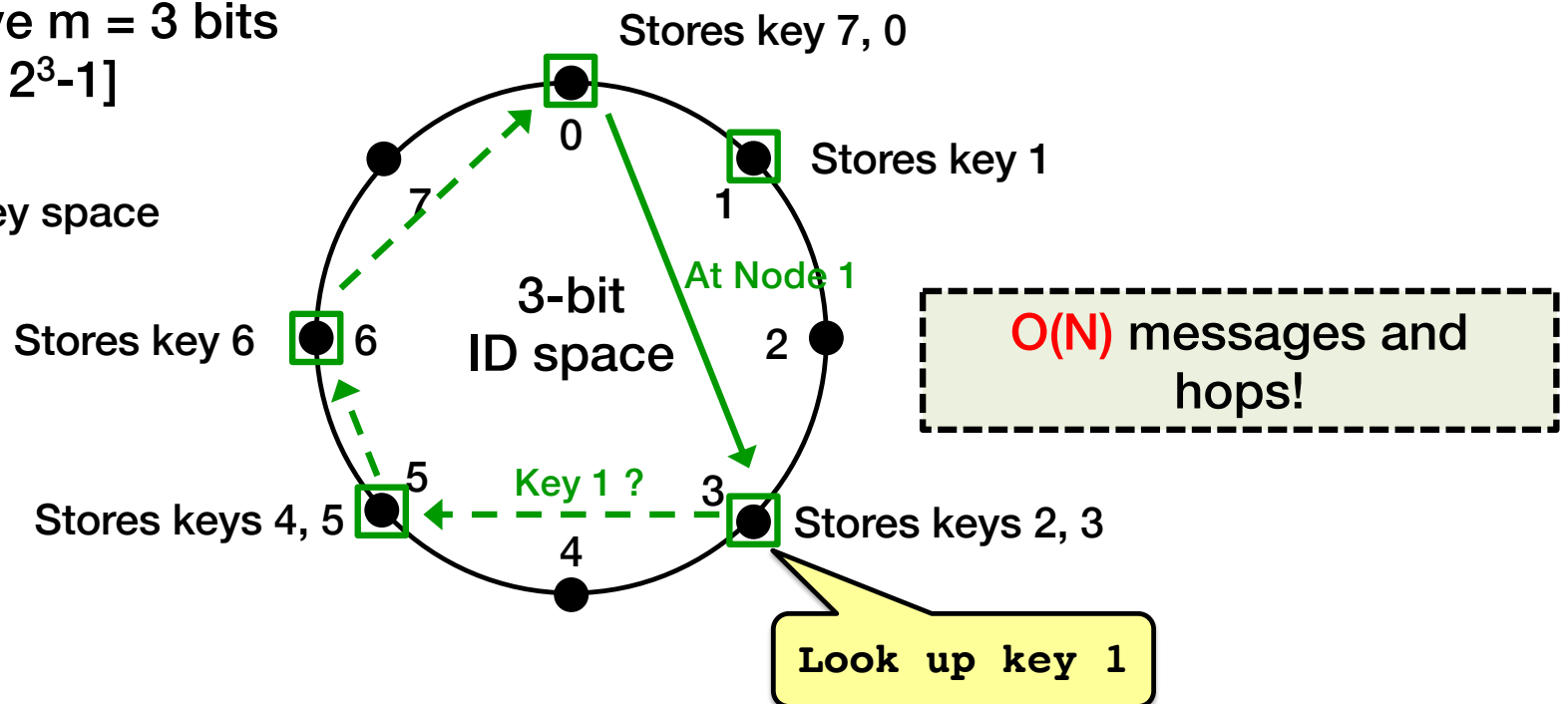Stores key 1

Stores key 6

Stores keys 4, 5

Stores keys 2, 3

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

# Consistent hashing [Karger '97] – basic lookup

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

●   Identifiers/key space

□   Node

--▸   Successor pointer

Stores key 7, 0

Stores key 1

At Node 1

3-bit ID space

Stores key 6

Stores keys 4, 5

Key 1 ?

Stores keys 2, 3

O(N) messages and hops!

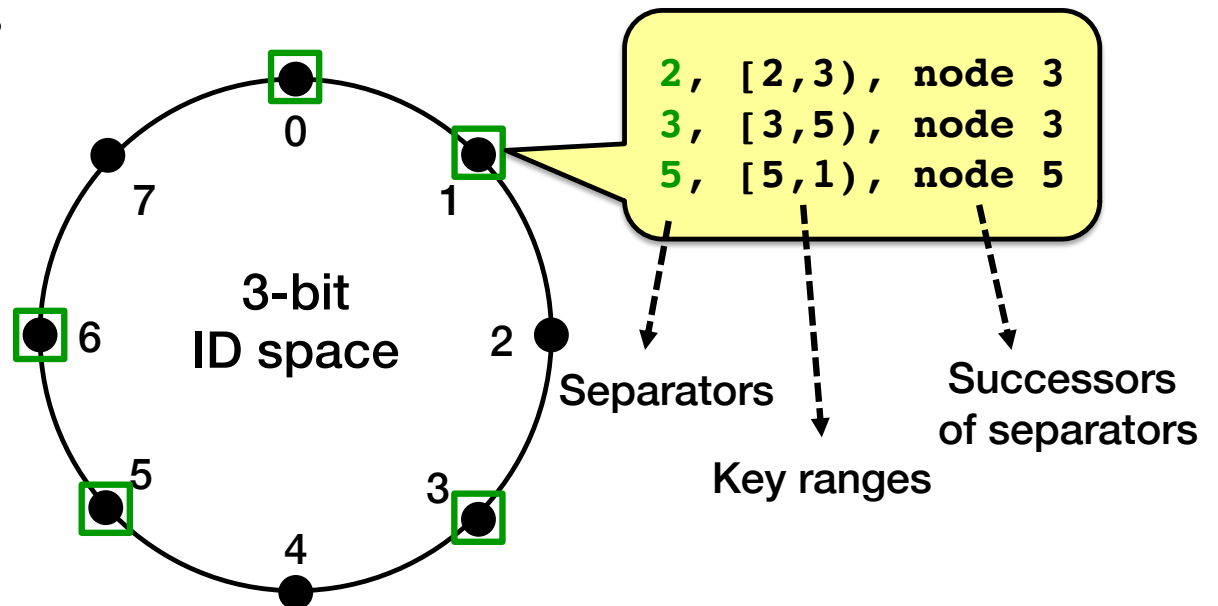Look up key 1

# 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
```
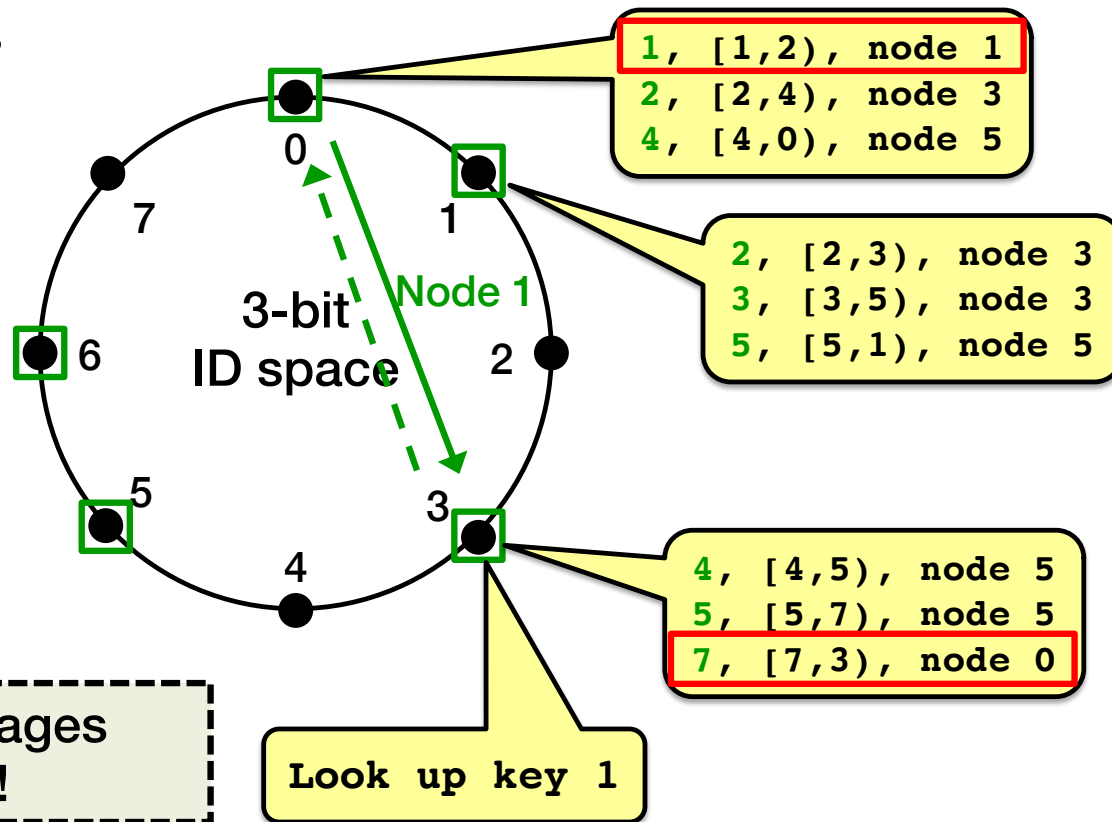
Separators

Key ranges

Successors
of separators

11

# 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
$(N+2^{k-1}) \mod 2^m, 1<=k<=m$

O(logN) messages
and hops!

3-bit
ID space

Node 1

Look up key 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
```

# 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
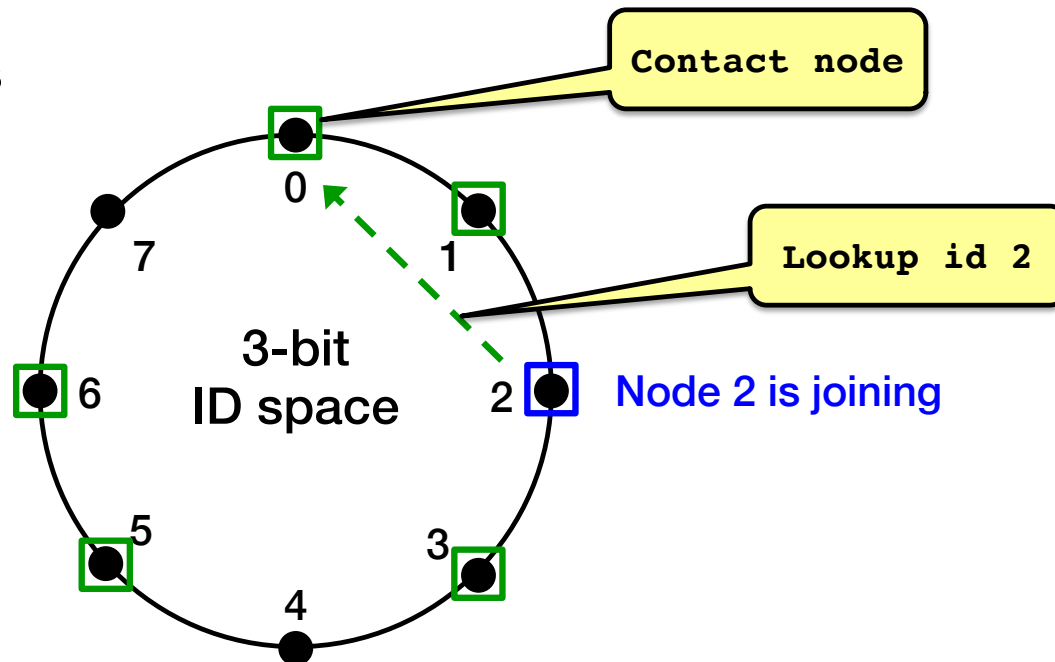    - No single point of failure
    - But a lot more state in total

# 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: $[0, 2^3-1]$

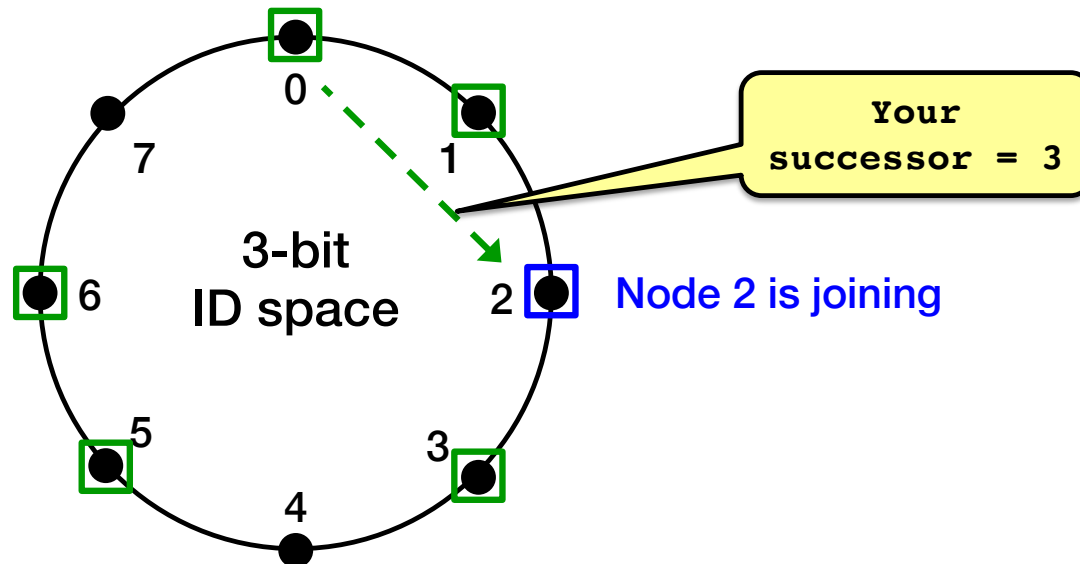●    Identifiers/key space

☐    Node

**Contact node**

**Lookup id 2**

3-bit
ID space

0
1
2
3
4
5
6
7

Node 2 is joining

# Chord – node joining

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

● Identifiers/key space

□ Node

3-bit ID space

Node 2 is joining

Your successor = 3

# Chord – node joining

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

●    Identifiers/key space

☐    Node



3-bit
ID space

Your
successor = 3

Node 2 is joining

Moves key 2 to node 2

18

# Chord – node joining

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

●    Identifiers/key space

☐    Node

3-bit
ID space

0
1
7
6
5
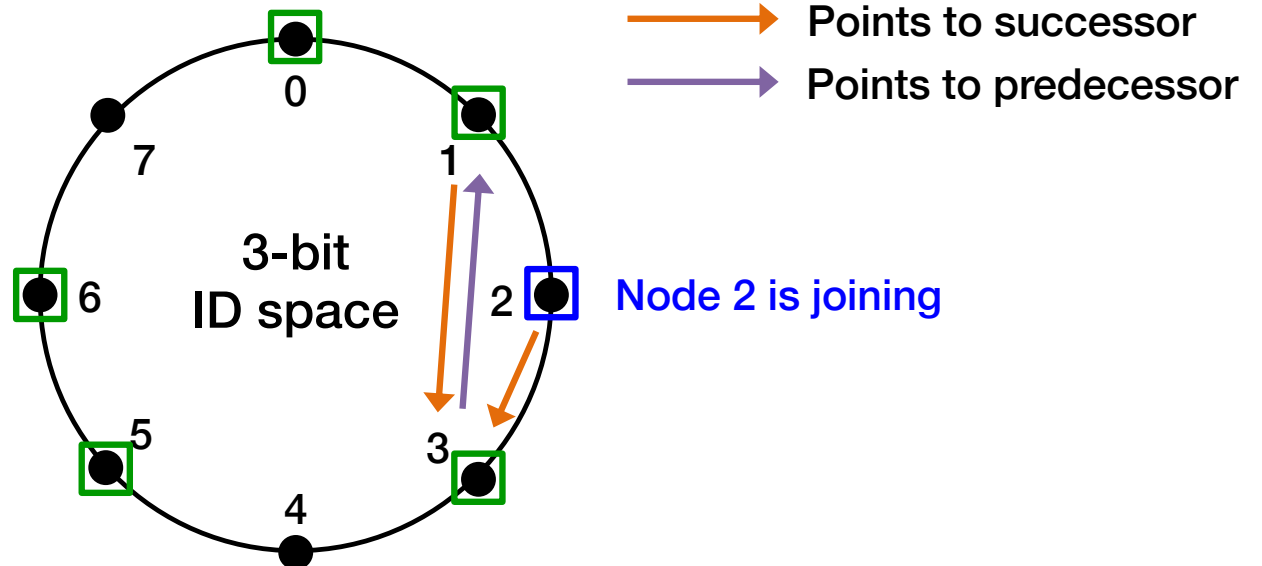4
3
2

→ Points to successor
→ Points to predecessor

Node 2 is joining

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

19

# Chord – node joining

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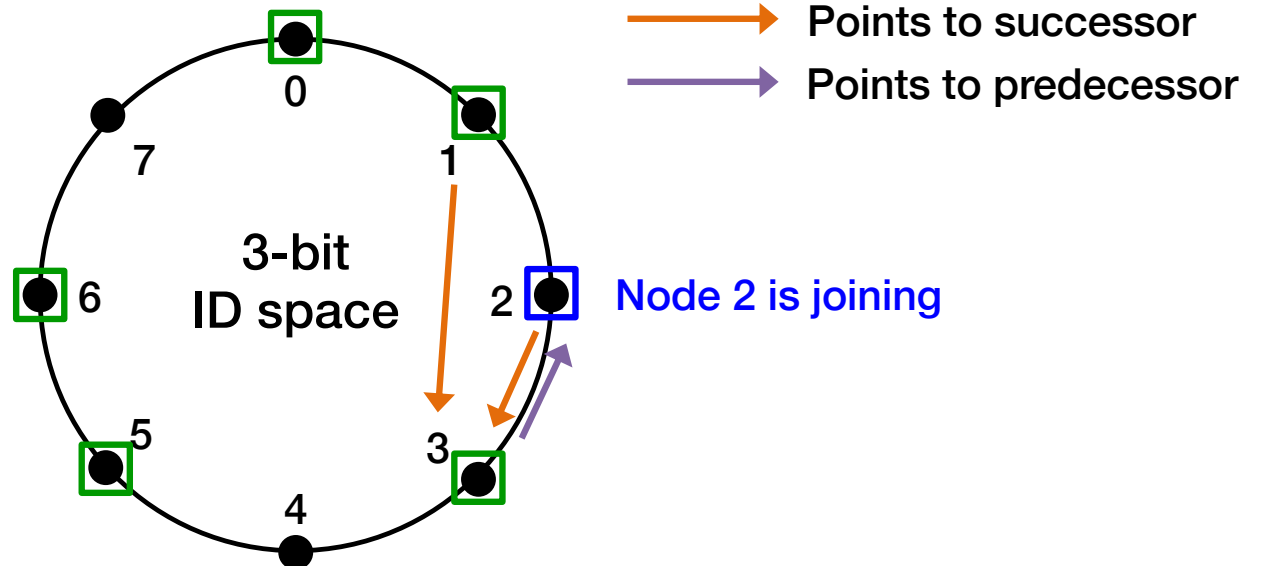Key space: $[0, 2^3-1]$

● Identifiers/key space

□ Node

3-bit
ID space

0
1
2  Node 2 is joining
3
4
5
6
7

→ Points to successor
→ Points to predecessor

# Chord – node joining

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

●    Identifiers/key space

☐    Node

→ Points to successor
→ Points to predecessor

3-bit
ID space

0
1
2    Node 2 is joining
3
4
5
6
7

21

# Chord – node joining

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

● Identifiers/key space

□ Node

3-bit
ID space

0
1
2  Node 2 is joining
3
4
5
6
7

→ Points to successor
→ Points to predecessor

22

# Chord – node joining

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

●    Identifiers/key space

□    Node

3-bit
ID space

0
1
2    Node 2 is joining
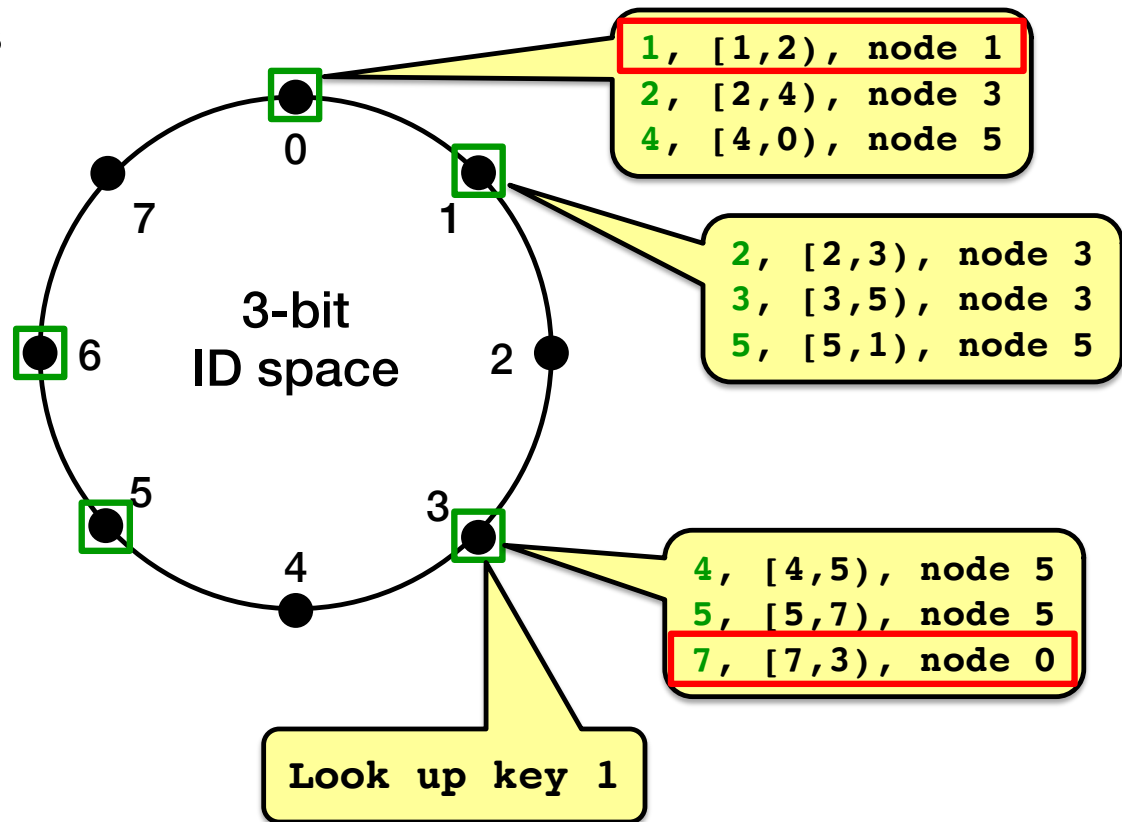3
4
5
6
7

→ Points to successor
→ Points to predecessor

23

# Chord – failures and successor list

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

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

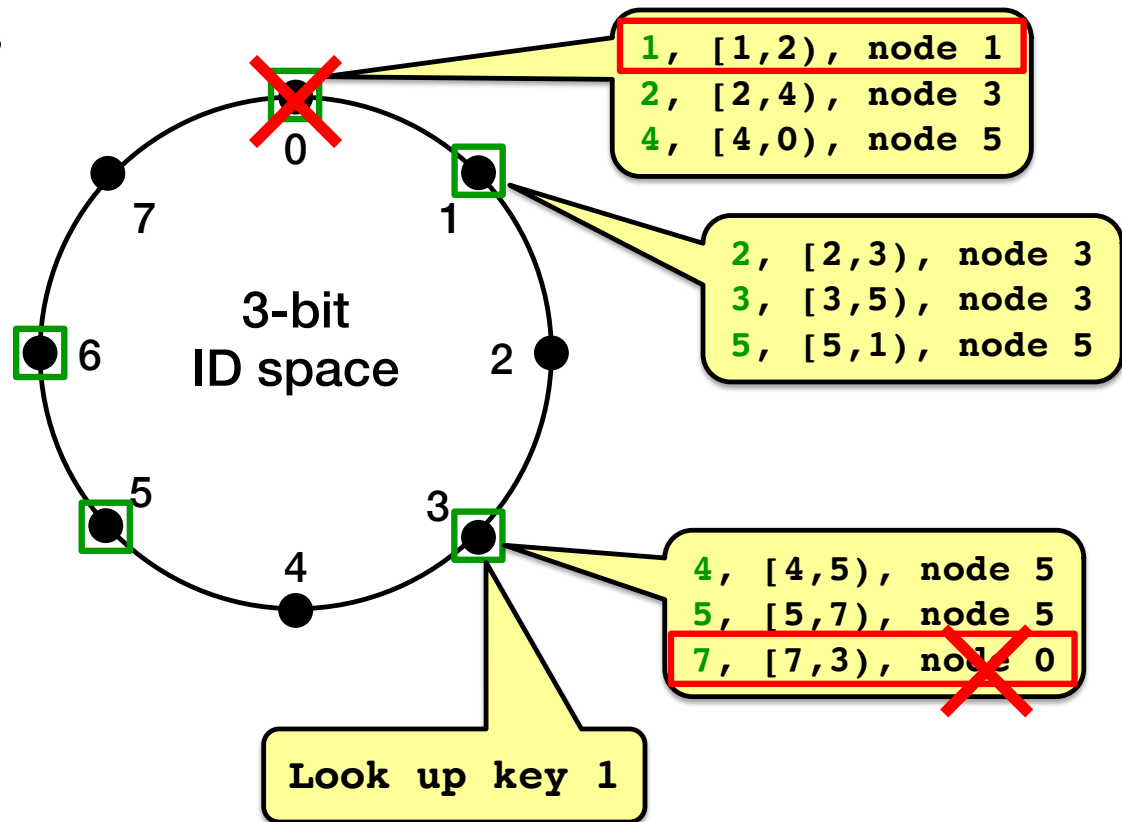Look up key 1

24

# Chord – failures and successor list

# 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, 2³-1]

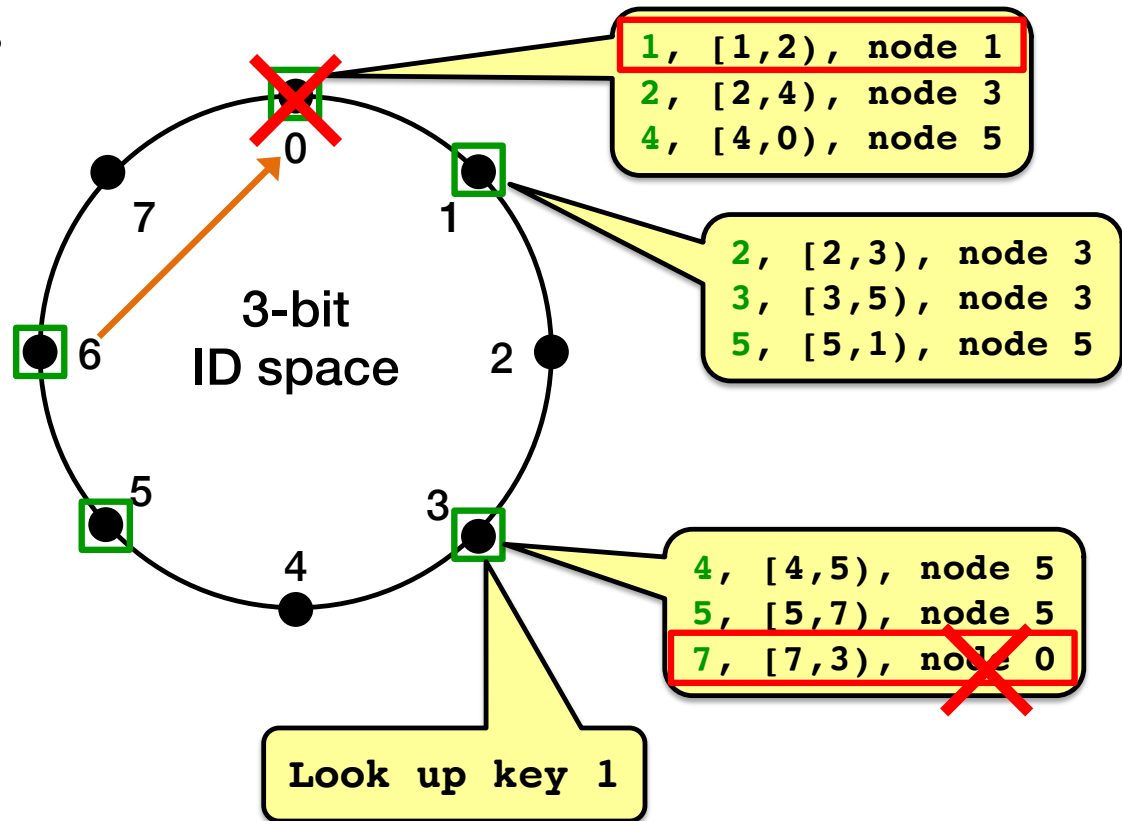● Identifiers/key space

☐ 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)

# Chord – failures and successor list

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

● Identifiers/key space
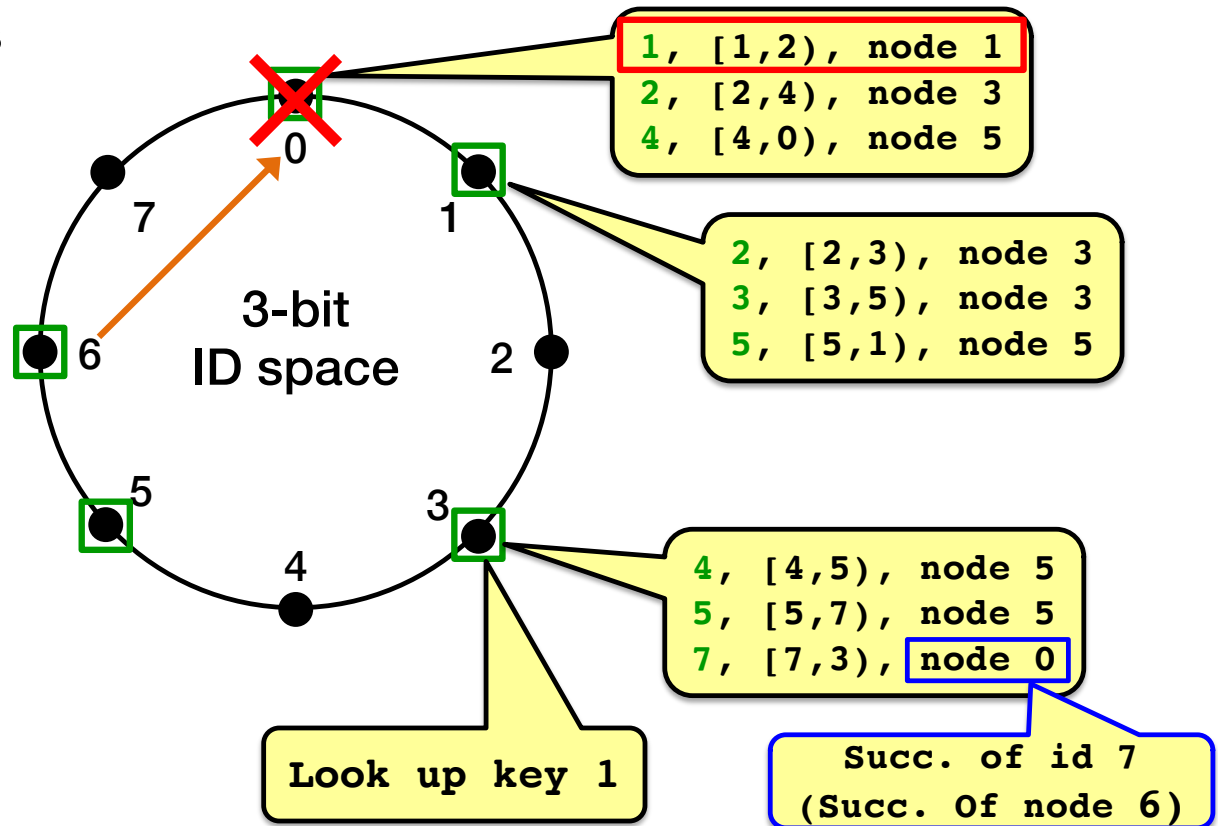
□ Node

→ Points to successor
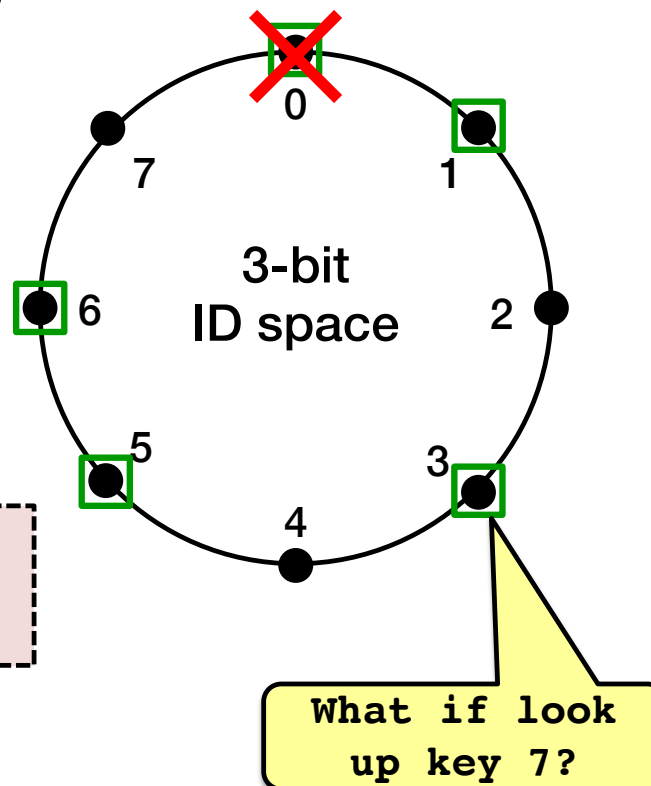
**r-nearest successors
(r = logN)**

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

28

# Chord – failures and successor list

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
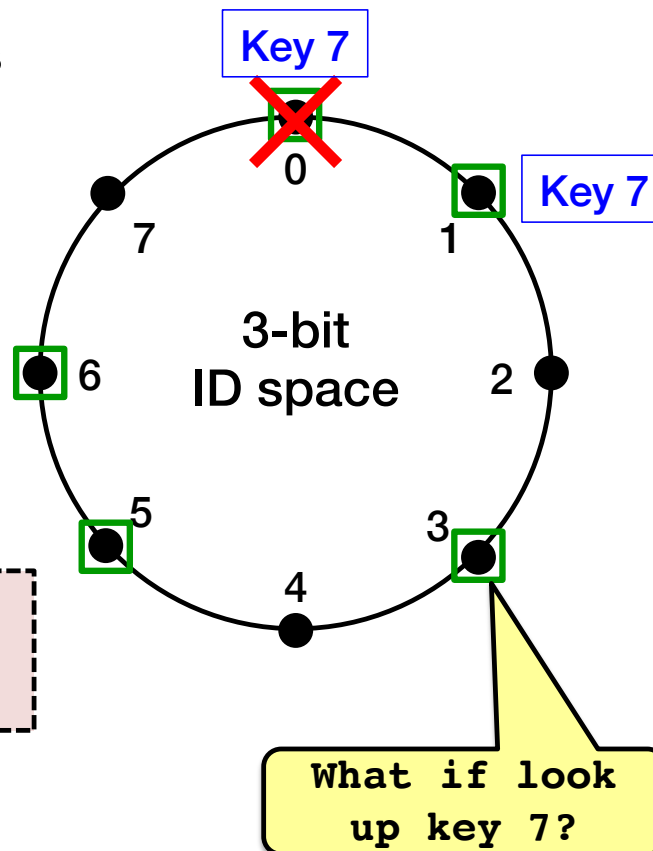
r-nearest successors
(r = logN)

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

"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?

30

# What DHTs got right

- **Consistent hashing**
  - Elegant way to divide a workload across machines
  - Very useful in clusters: actively used 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