## COMP 790-088: Networked & Distributed Systems

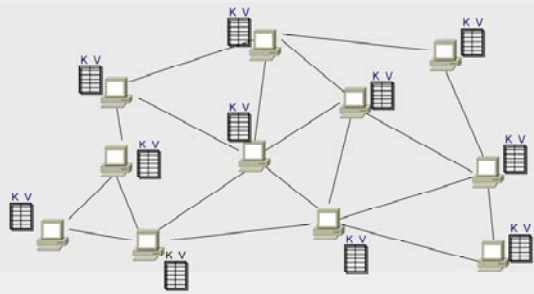# Distributed Hash Tables

*Jasleen Kaur*

November 11, 2009

1

---

## Distributed Hash Table (DHTs)

- Hash table: data structure that maps "keys" to "values"
  - essential building block in software systems

- Distributed Hash Table: similar, but spread across the Internet
  - Each node stores (key, value) pairs
  - Interface:
    - insert(key, value)
    - lookup(key)
    - Join/leave
  - Each DHT node in the overlay supports single operation:
    - given input key, route messages toward node holding key

- "Middleware" for building distributed systems
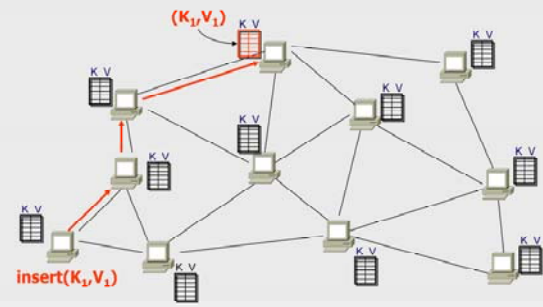  - DNS, File Systems, I3, Content distribution, ….

2

2

## DHT In Action



Operation: take *key* as input; route messages to node holding *key*

3

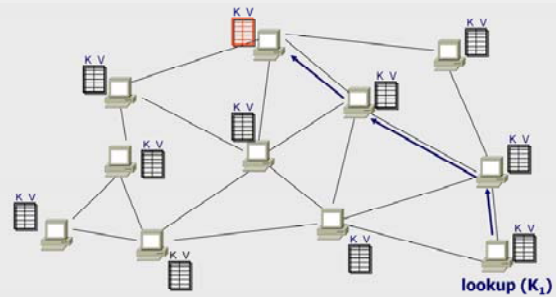## DHT In Action: insert()



(K₁,V₁)

insert(K₁,V₁)

Operation: take *key* as input; route messages to node holding *key*

4

3

4

DHT In Action: lookup()

**Operation:** take *key* as input; route messages to node holding *key*

lookup (K₁)

---

# DHT Design Goals

- An "overlay" network with:
  - » flexible mapping of keys to physical nodes
  - » small network diameter
  - » small degree
  - » local routing decisions

- A "storage" or "memory" mechanism with
  - » best-effort persistence (soft state)

- We'll look at two designs:
  - » Chord
  - » Pastry
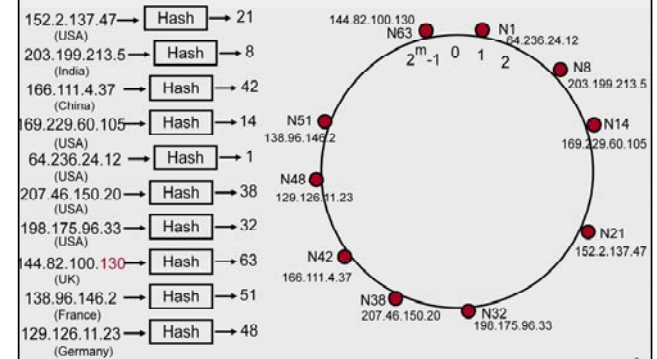
## Chord

- Based on logical *m*-bit identifiers
  - » 0 to $2^m - 1$ ordered in an identifier "circle" (modulo $2^m$)

- (Key, Value) pairs are stored/located by using a ***consistent hash*** function, *CH,* to map keys, *K,* onto a point, $\Phi$, on the circle
  - » $\Phi = CH(K)$

- System nodes are also mapped onto points, $N_i$, on the same identifier circle
  - » # Key values may be greater than # Nodes

- Node $N_i$ stores all (*K,V*) pairs where *K* maps to a point $\Phi$ such that $N_i$ is the <u>first</u> node where
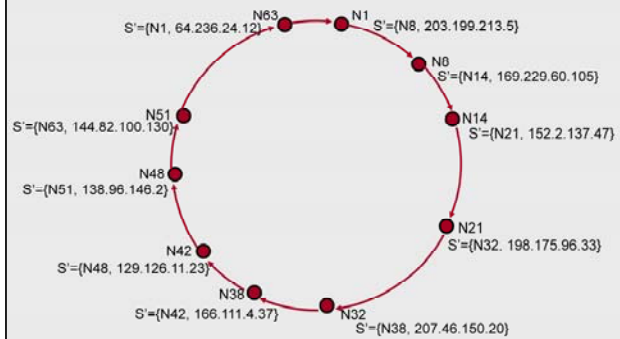  - » $\Phi \leq N_i$   ($N_i$ is the <u>*successor*</u> of $\Phi$)

7

## Hash IP address to Node ID (N) (example with m=6)



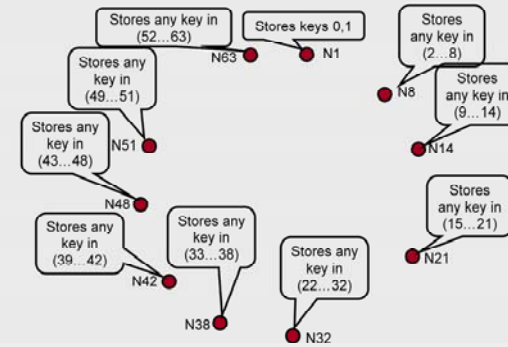| IP address | | Node ID |
|---|---|---|
| 152.2.137.47 (USA) | Hash | 21 |
| 203.199.213.5 (India) | Hash | 8 |
| 166.111.4.37 (China) | Hash | 42 |
| 169.229.60.105 (USA) | Hash | 14 |
| 64.236.24.12 (USA) | Hash | 1 |
| 207.46.150.20 (USA) | Hash | 38 |
| 198.175.96.33 (USA) | Hash | 32 |
| 144.82.100.130 (UK) | Hash | 63 |
| 138.96.146.2 (France) | Hash | 51 |
| 129.126.11.23 (Germany) | Hash | 48 |

8

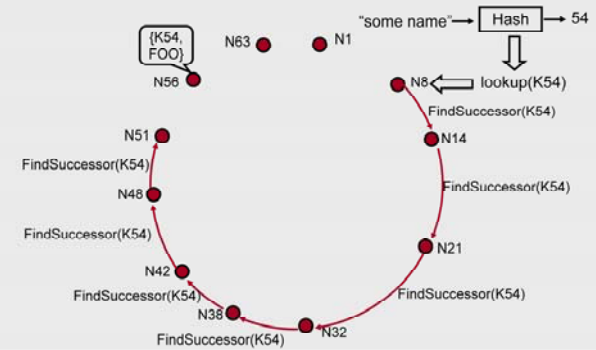Nodes Maintain Successor Pointer (S')



Key locations in example

## Chord

◆ DHT API:
  » Each node stores (key, value) pairs
  » Interface:
    ❖ insert(key, value)
    ❖ lookup(key)
    ❖ Join/leave
  » Each DHT node in the overlay supports single operation:
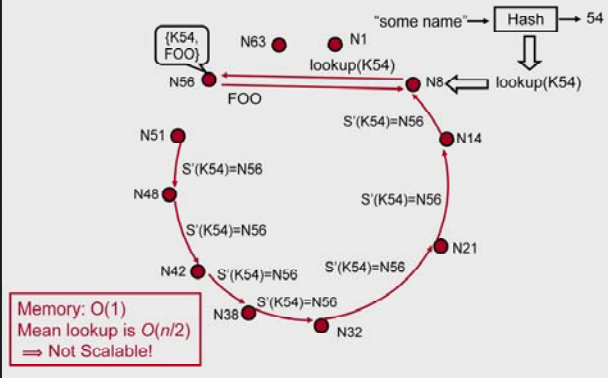    ❖ given input key, route messages toward node holding key

11

## Simple Lookup -- recursive mode (part one: find successor of key)
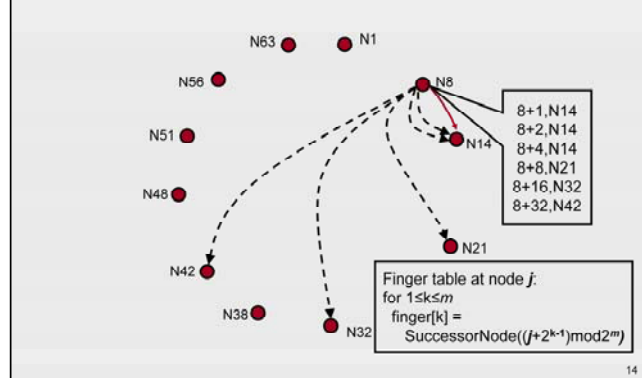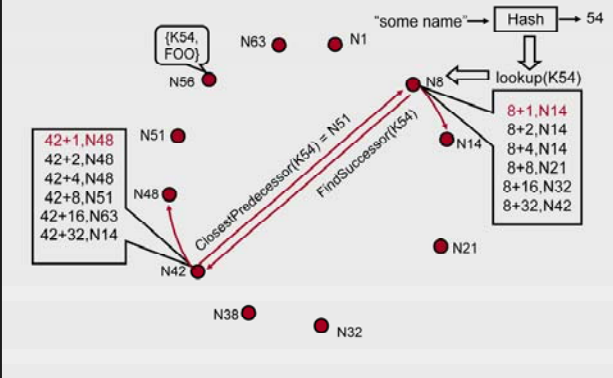


12

## Simple Lookup -- recursive mode
## (part two: return successor & send query)



"some name" → [Hash] → 54

{K54, FOO}

N63    N1

lookup(K54)

N56 ←→ N8 ⇐ lookup(K54)

FOO

S'(K54)=N56

N51    N14

S'(K54)=N56

N48

S'(K54)=N56    S'(K54)=N56

N42    S'(K54)=N56    S'(K54)=N56    N21

S'(K54)=N56

N38

N32

Memory: O(1)
Mean lookup is $O(n/2)$
⟹ Not Scalable!

13

## Scalable Lookup With Small Node State
## (part one: use local "finger table")



N63    N1

N56

N8

8+1,N14
8+2,N14
8+4,N14
8+8,N21
8+16,N32
8+32,N42

N51

N14

N48

N21

N42

Finger table at node $j$:
for $1 \leq k \leq m$
finger[k] =
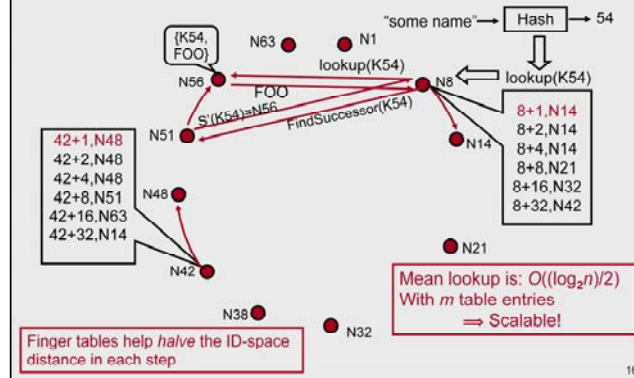SuccessorNode($(j+2^{k-1}) \bmod 2^m$)

N38

N32

14

Scalable Lookup With Small Node State (part two: use remote finger table data)



Scalable Lookup With Small Node State (part three: locate successor node)
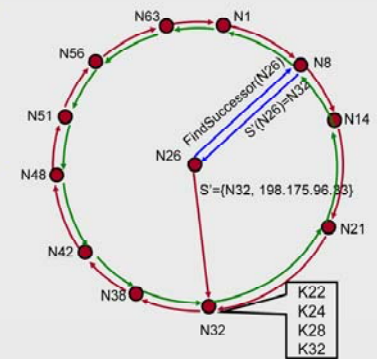
## Chord

- DHT API:
  - » Each node stores (key, value) pairs
  - » Interface:
    - ❖ lookup(key)
    - ❖ insert(key, value)
    - ❖ Join/leave
  - » Each DHT node in the overlay supports single operation:
    - ❖ given input key, route messages toward node holding key

17

## Node Join
## (example, Hash(128.250.6.182) = 26)



- Nodes also maintain a *predecessor* link (not used for search)

- (1) Joining node contacts any existing node to find successor

- (2) Successor link created from returned value.

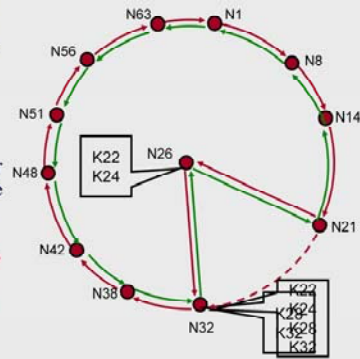FindSuccessor(N26)
S'(N26)=N32
S'={N32, 198.175.96.83}

K22
K24
K28
K32

18

## Node Join
## (example, Hash(128.250.6.182) = 26)

- (3) Successor *Notified* and data for keys < 26 moved and predecessor link made.

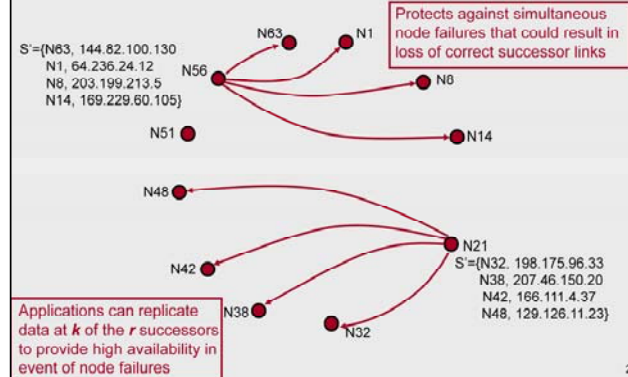- (4) Periodic *Stabilize* protocol run by all nodes updates successor link in predecessor node (N21) and predecessor link in new node; *Fix_Fingers* also run to fix finger tables (uses FindSuccessor() search)



19

## Replication & Robustness:
## Each node maintains list of *r* successors



Protects against simultaneous node failures that could result in loss of correct successor links

S'={N63, 144.82.100.130
N1, 64.236.24.12
N8, 203.199.213.5
N14, 169.229.60.105}

S'={N32, 198.175.96.33
N38, 207.46.150.20
N42, 166.111.4.37
N48, 129.126.11.23}

Applications can replicate data at *k* of the *r* successors to provide high availability in event of node failures

20

## The Chord Theorems

*Theorem IV.1:* For any set of $N$ nodes and $K$ keys, with high probability, the following is true.
1) Each node is responsible for at most $(1 + \epsilon)K/N$ keys.
2) When an $(N + 1)$th node joins or leaves the network, the responsibility for $O(K/N)$ keys changes hands (and only to or from the joining or leaving node).

*Theorem IV.2:* With high probability, the number of nodes that must be contacted to find a successor in an $N$-node network is $O(\log N)$.

*Theorem IV.3:* If any sequence of join operations is executed interleaved with stabilizations, then at some time after the last join the successor pointers will form a cycle on all the nodes in the network.

## The Chord Theorems (cont.)

*Theorem IV.4:* If we take a stable network with $N$ nodes with correct finger pointers, and another set of up to $N$ nodes joins the network, and all successor pointers (but perhaps not all finger pointers) are correct, then lookups will still take $O(\log N)$ time with high probability.

*Theorem IV.5:* If we use a successor list of length $r = \Omega(\log N)$ in a network that is initially stable, and then every node fails with probability 1/2, then with high probability $find\_successor$ returns the closest living successor to the query key.

*Theorem IV.6:* In a network that is initially stable, if every node then fails with probability 1/2, then the expected time to execute $find\_successor$ is $O(\log N)$.