

Distributed Hash Tables

15-415 (Fall 2010)

Adapted from a presentation by Jeff Pang in 15-744, Spring 2007

DHTs

- Like it sounds – a distributed hash table
- Put(Key, Value)
- Get(Key) → Value

Interface vs. Implementation

- Put/Get is an abstract interface
 - Very convenient to program to
 - Doesn't require a “DHT” in today's sense of the world.
 - e.g., Amazon's S³ storage service
 - /bucket-name/object-id → data
- We'll mostly focus on the back-end log(n) lookup systems like Chord
 - But researchers have proposed alternate architectures that may work better, depending on assumptions!

DHTs

- Two options:
 - lookup(key) → node ID
 - lookup(key) → data
- When you know the nodeID, you can ask it directly for the data, but specifying interface as → data provides more opportunities for caching and computation at intermediaries
- Different systems do either. We'll focus on the problem of *locating the node responsible for the data*. The solutions are basically the same.

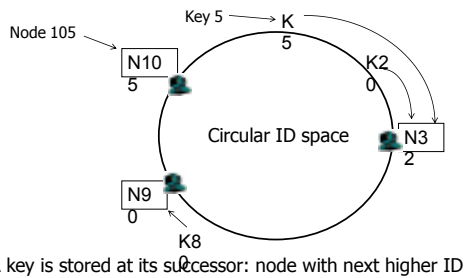
Algorithmic Requirements

- Every node can find the answer
- Keys are load-balanced among nodes
 - Note: We're not talking about *popularity* of keys, which may be wildly different. Addressing this is a further challenge...
- Routing tables must adapt to node failures and arrivals
- How many hops must lookups take?
 - Trade-off possible between state/maint. traffic and num lookups...

Consistent Hashing

- How can we map a key to a node?
- Consider ordinary hashing
 - $\text{func}(\text{key}) \% N \rightarrow \text{node ID}$
 - What happens if you add/remove a node?
- Consistent hashing:
 - Map node IDs to a (large) circular space
 - Map keys to same circular space
 - Key “belongs” to nearest node

DHT: Consistent Hashing



15-441 Spring 2004, Jeff Pang

7

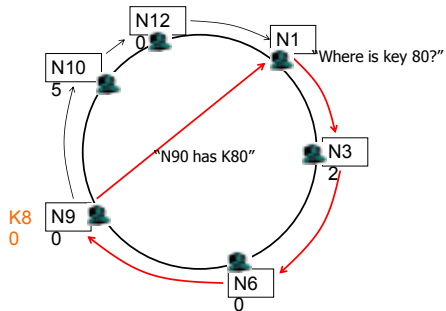
Consistent Hashing

- Very useful algorithmic trick outside of DHTs, etc.
 - Any time you want to not greatly change object distribution upon bucket arrival/departure
- Detail:
 - To have good load balance
 - Must represent each bucket by $\log(N)$ "virtual" buckets

15-441 Spring 2004, Jeff Pang

8

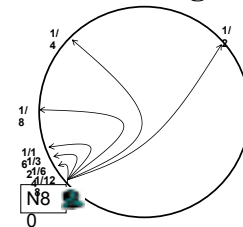
DHT: Chord Basic Lookup



15-441 Spring 2004, Jeff Pang

9

DHT: Chord "Finger Table"



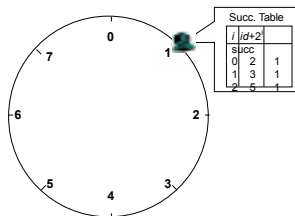
- Entry i in the finger table of node n is the first node that succeeds or equals $n + 2^i$
- In other words, the i th finger points $1/2^{i+1}$ way around the ring

15-441 Spring 2004, Jeff Pang

10

DHT: Chord Join

- Assume an identifier space $[0..8]$
- Node n_1 joins

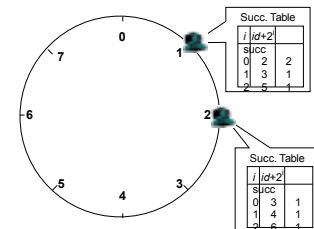


15-441 Spring 2004, Jeff Pang

11

DHT: Chord Join

- Node n_2 joins

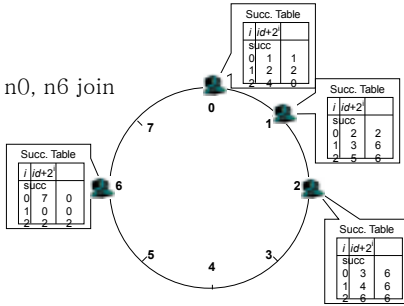


15-441 Spring 2004, Jeff Pang

12

DHT: Chord Join

- Nodes n_0, n_6 join

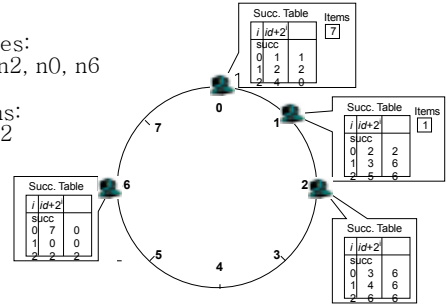


15-441 Spring 2004, Jeff Pang

13

DHT: Chord Join

- Nodes: n_1, n_2, n_0, n_6
- Items: f_7, f_2

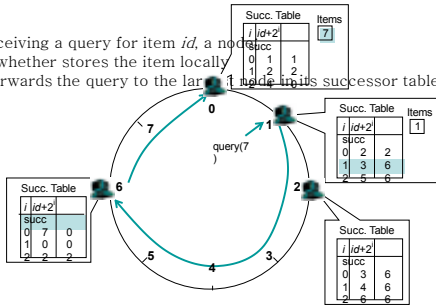


15-441 Spring 2004, Jeff Pang

14

DHT: Chord Routing

- Upon receiving a query for item id , a node
- Checks whether stores the item locally
- If not, forwards the query to the largest node in its successor table that does



15-441 Spring 2004, Jeff Pang

15

DHT: Chord Summary

- Routing table size?
 - Log N fingers
- Routing time?
 - Each hop expects to $1/2$ the distance to the desired

15-441 Spring 2004, Jeff Pang

16

LH*: A Distributed Linear Hash

- Just because we spoke about *Linear Hashing* earlier in the semester during our discussion of growable hashing schemes...
- It is easy to see that Linear Hashing can be distributed.
 - Each of the buckets is a host
 - The buckets can even be RAM-only
- A coordinator is invoked by the host of a bucket upon collision
 - The coordinator assigns a new host from a pool of available hosts
 - It then communicates with the two hosts to coordinate the split
 - After the split, the old hosts knows who it split with and can forward queries
- A retiring host is problematic.
 - Coordinator can supply replacement host to accept bucket of storage
 - Coordinator needs to inform all hosts that could have split with the retiring host over time, so that they can forward
 - Alternate approach: If unable to find a host, contact the coordinator to find its replacement
- One extension of Linear Hashing to the distributed environment is called *LH**

Cassandra and HBase

- Cassandra uses a Chord-based ring as its data store
- HBase is built above HDFS, the Hadoop File System.
 - Replicas go to (a) local node, (b) local rack, (c) some other rack, (d) random after that
 - NameNode knows the mapping - not a hash