Algorithmic Aspects of Distributed k-ary System

Seif Haridi Ali Ghodsi KTH/Royal Institute of Technology Swedish Institute of Computer Science (SICS)

DHT Introduction

- What is a DHT
- Chord: How to partition the data
- Chord: How to interconnect nodes
- Chord: How to speed up search
- Chord: How to maintain pointers

DKS Algorithms

- K-ary Search
- Atomic Join and Leaves
- Broadcast Algorithms
- Multicasting
- Bulk Operations
- Replication
- Topology Maintenance

What is a Distributed Hash Table (DHT)? 1/2

 An infrastructure that enables the distribution of an ordinary hash table onto a set of cooperating nodes



 The DHT provides a basic *lookup service*, which allows *any* node to find the value associated with a given key

• Example:

lookup("CS30"), at any node should return: "Distributed Sys."

Dynamo Workshop, Seif Haridi and Ali Ghodsi

What is a Distributed Hash Table (DHT)? 2/2

To provide the lookup service, the nodes must be interconnected



 Each node maintains a *routing table* with pointers to some other nodes such that lookup requests can be routed to the node storing the requested key/value-pair (a.k.a. item)

DHT Introduction

– What is a DHT

- Chord: How to partition the data

- Chord: How to interconnect nodes
- Chord: How to speed up search
- Chord: How to maintain pointers

DKS Algorithms

- K-ary Search
- Atomic Join and Leaves
- Broadcast Algorithms
- Multicasting
- Bulk Operations
- Replication

How do we distribute the hash table?

- Use a logical name space called the *identifier space*, consisting of identifiers {0,1,2,..., N-1}
- ${\ensuremath{\cdot}}$ The identifier space can be perceived as a logical ring modulo N
- Every node is assigned an identifier using a function H.
- Items are mapped to the identifier space using a function $H_{\rm r}$ every node knows H



- The items are stored at their *successor*, i.e. the first node encountered moving in the clockwise direction
 - Example: *N*=16, nodes {**a,b,c,d,e**}, and 5 items
 - Node **a** gets identifier 0 since $H_1(\mathbf{a})=0$, the other nodes
 - **b**, **c**, **d**, **e**, get identifiers 2, 5, 6, 11 the same way
 - Item ("cs15", "networking") is mapped to identifier 13 since H_2 ("cs15") = 13, other items are similarly to 15, 2, 4, 5

Dynamo Workshop, Seif Haridi and Ali Ghodsi

DHT Introduction

- What is a DHT
- Chord: How to partition the data

– Chord: How to interconnect nodes

- Chord: How to speed up search
- Chord: How to maintain pointers

DKS Algorithms

- K-ary Search
- Atomic Join and Leaves
- Broadcast Algorithms
- Multicasting
- Bulk Operations
- Replication
- Topology Maintenance

How do we interconnect the nodes?

 Each node maintains a *routing pointer* to the successor in the ring



10/27/06

 The successor of a **node** *n* is the first node met going in clockwise direction starting at n + 1

Successor of Node $0 \rightarrow \text{Node } 2$ Successor of Node $2 \rightarrow$ Node 5 Successor of Node $5 \rightarrow \text{Node } 6$ Successor of Node $6 \rightarrow$ Node 11 Successor of Node $11 \rightarrow Node 0$

Simple lookup Example



DHT Introduction

- What is a DHT
- Chord: How to partition the data
- Chord: How to interconnect nodes

- Chord: How to speed up search

– Chord: How to maintain pointers

DKS Algorithms

- K-ary Search
- Atomic Join and Leaves
- Broadcast Algorithms
- Multicasting
- Bulk Operations
- Replication
- Topology Maintenance

Speeding up lookups

 Each node, n, not only points to its successor but to the successors of

 $n+2^{1}, n+2^{2}, n+2^{3}, \dots, n+2^{L}$ (all arithmetic operations modulo N)

- At each step in the routing, the distance between the currrent node and destination is halved (worst case).
- Yields O(log₂N) hops at worst



DHT Introduction

- What is a DHT
- Chord: How to partition the data
- Chord: How to interconnect nodes
- Chord: How to speed up search
- Chord: How to maintain pointers

DKS Algorithms

- K-ary Search
- Atomic Join and Leaves
- Broadcast Algorithms
- Multicasting
- Bulk Operations
- Replication
- Topology Maintenance

Checking for new joins and departures

Each node periodically *stabilizes*

Algorithm 1 Chord's periodic stabilization protocol

- 1: procedure n.STABILIZE()
- 2: p := succ.GetPredecessor()
- 3: if $p \in (n, succ)$ then

```
4: succ := p
```

```
5: end if
```

```
6: succ.Notify(n)
```

7: end procedure

```
    procedure n.GETPREDECESSOR()
    return pred
```

```
10: end procedure
```

```
11: procedure n.NOTIFY(p)12: if p \in (pred, n] then13: pred := p14: end if15: end procedure
```



Stabilization: join

• A joining node, only notifies its successor about its existence



Stabilization: join



Handling dynamism: successor-lists

 Every node in the system maintains additional routing pointers to its *f* successors



Example: *f*=2

Every node knows its *f*=2 successors

Node **0** : Node 2 and Node 5

Node 2 : Node 5 and Node 6

Node 5 : Node 6 and Node 11

Node 6 : Node 11 and Node 0

Node 11 : Node 0 and Node 2

 If node n detect that its successor has failed, it replaces it with the first alive successor node it knows

Stabilization: failure

• A joining node, only notifies its successor about its existence



DHT Summary

- Completely decentralized
- Self-organizes as nodes join, leave, and fail
- Maximum $\log_2(N)$ number of hops to find items, where N is the number of hops
- Each node only stores a small amount of items, on average D/N (D is the number of items)
- Each node only maintains a small amount of routing information $\log_2(N)$
- Each join/leave/failure event requires D/N items to be reshuffled
- Each join/leave/failure event requires $(\log_2(N))^2$ messages to restore the routing state

DHT Introduction

- What is a DHT
- Chord: Distributing a hash table
- Chord: How to partition the data
- Chord: How to interconnect nodes
- Chord: How to speed up search
- Chord: How to maintain pointers

DKS Algorithms

- K-ary Search
- Atomic Join and Leaves
- Broadcast Algorithms
- Multicasting
- Bulk Operations
- Replication
- Topology Maintenance

K-ary Search: Intuition

• Goal:

- At most $\log_k(N)$ hops per lookup
 - k being a configurable parameter
 - ${\mbox{ \bullet }} N$ being the number of nodes
- Instead of only $\log_2(N)$

Achieving *log_kN* lookup effiency

- Each DKS node maintains log_kN levels, and each level contains k intervals with pointers to the first node encountered in the interval
- Example, k=4, N=64 (4³), node 0





Achieving *log_kN* lookup effiency

- Each DKS node maintains log_kN levels, and each level contains k intervals with pointers to the first node encountered in the interval
- Example, k=4, N=64 (4³), node 0



- level 1, 4 intervals
- level 2, 4 intervals

Achieving *log_kN* lookup effiency

- Each DKS node maintains log_kN levels, and each level contains k intervals with pointers to the first node encountered in the interval
- Example, k=4, N=64 (4³), node 0

48

10/27/06

- level 1, 4 intervals
- level 2, 4 intervals
- level 3, 4 intervals

16



Routing table

- For the desired based k, the identifier space N is a power of k
 - $\mathbf{N} = \mathbf{k}^{\mathsf{L}}$
- $L = log_k N$ is the number of levels
- RT is of size (k-1)L
- Views
- Level 1: V(1) = [n, n+N)
- Level 2: V(2) = [n, (n+N)/k)
- Level ℓ : V(I) = [n, (n+N)/k^{ℓ})
- Intervals: at any level l∈ _, V(l) is partitioned into k intervals I (l,i), for 0≤i≤k-1

Routing Table for node 21 Views, Levels (k = 4)



Routing Table for node 21 Views, Levels (k = 4)



Dynamo Workshop, Seif Haridi and Ali Ghodsi

(5,21)

(33, 48)

(24, 24)

Routing

- The routing table can be organized as a monotonically increasing set of pointers
- Each pointer refers to a node in corresponding interval [f(i), f(i+1))
- Start of each interval i, $1 \le i \le (k-1)\log_k N$: - f(i) = n + (1 + (i-1) mod (k-1)) k $\lfloor (i-1)/(k-1) \rfloor$



Simple Recursive Routing

Algorithm 11 Recursive Lookup Algorithm

- 1: procedure n.LOOKUP(i, OP)
- 2: if TERMINATE(i) then
- 3: $p := NEXT_HOP(i)$
- 4: res := p.OP()
- 5: return res
- 6: else
- 7: $m := \text{NEXT}_HOP(i)$
- 8: return m.LOOKUP(i)
- 9: end if

10: end procedure

▷ OP could carry parameters

Greedy Routing

Algorithm 14 Greedy Lookup

- 1: **procedure** *n*.TERMINATE(*i*)
- return $i \in (n, succ]$ 2:
- 3: end procedure
- 1: procedure *n*.NEXT_HOP(*i*)
- if TERMINATE(i) then 2:3:
 - return succ
- else 4: 5:

6:

7:

8:

9:

10:

11:

12:

```
r := succ
```

end for

end if

13: end procedure

return r

for j := 1 to K do

if $rt(j) \in (n,i)$ then

end if

r := rt(j)

 \triangleright Node has K pointers

32

Dynamo Workshop, Seif Haridi and Ali Ghodsi

Topology Maintenance

- Pointers in the routing table are monitored by an inaccurate failure detectors, and reset (nil) if a node is suspected
- Failed pointers leads to a lookup to the start of the corresponding entry
- The result is the successor and its successor list
- A suitable node is installed

DHT Introduction

- What is a DHT
- Chord: How to partition the data
- Chord: How to interconnect nodes
- Chord: How to speed up search
- Chord: How to maintain pointers

DKS Algorithms

- K-ary Search
- Atomic Join and Leaves
- Broadcast Algorithms
- Multicasting
- Bulk Operations
- Replication
- Topology Maintenance

Atomic Join and Leave

- Goal
 - Key or Lookup Consistency

Informally

- At any time, at most one node responsible for any key
- Joins/leaves should "not affect" functionality of lookups/inserts/updates/deletes

Lookup consistency is not guaranteed in traditional DHTs: lookup(4)



Dynamo Workshop, Seif Haridi and Ali Ghodsi

36

10/27/06

Locking to achieve atomicity

- Problem partly reduced to dining philosophers
 - Each node has a *fork* and a *lock queue*
 - Current node's fork and successor's fork acquired before modifying the ring
 - Avoiding deadlocks
 - Asymmetric Locking: one node acquires locks in reverse order
 - Probabilistic Locking: preempt locks



Locking to achieve atomicity

- Pitfalls:
 - Join adds a "philosopher"
 - Solution: some requests in the lock queue forwarded to to new node
 - Leave removes a "philosopher"
 - Problem:

if leaving node gives its lock queue to its successor, some nodes will get a worse position in queue *starvation*!

Correctness Properties

Safety

- Deadlock freedom
- Termination (trivial)

Liveness

- Livelock freedom (always some progress)
- Starvation freedom (every node makes progress)

Looking Consistency with Locking

- Lookup Consistency with Joins
 - Successor forwards requests to new node
- Lookup Consistency with Leaves
 - Leaving node forwards requests to successor
- Proving Lookup Consistency
 - Configuration is the pointers of all nodes
 - Show: any reachable configuration, only one node responsible for a key, starting lookup anywhere

DHT Introduction

- What is a DHT
- Chord: How to partition the data
- Chord: How to interconnect nodes
- Chord: How to speed up search
- Chord: How to maintain pointers

DKS Algorithms

- K-ary Search
- Atomic Join and Leaves
- Broadcast Algorithms
- Multicasting
- Bulk Operations
- Replication
- Topology Maintenance

Broadcasting

- Broadcast to all nodes in the DHT
- Fast dissemination
 - Broadcasting proceeds in parallel
 - Time complexity $O(\log N)$, N is number of nodes
- Efficient
 - No redundancy, full coverage
 - O(N) message complexity, N is number of nodes

Correctness of Broadcast

Safety

- Non-redundancy: A node never receives the same message more than once

Liveness

- **Termination:** every broadcast eventually terminates
- **Coverage:** every node eventually gets the message

Simple Best-Effort Broadcast

Algorithm 17 Simple Broadcast Algorithm

```
1: event n.STARTSIMPLEBCAST(msg) from m
```

sendto n.SIMPLEBCAST(msg, n) 2:

Local message to itself

3: end event

```
    event n.SIMPLEBCAST(msg, limit) from m
```

```
Deliver(msg)
2:
3:
```

```
Deliver msg to application
\triangleright Node has M unique pointers
```

```
for i := M downto 1 do
4:
```

```
if u(i) \in (n, limit) then
```

```
sendto u(i).SIMPLEBCAST(msg, limit)
```

```
limit := u(i)
```

```
end if
7:
```

```
end for
8:
```

```
9: end event
```

 5°

6:

DHT Introduction

- What is a DHT
- Chord: How to partition the data
- Chord: How to interconnect nodes
- Chord: How to speed up search
- Chord: How to maintain pointers

DKS Algorithms

- K-ary Search
- Atomic Join and Leaves
- Broadcast Algorithms

- Multicasting

- Bulk Operations
- Replication
- Topology Maintenance

Multicasting

•All nodes in the system are members of an instance of

•Group creation: –Create a *DKS* instance –Make it available in O –Joins the group

•Multicasting —Broadcasting within the group



DHT Introduction

- What is a DHT
- Chord: How to partition the data
- Chord: How to interconnect nodes
- Chord: How to speed up search
- Chord: How to maintain pointers

DKS Algorithms

- K-ary Search
- Atomic Join and Leaves
- Broadcast Algorithms
- Multicasting
- Bulk Operations
- Replication
- Topology Maintenance

Bulk Operations: motivation

- Background
 - Assume building filesystem on-top of DHT
 - 4mb file of 4kb block -> 4000 blocks
- Making 4000 lookups expensive
 - Marshaling/unmarshaling (XML?)

Bulk Operations

• Bulk(I)

- I is a set of identifiers
- Reach every node with identifier in I

Bulk_Owner(I)

- I is a set of identifiers
- Reach node responsible for every id in I

Bulk_Feedback(I)

Same as BULK, but gets feedback through the virtual dissemination tree

Bulk_Owner_Feedback(I)

- Same as BULK, but gets feedback through the virtual dissemination tree
- No redundant messages sent
- Max log(n) messages per node 10/27/06

Bulk Properties

- Extreme Case 1
 - I is all identifiers
 - N messages to reach N nodes
 - Completed in log(N) time
 - Identical to broadcast
- Extreme Case 2
 - I is a singleton with one identifier
 - log(n) messages to perform lookup
 - Completed in log(N) time
 - Identical to lookup

Dynamo Workshop, Seif Haridi and Ali Ghodsi

Bulk Applications

- Bulk lookup
 - Fetch all values associated with keys {x1, x2,...}
- Pseudo reliable broadcast
 - Broadcast with feedback, use time outs to avoid hanging, retransmit message with bulk to those not covered
- Range queries
 - Cover all nodes in range]i, j]
- Replication
 - Replicate item x to following responsible nodes {x1,x2,...}
 - Symmetric Replication
- Topology maintenance
 - Update all nodes with identifiers in range {x1,x2,...}
 - Correction-on-change

DHT Introduction

- What is a DHT
- Chord: How to partition the data
- Chord: How to interconnect nodes
- Chord: How to speed up search
- Chord: How to maintain pointers

DKS Algorithms

- K-ary Search
- Atomic Join and Leaves
- Broadcast Algorithms
- Multicasting
- Bulk Operations
- Replication
- Topology Maintenance

Background of Symmetric Replication

- DKS used to use Successor-lists, like Chord, Pastry, Koorde etcetera.
- This was abandoned in favor of Symmetric Replication because ...

Successor-lists and Leaf sets

- If a node joins or leaves
 - -f replicas need to be updated



Successor-lists and Leaf sets

- If a node joins or leaves
 - -f replicas need to be updated
 - Without central coordination epidemics are used f^2

Color represents a node's item

Node leaves

Yellow, blue, grey need to be re-distributed

Dynamo Workshop, Seif Haridi and Ali Ghodsi

10/27/06

Symmetric Replication

Goal of Symmetric Replication: – Simplicity!

- Enable concurrent requests
 - Do load-balancing
 - Increase robustness (high failure rates)
- Use with erasure codes
 - Given a k/n erasure used
 - Replicate n times
 - Fetch k replicas with bulk operations

Idea: - Partition the identifier space into *m* equivalence classes such that

The cardinality of each class is f
 m=N/f

 Each node replicates the equivalence class of every identifier it is responsible for

Symmetric replication Replication degree f=4, $Id=\{0,...,15\}$ **Congruence classes modulo 4:** 0, 4, 8, 12 1, 5, 9, 13 **Data**: 15, 0 2, 6, 10, 14 **Data**: 14, 13, 12, 11 0 15 1 3, 7, 11, 15 0 2 14 **Data**: 1, 2, 3 13 3 • 12 4 0 5 11 Data: 4, 5 \bigcirc 6 10 7 9 hop, Seif Haridi ar Ali Ghods Data: 6, 7, 8, 9, 10 10/27/06 58





DHT Introduction

- What is a DHT
- Chord: How to partition the data
- Chord: How to interconnect nodes
- Chord: How to speed up search
- Chord: How to maintain pointers

DKS Algorithms

- K-ary Search
- Atomic Join and Leaves
- Broadcast Algorithms
- Multicasting
- Bulk Operations
- Replication
- Topology Maintenance

Applications ontop of DKS

- MyriadStore
 - Distributed Backup
 - Use DKS for metadata storage
- Keso
 - Distributed File System
- DOH
 - Replicated Web Servers

• Delegent

Decentralized Trust Management Systems

Dynamo Workshop, Seif Haridi and Ali Ghodsi

10/27/06



Thank You!

http://dks.sics.se

http://www.sics.se/~ali/dissert.pdf

Dynamo Workshop, Seif Haridi and Ali Ghodsi

63

10/27/06

DHT Introduction

- What is a DHT
- Chord: How to partition the data
- Chord: How to interconnect nodes
- Chord: How to speed up search
- Chord: How to maintain pointers

DKS Algorithms

- K-ary Search
- Atomic Join and Leaves
- Broadcast Algorithms
- Multicasting
- Bulk Operations
- Replication

IP multicast vs Overlay Multicast

Idea

- Utilize IP multicast where available

- Before joining
 - IP multicast to discover *local* nodes
 - Use same identifier as local nodes
 - Otherwise join as usual
- Whenever receive a message
 IP multicast to all local nodes

2 5

Topology Maintenance

- Chord:
 - Extra pointers to *successor* of intervals

 $n+2^{1}, n+2^{2}, n+2^{3}, ..., n+2^{L}$ (arithmetic modulo N)

• DKS:

- Extra pointers to *predecessor* of intervals
- Ring kept as Chord
- Responsibility as Chord
 - Items stored on successor
- Consequence:
 - DKS: Cost of maintenance O(1) if no churn
 - Chord Cost of maintenance O(log N) if no churn

