

# Distributed Hash Tables

## Pastry

Smruti R. Sarangi

Department of Computer Science  
Indian Institute of Technology  
New Delhi, India

# Outline

- 1 Distributed Hash Tables
- 2 Pastry
  - Overview
  - Operation
  - Arrival, Departure, and Locality
  - Results

# Normal Hashtables

- **Hashtable** : Contains a set of key-value pairs. If the user supplies the key, the hashtable returns the value.
- Basic operations.

**insert(key,value)** Inserts the key,value pair into the hashtable.

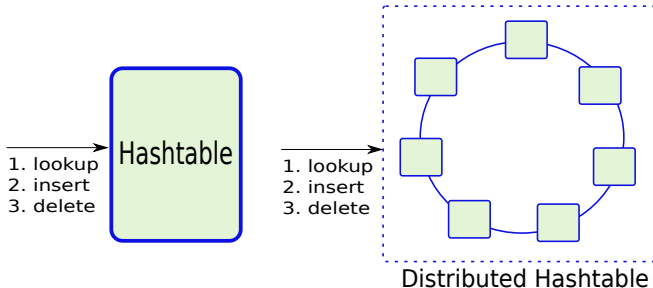
**lookup(key)** Returns the value, or **null** if there is no value.

**delete(key)** Deletes the key

**Time Complexity** Approximately,  $\theta(1)$

- Need a sophisticated hash function to map keys to unique locations.
- Need to resolve collisions through chaining or rehashing.

# Distributed Hashtables(DHT)



## Salient Points of DHTs

- They can store more data than centralized databases.
- DHTs are the only feasible options for web-scale data: facebook, linkedin, google
  - Assume that a bank has 10 crore customers (0.1 billion)
  - Each customer requires storage equivalent to the size of this latex file.
  - Total storage requirement:  $8 \text{ KB} \times 0.1 \text{ billion} = 0.8 \text{ TB}$

## Salient Points of DHTs

- They can store more data than centralized databases.
  - DHTs are the only feasible options for web-scale data: facebook, linkedin, google
    - Assume that a bank has 10 crore customers (0.1 billion)
    - Each customer requires storage equivalent to the size of this latex file.
    - Total storage requirement:  $8 \text{ KB} \times 0.1 \text{ billion} = 0.8 \text{ TB}$
- A user is sharing 100 songs : 500 MB/user
  - There are 10 crore(0.1 billion) users
  - Storage: 50 PB (petabytes)

## Salient Points of DHTs

- They can store more data than centralized databases.
- DHTs are the only feasible options for web-scale data: facebook, linkedin, google
  - Assume that a bank has 10 crore customers (0.1 billion)
  - Each customer requires storage equivalent to the size of this latex file.
  - Total storage requirement:  $8 \text{ KB} \times 0.1 \text{ billion} = 0.8 \text{ TB}$

- A user is sharing 100 songs : 500 MB/user
- There are 10 crore(0.1 billion) users
- Storage: 50 PB (petabytes)

**There is a difference of an order of magnitude !!!**

# Advantages of DHTs

- DHTs scale, and are ideal candidates for web scale storage.
- They are more immune to node failures. They use extensive data replication.
- DHTs also scale in terms of number of users. Different users are redirected to different nodes based on their keys. ( [better load balancing](#) ).
- In the case of torrent applications: they reduce the legal liability since there is no dedicated central server. 😊

## Major Proposals

Pastry, Chord, Tapestry, CAN, Fawn



# Outline

## 1 Distributed Hash Tables

## 2 Pastry

- Overview
- Operation
- Arrival, Departure, and Locality
- Results

# Salient Points

- Scalable distributed object location service.
- Uses a ring based **overlay network**.
- The overlay network takes into account network locality.
- Automatically adapts to the arrival, departure, and failure or nodes.
- Pastry has been used as the substrate to make large storage services (PAST) and scalable publish/subscribe system (SCRIBE).
  - PAST is a large scale file storage service.
  - SCRIBE stores a massive number of topics in the DHT. When a topic changes, the list of subscribers are notified.

# Design of Pastry

- **Node** → has a unique 128 bit nodeId.
- The nodes are conceptually organized as a ring, arranged in ascending order of nodeIds.
- nodeIds are generated by computing a cryptographic hash of the node's IP address or public key.
- Basic idea of routing:
  - Given a key, find its 128 bit hash.
  - Find the node whose nodeId is numerically closest to the hash-value.
  - Send the request to that node.

# Advantages of Pastry

- Pastry can route a request to the right node in less than  $\lceil \log_{2^b}(N) \rceil$  steps.  $b$  is typically 4.
- Eventual delivery is guaranteed unless  $L/2$  nodes with adjacent nodeIds fail simultaneously.  $L = 16$  or  $32$ .
- Both nodeIds and keys are thought to be base  $2^b$  numbers. If we assume that  $b = 4$ , then these are hexadecimal numbers.
- In each step, the request is forwarded to a node whose shared prefix with the key is at least 1 digit ( $b$  bits) more than the length of the shared prefix with the current node's nodeId.
- If no such node is found, then the request is forwarded to a node, which has the same length of the prefix but is numerically closer to the key.
- If no such node is found, then the request is forwarded to a node, which has the same length of the prefix but is numerically closer to the key.

# Outline

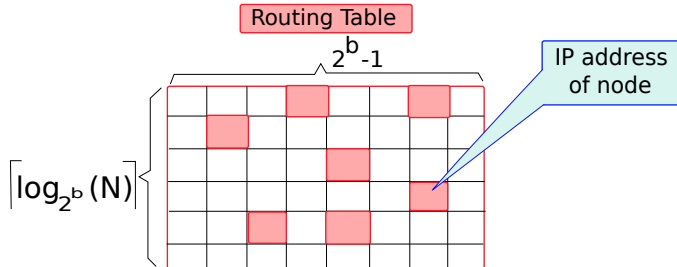
- 1 Distributed Hash Tables
- 2 **Pastry**
  - Overview
  - **Operation**
  - Arrival, Departure, and Locality
  - Results

# Pastry Node

## Structure of a Pastry Node

It contains a **routing table**, **neighborhood table**, and a **leaf set**

### Structure of the routing table



# Routing Table ( $\mathcal{R}$ )

- The routing table contains  $\lceil \log_{2^b}(N) \rceil$  rows.
- The entries at row  $i$  point to nodes that share the first  $i$  digits of the prefix with the key.
- Each row contains  $2^b - 1$  columns.
- Each cell refers to a base  $2^b$  digit.
- If the digit associated with a cell matches the  $(i + 1)^{th}$  digit of the key, then we have a node that matches the key with a longer prefix.
- We should route the request to that node.

# Neighborhood Set and Leaf Set

## Neighborhood Set ( $\mathcal{M}$ )

- It contains  $M$  nodes that are closest to the node according to a proximity metric.
- Typically contains  $2^{b+1}$  entries.

## Leaf Set ( $\mathcal{L}$ )

- Contains  $L/2$  nodes with a numerically closest larger nodeIds.
- Contains  $L/2$  nodes with a numerically closest smaller nodeIds.
- Typically  $2^b$ .



# Routing Algorithm

**Algorithm 1:** Routing algorithm

1 **Input:** key  $K$ , routing table  $\mathcal{R}$

**Output:** Value  $V$

**if**  $\mathcal{L}_{-L/2} \leq D \leq \mathcal{L}_{L/2}$  **then**

    /\*  $K$  is within the range of the leaf set      \*/

2     forward  $K$  to  $L_i$  such that  $|L_i - K|$  is minimal

3 **end**

4 **else**

5      $I \leftarrow \text{prefix}(K, \text{nodeId})$

**if**  $R_i^{D_i} \neq$  **then**

6         forward to  $R_i^{D_i}$

7     **end**

8     **else**

9         forward to  $(T \in \mathcal{L} \cup \mathcal{R} \cup \mathcal{M})$  such that

$\text{prefix}(T, K) \geq I$

$|T - K| < |\text{nodeId} - K|$

# Explanation

- The node first checks to find if the key is within the leaf set. If so, it forwards the messages to the closest node (by `nodeId`) in the leaf set.
- Otherwise, Pastry forwards the message to a node with one more matching digit in the common prefix.
- In the rare case, when we are not able to find a node that matches the first two criteria, we forward the request to any node that is closer to the key than the current `nodeId`. Note that it still needs to have a match of at least  $l$  digits.

# Performance and Reliability

- If  $L/2$  nodes in the leaf set are alive then the message can always be passed on to some other node.
- At every step, we are ( **with high probability** ) either searching in the leaf set or moving to another node.
- If the key is within the range of the leaf set  $\rightarrow$  it is at most **one hop** away
- Otherwise, at every step we increase the length of the matched prefix by  $2^b$ .

## Routing Time Complexity

The average routing time is thus  $O(\log_{2^b}(N))$ .

# Outline

## 1 Distributed Hash Tables

## 2 Pastry

- Overview
- Operation
- **Arrival, Departure, and Locality**
- Results

# Node Arrival

- Assume that node  $X$  wants to join the network.
- It nodes another nearby node  $A$ .
- $A$  can also be found with **expanding ring multicast**.
- $A$  forwards the request to the numerically closest node –  $Z$ .
- Nodes,  $A$ ,  $Z$ , and all the nodes in the path from  $A$  to  $Z$  send their routing tables to  $X$ .
- $X$  uses all of this information to initialize its tables.

# Table Initialization

## Neighborhood Set

Since  $A$  is close to  $X$ .  $X$  copies  $A$ 's neighborhood set.

## Leaf Set

$Z$  is the closest to  $X$  (nodeId).  $X$  uses  $Z$ 's leaf set to form its own leaf set.

# Table Initialization - II

## Routing Table

- Assume that  $A$  and  $X$  do not share any digits in the prefix ( **General Case** ).
- The first row of the routing table are independent of the node id ( **No Match** ).  $X$  can use it to initialize its first row.
- Every node in the path from  $A$  to  $Z$  has one additional digit matching with  $X$ . Let  $B_i$  be the  $i^{th}$  entry in the path from  $A$  to  $Z$ .
- **Observation** :  $B_i$  shares  $i$  bits of its prefix with  $X$ . Use its  $i^{th}$  row in its routing table to populate the  $i^{th}$  row of the routing table of  $X$ .
- **Finally**  $X$  transmits its state to all the nodes in  $\mathcal{L} \cup \mathcal{R} \cup \mathcal{M}$ .

# Node Departure – Leaf Set

- Nodes might just fail or leave without notifying.

## Repairing the leaf set

- Assume that a leaf  $\mathcal{L}_{-k}$  fails. ( $-L/2 < k < 0$ ).
- In this case, the node contacts  $\mathcal{L}_{-L/2}$ .
- It gets its leaf set and merges it with its leaf set.
- For any new nodes added, it verifies their existence by pinging them.



# Node Departure – Routing Table and Neighborhood Set

## Repairing the Routing Table

- Assume that  $R_i^d$  fails.
- Try to get a replacement for it by contacting  $R_i^j$  ( $d \neq i$ )
- If it is not able to find a candidate, it casts a wider net by asking  $R_{i+1}^j$  ( $i \neq d$ )

## Repairing the Neighborhood Set ( $\mathcal{M}$ )

- A node periodically pings its neighbors.
- If a neighbor is found to be dead, it gets  $\mathcal{M}$  from its neighbors and repairs its state.

# Maintaining Locality

- Assume that before node  $X$  is added, there is good locality.
  - All the nodes in the routing table point to nearby nodes.
- We add a new node  $X$ 
  - We start with a nearby node  $A$  and move towards  $Z$  (closest by nodeId).
  - Let  $B_i$  be the  $i^{th}$  node in the path. ( **Induction assumption:  $B_i$  has locality** )
  - $B_i$  is fairly close to  $X$  because it is fairly close to  $A$ .
  - Since we get the  $i^{th}$  row of the routing table from  $B_i$ , and  $B_i$  has locality, the  $i^{th}$  row of the routing table of  $X$ .
  - Thus,  $X$  has locality.

# Maintaining Locality

- Assume that before node  $X$  is added, there is good locality.
  - All the nodes in the routing table point to nearby nodes.
- We add a new node  $X$ 
  - We start with a nearby node  $A$  and move towards  $Z$  (closest by nodeId).
  - Let  $B_i$  be the  $i^{th}$  node in the path. ( **Induction assumption:  $B_i$  has locality** )
  - $B_i$  is fairly close to  $X$  because it is fairly close to  $A$ .
  - Since we get the  $i^{th}$  row of the routing table from  $B_i$ , and  $B_i$  has locality, the  $i^{th}$  row of the routing table of  $X$ .
  - Thus,  $X$  has locality.

Induction hypothesis proved

# Tolerating Byzantine Failures

- Have multiple entries in each cell of the routing table.
- Randomize the routing strategy.
- Periodically send IP broadcasts and multicasts (expanding ring) to connect disconnected networks.

# Outline

## 1 Distributed Hash Tables

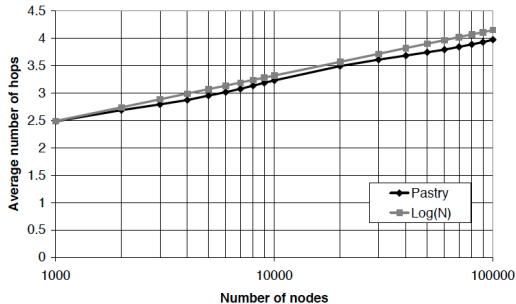
## 2 Pastry

- Overview
- Operation
- Arrival, Departure, and Locality
- Results

# Setup

- 100,000 nodes
- Each node runs a Java based VM
- Each node is assigned a co-ordinate in an Euclidean plane.

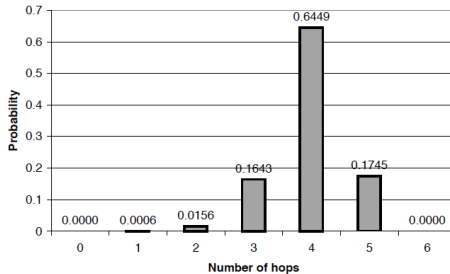
# Results - I



**Fig. 4.** Average number of routing hops versus number of Pastry nodes,  $b = 4$ ,  $|L| = 16$ ,  $|M| = 32$  and 200,000 lookups.

Source [1]

## Results - II

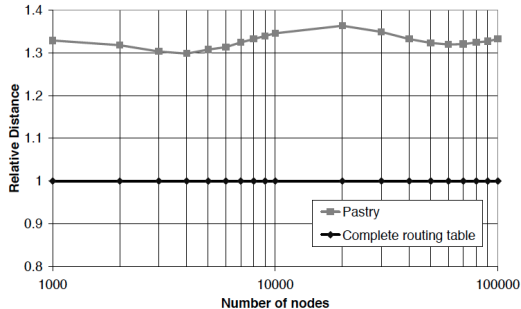


**Fig. 5.** Probability versus number of routing hops,  $b = 4$ ,  $|L| = 16$ ,  $|M| = 32$ ,  $N = 100,000$  and 200,000 lookups.

Source [1]



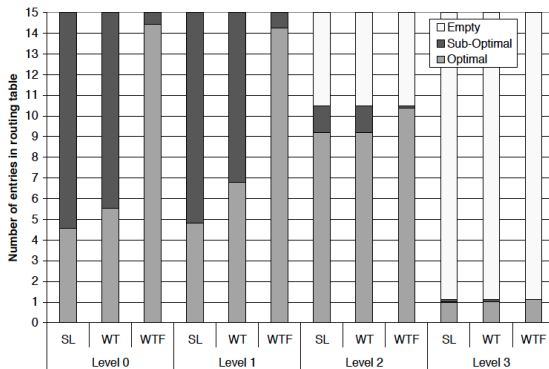
## Results - III



**Fig. 6.** Route distance versus number of Pastry nodes,  $b = 4$ ,  $|L| = 16$ ,  $|M| = 32$ , and 200,000 lookups.

Source [1]

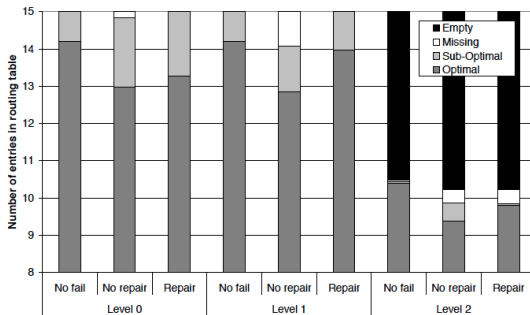
## Results - IV



**Fig.7.** Quality of routing tables (locality),  $b = 4$ ,  $|L| = 16$ ,  $|M| = 32$  and 5,000 nodes.

Source [1]

# Results - V



**Fig. 9.** Quality of routing tables before and after 500 node failures,  $b = 4$ ,  $|L| = 16$ ,  $|M| = 32$  and 5,000 starting nodes.

Source [1]



Pastry: Scalable, decentralized object location and routing for large-scale peer-to-peer systems by Antony Rowstron and Peter Druschel in Middleware 2001