# Distributed Hash Table Final Report

XUZHENG LU

G34363475
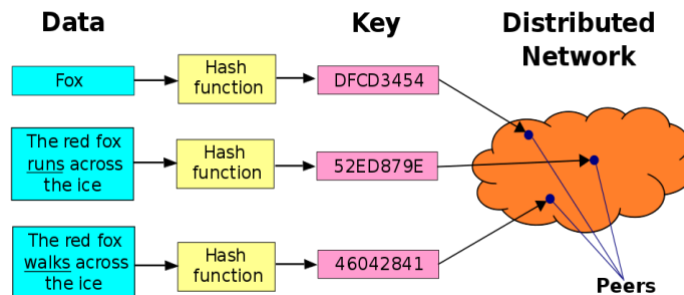
ZETIAN ZHENG

G36558368

XIAOYU SHEN

G33150730

GUANGYUAN SHEN

G48268794

## 1  INTRODUCTION

Peer-to-peer (P2P) computing or networking is a distributed application architecture that partitions tasks or workloads between peers. (Wikipedia: peer to peer)

In this kind of distributed system, "peer to peer" means that the function of each node is equal, and each node is connected with others via the internet. In this network, files are transported between nodes. There are many features for peer-to-peer systems, for example, anonymity, hierarchical naming, Lookup, and data storage. The most important feature is the efficient location of data items (Stoica et al., 2003).



A distributed hash table (DHT) is a protocol for structured peer-to-peer (p2p) networks used. Each node in the network follows a specific rule to divide and store the large file index hash table. The DHT entries should be the form: (key, value). The key is the hash value of the file, and the value is the IP address where the file is stored. It is efficient to find and return value/address from the query node.

In this distributed hash table, the joined nodes will automatically form a decentralized structure without the scheduling from the master node. In this way, there will be no system crash when the master node is failed or busy. Also, when a large number of nodes are flooded into the system, there may be a short-term load imbalance. The system will randomly allocate the same amount of storage to each node. Since the hash table can have node exit or sudden node failure, we also need to address the issue of exiting the nodes to ensure query and data stability.

## 2  RELATED WORK

### 2.1  CAN DHT

Sylvia Paul Ratnasamy et al. proposed the Content-Addressable Network (CAN DHT) (Ratnasamy et al., 2001). CAN was designed as a scalable, fault-tolerant, and self-organizing distributed hash table. It is similar to a multi-

dimensional Cartesian coordinate space, and each node in the network can be identified in the space. Each CAN node maintains a routing table that saves each neighbor's IP address and the coordinate area in the virtual space. The node can route the message to the target area in the coordinate space and organize the entire network.

## 2.2 Chord DHT

Ion Stoica et al. proposed Chord DHT (Stoica et al., 2001). Chord covers the entire network by stringing the nodes' IDs from small to large to form a ring, and each data is stored at the backward node closest to the data key. Chord uses the similar dichotomy search method, a non-linear search algorithm, and the convergence speed can be very fast. The routing complexity can be reduced to the logarithm of the number of nodes.

## 2.3 Pastry DHT

Antony Rowstron and Peter Drusche proposed Pastry DHT (Rowstron et al. 2001). Like Chord, Pastry also uses a ring network topology, concatenating each node's hash identifiers into a ring. The difference from Chord is that Pastry introduces the idea of grading, and it does not directly use node hash values to construct a ring but only take the first 'a' bit of the node's hash value and the back 'b' bit as the leaf node of a node on the ring. The advantage of this is that it greatly reduces the ring's space size and can reduce the routing time. Another difference is that Pastry's data is stored on the node closest to the key-value, not just in the backward node.
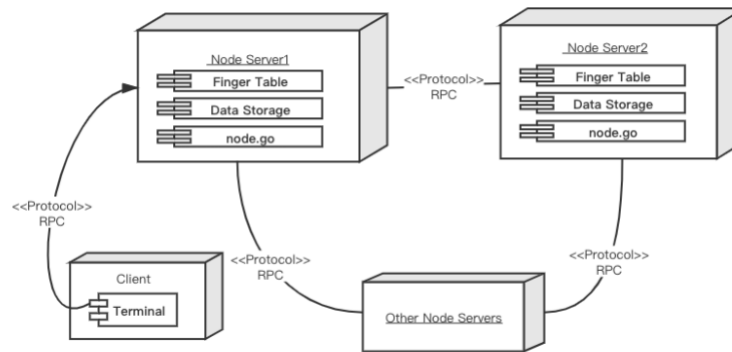
## 2.4 Tapestry DHT

Ben Y. Zhao et al. proposed Tapestry DHT (Zhao et al., 2004). In the Tapestry network, each node will assign itself a globally unique node ID, and at the same time, there is a unique root node ID in the entire system. All nodes form a tree topology in the network and data is saved on the node closest to the key value of the data.
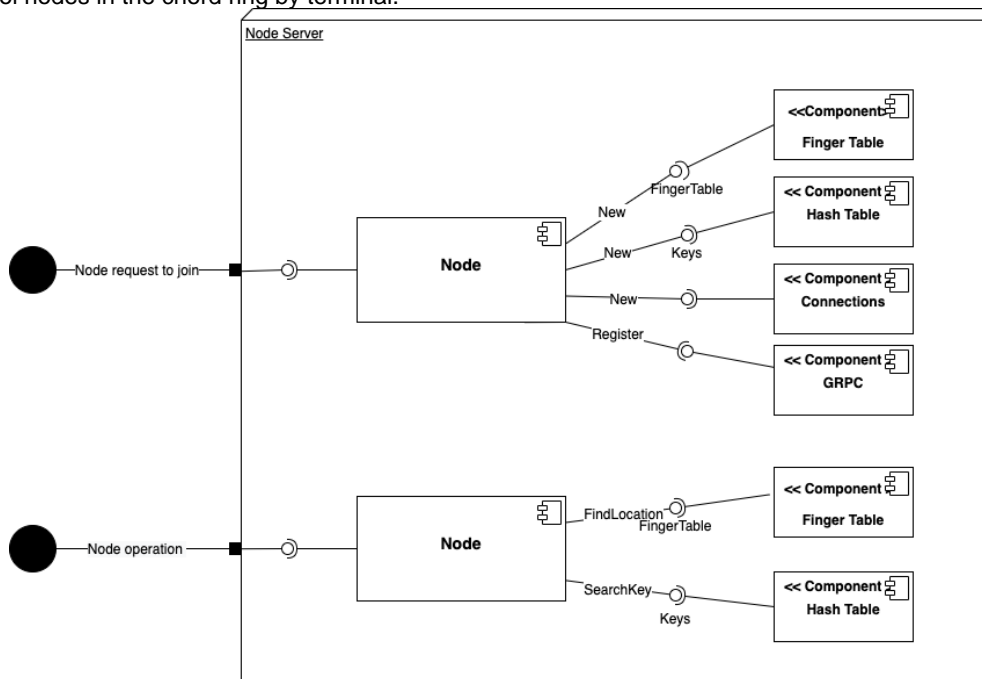
## 3  PROJECT DESIGN

### 3.1  CAN DHT

The main idea of DHT are as follows:

1. Each file index is represented as a (K, V) pair. K is the keyword, which can be the hash value of the file name (or other information of the file), and V is the IP address of the node that stores the file (or other information of the node).
2. All the file index entries (i.e. (K, V) pairs) form a large hash table. We can find the node address that stores a file as long as we input the K value.
3. The large file hash table will be divided into many local small blocks. The local hash tables are distributed to all nodes in the system according to specific rules so that each node is responsible for maintaining one of the blocks.
4. When a node queries a file, it only needs to route the query message to the corresponding node, where the hash table block maintained the (K, V) pair to be searched.

The deployment diagram is shown in this graph. Each node is a node as a container that it holds artifacts: finger table, data storage which is keys and node.go which include the basic function of the node. RPC is a communication protocol between each instance. The user can do some actions which show in use case diagram to control nodes in the chord ring by terminal.
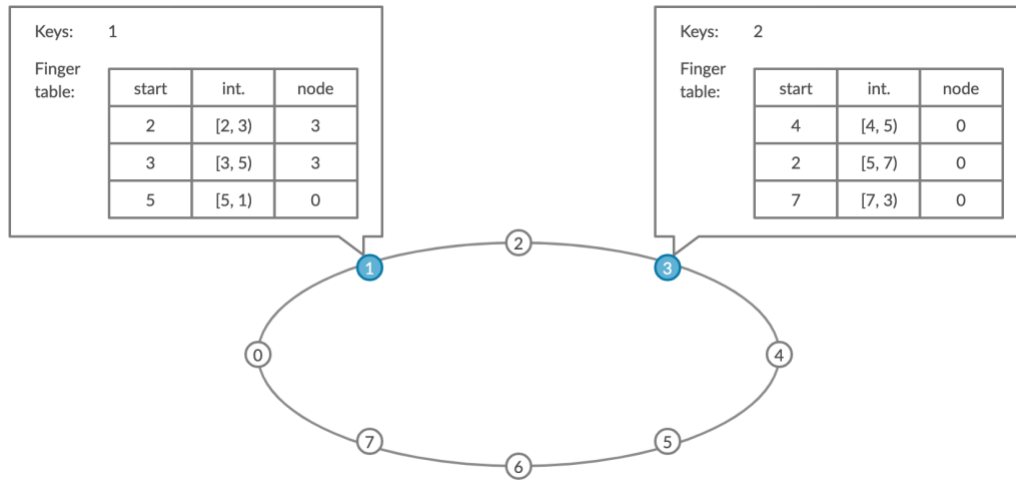


This is a component diagram. When the node starts and joins the chord, the node sets the preset parameters locally and then enters the node through a port call, then creates a new finger table, a new local hash table, a connection and a GRPC in the node. when the node performs an operation, the node is called through the port and then first finds the precedence and successor based on the finger table node, and then find the corresponding hash value based on the node location.

### 3.2 Architecture

We will use the architecture similar to Chord DHT proposed by Ion Stoica et al. in 2003 (Stoica et al., 2003). Chord covers the entire network by stringing the nodes' IDs from small to large to form a ring, and each data is stored at the backward node closest to the data key. Chord uses the similar dichotomy search method, a non-linear search algorithm, and the convergence speed can be very fast. The routing complexity can be reduced to the logarithm of the number of nodes.

The structure of a basic Chord-based system is shown as follows:

Keys: 1

Finger table:

| start | int. | node |
|---|---|---|
| 2 | [2, 3) | 3 |
| 3 | [3, 5) | 3 |
| 5 | [5, 1) | 0 |

Keys: 2

Finger table:

| start | int. | node |
|---|---|---|
| 4 | [4, 5) | 0 |
| 2 | [5, 7) | 0 |
| 7 | [7, 3) | 0 |

1. **NID (Node Identifier):** an m-bits number representing a physical machine. m must be large enough to make the probability of sharing the same NID be negligible.
2. **KID (Key Identifier):** an m-bits number representing a key (or resource). Key is hash-bound to a resource so that the KID can be obtained by the Key through the hash operation. m must be large enough to make the probability of sharing the same NID be negligible.
3. **Constant hash function:** Compared with the general hash function, the joining and leaving of nodes have the least impact on the entire system using the constant hash function. SHA-1 is used in Chord for constant hash calculation.
4. **Chord Ring:** NID and KID are allocated on a ring with the size of $2^m$, for key distribution, node distribution, and key location. The key is allocated to the node with NID>=KID. This node becomes the successor node of k, the first node in the clockwise direction from k on the ring, denoted as successor(k). The nodes are placed on the ring from large to small in the clockwise direction.

### 3.3  Core Functions Design



The function design is shown in this graph. In our system, there are three major parts: node management, keys management and node stabilization, the core functions are: find keys (get key location), add new node and node stabilization.

1. Find Keys (Key Location)
    a. The key location is the core function of the Chord system. When node N is looking for a key with KID=id, it will ask whether the key is on the next node. If the KID is between the NID of the node and the NID of the next node, then the key is allocated to the next node. Otherwise, the same query will be asked on the next node, recursively. Finally, the location of the key will be found.
    b. The query time complexity is O(N), where N is the number of nodes in the ring. To accelerate the searching process, Chord uses finger tables to maintain up to m entries on each node, where m is the number of bits in the hash key.
2. Add New Node

    When there is a new node added into the system, we must maintain three invariants:
    a. The successors of every node should point to their successors correctly.
    b. Every key should be correctly stored in the successor(k).
    c. The finger tables on every node should be updated correctly.
    d. Therefore,  there are some tasks that need to be done: the first is Initialize node n with its predecessor and the finger table. And then the new node notifies other nodes to update their predecessors and finger tables. Last, the new node takes over its responsible keys from its successor.
3. Node Stabilization:

    To ensure the reliability of our system, all nodes need to automatically update their status, therefore, this stabilization function is running periodically in the background. In our design there are some functions to maintain node stabilization:
    a. Stabilize(): In the case p recently joined the system,  n will ask its successor for its predecessor p and decides whether p should be n's successor instead

5

b. Inform(): which notifies node n's successor of its existence, and change its predecessor to n
c. findNextFinger(): which updates finger tables
d. checkPreNode(): Periodically checks if its' predecessor is alive

## 4  DISTRIBUTED SYSTEMS CHALLENGES

There are several challenges in implementing a distributed hash table:

### 4.1  Fault Tolerance

We have considered fault-tolerant solutions during the project implementation. We create a corresponding channel for each node when it starts, and when the node exits normally or unexpectedly, it will broadcast the message of node exit to other nodes through the channel, and the other nodes will remove the existing node from the finger table in the next data operation. And each node will also modify its own predecessor and successor nodes according to the nodes before and after.

But we didn't consider the nodes starting at the same time when we started them. So we need to wait for one node to finish starting before we can start the next node when it starts. But in actual operation, a large distributed hash table cannot wait for other nodes to finish starting.

### 4.2  Scalability

Our design has a good scalability. When a new node is added, the node is still constantly checking the node before and after, updating the routing table if it is a new node, and keeping it unchanged if it is the original node. Besides, other nodes need to be notified of the addition of new nodes. Last, the new node will take over responsible data from its successor.

### 4.3  Heterogeneity

In our project, we used Protocol Buffers which is a technology developed by Google as our data serialized method. This technology can help to transform the data structure to a storing or transmission form and the protoBuf can be used in any platform with any language: C++, C#, Dart, Go, Java, Python…. So in our design, the data storing and transmission is compatible with different machines.

### 4.4  Openness

Our project follows a standard rule to offer services. GRPC is a communication protocol to ensure connections between nodes. GRPC uses Protocol Buffers to encode the data, in this way, this protocol is more strict. Another feature of GRPC is authentication，it provides two types of credentials; channel credentials and call credentials, therefore it ensures openness of our project.

### 4.5  Security

Our project is based on the DHT system, it has more ability against hostile attacks from outside than a centralized system. Because our data storage is distributed, each node only maintains partial data, therefore when the system is being attacked, it is difficult to get the whole data from the system.

### 4.6  Concurrency

In our project, all node operations add, delete, and query support concurrent operations, which are processed sequentially according to the order of requests, and we add control locks to each operation to prevent read and write inconsistencies.We also create channels on each node so that also our programs can be executed concurrently and the order in which the messages arrive determines the order of execution.

### 4.7 Quality of Service

Based on our local testing of the program, we found that all data and operations were delivered accurately, although the project would experience some timeouts. Based on the local quality of service, our project can perform all the functions of the same type of hash table.

### 4.8 Transparency

Our project has transparency, when any user gets the node code, they can do any operation on the hash table as long as they join the network. And there is no master node in the project, when the first node exits, our project can still run.

## 5 EXPERIMENTS AND RESULTS

The node instance will be running locally, and to identify different nodes, we assigned a different port number to each node.

### 5.1 Baseline DHT System

Before testing the functionality, we built a baseline DHT system, which consists of only three nodes, following steps shown below:

Step1: Set the environment for each node

Step2: Initialize the first node

```
(base) → project-schrodinger git:(main) ✗ source setEnv.sh
(base) → project-schrodinger git:(main) ✗ cd src/dht/node_client
s
(base) → node_clients git:(main) ✗ go run node1.go
--------------------------------------------------------------
[Log] Created a new node
        Address: 0.0.0.0:8001
        Initial Node ID: 1
        New Node ID: 30494258244493662932569936375743582007759025
9883
[Log] Created finger table for the node.
--------------------------------------------------------------
[Log] Added a new node!
My Node ID:
        30494258244493662932569936375743582007759025983
--------------------------------------------------------------
[Log] Set predecessor
        of: 30494258244493662932569936375743582007759025983
        to be: 30494258244493662932569936375743582007759025983
```

As shown above, the node got a new ID, which is the hash of the initial ID that we assigned. And it set itself to be the predecessor.

Step3: Add new nodes to the node that already in the DHT system

7

```
(base) → node_clients git:(main) x go run node2.go
────────────────────────────────────────────────────────────
[Log] Created a new node
        Address: 0.0.0.0:8002
        Initial Node ID: 4
        New Node ID: 1563801023189659902646662860181919005906589052
10
[Log] Created finger table for the node.
────────────────────────────────────────────────────────────
[Log] Added a new node!
My Node ID:
        15638010231896599026466628601819190059065890521O
────────────────────────────────────────────────────────────
[Log] Set successor
        of: 15638010231896599026466628601819190059065890521O
        to be: 304942582444493662932569936375743582OO77590259883
────────────────────────────────────────────────────────────
[Log] Set successor
        of: 15638010231896599026466628601819190059065890521O
        to be: 304942582444493662932569936375743582OO77590259883
────────────────────────────────────────────────────────────
[Log] Set predecessor
        of: 304942582444493662932569936375743582OO77590259883
        to be: 15638010231896599026466628601819190059065890521O
```

```
(base) → node_clients git:(main) x go run node3.go
────────────────────────────────────────────────────────────
[Log] Created a new node
        Address: 0.0.0.0:8003
        Initial Node ID: 8
        New Node ID: 14521739854087502033184751171896360629112140
2143
[Log] Created finger table for the node.
────────────────────────────────────────────────────────────
[Log] Added a new node!
My Node ID:
        14521739854087502033184751171896360629112140421
────────────────────────────────────────────────────────────
[Log] Received a new (key, value) pair:
        (27, "Hello, this is node2!")
My Node ID:
        14521739854087502033184751171896360629112140421
────────────────────────────────────────────────────────────
[Log] Set predecessor
        of: 304942582444493662932569936375743582OO77590259883
        to be: 14521739854087502033184751171896360629112140421
```

As shown above, when node 2 and 3 are initialized, the connections were updated on all the nodes.

Step4: Add keys to the system



We let the second node automatically add a key-value pair to the system every second. And we can find out which node finally received the key by checking the received log. As shown above, all of the three nodes received some keys.

## 5.2 Functionality Testing

5.3 **We implemented three interactive nodes, in which we can execute some commands manually, such as add the key, get the value of the key, find the location of a key, delete a key, or print out the finger table.**

5.4 **Print out Finger Table and Node ID**



**We initialized three interactive nodes, and then printed out the finger table and the Node ID on each node. As shown above, the finger tables are updated correctly.**

### b. Add keys



**After a key is added, it will be assigned and stored in a node on the ring.**

```
[Log] Received a new (key, value) pair:
        (1, "2")
My Node ID:
        44002260241957126155656285565789015079246304075l
```

### c. Get the value of a given key

**We can get the value of the key by checking the node that stored the key.**

```
get 3
--------------------------------------------------------------------------
[Command] Get the value of a key: 3
        HashKey: 68432980133622366135695254607826988903893870277l
        From: 7040345787956153021819957052950481364180613870l
        Value: 4
Output: 4
```

### d. Get the location of a given key

**We can also get the location of the key by calculating the hash of the key.**

```
Input: getloc 3
--------------------------------------------------------------------------
[Command] Get the location of a key: 3
        HashKey: 68432980133622366135695254607826988903893870277l
        Location: 7040345787956153021819957052950481364180613870l
Output: 7040345787956153021819957052950481364180613870l
```

### e. Delete a key

**The node that stored the key will delete the key.**

```
Input: del 3
[Command] Delete a key: 3
        HashKey: 68432980133622366135695254607826988903893870277l
        From: 7040345787956153021819957052950481364180613870l
        State: Succeed
```

**After the key was deleted, we could not find the value nor location of the key.**

```
Input: getloc 3
--------------------------------------------------------------------------
[Command] Get the location of a key: 3
        HashKey: 68432980133622366135695254607826988903893870277l
        Location: Key not found
Output: Key not found
```

## 3. Fault Tolerance Testing

To test the fault tolerance of the system, we could simply shut down a node and to test whether the system still works, and reconnect the node again and test the performance.

### Step1: Shut down a node

```
Input: get 3
----------------------------------------------------------------
[Command] Get the value of a key: 3
        HashKey: 6843298013362236613569525460782698890389387027779
        From: 90439188749661471240920510780665347945322495176
        Value: Key not found
Output: Key not found
----------------------------------------------------------------
```

After a node is shut down, we cannot find the value of the key that is stored on the node.

```
[Command] Add a new key:
        (5, "6")
        HashKey: 98311657783177760831276567051553810276466389276
        From: 90439188749661471240920510780665347945322495176
        To: 440022602419571261556562855657890150792463040751
----------------------------------------------------------------
Input: get 5
----------------------------------------------------------------
[Command] Get the value of a key: 5
        HashKey: 98311657783177760831276567051553810276466389276
        From: 440022602419571261556562855657890150792463040751
        Value: 6
Output: 6
----------------------------------------------------------------
Input: getloc 5
----------------------------------------------------------------
[Command] Get the location of a key: 5
        HashKey: 98311657783177760831276567051553810276466389276
        Location: 440022602419571261556562855657890150792463040751
Output: 440022602419571261556562855657890150792463040751
```

Nonetheless, the rest part of the system could still work as normal. As shown above, we can add and find a key on the node that is still working.

### Step2: Reconnect a node

11

**Then we tried to reconnect the node to the system. As shown below, after being added, the node could work as normal.**

```
(base) → node_clients git:(main) ✗ go run interactive_node3.go
————————————————————————————————————————————————————————————
[Log] Created a new node
        Address: 0.0.0.0:8006
        Initial Node ID: 24
        New Node ID: 440022602419571261556562855657890150792463040751
[Log] Created finger table for the node.
————————————————————————————————————————————————————————————
[Log] Added a new node!
My Node ID:
        440022602419571261556562855657890150792463040751
————————————————————————————————————————————————————————————
[Log] Set successor
        of: 440022602419571261556562855657890150792463040751
        to be: 704034578795615302181995705295048136418061387092
————————————————————————————————————————————————————————————
[Log] Set predecessor
        of: 90439188749661471240920510780665347945322495417 6
        to be: 440022602419571261556562855657890150792463040751
————————————————————————————————————————————————————————————
Input: ——————————————————————————————————————————————————————
[Log] Received a new (key, value) pair:
        (1, "2")
My Node ID:
        440022602419571261556562855657890150792463040751
————————————————————————————————————————————————————————————
[Log] Received a new (key, value) pair:
        (3, "4")
My Node ID:
        440022602419571261556562855657890150792463040751
^CErrors occurred and stop:  <nil> <nil>
(base) → node_clients git:(main) ✗ go run interactive_node3.go
————————————————————————————————————————————————————————————
[Log] Created a new node
        Address: 0.0.0.0:8006
        Initial Node ID: 24
        New Node ID: 440022602419571261556562855657890150792463040751
[Log] Created finger table for the node.
————————————————————————————————————————————————————————————
[Log] Added a new node!
My Node ID:
        440022602419571261556562855657890150792463040751
————————————————————————————————————————————————————————————
[Log] Set successor
        of: 440022602419571261556562855657890150792463040751
        to be: 704034578795615302181995705295048136418061387092
```

12

**The finger tables and Node IDs are shown as follows, and we could find out that the nodes were reconnected correctly.**



1. **Task 1: Basic query**

   This task includes some basic functions to realize the query function like the "get" of the HashMap in java.

2. **Task 2: Finger Table**

   This task is an implementation of the figure table, which includes designing a new data struct for figure table and basic functions for it.

3. **Task 3: Add New Node**

   This task is meant for implementing a new data struct for Nodes and the related changes when a new node is added into the chord.

4. **Task 4: System stability**

   - To make sure that each query is correct, update the successor and predecessor of each node.
   - Consider more factors that may influence the system like memory usage, CPU and network and so on. Try to enhance the robustness of the system and shorten the response time.
   - Tests and debug.

The basic schedule of these tasks is shown in the figure below:



**5.5 Process Checking**

We plan to have a regular meeting every Friday night to check the processing. But we will track and manage our work every day. And We design a table for recording the weekly task. Here is a simple template.

|  | Task | Work Log |
|---|---|---|
| Week 1 | Task 1 |  |

13

| | | |
|---|---|---|
| | Task 2 | |
| Week 2 | Task 2 | |
| | Task 3 | |
| Week 3 | Taks 3 | |
| Week 4 | Task 4 | |

## 5.6  Issues Recording

We also create an issue table for recording our problem status. The table below is a simple prototype.

| Issue ID | Date | Task | Editor | Problem | Solved or not |
|---|---|---|---|---|---|
| 1 | 11/09 | Task 1 | Guangyuan Shen | Cannot build the findID function | yes |
| 2 | 11/10 | Task 1 | Xiaoyu Shen | Solve Issue 1 | - |
| 3 | 11/12 | Task 2 | Zetian Zheng | Some test samples cannot pass | no |
| 4 | 11/14 | Task 1 | XuZheng Lu | Update Issue 2 | - |

**REFERENCES**

[1]  Kaija, Kasperi. "The implementation and performance of Chord." (2018).

[2]  Stoica, Ion, et al. "Chord: a scalable peer-to-peer lookup protocol for internet applications." IEEE/ACM Transactions on networking 11.1 (2003): 17-32.

[3]  Karger, David R., and Matthias Ruhl. "Simple efficient load balancing algorithms for peer-to-peer systems." Proceedings of the sixteenth annual ACM symposium on Parallelism in algorithms and architectures. 2004.

[4]  Naor, Moni, and Udi Wieder. "A simple fault tolerant distributed hash table." International Workshop on Peer-to-Peer Systems. Springer, Berlin, Heidelberg, 2003.

[5]  https://medium.com/techlog/chord-building-a-dht-distributed-hash-table-in-golang-67c3ce17417b

[6]  https://www.cnblogs.com/gnuhpc/archive/2012/01/13/2321476.html