



DISTRIBUTED

HASH

TABLE



DHT ?

THIS PART INTRODUCE THE THEORY AND IMPLEMENTATION OF
DISTRIBUTED HASH TABLE



A DISTRIBUTED HASH TABLE IS A CLASS OF
A DECENTRALIZED DISTRIBUTED SYSTEM
THAT PROVIDES A LOOKUP SERVICE SIMILAR
TO A HASH TABLE: (KEY, VALUE)

--WIKIPEDIA

Three hand-drawn white arrows on a black background. One arrow points diagonally down towards the top of the word 'DECENTRALIZED'. A second arrow starts from the left, curves around the bottom of the word, and points towards the letter 'D'. A third arrow starts from the right, curves around the top of the word, and points towards the letter 'D'.

DECENTRALIZED

Decentralized is a vital character of Distributed Hash Table



WHY AND HOW TO

Why Decentralized

- X Too much data & high expense
- X Convenient
- X Anonymous
- X Security

How to become Decentralized

- X Smart approach to organize computers in network
- X Efficient hash function
- X Efficient strategy to look up values via key

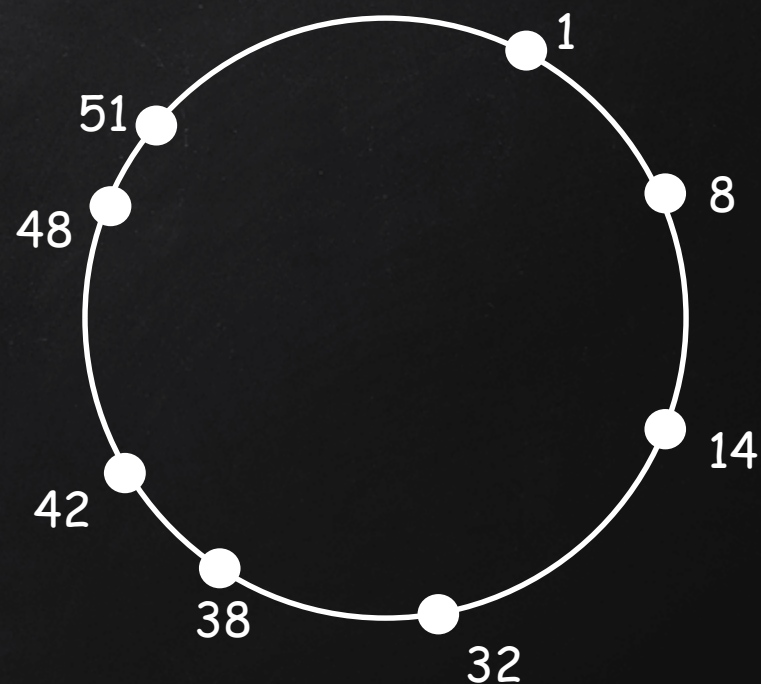
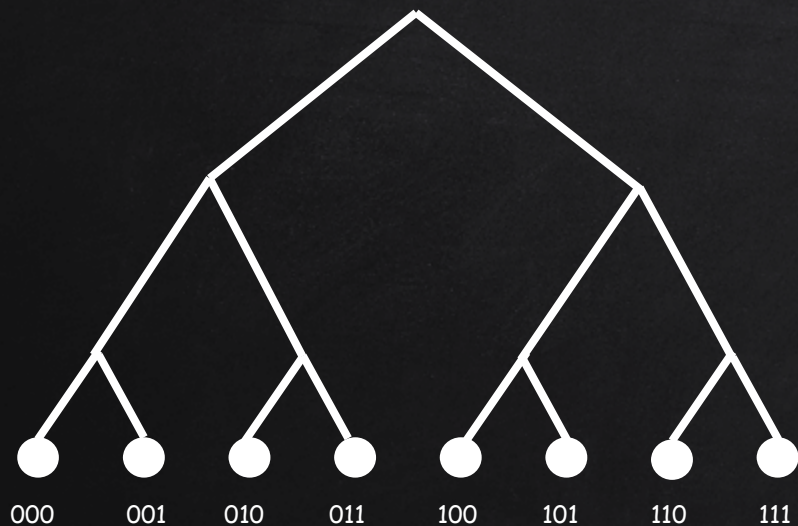


IMPLEMENTATION

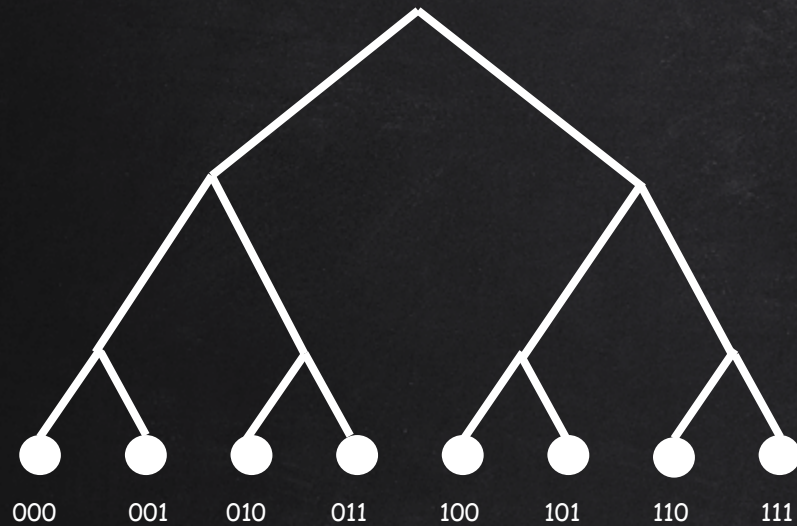
KADEMLIA / CHORD / PASTRY / CAN / KOORDES...

Kademlia – The most widely used

Chord – What we implemented in PPCA



Kademlia - The most widely used

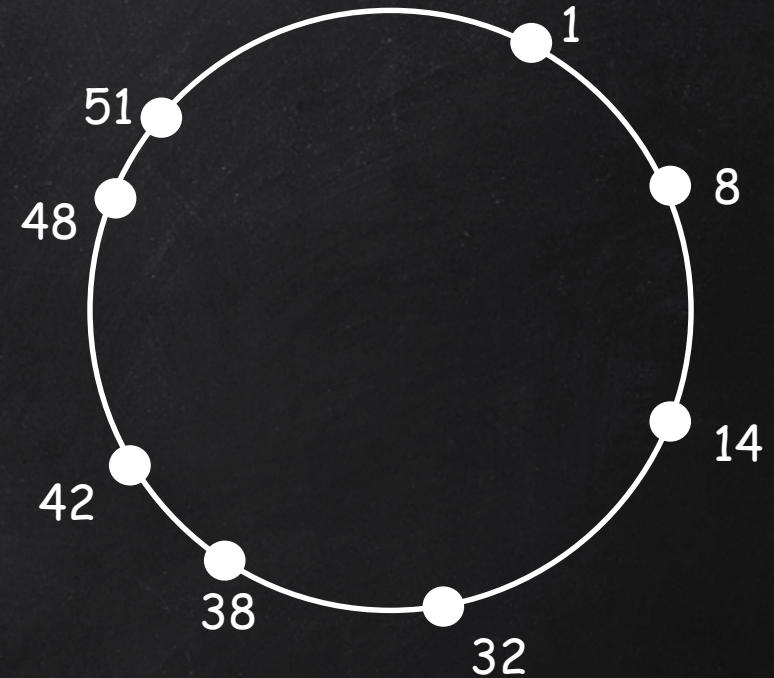


Key := "Fan Zhou"

Node ID := "114.51.41.91:810"

SHA-1

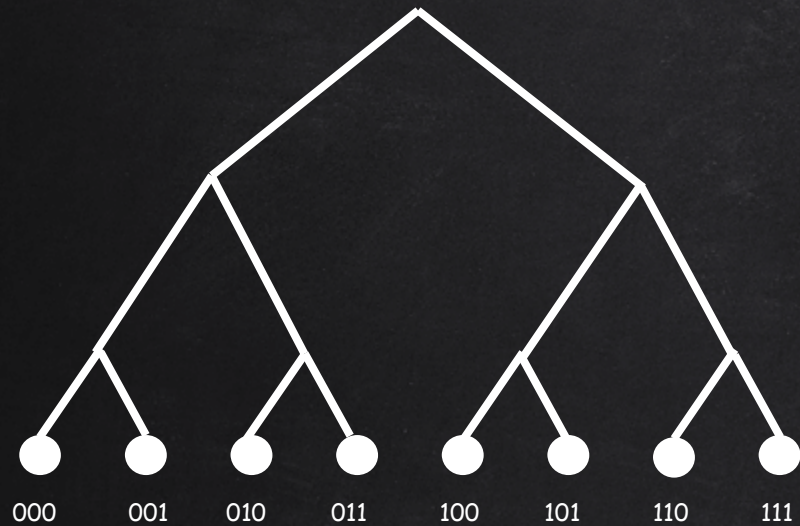
Chord - What we implemented in PPCA



Hashed Id (hex):

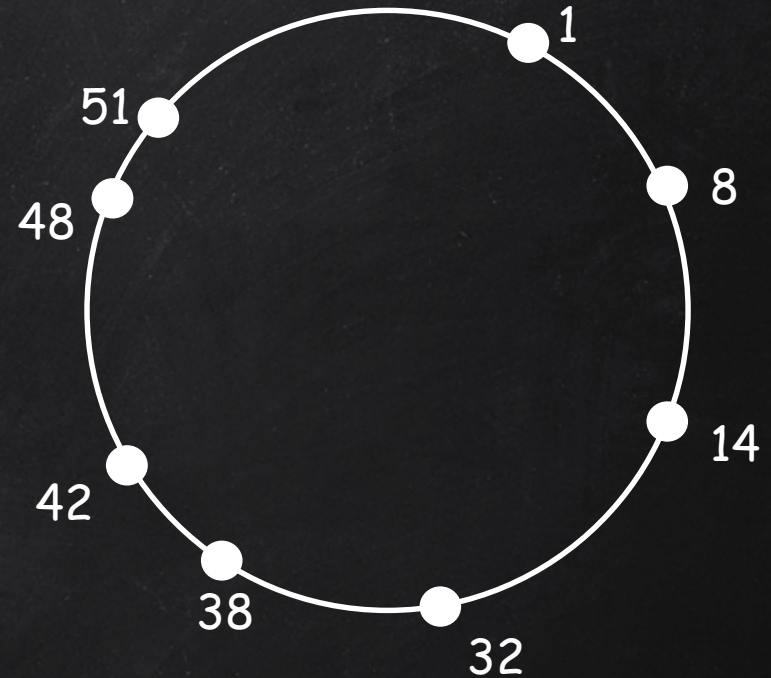
2486a1e52703c5b049d1
61c0cec2137a9ca97f73
965c4e5d7803f1f3f65c
bea68888d86e84ca5950

Kademlia - The most widely used



$\text{dis}(a, b) :=$
 $a \text{ xor } b == b \text{ xor } a$

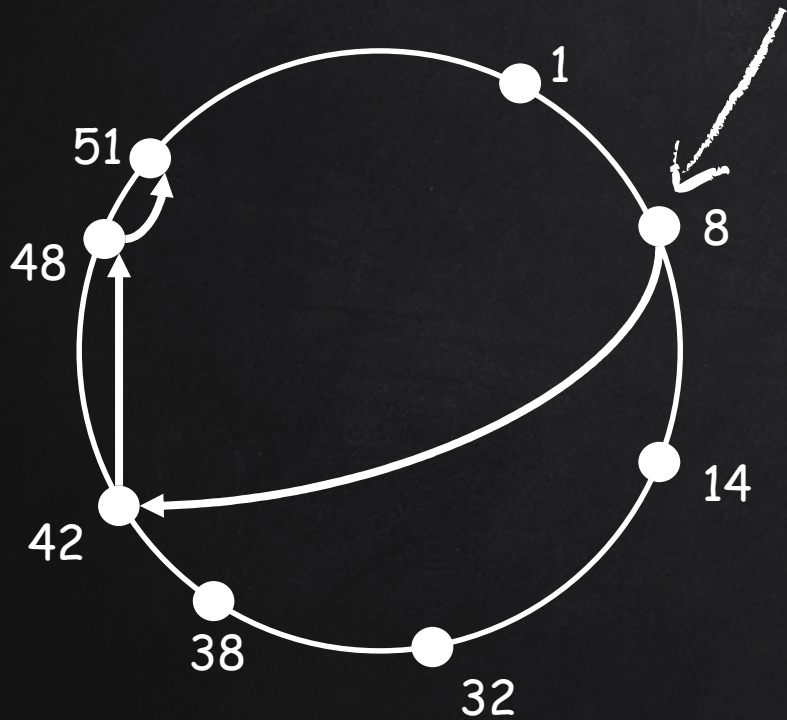
Chord - What we implemented in PPCA



$\text{dis}(a, b) :=$
 $(\text{int}(b) - \text{int}(a) + 2^{160}) \bmod 2^{160}$

The $\langle \text{KEY}, \text{VALUE} \rangle$ will be given to the closest node

CHORD - LOOK UP

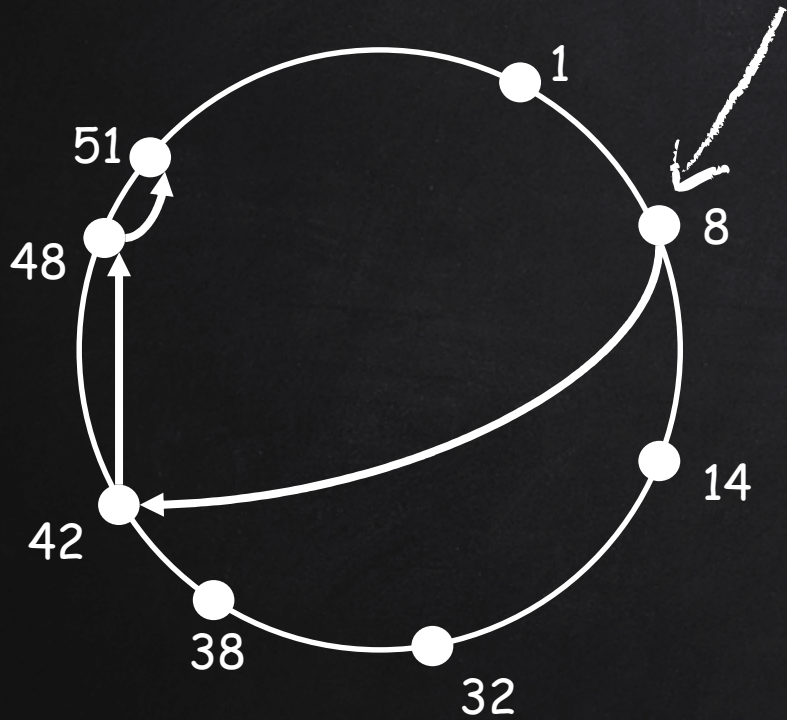


ID	Node
$8 + 2^0$	#14
$8 + 2^1$	#14
$8 + 2^2$	#14
$8 + 2^3$	#32
$8 + 2^4$	#32
$8 + 2^5$	#42

Steps of look up:

- X 1. Node #8 want to look up who has key 49.
- X 2. Node #8 check whether his successor 14 is the nearest node of 49, obviously not.
- X 3. Node #8 iterate its finger table from bottom to top for the available closest preceding node of 49 he knows, then he find Node #42.
- X 4. Node #8 ask node Node #42 to repeat these steps until the closest node of 49, Node #51 is found.

CHORD - LOOK UP



ID	Node
$8 + 2^0$	#14
$8 + 2^1$	#14
$8 + 2^2$	#14
$8 + 2^3$	#32
$8 + 2^4$	#32
$8 + 2^5$	#42

$O(\log N)$

CHORD - LOOK UP

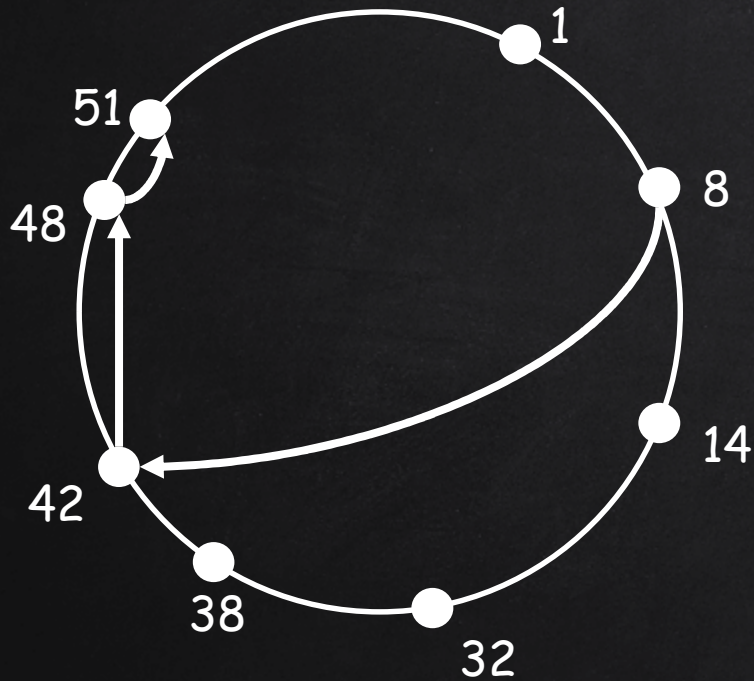
Available operation for any computer in the network:

X 1. Put $\langle \text{Key}, \text{Value} \rangle$

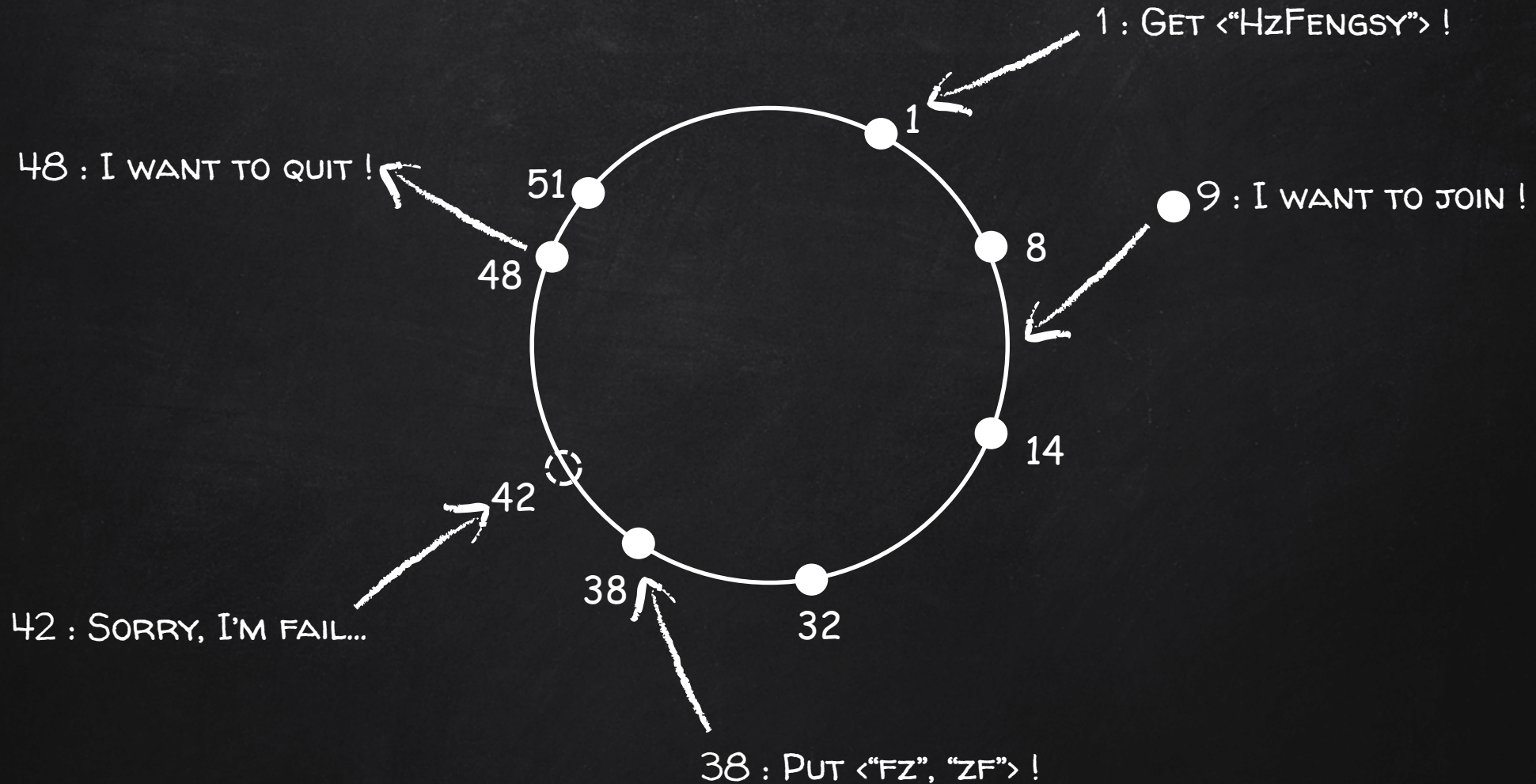
X 2. Look up $\langle \text{Key} \rangle$ for $\langle \text{Value} \rangle$

X 3. Remove $\langle \text{Key}, \text{Value} \rangle$ from network

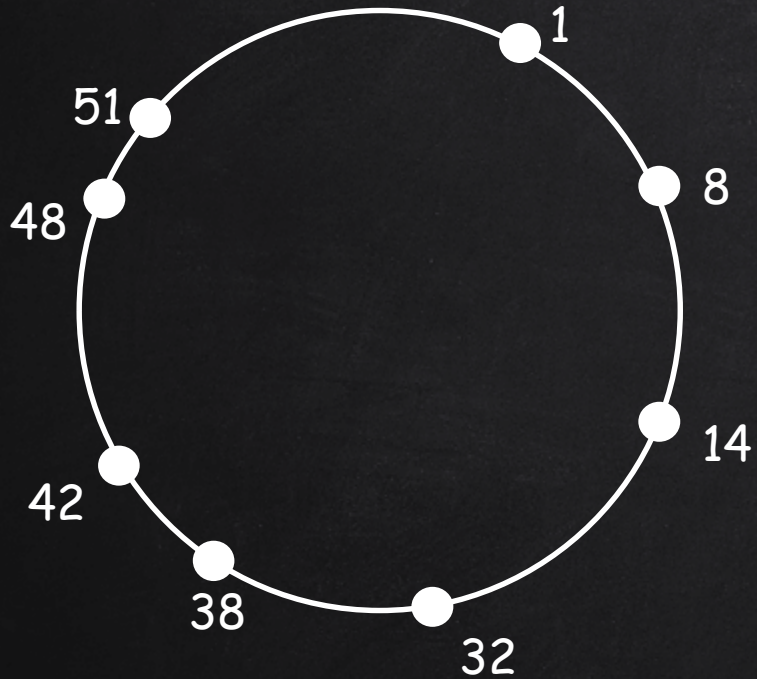
X 4. Modify $\langle \text{Key}, \text{Value} \rangle$ to $\langle \text{Key}, \text{New Value} \rangle$



CHORD - ROBUSTNESS



CHORD - ROBUSTNESS



SACRIFICE TIME
FOR
STABILITY AND CORRECTNESS

CHORD – ROBUSTNESS

Strategies to maintaining the network:

- X 1. Maintain a successor list via copying successor's SUCCESSION LIST, removing the last entry of it and adding the successor as the first entry of the list.
- X 2. Check whether successor changed and notify the new successor about himself if changed periodically.
- X 3. Tell the successor about himself periodically so that the successor can update his predecessor
- X 4. Fix one finger table entry per fix interval
- X 5. Notify the predecessor and successor when a node want to quit
- X 6. Check whether predecessor is failed periodically so that it can be updated to new available predecessor

CHORD – ROBUSTNESS

Strategies to maintaining the network:

- X 1. Maintain a successor list via copying successor's SUCCESSOR LIST, removing the last entry of it and adding the successor as the first entry of the list.
- X 2. Check whether successor changed and notify the new successor about himself if changed periodically.
- X 3. Tell the successor about himself periodically so that the successor can update his predecessor
- X 4. Fix one finger table entry per fix interval
- X 5. Notify the predecessor and successor when a node want to quit
- X 6. Check whether predecessor is failed periodically so that it can be updated to new available predecessor

For 1,2 and 3 we use stabilize() and notify()

n.stabilize()

x := n.succ.pre

if x between n and n.succ

n.succ = x

n.succList = n.succ.succList

n.succList.pop_back()

n.succList.push_front(n.succ)

n.succ.notify(n)

n.notify(Node n')

if n.pre == nil or n' is between n.pre and n

n.pre = n'

CHORD – ROBUSTNESS

Strategies to maintaining the network:

- X 1. Maintain a successor list via copying successor's SUCCESSOR LIST, removing the last entry of it and adding the successor as the first entry of the list.
- X 2. Check whether successor changed and notify the new successor about himself if changed periodically.
- X 3. Tell the successor about himself periodically so that the successor can update his predecessor
- X 4. Fix one finger table entry per fix interval
- X 5. Notify the predecessor and successor when a node want to quit
- X 6. Check whether predecessor is failed periodically so that it can be updated to new available predecessor

For 4 we use function fix_finger()

```
n.fix_finger()
```

```
n.finger[next] = n.find_succ(n.id + 2next)  
next = next + 1
```

```
if next > 159
```

```
    next = 0
```

//n.finger is a array with length of 160, and here we use C-Style indexing

//find_succ() is a function that find which node the specific key belonged to.

CHORD – ROBUSTNESS

Strategies to maintaining the network:

- X 1. Maintain a successor list via copying successor's SUCCESSOR LIST, removing the last entry of it and adding the successor as the first entry of the list.
- X 2. Check whether successor changed and notify the new successor about himself if changed periodically.
- X 3. Tell the successor about himself periodically so that the successor can update his predecessor
- X 4. Fix one finger table entry per fix interval
- X 5. Notify the predecessor and successor when a node want to quit
- X 6. Check whether predecessor is failed periodically so that it can be updated to new available predecessor

For 5 we use quit()
n.quit()

```
n.pre.succList.remove(n)
```

```
last := len(n.succList)
```

```
n.pre.succList.append(n.succList[len - 1])
```

```
n.succ.pre = n.pre
```

// data also should be transferred here, this will be discuss in next part

CHORD – ROBUSTNESS

Strategies to maintaining the network:

- X 1. Maintain a successor list via copying successor's SUCCESSOR LIST, removing the last entry of it and adding the successor as the first entry of the list.
- X 2. Check whether successor changed and notify the new successor about himself if changed periodically.
- X 3. Tell the successor about himself periodically so that the successor can update his predecessor
- X 4. Fix one finger table entry per fix interval
- X 5. Notify the predecessor and successor when a node want to quit
- X 6. Check whether predecessor is failed periodically so that it can be updated to new available predecessor

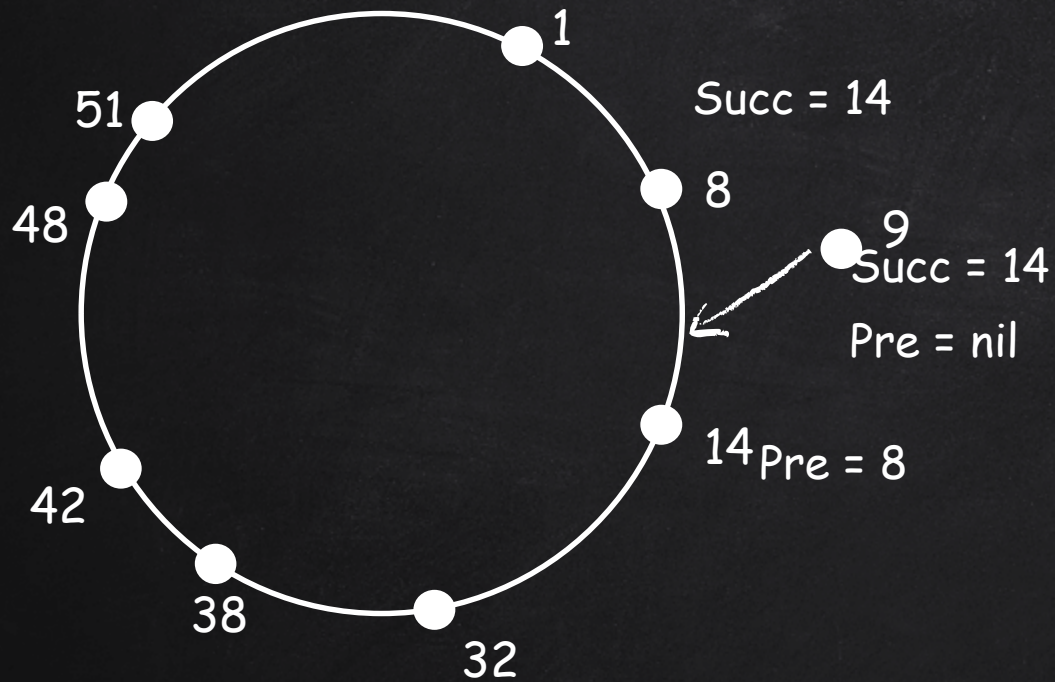
For 6 we use check_predecessor()
n.check_predecessor()

if n.pre is failed

n.pre = nil

//we can use function like ping to check whether the predecessor is failed

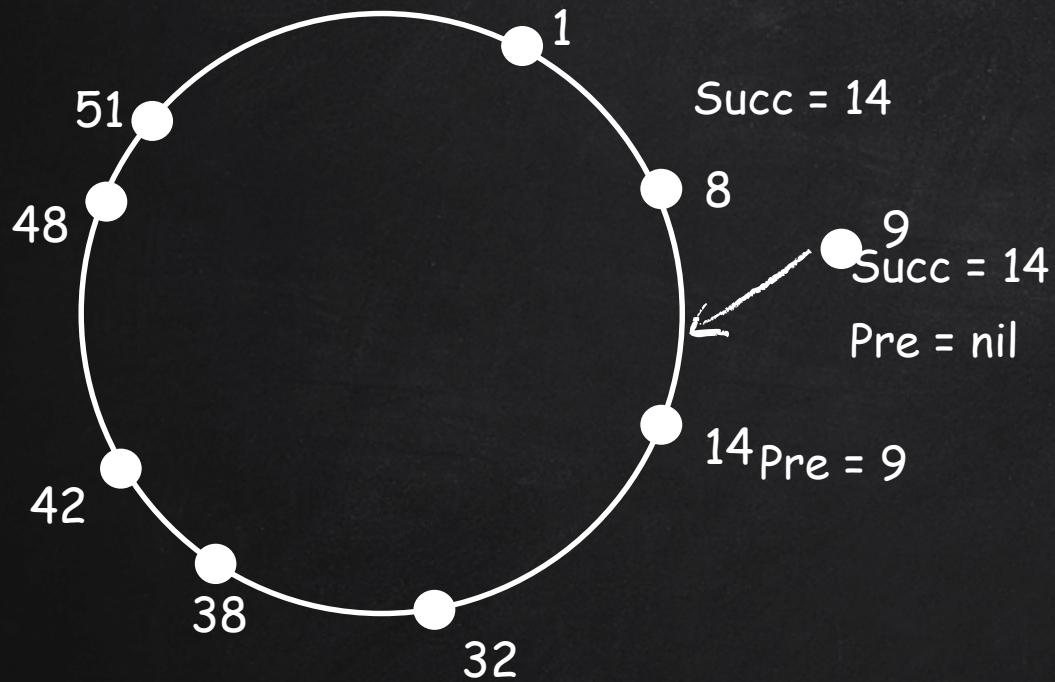
CHORD - ROBUSTNESS - SIMULATION



Condition: a node with id 9 join the network via node #51.

X node #9 call node #51's find_succ() and it get #14 as successor

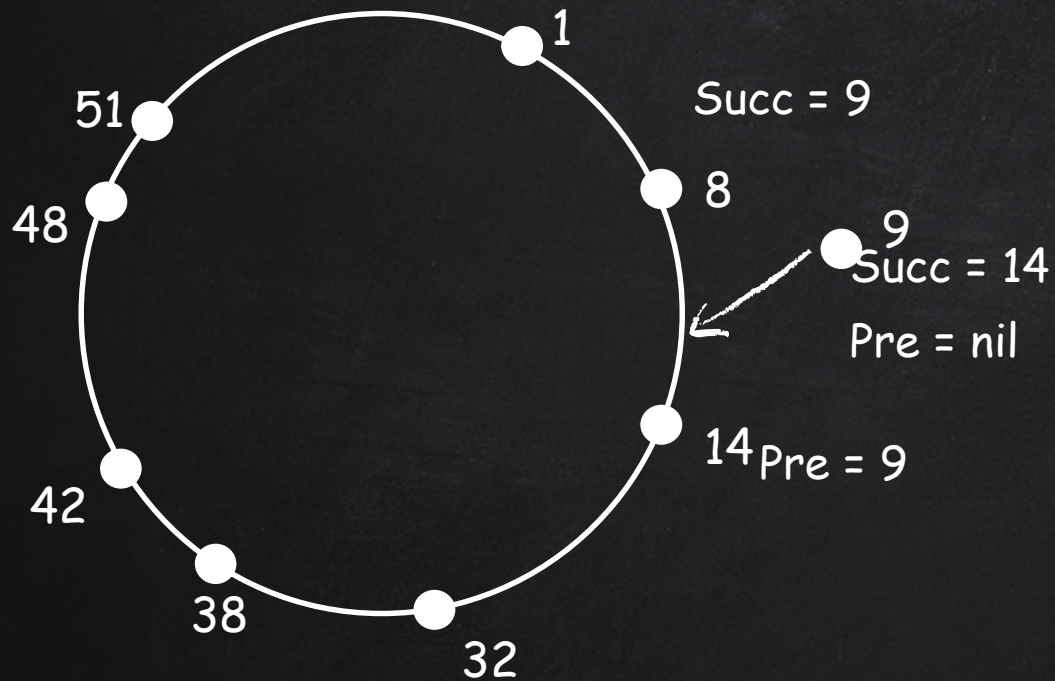
CHORD - ROBUSTNESS - SIMULATION



Condition: a node with id 9 join the network via node #51.

- X node #9 call node #51's find_succ() and it get #14 as successor
- X node #9 runs notify() in stabilize() and node #14 update his predecessor to 9

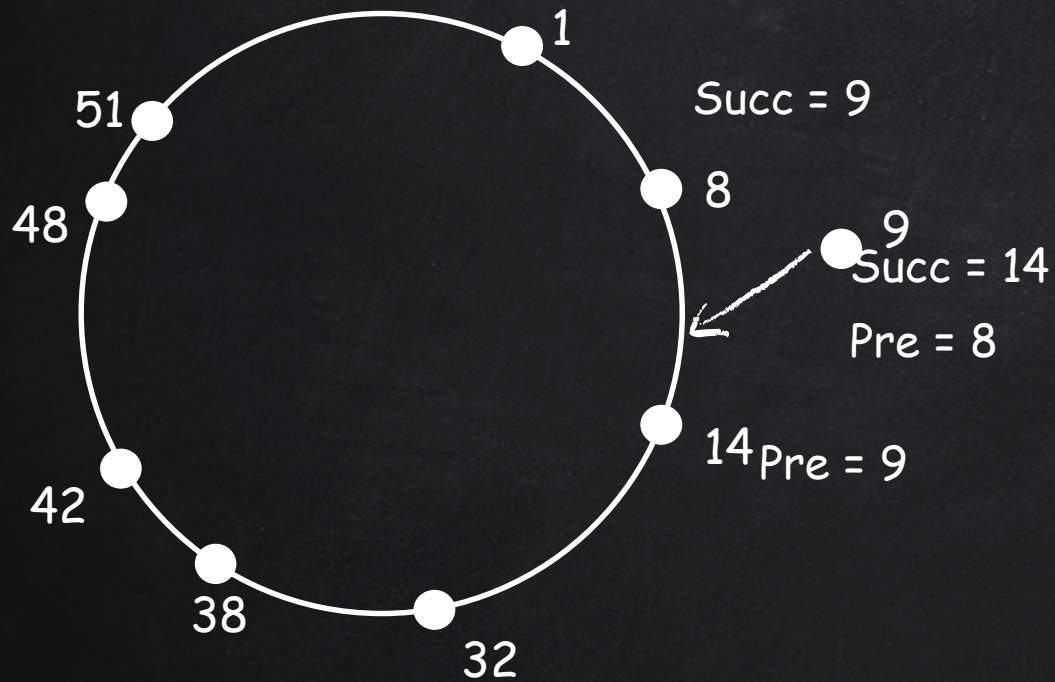
CHORD - ROBUSTNESS - SIMULATION



Condition: a node with id 9 join the network via node #51.

- X node #9 call node #51's find_succ() and it get #14 as successor
- X node #9 runs notify() in stabilize() and node #14 update his predecessor to 9
- X node #8 runs stabilize() and update his succ to 9

CHORD - ROBUSTNESS - SIMULATION



Condition: a node with id 9 join the network via node #51.

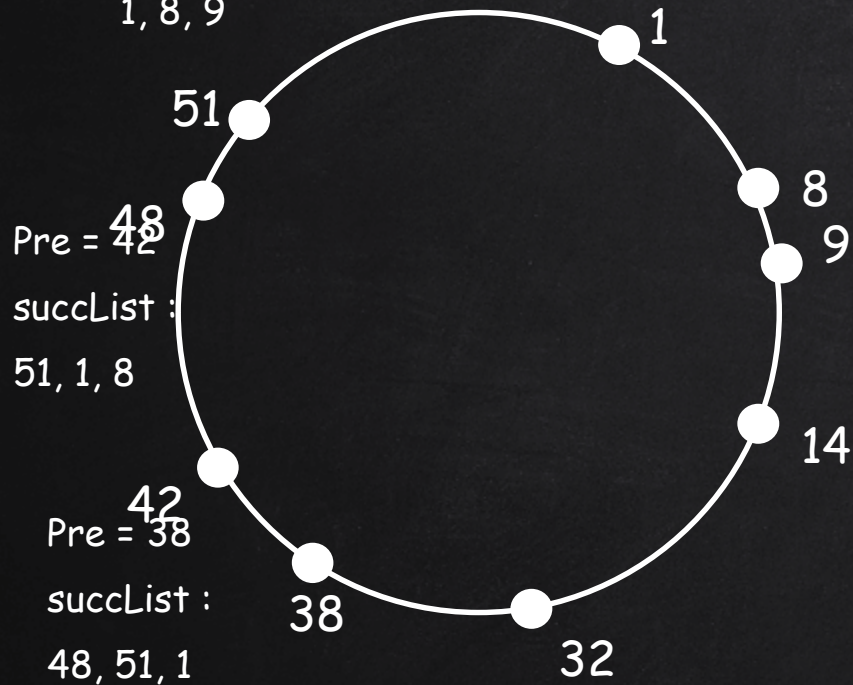
- X node #9 call node #51's find_succ() and it get #14 as successor
- X node #9 runs notify() in stabilize() and node #14 update his predecessor to 9
- X node #8 runs stabilize() and update his succ to 9
- X node #8 runs notify() in stabilize() and node #9 update his pre to 8
- X As the time passed, finger table of each node will be fixed gradually by fix_finger() function, then the network become stable again

CHORD - ROBUSTNESS - SIMULATION

Pre = 48

succList :

1, 8, 9



Condition: node #48 quit

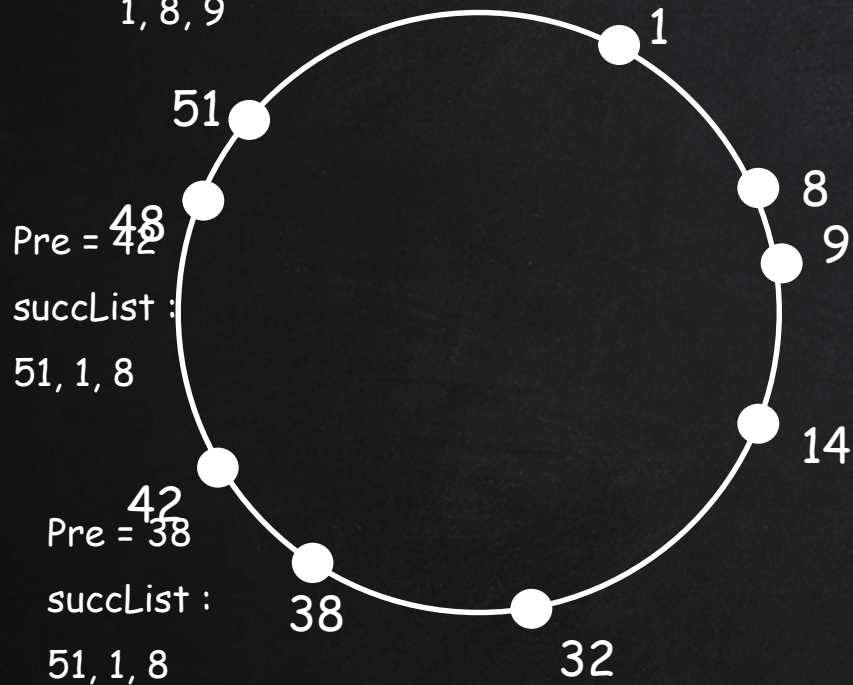
X node #48 runs function quit()

CHORD - ROBUSTNESS - SIMULATION

Pre = 48

succList :

1, 8, 9

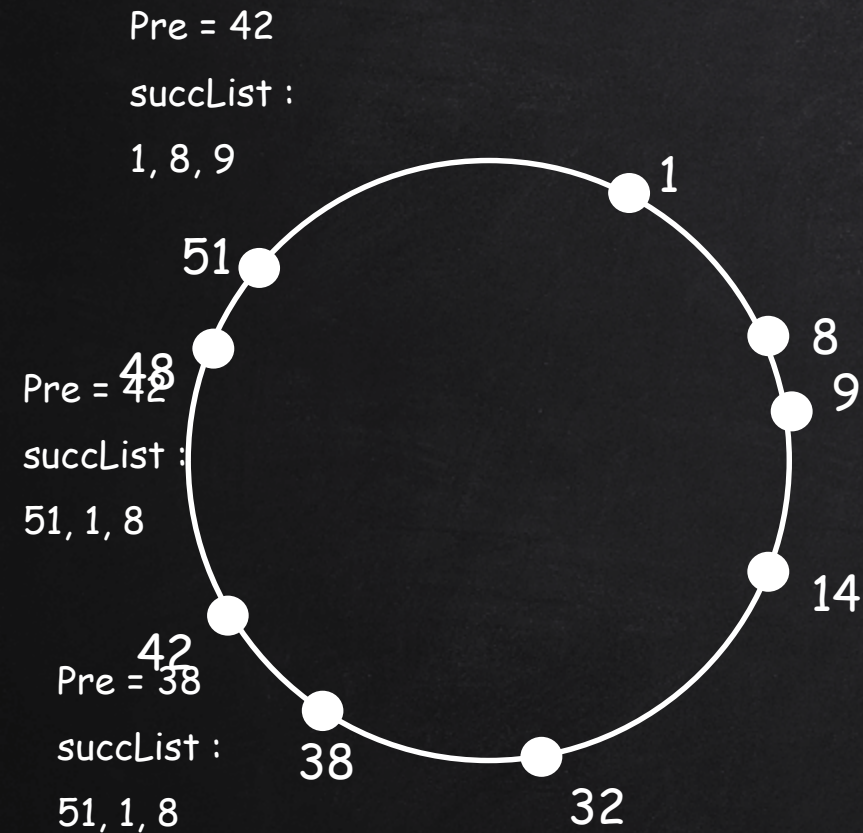


Condition: node #48 quit

X node #48 runs function quit()

X node #42 is called to remove 48 from his succList and add 8 to the end of his succList

CHORD - ROBUSTNESS - SIMULATION



Condition: node #48 quit

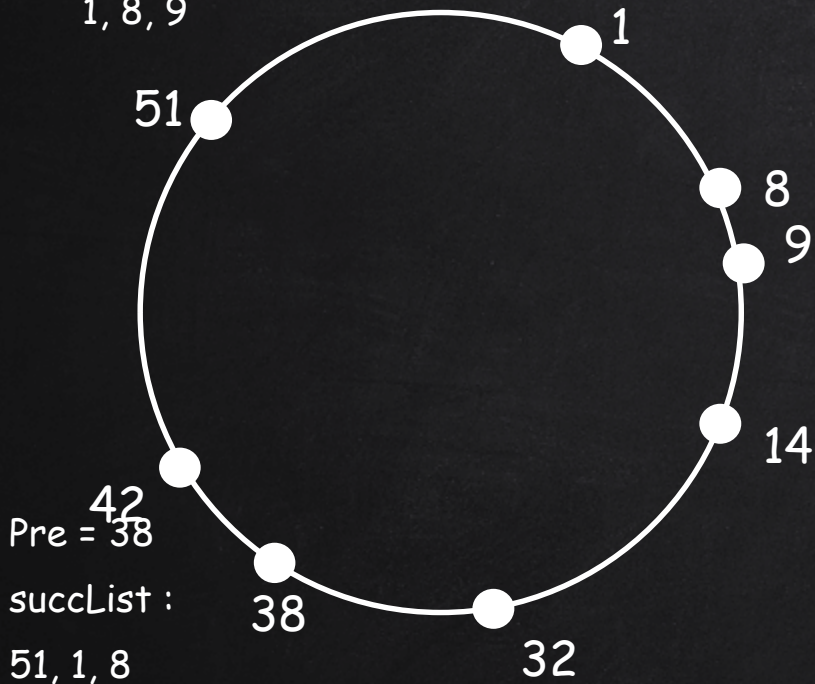
- X node #48 runs function quit()
- X node #42 is called to remove 48 from his succList and add 8 to the end of his succList
- X node #51 is called to update his predecessor to 42

CHORD - ROBUSTNESS - SIMULATION

Pre = 42

succList :

1, 8, 9



Condition: node #48 quit

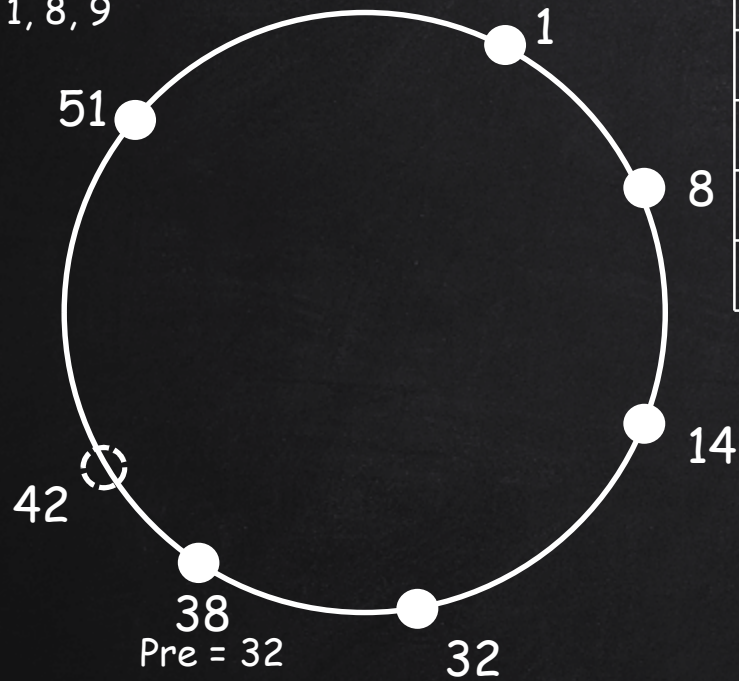
- X node #48 runs function quit()
- X node #42 is called to remove 48 from his succList and add 8 to the end of his succList
- X node #51 is called to update his predecessor to 42
- X then node #48 quit successfully
- X As the time passed, finger table of each node will be fixed gradually by fix_finger() function, then the network become stable again

CHORD - ROBUSTNESS - SIMULATION

Pre = nil

succList :

1, 8, 9



succList :

42, 51, 1

ID	Node
$8 + 2^0$	#14
$8 + 2^1$	#14
$8 + 2^2$	#14
$8 + 2^3$	#32
$8 + 2^4$	#32
$8 + 2^5$	#42

Condition: node #42 failed

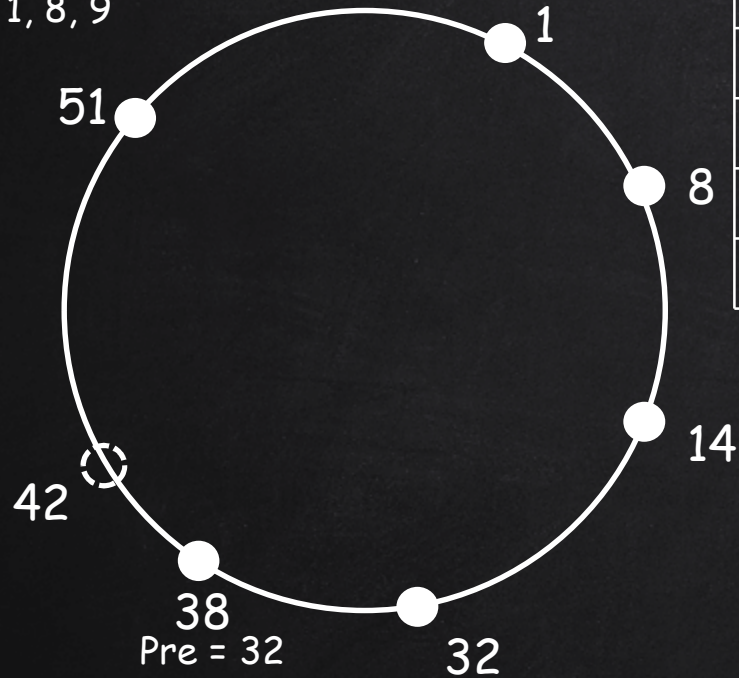
X node #51 runs check_predecessor and find his predecessor failed, so he set to nil

CHORD - ROBUSTNESS - SIMULATION

Pre = 42

succList :

1, 8, 9



succList :

51, 1

ID	Node
$8 + 2^0$	#14
$8 + 2^1$	#14
$8 + 2^2$	#14
$8 + 2^3$	#32
$8 + 2^4$	#32
$8 + 2^5$	#42

Condition: node #42 failed

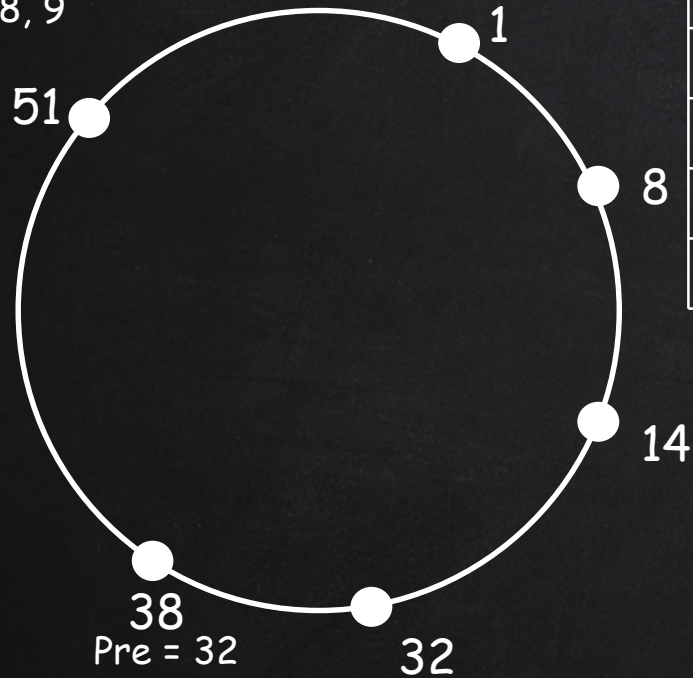
- X node #51 runs check_predecessor and find his predecessor failed, so he set to nil
- X node #38 runs stabilize() and find his successor failed, so he remove 42 from succList and choose 51 as new successor.

CHORD - ROBUSTNESS - SIMULATION

Pre = 38

succList :

1, 8, 9



succList :

51, 1

ID	Node
$8 + 2^0$	#14
$8 + 2^1$	#14
$8 + 2^2$	#14
$8 + 2^3$	#32
$8 + 2^4$	#32
$8 + 2^5$	#42

Condition: node #42 failed

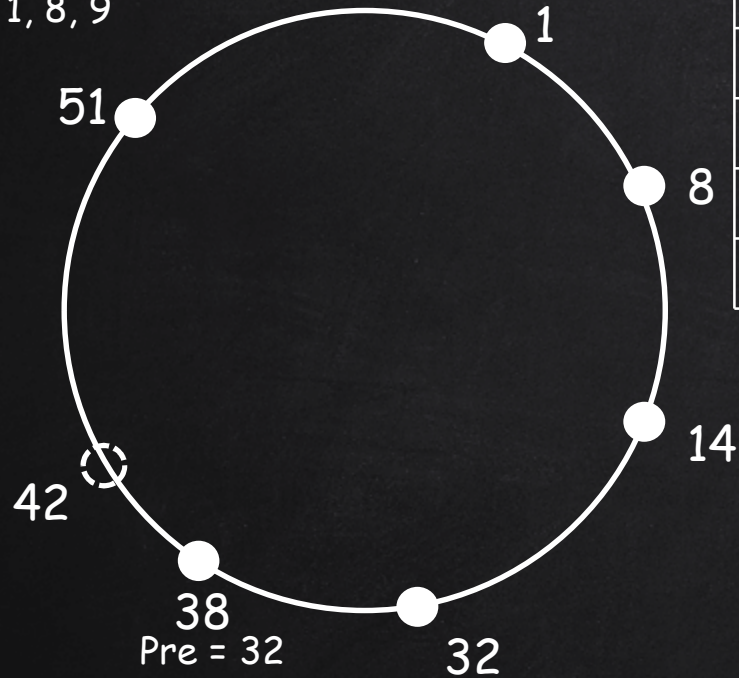
- X node #51 runs check_predecessor and find his predecessor failed, so he set to nil
- X node #38 runs stabilize() and find his successor failed, so he remove 42 from succList and choose 51 as new successor.
- X then node #38 notify his new successor for his existence, and node #51's predecessor is nil, so node #51 set his predecessor to 38
- X then node impact of the failure of node #42 is now eliminated

CHORD - ROBUSTNESS - SIMULATION

Pre = 38

succList :

1, 8, 9



Pre = 32

succList :

51, 1

ID	Node
$8 + 2^0$	#14
$8 + 2^1$	#14
$8 + 2^2$	#14
$8 + 2^3$	#32
$8 + 2^4$	#32
$8 + 2^5$	#42

Condition: node #8 fix finger table

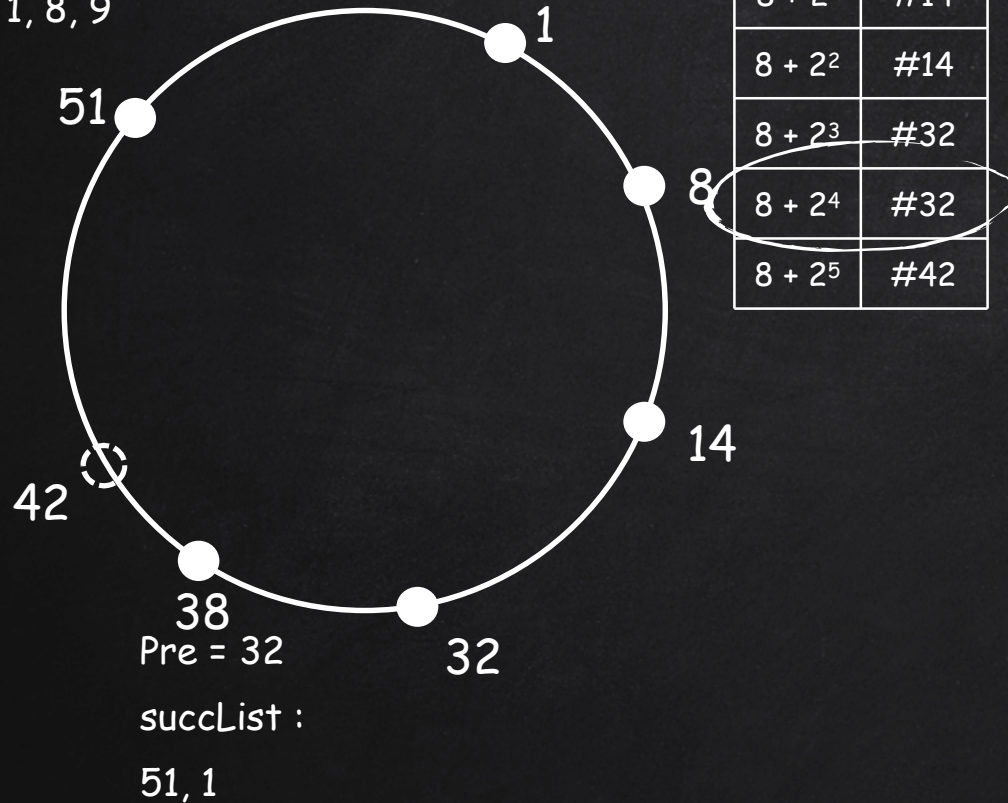
X take the previous condition, and node #8 want to fix the 5th entry of his finger table

CHORD - ROBUSTNESS - SIMULATION

Pre = 38

succList :

1, 8, 9



ID	Node
$8 + 2^0$	#14
$8 + 2^1$	#14
$8 + 2^2$	#14
$8 + 2^3$	#32
$8 + 2^4$	#32
$8 + 2^5$	#42

Condition: node #8 fix finger table

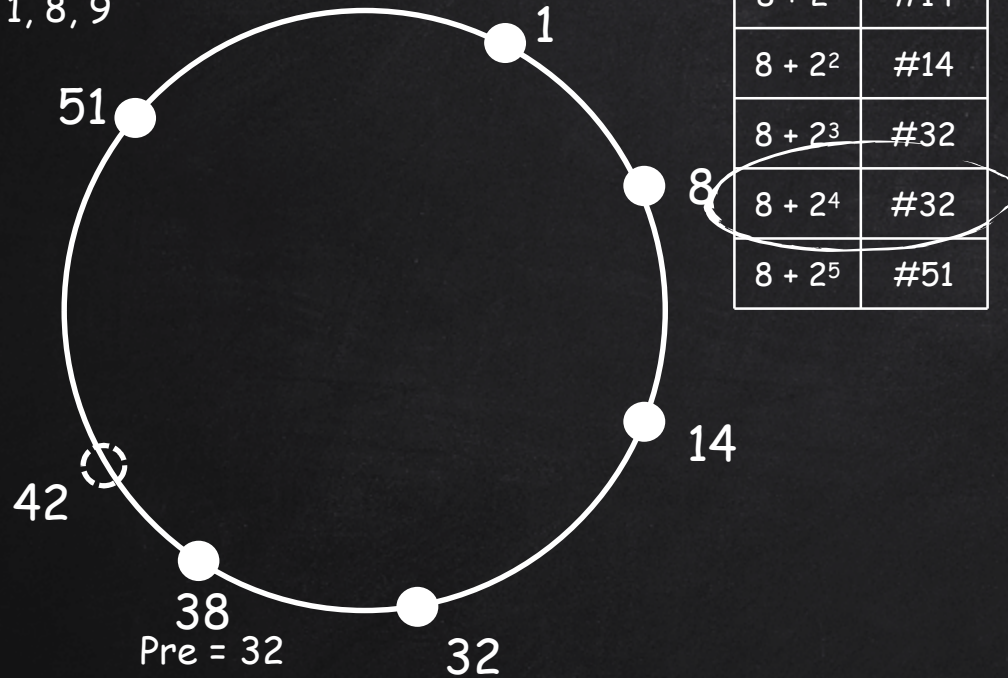
- X take the previous condition, and node #8 want to fix the 5th entry of his finger table
- X he find in his for first valid predecessor of id $40 = 8 + 2^5$, and he get 32

CHORD - ROBUSTNESS - SIMULATION

Pre = 38

succList :

1, 8, 9



Pre = 32

succList :

51, 1

ID	Node
$8 + 2^0$	#14
$8 + 2^1$	#14
$8 + 2^2$	#14
$8 + 2^3$	#32
$8 + 2^4$	#32
$8 + 2^5$	#51

Condition: node #8 fix finger table

- X take the previous condition, and node #8 want to fix the 5th entry of his finger table
- X he find in his for first valid predecessor of id $40 = 8 + 2^5$, and he get 32
- X he asks node #32 to find the successor of 40, then it become a normal look up. And finally the node #8 get 51 and update the 5th entry of his finger table to 51

CHORD - DATA

Strategies to maintaining data in Chord

- X 1. When a node joined, some keys has assigned to his successor should be moved to him
- X 2. When a node quit, all his data should be move to his successor
- X 3. (Personal) When predecessor changed, some key assigned to him should be moved to the predecessor

For 1, we use a modified join() function
n.join(n')

n.pre = nil

n.succ = n'.find_predecessor()

for k, v in n.succ.data

if k is not between n and n.succ

n.data.append(k, v)

n.succ.data.remove(k)

CHORD - DATA

Strategies to maintaining data in Chord

- X 1. When a node joined, some keys has assigned to his successor should be moved to him
- X 2. When a node quit, all his data should be move to his successor
- X 3. (Personal) When predecessor changed, some key assigned to him should be moved to the predecessor

For 2, we use a modified quit() function
n.quit()

```
n.pre.succList.remove(n)
last := len(n.succList)
n.pre.succList.append(n.succList[len - 1])
n.succ.pre = n.pre
for k, v in n.data
    n.succ.data.append(k ,v)
```


CHORD - DATA

Strategies to maintaining data in Chord

- X 1. When a node joined, some keys has assigned to his successor should be moved to him
- X 2. When a node quit, all his data should be move to his successor
- X 3. (Personal) When predecessor changed, some key assigned to him should be moved to the predecessor

For 3, we use a modified notify() function

```
n.notify(n')
```

```
    if n.pre == nil or n' is between n.pre and n
```

```
        n.pre = n'
```

```
        for k, v in n.data
```

```
            if k not between n.pre and n
```

```
                n.pre.data.append(k, v)
```

```
            n.data.remove(k)
```


CHORD - DATA

WHAT ABOUT SOME NODE FAILS?

CHORD - DATA

WHAT ABOUT SOME NODE FAILS?





APPLICATIONS ON DHT

THIS PART INTRODUCE SOME APPROACHES THAT MAKE DHT SAFER AND MY APPLICATION ON DHT



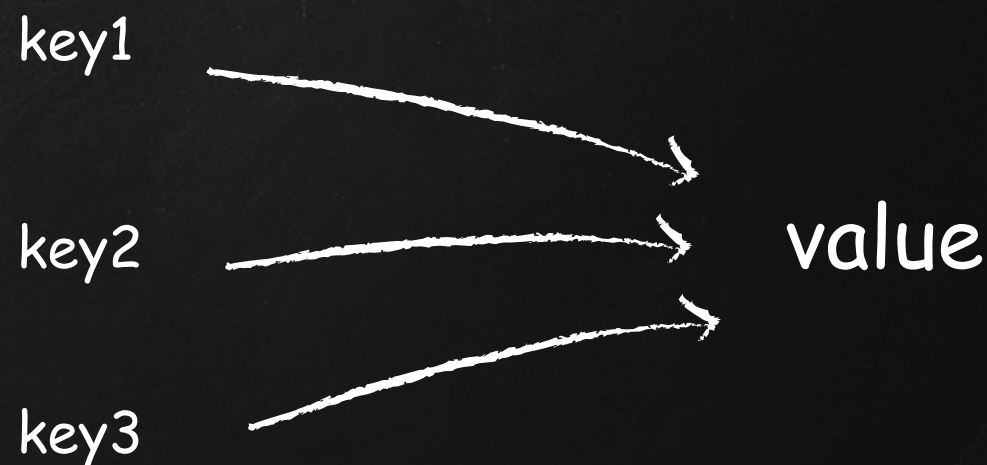
MAKE DHT SAFER

Create local buffer

Dump data in the memory to hard disk periodically in case the node fails



Multi key to a value





MY APPLICATION BASED ON DHT

magnet:?xt=urn:btih:

e04f646aa669375c7a06b6569b70657a9bd82a0c



MY APPLICATION BASED ON DHT

*Stupid Chord Peer to peer File
Sharing System
"SCPFSS"*



MY APPLICATION BASED ON DHT

"SCPFSS"

A BitTorrent like system



STRUCTUE OF SCPFSS

Basement: DHT

Put file hash as key and a LIST of the address of node who share this file as value

Middle: File Server RPC

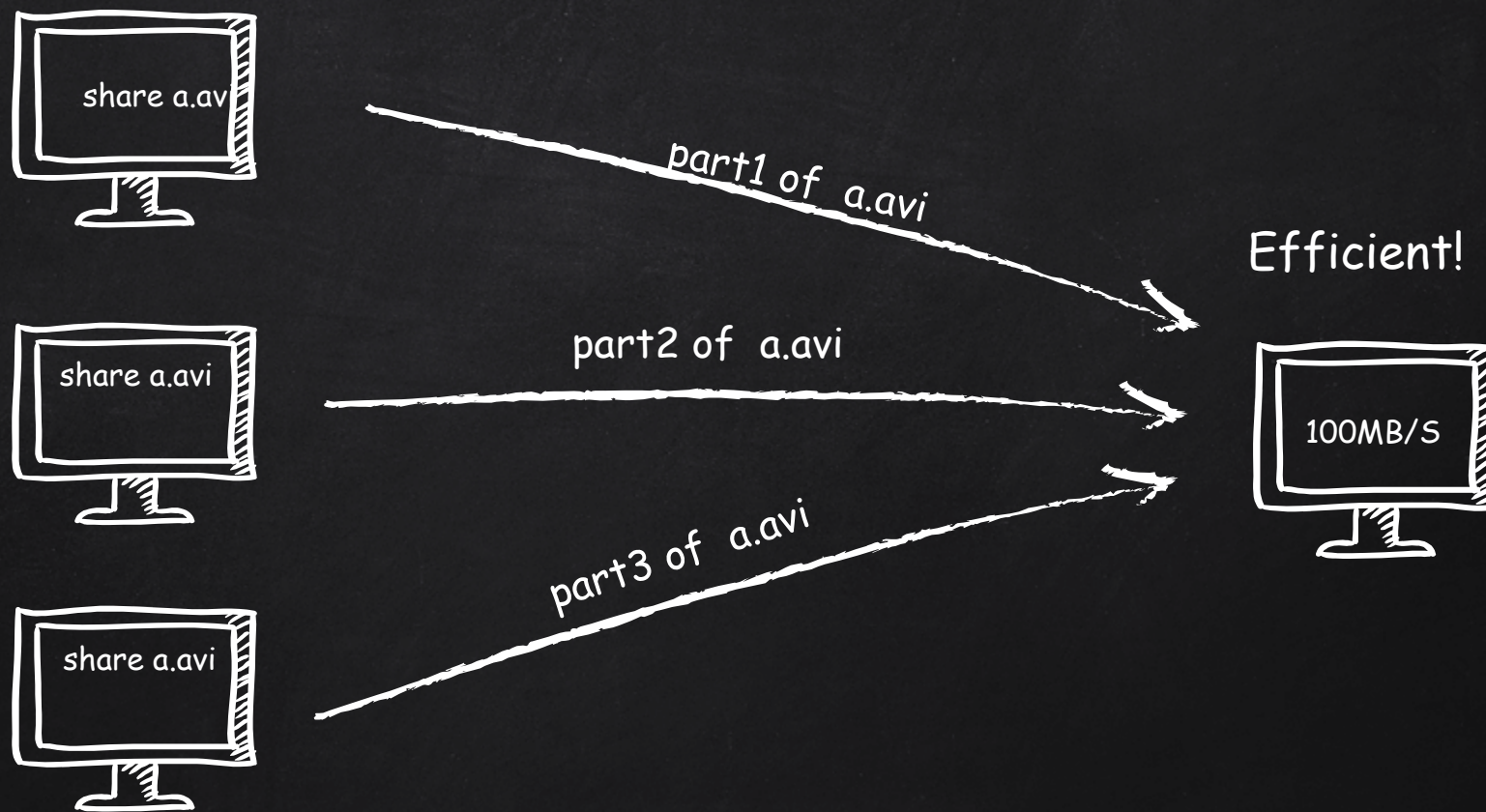
Providing secive of looking up file and send file information

Upper: TCP File Server

Send (part of) files to other peers



CHARACTER OF SCPFSS—MULTI THREADING





THANKS!

Any questions?