

GSCI 1380 : Day 9 (Tapestry Continued)

---

---

---

---



## Last class

- \* Consistent Hash Rep

\* Scaling consistent hash

$$\text{Chord (Table size} = \log_2(N)$$

$$\text{Tapestry (Table size} = \boxed{\log_B(N)} * B$$

MaxHop

load imbalance  $\Rightarrow$  Virtual nodes

route table size (# of servers)

## This Class

- \* Tapestry In Detail

① Routing Table Creation (closest node)

flag #3  
in paper

② Key lookup (Root Node, Better Choice, MaxHop)

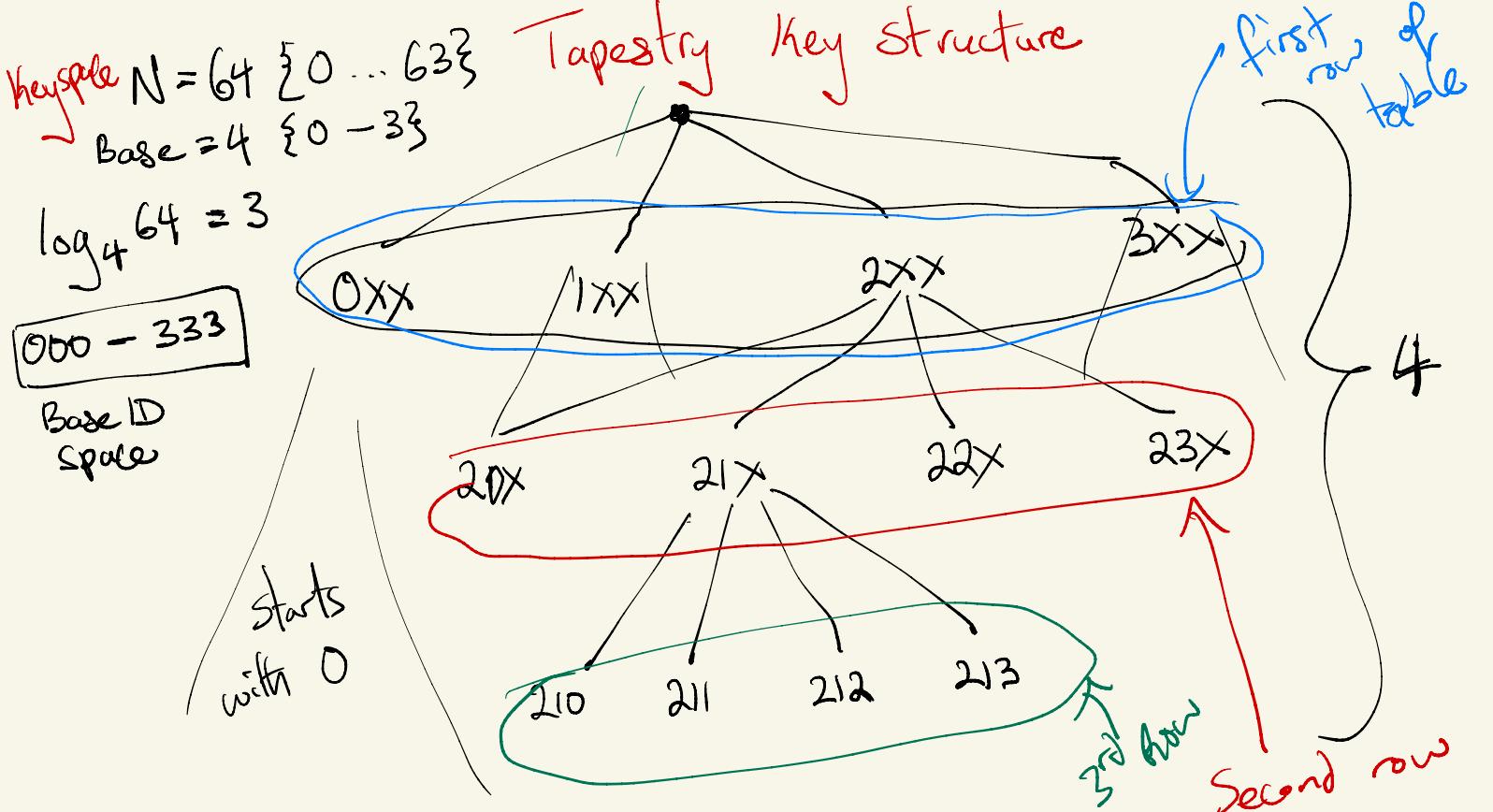
③ Overview (back pointers, object store, Publishing)

④ Deleting Nodes (Graceful/planned V. unGraceful/updated)

⑤ Adding Nodes (AddNode Multicast, Need To Know Nodes)

Node = Server

Host



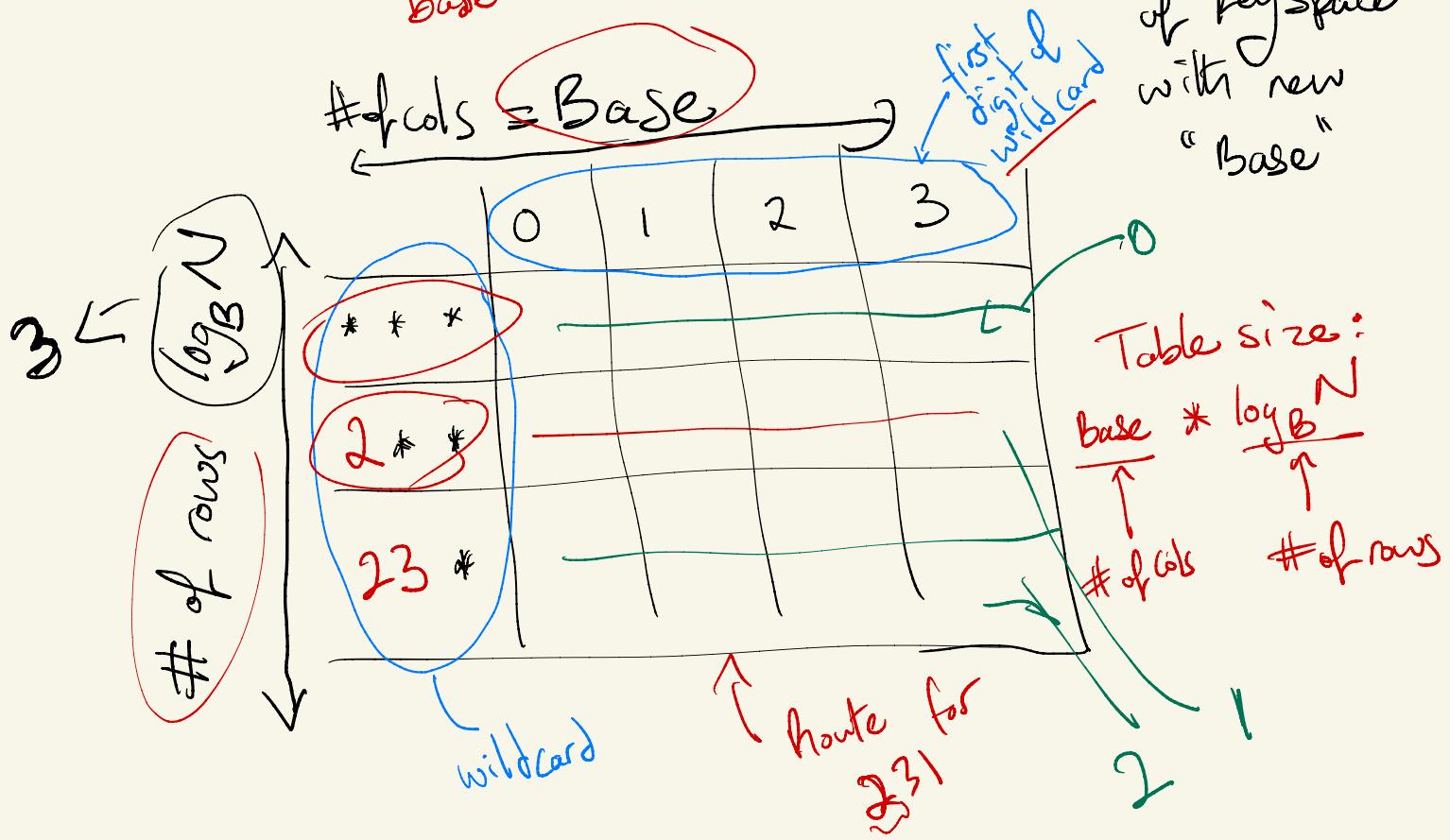
Key ID Space =  $0 \dots 63$

$N=64$

base 10

Pick a Base here 4  $\Rightarrow 000 - 333$

new representation of key space with new "base"



Node IDs : 011, 211, 222, 231, 201, 300, 333

211			
0	1	2	3
011	211	300	
201	211	222	231
210	211	222	233

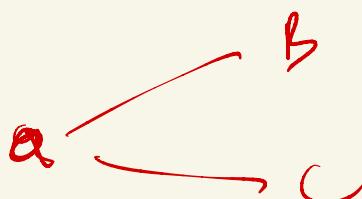
011			
0	1	2	3
*	011	-	222 333
0	*	011	-
01	*	011	-

## Options for closeness

- ① random
- ② no hop (minimize latency)
- ③  $\min(|A - B|)$ 
  - = absolute distance / difference in IDs
  - ↳ project

0	1	2	3
*	*		
*	*		
*			

closeness

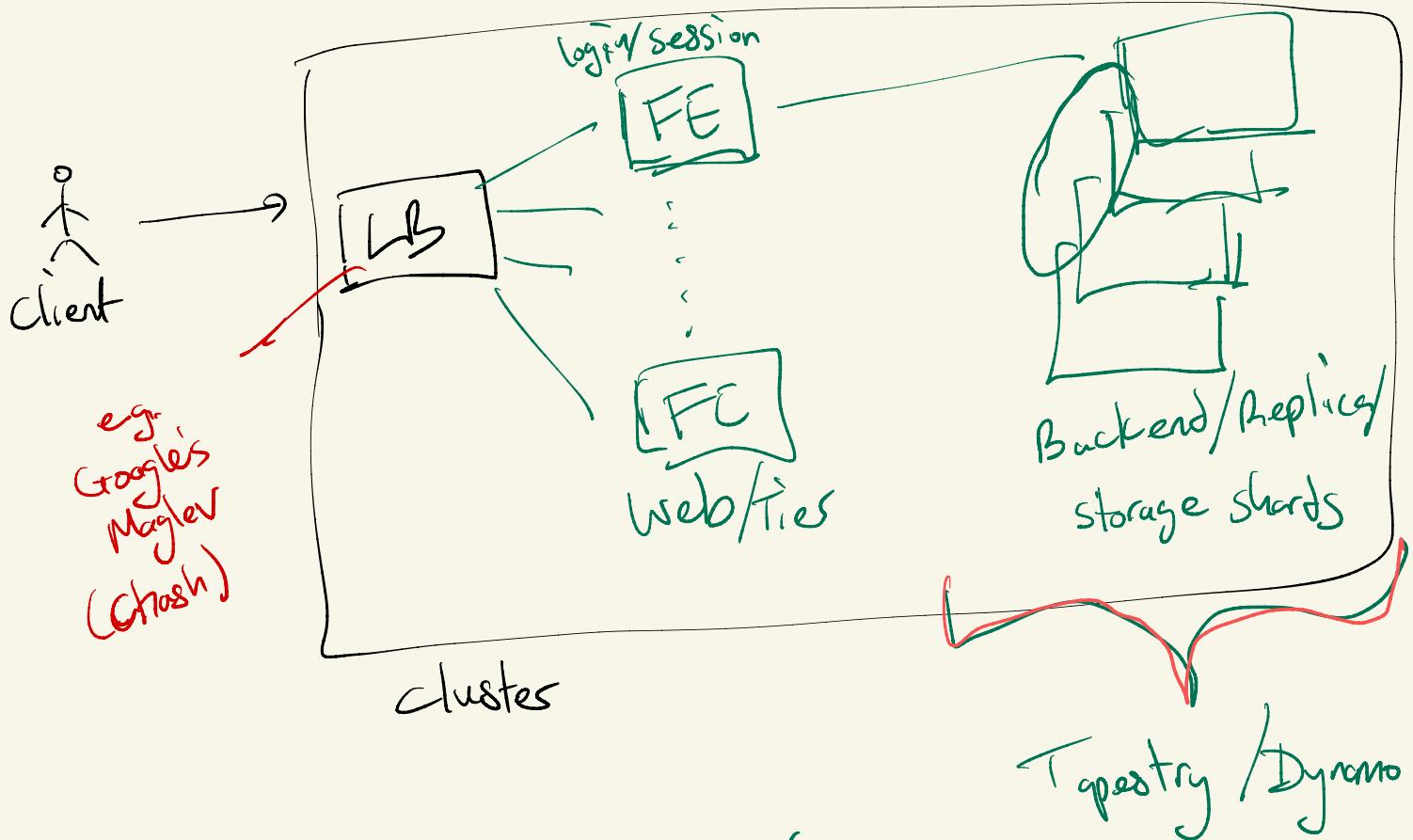


$$\min(|a - b|, |a - c|)$$

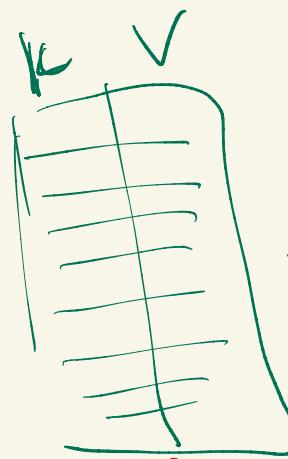
Points about table creation

- ① follow the rules (constraints placed by columns & rows)
- ② if you have several options pick based on closeness
- ③ there may be empty entries
- ④ always one entry per row

# Context / overview



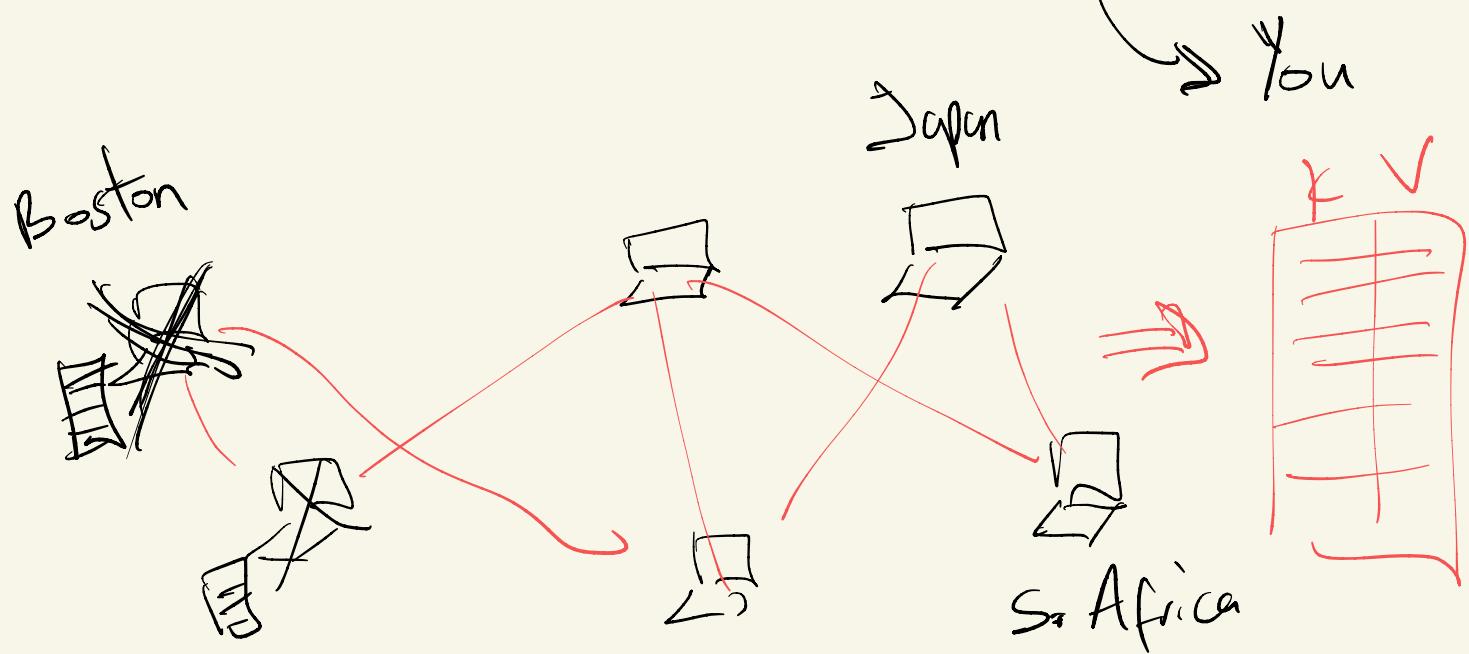
abstraction  
of a large  
Key Value store  
(distributed hash table)



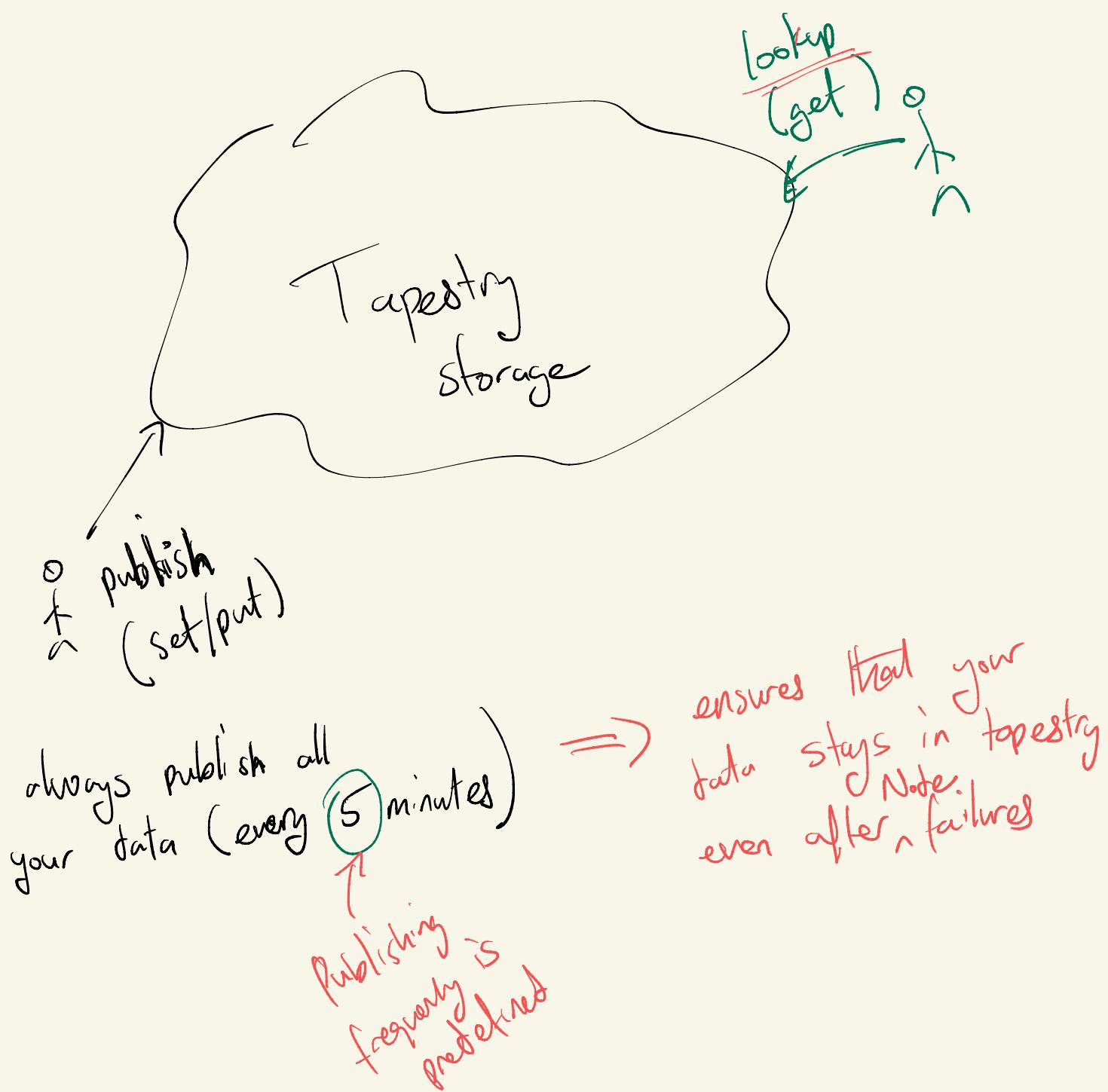
set(KeyID, "in")  
get(KeyID)  
move  
email  
shopping

Storage servers  
(shards)

# P2P (Peer to Peer)



- ① H/w ← different resources
- ② nodes come and go at will
- ③ Geo distributed



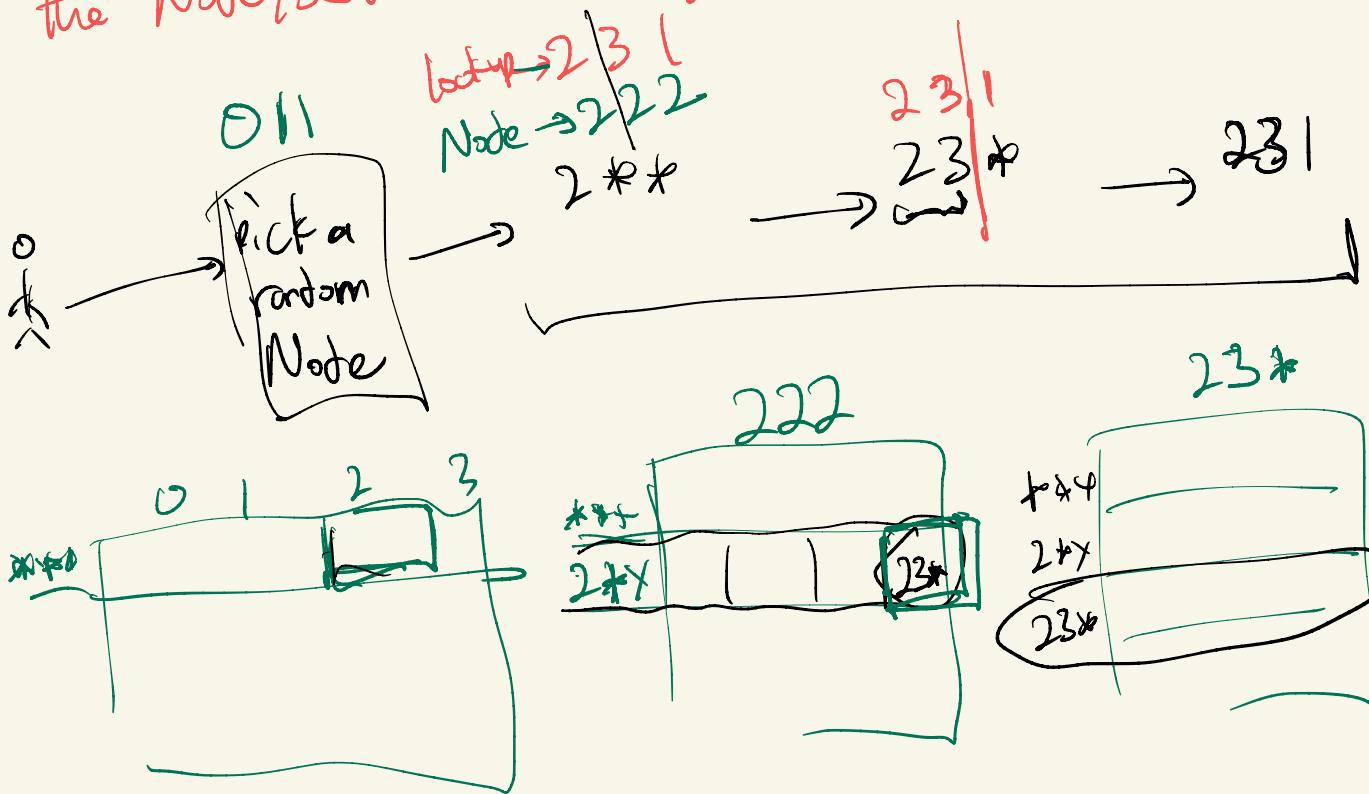
# Lookups (in Tapestry)

① Goal : find "Root"

(Node)  
Server in charge of some set of IDs

Root = Lookup (231)

AKA the Node/SERVER in charge of that ID



$$O(\log_B N)$$

## Prefix Routing

- \* every node
- \* tries to find an entry in its table
- \* that has the largest overlap
- \* then forwards you to this table

Node IDs : 011, 211, 222, 231, 201, 300, 333

211				300
0	1	2	3	
011	—	211	300	
—	211	222	231	
—	211	—	—	

0	1	2	3
011	—	222	333
300	—	—	333
300	—	—	—

Lookup (KeyID, step) {

row = step

col = KeyID.at(step)

nextNode = Table[col, row]

step++

nextNode.lookup(KeyID, step)

333				
0	1	2	3	
011	—	211	300	
300	—	—	—	333
—	—	—	—	333

3

start @ 0

lookup (333, 0) → 211

lookup (333, 1) → 300

lookup (333, 2) → 333

Node IDs : 011, 211, 222, 231, 201, 300, 333

211				300
0	1	2	3	
011	—	211	300	
—	211	222	231	
—	211	—	—	

0	1	2	3
011	—	222	333
300	—	—	333
300	—	—	—

Lookup (KeyID, step) {

row = step

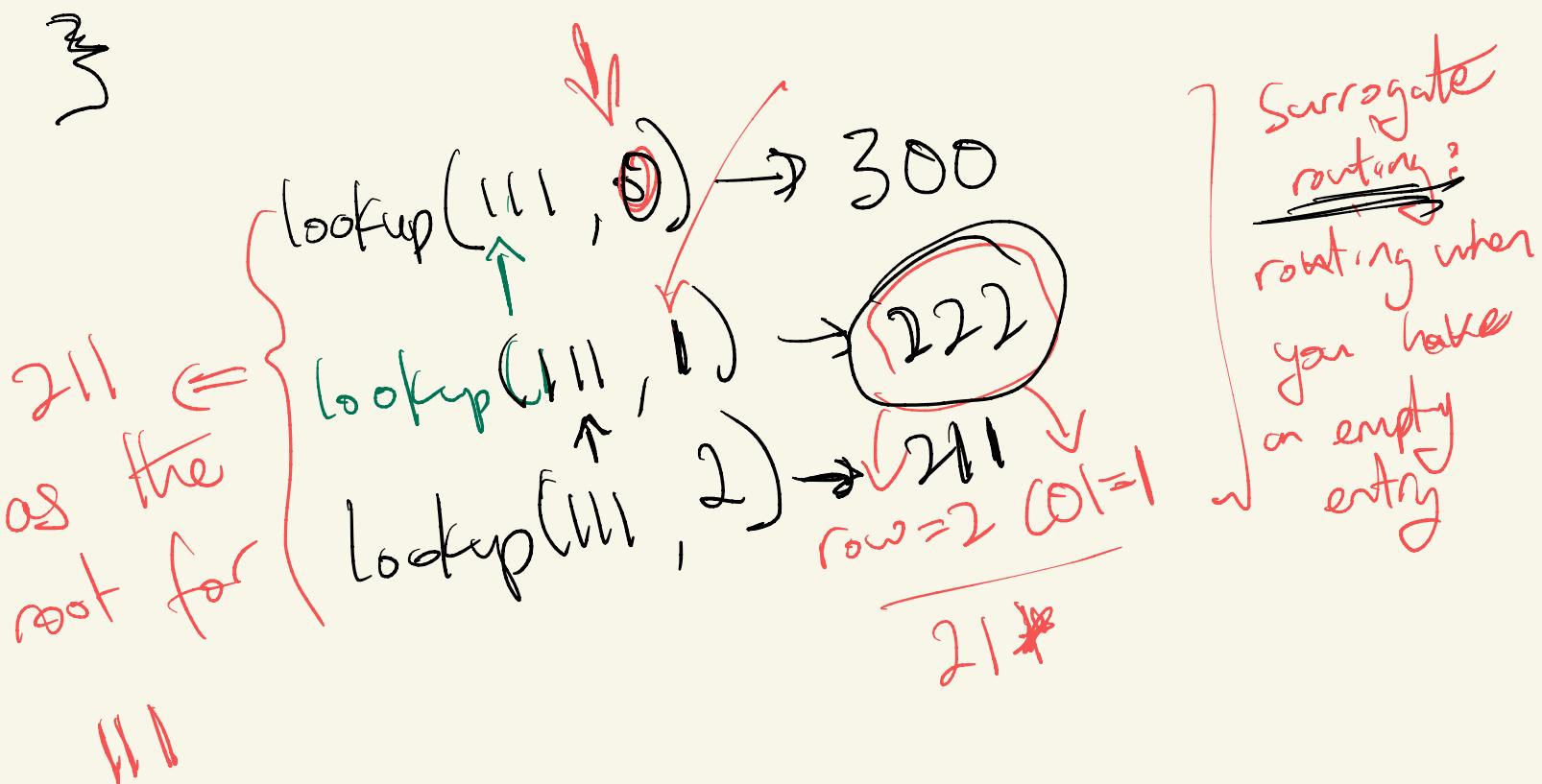
col = KeyID.at(step)

nextNode = Table[col, row]

step++

nextNode.lookup(KeyID, step)

333				
0	1	2	3	
011	—	211	300	
300	—	—	—	333
—	—	—	—	333



# Node failure / removal

## Graceful

- \* Say Good bye to nodes & hand over your data

① Who to give your data to?

New root

② how to update broken routing tables

need to figure out which tables are impacted by exit!!!

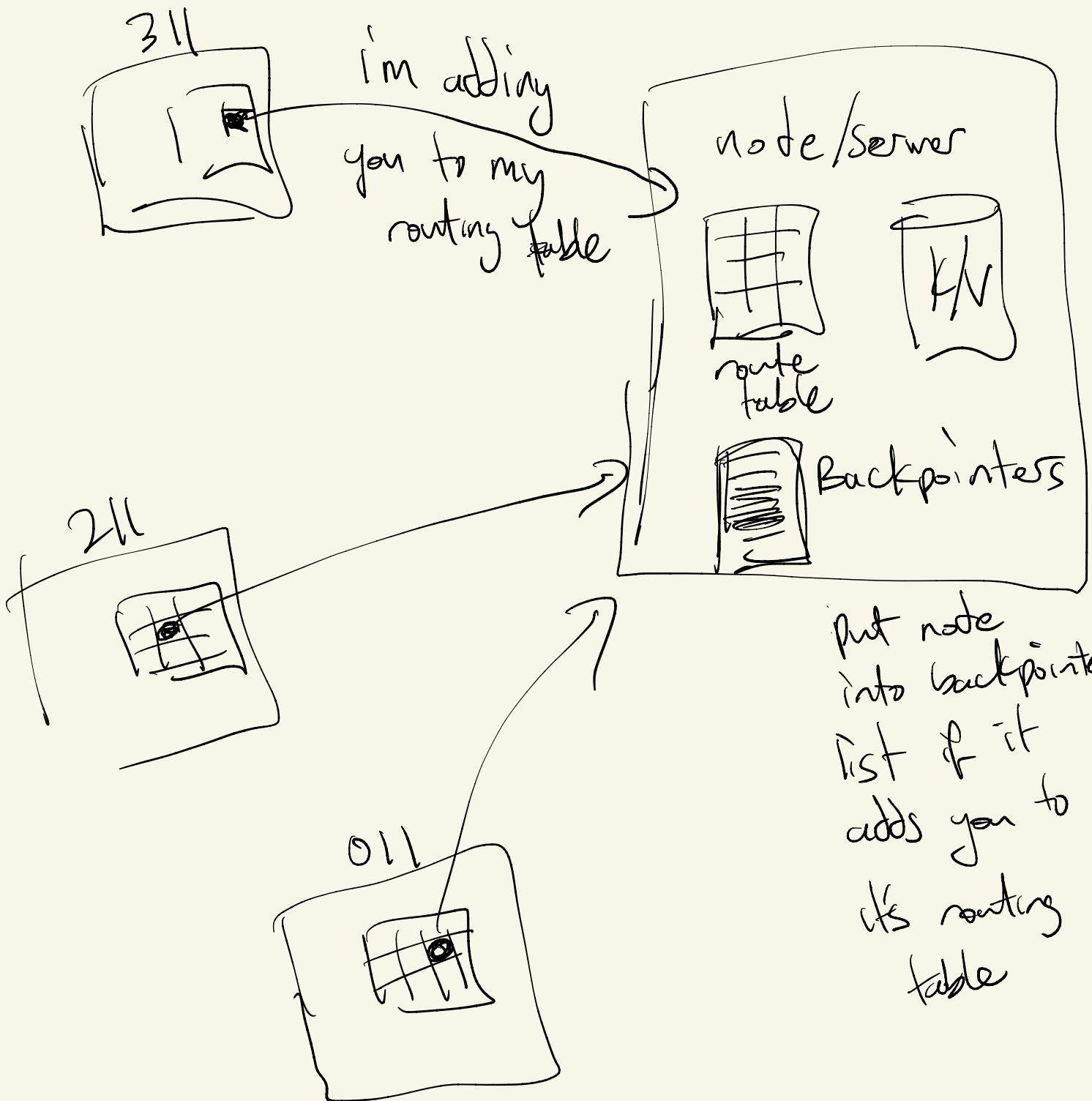
Use Backpointers to determine nodes to notify

## not Graceful

Just disconnect

- \* failure is detected w/ heartbeats

\* data persists because clients republish the data



# Summary

---

- ① Overview / context for Tapestry
- ② Building a route table
- ③ Joining in Tapestry (lookups)
- ④ Exiting the Tapestry P2P network
  - graceful & ungraceful

- 
- ⑤ How to add nodes into the network

Node IDs : 011, 211, 222, 231, 201, 300, 333

211			
0	1	2	3
*	*	*	
2	*	*	
2	1	*	

0	1	2	3
*	*	*	
*	*		
	*		

0	1	2	3
*	*	*	
*	*		
	*		

Node IDs : 011, 211, 222, 231, 201, 300, 333