

CSX55: DISTRIBUTED SYSTEMS [DHTs]

Routing in DHTs

So many, many
nodes and items
But the mapping's unambiguous
Deterministic to boot

Each node in the know only
about a few others
Messages relayed closer and closer
In a few bounded hops

Shrideep Pallickara
Computer Science
Colorado State University

COMPUTER SCIENCE DEPARTMENT



COLORADO STATE UNIVERSITY

1

Frequently asked questions from the previous class survey

- Is there ever a problem comparing such large numbers?
- N-fold replication of content in routers vs overlays
- Clockwise vs anti-clockwise traversals
- How can a node hierarchy be built up to optimize lookups based on text?



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALICKARA
COMPUTER SCIENCE DEPARTMENT

PEER-TO-PEER SYSTEMS

L15.2

2

Topics covered in this lecture

- Distributed Hash Tables
- Chord



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALICKARA
COMPUTER SCIENCE DEPARTMENT

PEER-TO-PEER SYSTEMS

L15.3

3

Distributed hash tables

- Few constraints on the structure of the keys
- REQUIREMENTS
 - ▣ Data identified using numeric **keys**
 - ▣ Nodes must be willing to store keys **for each other**



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALICKARA
COMPUTER SCIENCE DEPARTMENT

PEER-TO-PEER SYSTEMS

L15.4

4

Storage and retrieval in distributed hash tables

- Data items are *inserted* and *found* by specifying a unique **key** for the data
- Underlying algorithm must determine *which node* is responsible for storing the data



Distributed Storage using DHTs: Publishing a file

- **Convert** file-name to numeric key
 - Using one-way hash functions like MD5 or SHA-1
- Call **lookup (key)**
 - Returns IP address of node responsible for key
- **Send file** to be stored at node returned by lookup



Distributed Storage using DHTs: Retrieving a file

- ① Obtain name of file
- ② Convert it to a key using one-way hash function
- ③ Call lookup (key)
- ④ Ask resulting node, from (3), for a copy of the file



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

PEER-TO-PEER SYSTEMS

L15.7

7



8

Implementing DHTs: 3 core elements

- **Mapping** keys to nodes
- **Forwarding** a lookup for a key to the appropriate node
- Building **routing tables**



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

PEER-TO-PEER SYSTEMS

L15.9

9

Implementing DHTs: Mapping keys to nodes

- Must be load balanced
- Done using one-way hash functions
 - MD5 (128-bit) or SHA-1 (160-bit)
- Ensures that content is distributed **uniformly**



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

PEER-TO-PEER SYSTEMS

L15.10

10

Implementing DHTs

Forwarding lookups

- Any node that receives query for key
 - ▣ Must forward it to a node whose ID is **closer** to the key
- Above rule guarantees that query **eventually arrives** at the closest node
- For e.g.:
 - ▣ Node has ID 346, and key has ID 542
 - ▣ Forwarding to node 495 gets it numerically closer



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

PEER-TO-PEER SYSTEMS

L15.11

11

Implementing DHTs:

Building routing tables

- Multiple nodes participate in locating content
- Each node must know about **some other** nodes
 - ▣ To forward lookup requests
 - ▣ SUCCESSOR
 - The node with the **closest succeeding** ID
 - ▣ Other nodes
 - For efficiency in routing



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

PEER-TO-PEER SYSTEMS

L15.12

12

Distributed hash tables: Identifiers

- Data items are assigned an identifier from a large random space
 - ▣ 128-bit UUIDs or 160-bit SHA1 digests
- **Nodes are also assigned a number from the same identifier space**



Crux of the DHT problem

- Implement an efficient, **deterministic** scheme to
 - ▣ Map data items to node
- When you **look up** a data item
 - ▣ Network address of node holding the data is returned





15

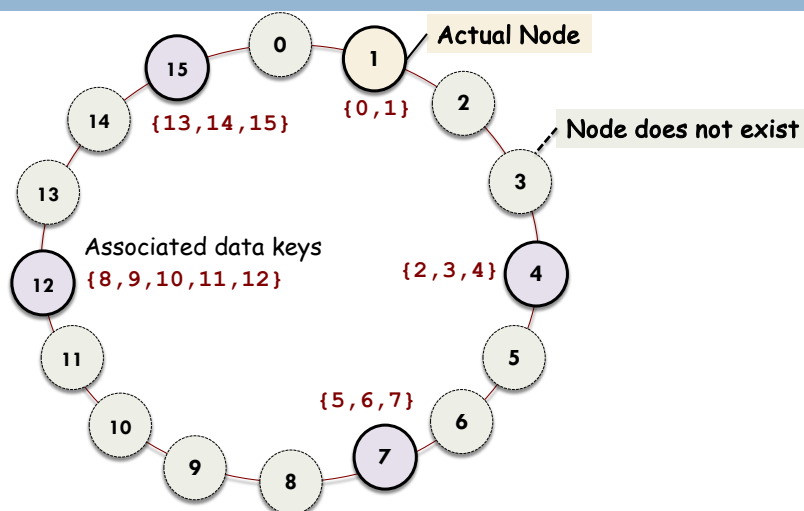
The Chord System

- Assigns IDs to keys and nodes from the same 1-dimensional ID space
- Nodes are organized into a **ring**
- Data item with key k is mapped to a node with the **smallest** $id \geq k$
 - ▣ Also referred to as `successor(k)`



16

Mapping of data items to nodes in Chord



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALICKARA
COMPUTER SCIENCE DEPARTMENT

PEER-TO-PEER SYSTEMS

L15.17

17

Chord lookups

- N is the number of possible nodes in the system
- Each node maintains a **finger table**
 - With $\log N$ entries
 - Entries contains IP addresses of nodes
 - Half-way around the ID space from it
 - $1/4^{\text{th}}$, $1/8^{\text{th}}$, ... **in powers of two**
 - Ensures node can forward lookup query to at least $1/2$ of the remaining ID-space distance to key
 - Lookups in $O(\log N)$



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALICKARA
COMPUTER SCIENCE DEPARTMENT

PEER-TO-PEER SYSTEMS

L15.18

18

Storing keys and forwarding lookups

- An entity with key k falls under the **jurisdiction** of node with the **smallest identifier** id
 - $id \geq k$
 - Referred to as the successor of k or $succ(k)$
- A node **forwards** query for key k to node (in its FT) with highest ID $\leq k$
 - The exception is ONLY when the first entry is greater than k
 - In this case, that node is responsible for storing that element



COLORADO STATE UNIVERSITY

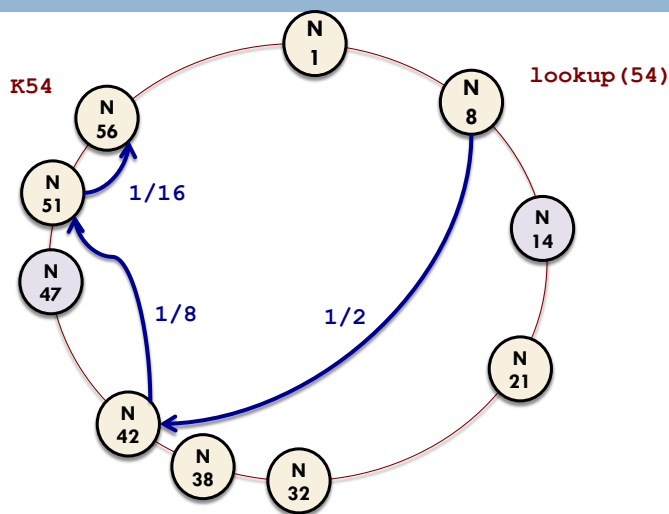
Professor: SHRIDEEP PALICKARA
COMPUTER SCIENCE DEPARTMENT

PEER-TO-PEER SYSTEMS

L15.19

19

Chord lookup example for $k=54$



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALICKARA
COMPUTER SCIENCE DEPARTMENT

PEER-TO-PEER SYSTEMS

L15.20

20

When a node wants to join

- Generate a random id
 - ▣ Probability of collisions is low
- **lookup (id)**
 - ▣ Will return $\text{successor}(\text{id})$
- Contact $\text{successor}(\text{id})$ and its predecessor
 - ▣ Insert self in the ring
 - ▣ **Transfer** data items
 - All keys must be fetched from the new node's successor



COLORADO STATE UNIVERSITY

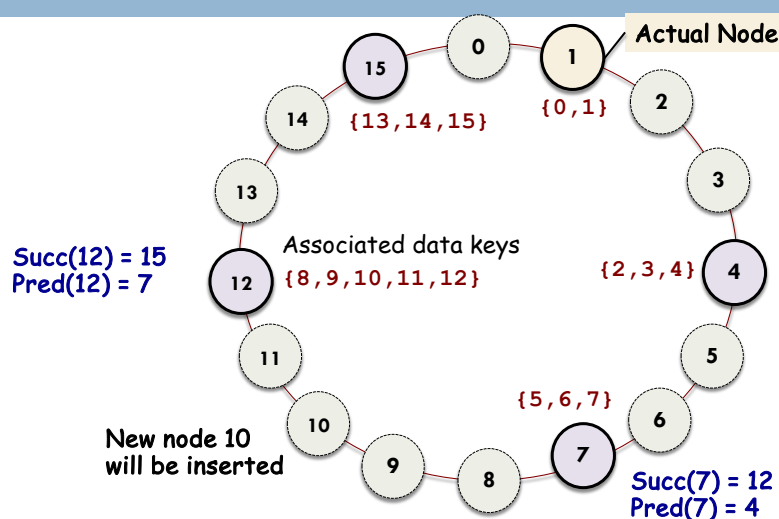
Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

PEER-TO-PEER SYSTEMS

L15.21

21

An example of inserting a new node



COLORADO STATE UNIVERSITY

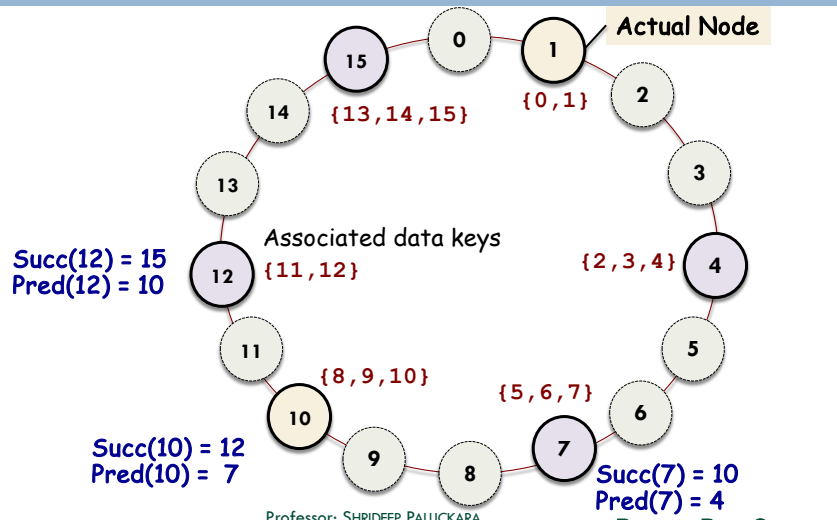
Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

PEER-TO-PEER SYSTEMS

L15.22

22

An example of inserting a new node



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALLICKARA
 COMPUTER SCIENCE DEPARTMENT

PEER-TO-PEER SYSTEMS

L15.23

23

Finger Table in Chord

- Chord uses an m -bit identifier space
 - 2^m possible peers
- Each node, p , in Chord maintains a Finger Table with m -entries

■ $FT_p[i] = \text{succ}(p + 2^{i-1})$

Note: This is when you count your indices from 1.
 When you code, and we are counting from 0 this would be
 ■ $FT_p[i] = \text{succ}(p + 2^i)$



COLORADO STATE UNIVERSITY

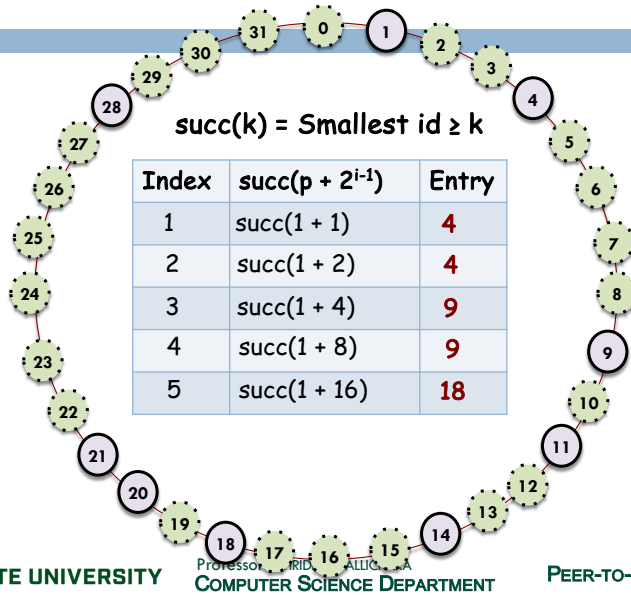
Professor: SHRIDEEP PALLICKARA
 COMPUTER SCIENCE DEPARTMENT

PEER-TO-PEER SYSTEMS

L15.24

24

Constructing the Finger Table: Node 1



We are looking at a 5-bit ID space.
 IDs go from 0 through $(2^5 - 1)$



COLORADO STATE UNIVERSITY

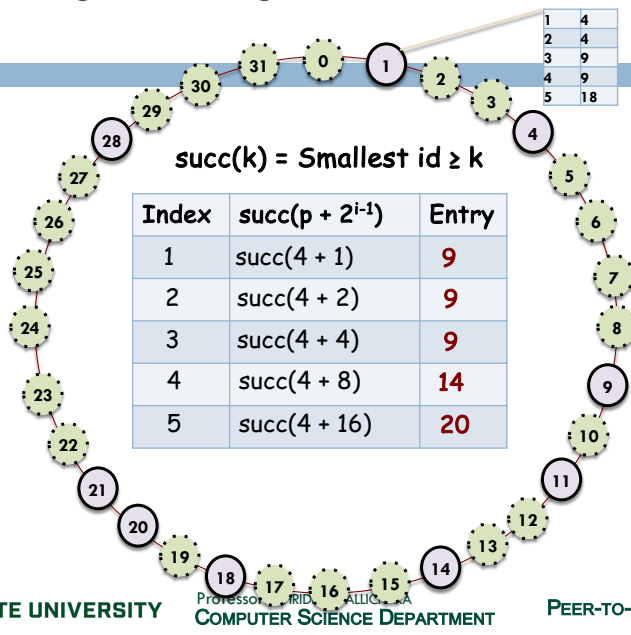
Professor SHRIDEEP PALICKARA
 COMPUTER SCIENCE DEPARTMENT

PEER-TO-PEER SYSTEMS

L15.25

25

Constructing the Finger Table: Node 4



1	4
2	4
3	9
4	9
5	18



COLORADO STATE UNIVERSITY

Professor SHRIDEEP PALICKARA
 COMPUTER SCIENCE DEPARTMENT

PEER-TO-PEER SYSTEMS

L15.26

26

Constructing the Finger Table: Node 9

$\text{succ}(k) = \text{Smallest id} \geq k$

Index	$\text{succ}(p + 2^{i-1})$	Entry
1	$\text{succ}(9 + 1)$	11
2	$\text{succ}(9 + 2)$	11
3	$\text{succ}(9 + 4)$	14
4	$\text{succ}(9 + 8)$	18
5	$\text{succ}(9 + 16)$	28

```

if (val ≥ 2m) {
    val = val (mod 2m)
}
        
```

COLORADO STATE UNIVERSITY
Professor SHRIDEEP PALICKARA
COMPUTER SCIENCE DEPARTMENT

PEER-TO-PEER SYSTEMS

L15.27

27

Constructing the Finger Table: Node 28

$\text{succ}(k) = \text{Smallest id} \geq k$

Index	$\text{succ}(p + 2^{i-1})$	Entry
1	$\text{succ}(28 + 1)$	1
2	$\text{succ}(28 + 2)$	1
3	$\text{succ}(28 + 4)$	1
4	$\text{succ}(28 + 8)$	4
5	$\text{succ}(28 + 16)$	14

```

if (val ≥ 2m) {
    val = val (mod 2m)
}
        
```

COLORADO STATE UNIVERSITY
Professor SHRIDEEP PALICKARA
COMPUTER SCIENCE DEPARTMENT

PEER-TO-PEER SYSTEMS

L15.28

28

Using the finger table to route queries: Make sure you don't overshoot

- To lookup a key k , node p will forward query to node q with index j in p 's FT where:

Node with
greatest ID less than or equal to k

$$q = FT_p[j] \leq k < FT_p[j+1]$$

OR

$$q = FT_p[1] \text{ when } p < k < FT_p[1]$$

First entry ONLY if its ID is greater than k



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

PEER-TO-PEER SYSTEMS

L15.29

29

Stop forwarding the query when you are the target node

- A node is **responsible** for keys that fall in the range
 $key > predecessor$
 $key \leq self$



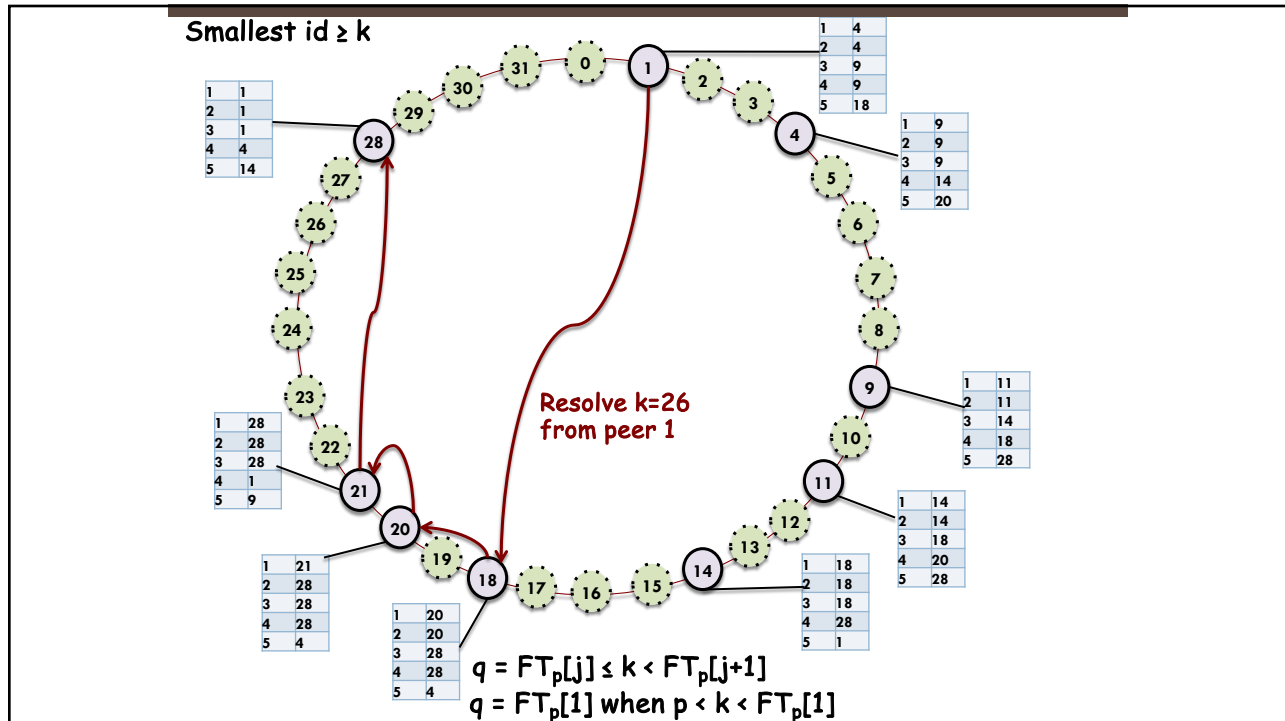
COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

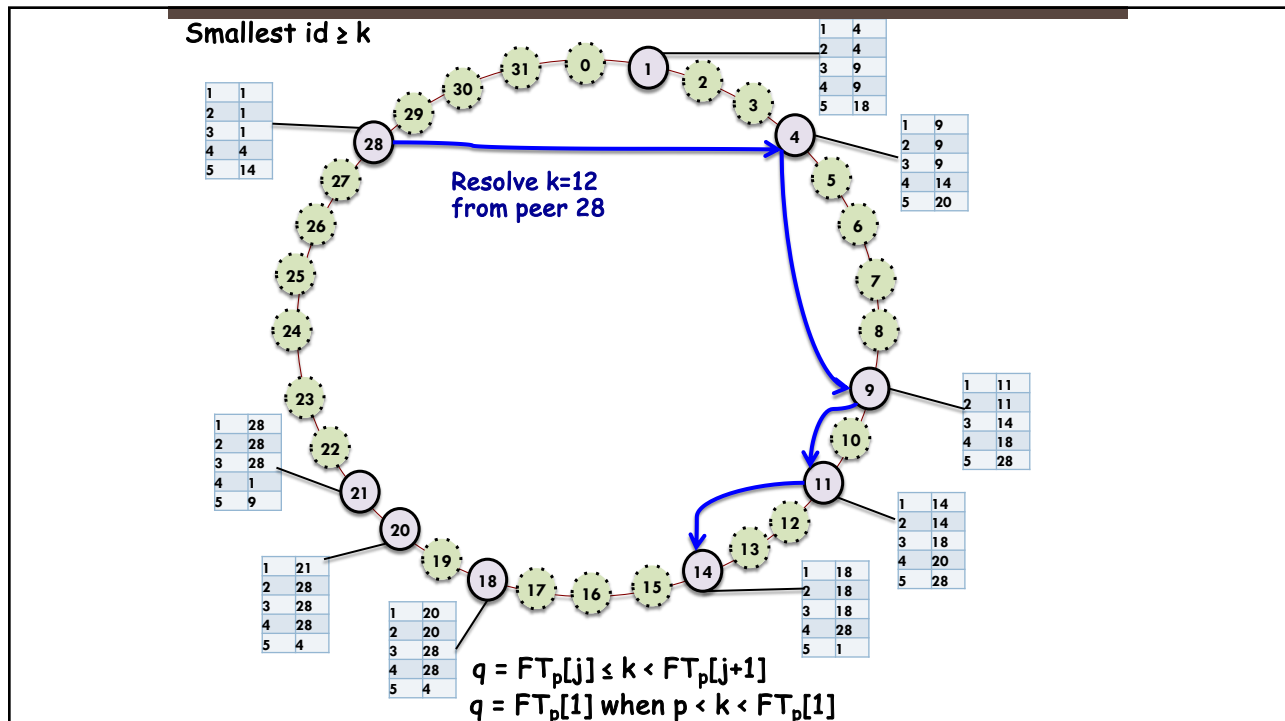
PEER-TO-PEER SYSTEMS

L15.30

30



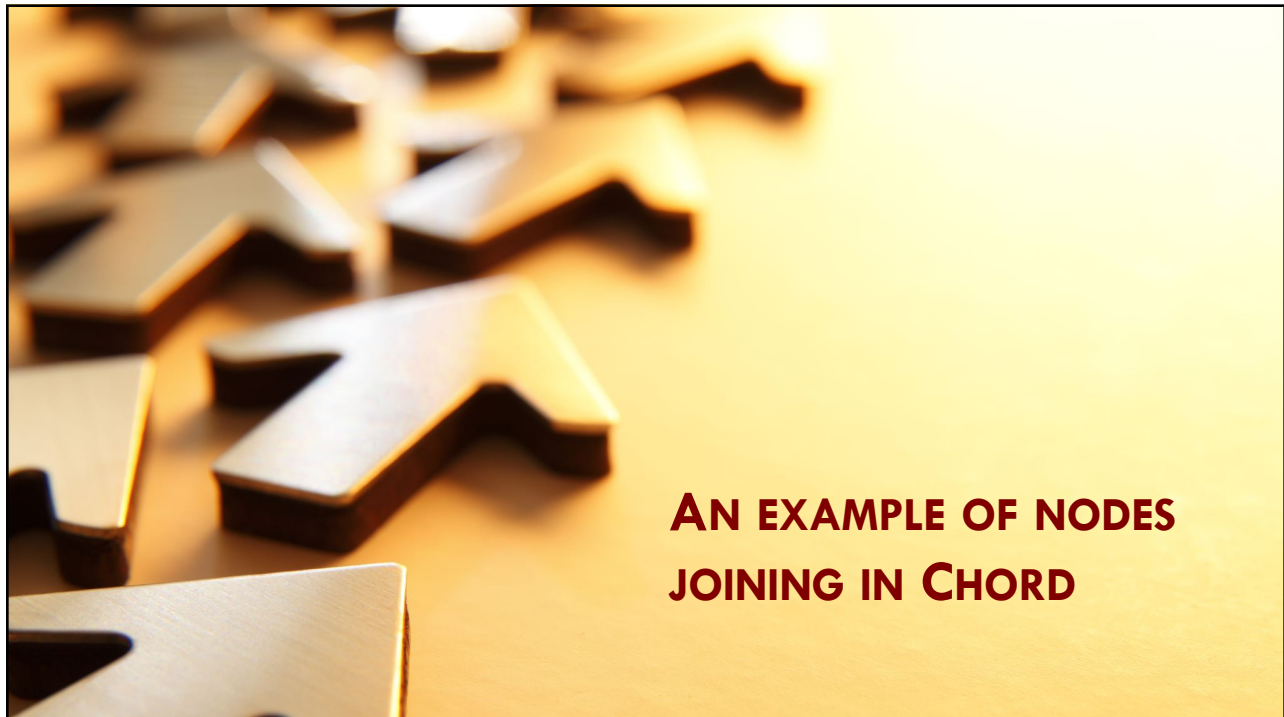
31



32

Keeping the finger table up-to-date: At node q , $FT_q[1]$ must be accurate

- ① Contact $succ(q+1)$ {This is $FT_q[1]$ }
 - ▣ Have it return its predecessor
- ② If $q = pred(succ(q+1))$
 - ▣ Everything is fine
- ③ Otherwise:
 - ▣ There is a new node p such that $q < p \leq succ(q+1)$
 - ▣ $FT_q[1] = p$
 - ▣ Check if p has recorded q as its predecessor
No? Go to step (1)



**AN EXAMPLE OF NODES
JOINING IN CHORD**

Updating the FT at N-1

1	4
2	4
3	1

Pred(1) = 4
Succ(1) = 4

1	1
2	1
3	1

Succ(4) = 1
Pred(4) = 1

COLORADO STATE UNIVERSITY
 Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

PEER-TO-PEER SYSTEMS L15.37

37

An example of inserting a new node N-7: N-7 contacts N-1 for filling its FT

1	4
2	4
3	1

Pred(1) = 4
Succ(1) = 4

1	1
2	1
3	1

Succ(4) = 1
Pred(4) = 1

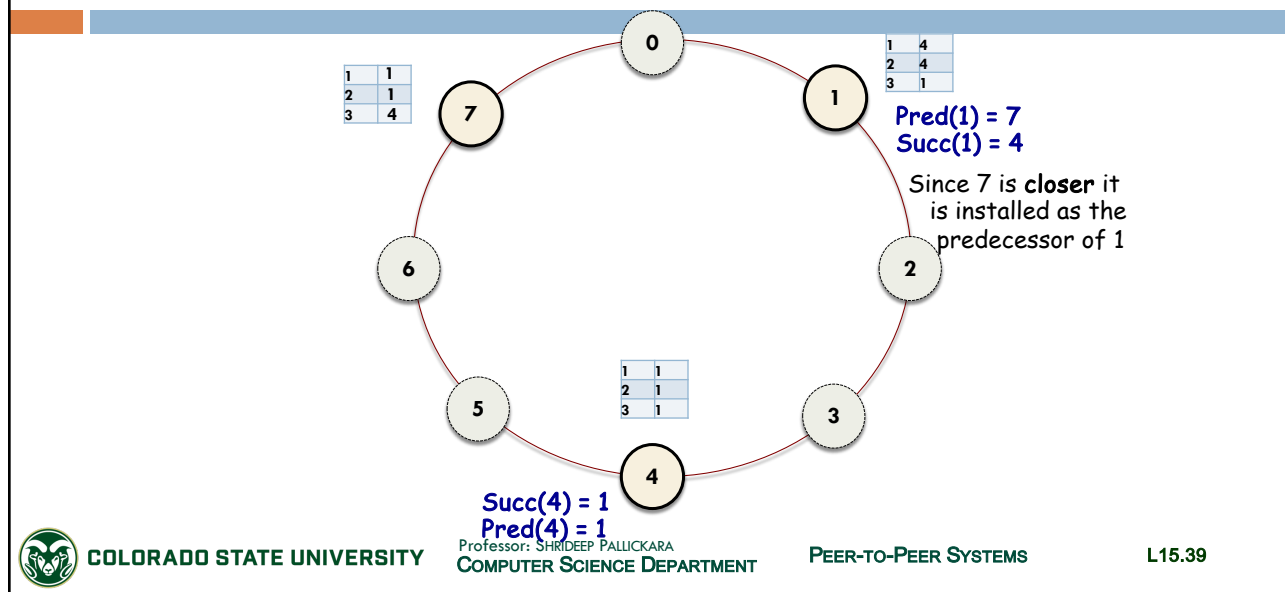
1	1
2	1
3	4

COLORADO STATE UNIVERSITY
 Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

PEER-TO-PEER SYSTEMS L15.38

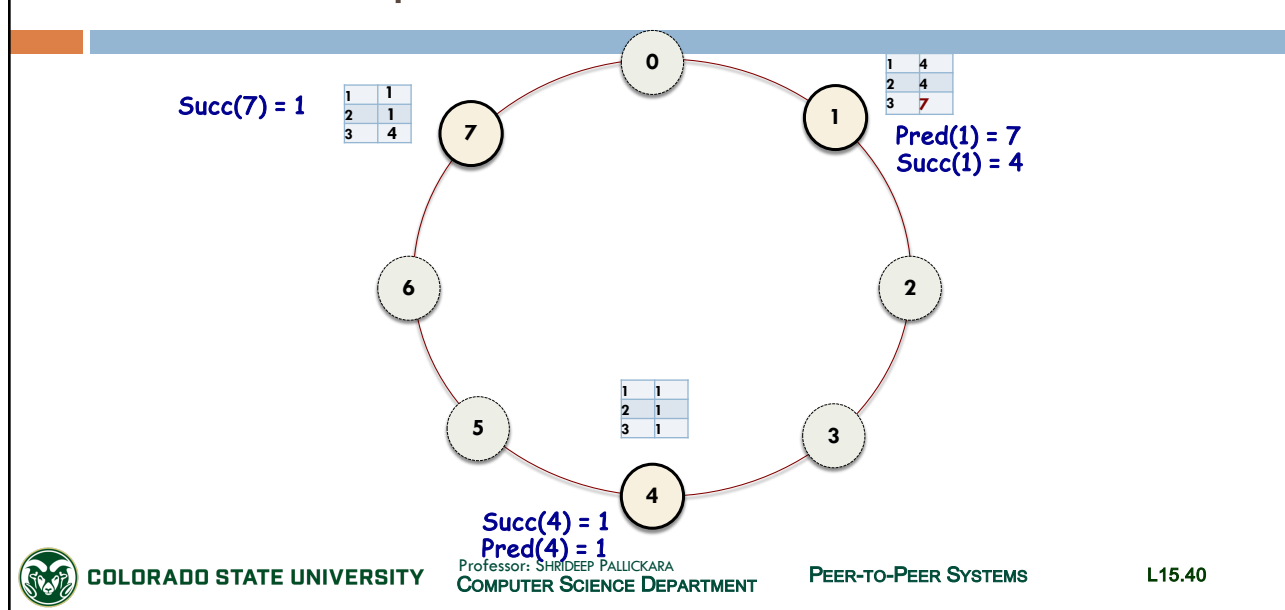
38

N-7 informs N-1 that it (N-7) is now N-1's predecessor

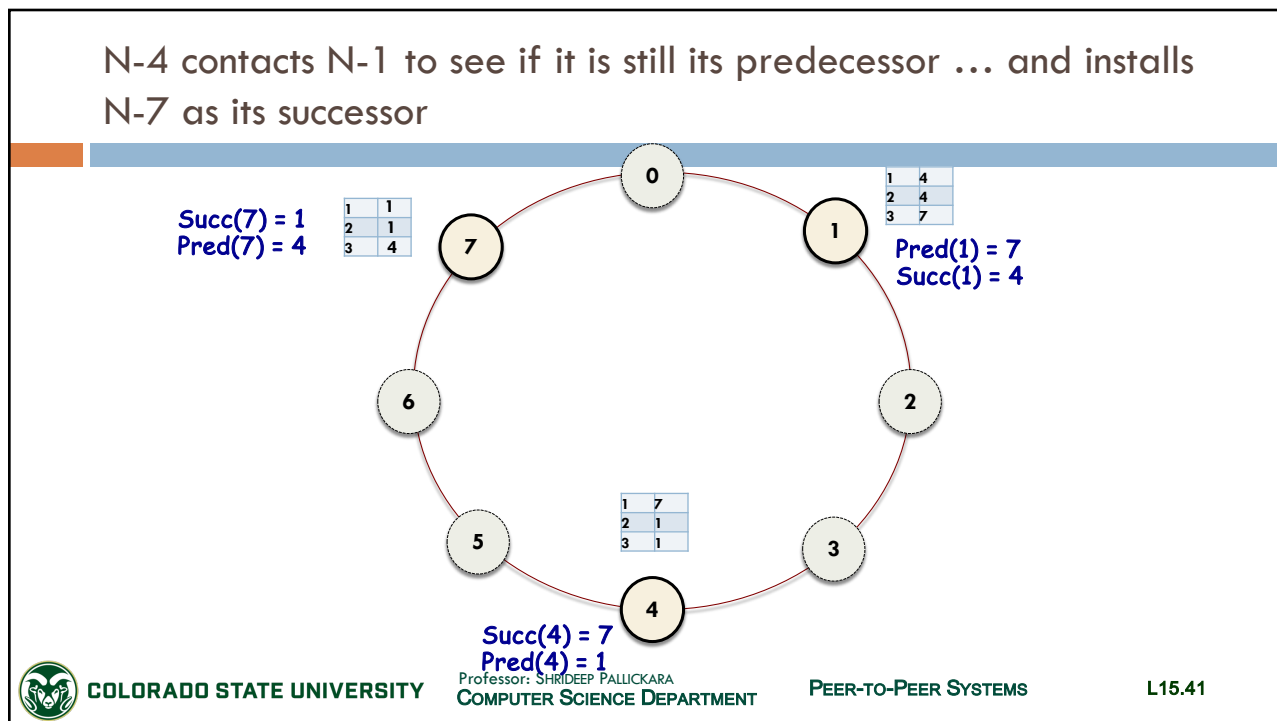


39

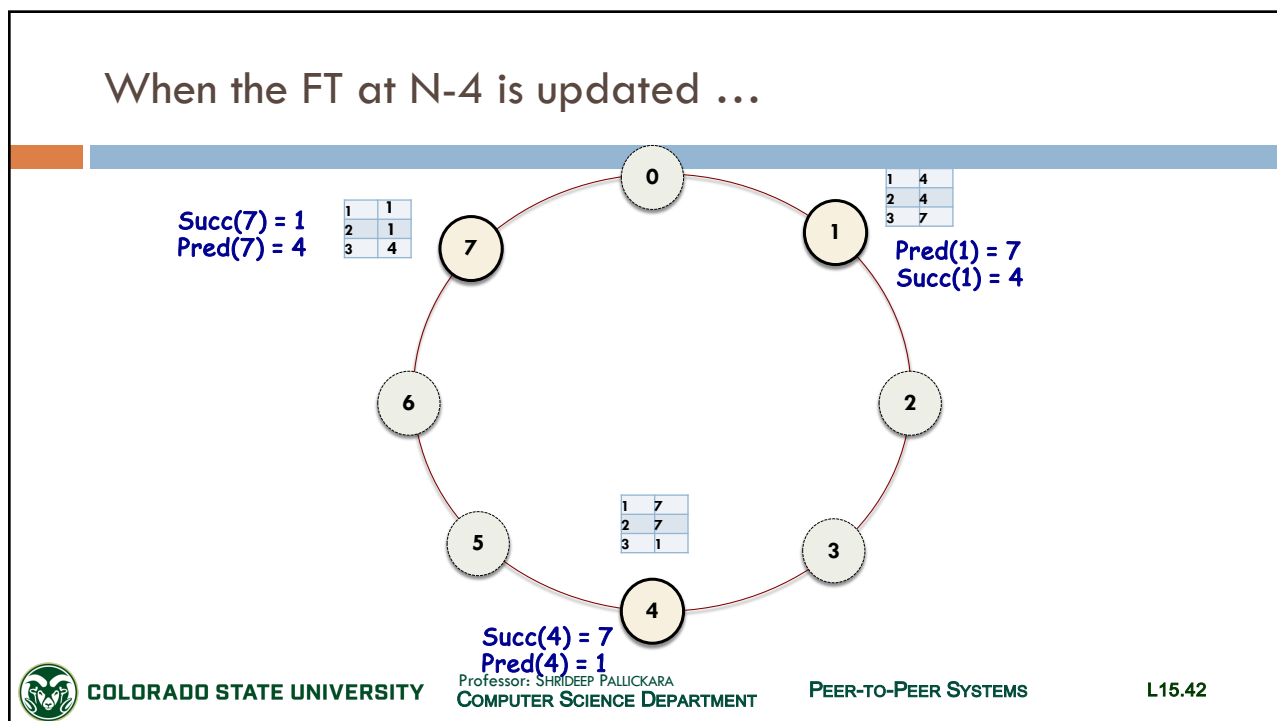
When N-1 updates its FT later on ...



40



41



42

The contents of this slide-set are based on the following references

- *Distributed Systems: Principles and Paradigms*. Andrew S. Tanenbaum and Maarten Van der Steen. 2nd Edition. Prentice Hall. ISBN: 0132392275/978-0132392273. [Chapter 5]
- *Distributed Systems: Concepts and Design*. George Coulouris, Jean Dollimore, Tim Kindberg, Gordon Blair. 5th Edition. Addison Wesley. ISBN: 978-0132143011. [Chapter 10]

