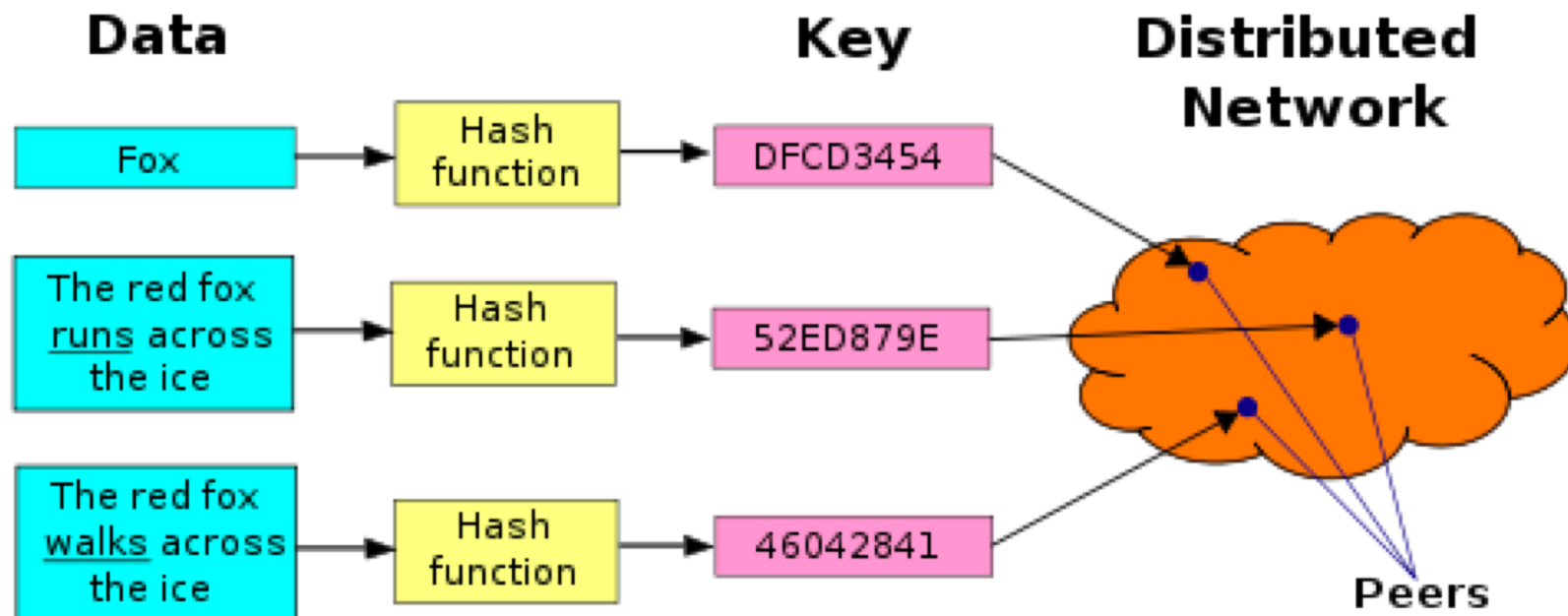


# Distributed Hash Table

-In Algorithm Perspective-

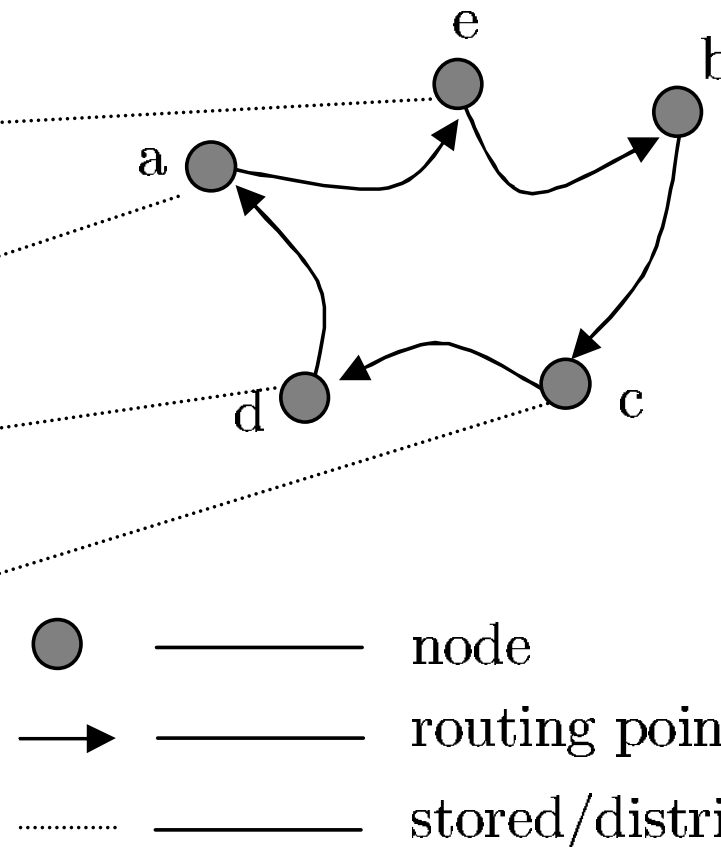
Wenguang Liu

# Distributed Hash Table



# Distributed Hash Table

Key	Value
"abc.txt"	<url <sub>1</sub> >
"me.jpg"	<url <sub>2</sub> >
"music.mp3"	<url <sub>3</sub> >
"piano.mp3"	<url <sub>4</sub> >
"cv.doc"	<url <sub>5</sub> >
"source.zip"	<url <sub>6</sub> >



## Properties:

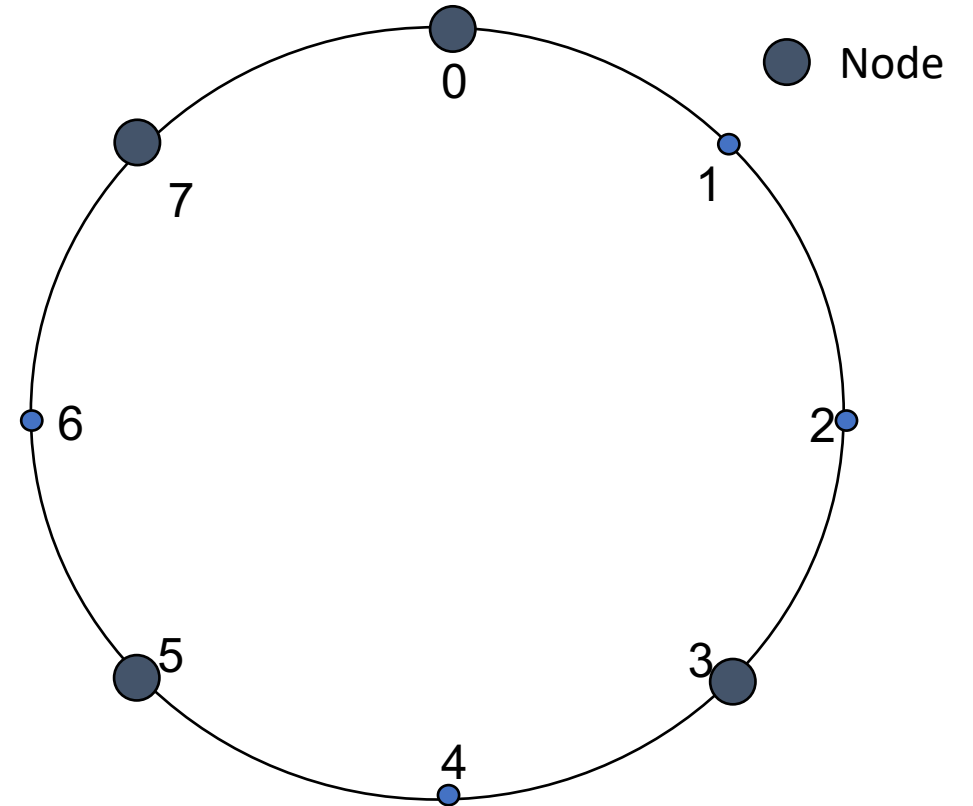
- Scalable
- Dispersed items
- Scales with dynamism
- Self-manage
  - Node JOIN,
  - Node LEAVE,
  - Node FAIL

## Operations:

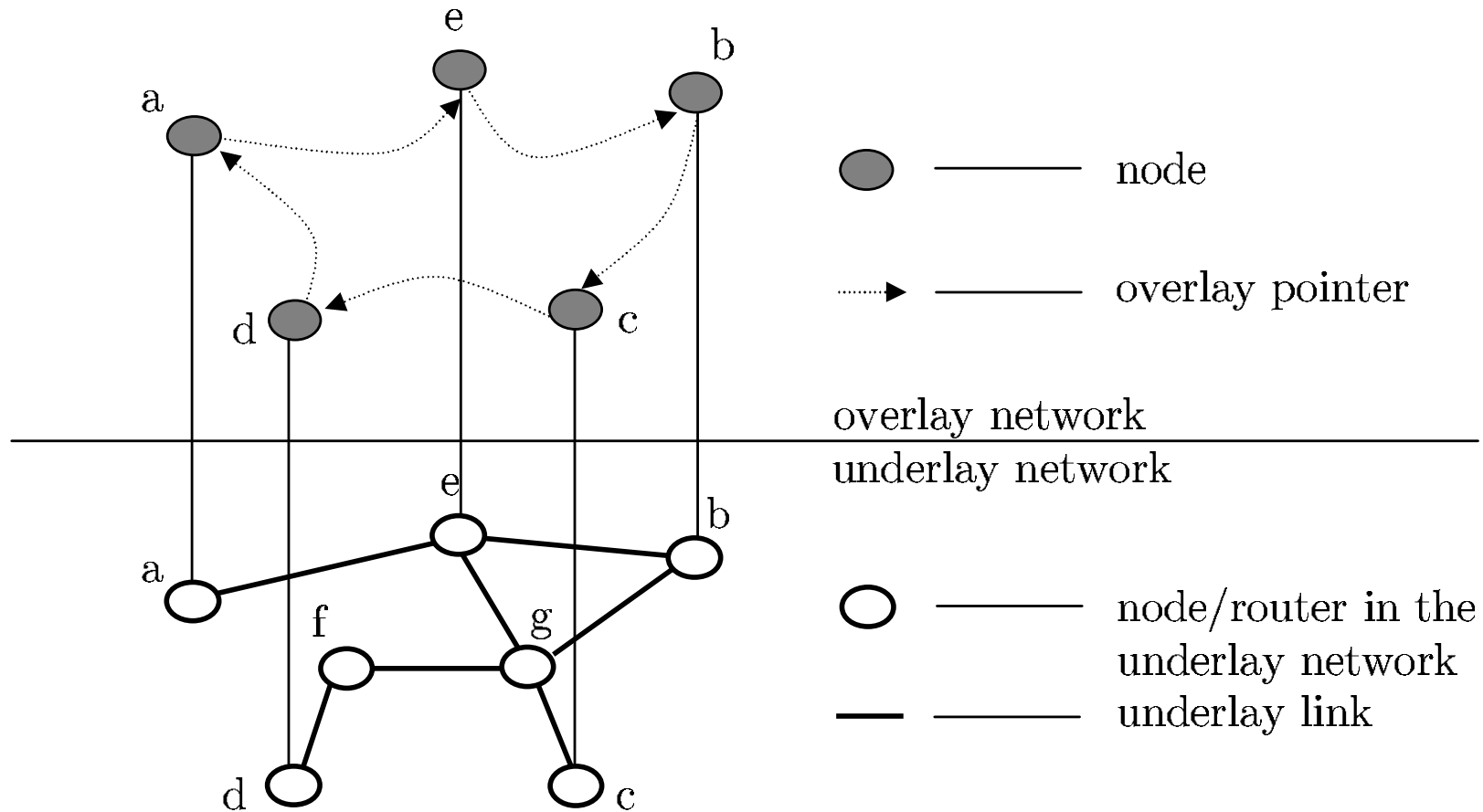
- Node management
- K/V management
- Range query
- Group communication

# Distributed Hash Table

- Structure
  - Centralized
  - Decentralized/P2P
- Concept
  - Keyspace, e.g.  $[0, 2^N)$
  - Keyspace Partition, e.g. Consistent Hashing



# Overlay Networks



# Topic

- Atomic Ring Maintenance
- Routing Maintenance
- Group Communication
- Replication
- Applications

# Atomic Ring Maintenance

# Atomic Ring Maintenance

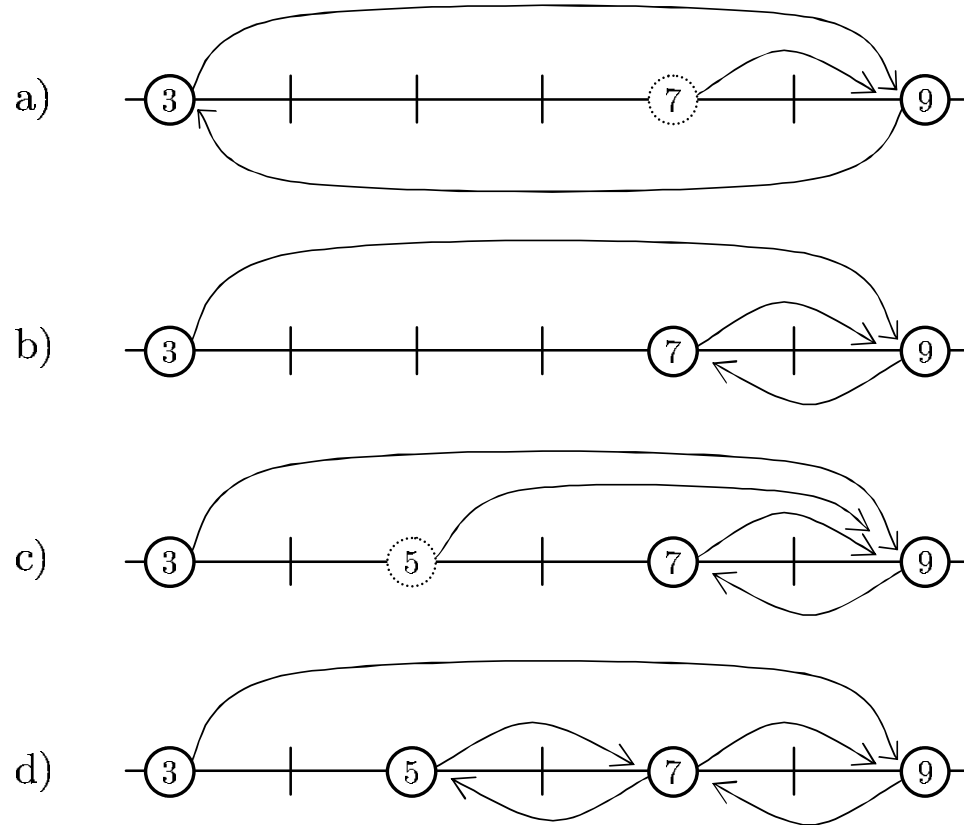


Figure 3.1: Example of inconsistent stabilization.

## Algorithm 1 Chord's periodic stabilization protocol

```

1: procedure  $n$ .STABILIZE()
2:    $p := succ.GetPredecessor()$ 
3:   if  $p \in (n, succ)$  then
4:      $succ := p$ 
5:   end if
6:    $succ.Notify(n)$ 
7: end procedure

8: procedure  $n$ .GETPREDECESSOR()
9:   return  $pred$ 
10: end procedure

11: procedure  $n$ .NOTIFY( $p$ )
12:   if  $p \in (pred, n]$  then
13:      $pred := p$ 
14:   end if
15: end procedure

```

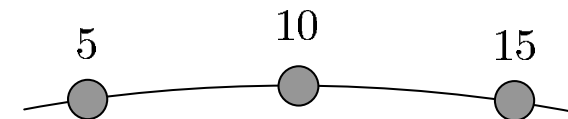


Figure 3.2: Perfect system state before a leave operation.  
10 and 15 leave concurrently



# Atomic Ring Maintenance

- Concurrency Control strategies:
  - Lock the whole ring;
  - Three locks: predecessor's lock, successor's lock, joining/leaving node's lock;
  - Two locks: Joining/leaving node's own lock, and its' successor's lock;
    - Used in this paper
    - Suffer with Dining philosophers
- Safety
- Liveness
  - Asymmetric locking
  - Randomized locking: release all locks when timeout.

# Atomic Ring Maintenance

- Asymmetric Locking: *LockQueue*

---

**Algorithm 2** Asymmetric locking with forwarding

---

1: <b>procedure</b> <i>n</i> .JOIN( <i>succ</i> )	▷ Join the ring with <i>succ</i> as successor	
2: <i>Leaving</i> := false	▷ Initialize variable	
3: <i>LockQueue</i> .ENQUEUE( <i>n</i> )	▷ Enqueue request to local lock	
4: <i>slock</i> := GETSUCCLOCK()		
5: <i>pred</i> := <i>succ</i> . <i>pred</i>		
6: <i>pred</i> . <i>succ</i> := <i>n</i>		
7: <i>succ</i> . <i>pred</i> := <i>n</i>		
8: <i>LockQueue</i> := <i>succ</i> . <i>LockQueue</i>	▷ Copy successor's queue	
9: <i>LockQueue</i> .FILTER( <i>((pred, n])</i> )	▷ Keep requests in the range	
10: <i>succ</i> . <i>LockQueue</i> .FILTER( <i>((n, pred])</i> )	▷ Keep requests in the range	
11: <i>LockQueue</i> .DEQUEUE()	▷ Remove local request	
12:   RELEASELOCK( <i>slock</i> )		
13: <b>end procedure</b>		

14: <b>procedure</b> <i>n</i> .LEAVE()	▷ Leave the ring
15: <b>if</b> <i>n</i> > <i>succ</i> <b>then</b>	▷ Asymmetric Locking
16: <i>slock</i> := GETSUCCLOCK()	
17: <i>Leaving</i> := true	▷ Enable forwarding
18: <i>LockQueue</i> .ENQUEUE( <i>n</i> )	▷ Enqueue request to local lock
19: <b>else</b>	
20: <i>Leaving</i> := true	▷ Enable forwarding
21: <i>LockQueue</i> .ENQUEUE( <i>n</i> )	▷ Enqueue request to local lock
22: <i>slock</i> := GETSUCCLOCK()	
23: <b>end if</b>	
24: <i>pred</i> . <i>succ</i> := <i>succ</i>	
25: <i>succ</i> . <i>pred</i> := <i>pred</i>	
26: <i>LockQueue</i> .DEQUEUE()	▷ Remove local request
27:   RELEASELOCK( <i>slock</i> )	
28: <b>end procedure</b>	

# Atomic Ring Maintenance

- Lookup consistency in the presence of JOINs: *JoinForward*

---

**Algorithm 4** Pointer updates during joins

---

```
1: event  $n$ .UPDATEJOIN() from  $n$  ▷ Assuming  $succ$  is correct
2:   sendto  $succ$ .UPDATEPRED()
3: end event

4: event  $n$ .UPDATEPRED() from  $m$ 
5:    $JoinForward := \text{true}$  ▷ Forwarding Enabled
6:   sendto  $m$ .JOINPOINT( $pred$ ) ▷ Join Point
7:    $oldpred := pred$ 
8:    $pred := m$ 
9: end event
```

```
10: event  $n$ .JOINPOINT( $p$ ) from  $m$ 
11:    $pred := p$ 
12:    $succ := m$ 
13:   sendto  $pred$ .UPDATESUCC()
14: end event
```

```
15: event  $n$ .UPDATESUCC() from  $m$ 
16:   sendto  $succ$ .STOPFORWARDING()
17:    $succ := m$ 
18: end event
```

```
19: event  $n$ .STOPFORWARDING() from  $m$ 
20:    $JoinForward := \text{false}$ 
21:   sendto  $pred$ .FINISH()
22: end event
```

---

# Atomic Ring Maintenance

- Lookup consistency in the presence of Leaves: *LeaveForward*

---

**Algorithm 5** Pointer updates during leaves

---

```
1: event n.UPDATELEAVE() from n
2:   LeaveForward := true                                ▷ Forwarding Enabled
3:   sendto succ.LEAVEPOINT(pred)
4: end event

5: event n.LEAVEPOINT(p) from m
6:   pred := p
7:   sendto pred.UPDATESUCC()
8: end event

9: event n.UPDATESUCC() from m
10:  sendto succ.STOPFORWARDING()
11:  succ := m
12: end event

13: event n.STOPFORWARDING() from m
14:   LeaveForward := false                                ▷ Forwarding Disabled
15: end event
```

---

# Atomic Ring Maintenance

- Combine together

---

## Algorithm 6 Lookup algorithm

---

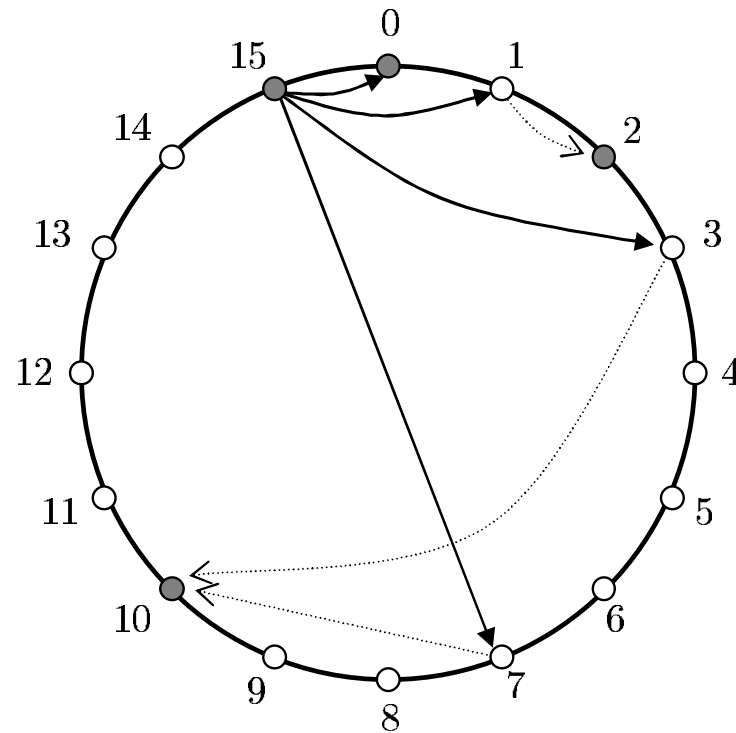
```
1: event  $n$ .LOOKUP( $id, src$ ) from  $m$ 
2:   if  $JoinForward = \text{true}$  and  $m = oldpred$  then
3:     sendto  $pred$ .LOOKUP( $id, src$ )           ▷ Redirect Message
4:   else if  $LeaveForward = \text{true}$  then
5:     sendto  $succ$ .LOOKUP( $id, src$ )           ▷ Redirect Message
6:   else if  $pred \neq \text{nil}$  and  $id \in (pred, n]$  then
7:     sendto  $src$ .LOOKUPDONE( $n$ )
8:   else
9:     sendto  $succ$ .LOOKUP( $id, src$ )
10:  end if
11: end event
```

---

# Routing Maintenance

# Routing Maintenance

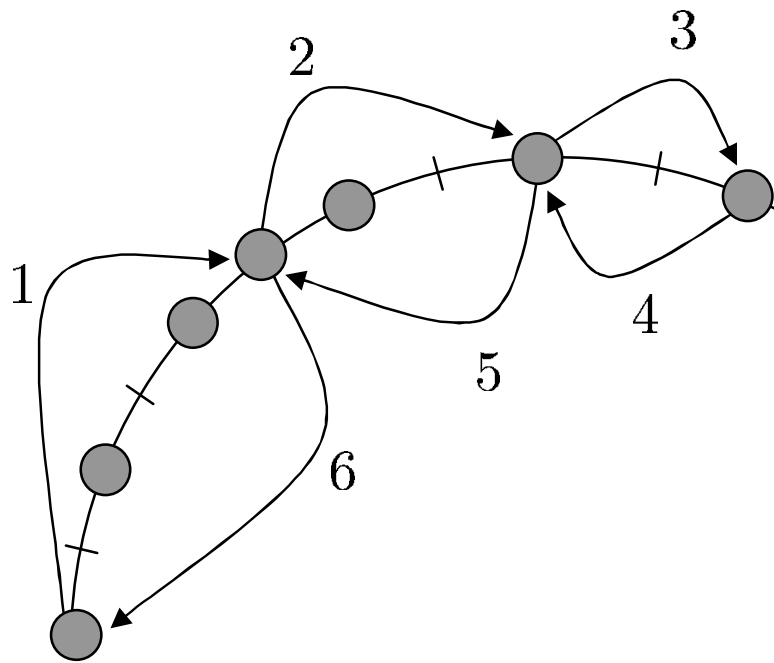
- Extension to the ring: pred, succ, successor-list



$$p \oplus 2^{L-1},$$

# Routing Maintenance

- Lookup: strategy-Recursive Lookup



---

**Algorithm 12** Recursive lookup algorithm

---

```
1: procedure  $n$ .LOOKUP( $i$ , OP)
2:   if TERMINATE( $i$ ) then
3:      $p := \text{NEXT\_HOP}(i)$ 
4:      $res := p.OP(i)$ 
5:     return  $res$ 
6:   else
7:      $m := \text{NEXT\_HOP}(i)$ 
8:     return  $m$ .LOOKUP( $i$ , OP)
9:   end if
10: end procedure
```

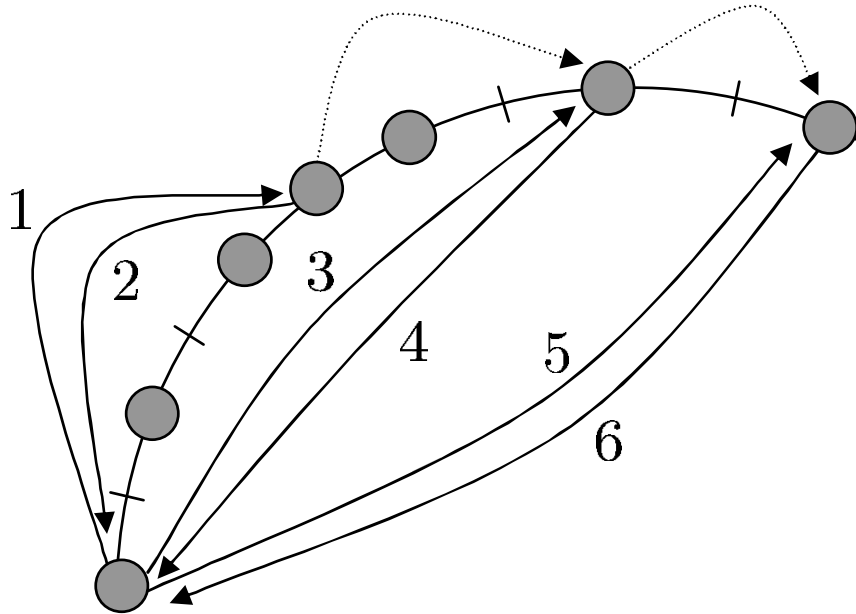
---

▷ OP could carry parameters



# Routing Maintenance

- Lookup: strategy-Iterative Lookup



---

## Algorithm 13 Iterative lookup algorithm

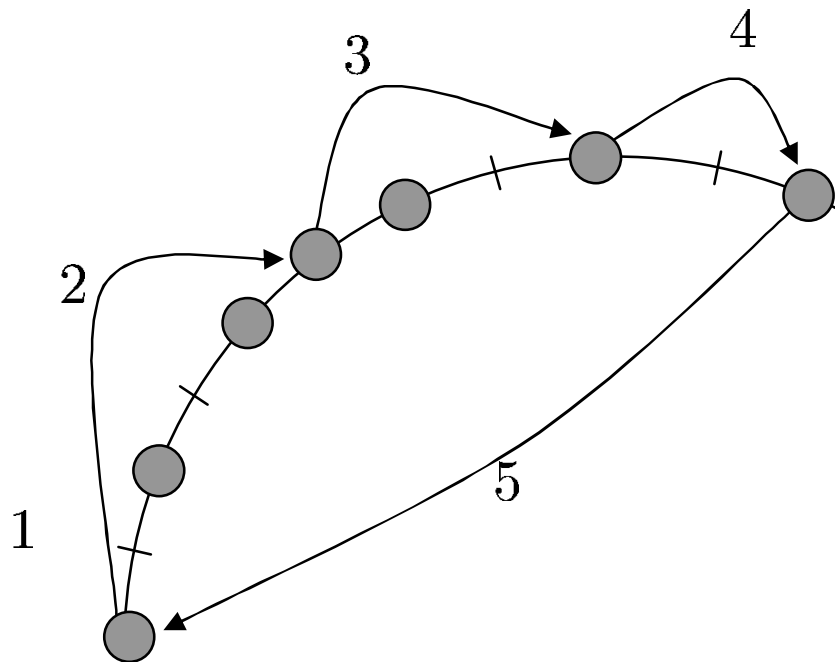
---

```
1: procedure  $n$ .LOOKUP( $i$ , OP)
2:    $m := n$ 
3:   while not  $m$ .TERMINATE( $i$ ) do
4:      $m := m$ .NEXT_HOP( $i$ )
5:   end while
6:    $p := m$ .NEXT_HOP( $i$ )
7:   return  $p$ .OP( $i$ )
8: end procedure
```

---

# Routing Maintenance

- Lookup: strategy-Transitive Lookup



---

**Algorithm 14** Transitive lookup algorithm

---

```
1: procedure  $n$ .LOOKUP( $i$ , OP)
2:   sendto  $n$ .LOOKUP_AUX( $n$ ,  $i$ , OP)
3:   receive LOOKUP_RES( $r$ ) from  $q$ 
4:   return  $r$ 
5: end procedure

6: event  $n$ .LOOKUP_AUX( $q$ ,  $i$ , OP) from  $m$ 
7:   if TERMINATE( $i$ ) then
8:      $p := \text{NEXT\_HOP}(i)$ 
9:     sendto  $p$ .LOOKUP_FIN( $q$ ,  $i$ , OP)
10:  else
11:     $p := \text{NEXT\_HOP}(i)$ 
12:    sendto  $p$ .LOOKUP_AUX( $q$ ,  $i$ , OP)
13:  end if
14: end event

15: event  $n$ .LOOKUP_FIN( $q$ ,  $i$ , OP) from  $m$ 
16:    $r := \text{OP}(i)$ 
17:   sendto  $q$ .LOOKUP_RES( $r$ )
18: end event
```

---

# Routing Maintenance

- Greedy Lookup Algorithm
  - $rt(i)$ : successors sorted by routing distance asc.

---

**Algorithm 15** Greedy lookup

---

```
1: procedure  $n.TERMINATE(i)$ 
2:   return  $i \in (n, succ]$ 
3: end procedure

1: procedure  $n.NEXT\_HOP(i)$ 
2:   if  $TERMINATE(i)$  then
3:     return  $succ$ 
4:   else
5:      $r := succ$ 
6:     for  $j := 1$  to  $K$  do
7:       if  $rt(j) \in (n, i)$  then
8:          $r := rt(j)$ 
9:       end if
10:    end for
11:    return  $r$ 
12:   end if
13: end procedure
```

# Group Communication

# Group Communication

- Motivation
  - Exactly match → Wildcard expression
- Desirable properties
  - Termination
  - Coverage,
    - All the *designated* nodes that are *reachable* should receive the message
  - Non-redundancy
    - Never receive a message more than once.

# Group Communication

- Broadcast algorithms – simple broadcast

---

**Algorithm 19** Simple broadcast algorithm

---

```
1: event  $n$ .STARTSIMPLEBCAST( $msg$ ) from  $app$ 
2:   sendto  $n$ .SIMPLEBCAST( $msg$ ,  $n$ )           ▷ Local message to itself
3: end event
```

```
1: event  $n$ .SIMPLEBCAST( $msg$ ,  $limit$ ) from  $m$ 
2:   Deliver( $msg$ )                               ▷ Deliver  $msg$  to application
3:   for  $i := M$  downto 1 do                       ▷ Node has  $M$  unique pointers
4:     if  $u(i) \in (n, limit)$  then
5:       sendto  $u(i)$ .SIMPLEBCAST( $msg$ ,  $limit$ )
6:        $limit := u(i)$ 
7:     end if
8:   end for
9: end event
```

---

# Group Communication

- Broadcast algorithms – simple broadcast with feedback

---

**Algorithm 20** Simple broadcast with feedback algorithm

---

```
1: event  $n$ .STARTBCAST( $msg$ ) from  $app$ 
2:   sendto  $n$ .BCAST( $msg$ ,  $n$ )           ▷ Local message to itself
3: end event

1: event  $n$ .BCAST( $msg$ ,  $limit$ ) from  $m$ 
2:    $FB := \text{Deliver}(msg)$            ▷ Deliver  $msg$  and get set of feedback
3:    $par := n$ 
4:    $Ack := \emptyset$ 
5:   for  $i := M$  downto 1 do           ▷ Node has  $M$  unique pointers
6:     if  $u(i) \in (n, limit)$  then
7:       sendto  $u(i)$ .BCAST( $msg$ ,  $limit$ )
8:        $Ack := Ack \cup \{u(i)\}$ 
9:        $limit := u(i)$ 
10:    end if
11:  end for
12:  if  $Ack = \emptyset$  then
13:    sendto  $par$ .BCASTRESP( $FB$ )
14:  end if
15: end event
```

```
1: event  $n$ .BCASTRESP( $F$ ) from  $m$ 
2:   if  $m = n$  then
3:     sendto  $app$ .BCASTTERM( $FB$ )
4:   else
5:      $Ack := Ack - \{m\}$ 
6:      $FB := FB \cup F$ 
7:     if  $Ack = \emptyset$  then
8:       sendto  $par$ .BCASTRESP( $FB$ )
9:     end if
10:  end if
11: end event
```

---

# Group Communication

- Bulk Operations,
  - To designated nodes with identifier in **Bulk Set  $I$** .

---

**Algorithm 21** Bulk operation algorithm

---

```
1: event  $n$ .BULK( $I$ ,  $msg$ ) from  $m$ 
2:   if  $n \in I$  then
3:     Deliver( $msg$ )                                ▷ Deliver  $msg$  to application
4:   end if
5:    $limit := n$ 
6:   for  $i := M$  downto 1 do                        ▷ Node has  $M$  unique pointers
7:      $J := [u(i), limit)$ 
8:     if  $I \cap J \neq \emptyset$  then
9:       sendto  $u(i)$ .BULK( $I \cap J$ ,  $msg$ )
10:       $I := I - J$                                 ▷ Same as  $I := I - (I \cap J)$ 
11:       $limit := u(i)$ 
12:     end if
13:   end for
14: end event
```

---



# Group Communication

- Bulk Operations with feedback

---

**Algorithm 22** Bulk operation with feedback algorithm

---

```
1: event  $n$ .BULKFEED( $I, msg$ ) from  $m$ 
2:   if  $n \in I$  then
3:      $FB := \text{Deliver}(msg)$  ▷ Deliver and get set of feedback
4:   else
5:      $FB := \emptyset$  ▷ No feedback
6:   end if
7:    $par := m$ 
8:    $Ack := \emptyset$ 
9:    $limit := n$ 
10:  for  $i := M$  downto 1 do ▷ Node has  $M$  unique pointers
11:     $J := [u(i), limit)$ 
12:    if  $I \cap J \neq \emptyset$  then
13:      sendto  $u(i)$ .BULKFEED( $I \cap J, msg$ )
14:       $I := I - J$  ▷ Same as  $I := I - (I \cap J)$ 
15:       $Ack := Ack \cup \{u(i)\}$ 
16:       $limit := u(i)$ 
17:    end if
18:  end for
19:  if  $Ack = \emptyset$  then
20:    sendto  $par$ .BULKRESP( $FB$ )
21:  end if
22: end event
```

```
1: event  $n$ .BULKRESP( $F$ ) from  $m$ 
2:   if  $m = n$  then
3:     sendto  $app$ .BULKFEEDTERM( $FB$ )
4:   else
5:      $Ack := Ack - \{m\}$ 
6:      $FB := FB \cup F$ 
7:     if  $Ack = \emptyset$  then
8:       sendto  $par$ .BULKRESP( $FB$ )
9:     end if
10:  end if
11: end event
```

---

# Group Communication

- Bulk Owner Operations

---

**Algorithm 24** Bulk owner operation algorithm

---

```
1: event  $n$ .STARTBULKOWN( $I, msg$ ) from  $m$ 
2:   sendto  $n$ .BULKOWN( $I, I, n, msg$ )      ▷ Local message to itself
3: end event
```

```
1: event  $n$ .BULKOWN( $I, R, next, msg$ ) from  $m$ 
2:    $MS := R \cap (u(M), n]$                   ▷  $u(M)$  is same as  $pred$ 
3:   if  $MS \neq \emptyset$  then
4:     Deliver( $msg, MS$ )                      ▷ App is responsible for ids in  $MS$ 
5:   end if
6:    $limit := n$ 
7:    $lnext := next$ 
8:    $sentsucc := \text{false}$ 
9:   for  $i := M$  downto 1 do                  ▷ Node has  $M$  unique pointers
10:     $J := (u(i), limit]$ 
11:    if  $I \cap J \neq \emptyset$  then
12:       $K := (u(i-1), u(i)]$ 
13:      sendto  $u(i)$ .BULKOWN( $I \cap J, I \cap K, lnext, msg$ )
14:       $I := I - J$                           ▷ Same as  $I := I - (I \cap J)$ 
15:       $limit := u(i)$ 
16:       $lnext := u(i)$ 
17:      if  $i = 1$  then
18:         $sentsucc := \text{true}$ 
19:      end if
20:    end if
21:  end for
22:   $J := (n, u(1)]$ 
23:  if  $I \cap J \neq \emptyset$  and  $sentsucc = \text{false}$  and  $next \neq u(1)$  then
24:    sendto  $u(1)$ .BULKOWN( $\emptyset, I \cap J, limit, msg$ )
25:  end if
26: end event
```

---

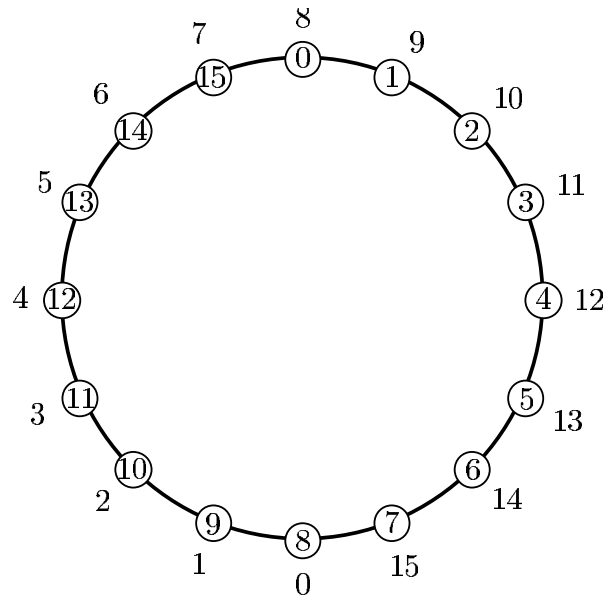
# Group Communication

- Fault-tolerance
  - Use timeouts to detect node failure.
- Efficient overlay multicast
  - Multicast-Group: creating, joining, leaving.
  - IP multicast integration

Replication

# Replication

- Symmetric replication scheme
  - the identifier space is partitioned into  $N/f$  equivalence classes,
    - identifiers in an equivalence class are all associated with each other
  - if the identifier  $i$  is associated with the set of identifiers  $r1, \dots, rf$ , then the identifier  $rx$ , for  $1 \leq x \leq f$ , is associated with the identifiers  $r1, \dots, rf$  as well.



$$r(i, x) = i \oplus (x - 1) \frac{N}{f}$$

$$1 \leq x \leq f, \quad f = 2.$$

$$1 \equiv 9 \pmod{8}.$$

# Replication

- Join and Leave Algo

---

**Algorithm 25** Symmetric replication for joins and leaves

---

```
1: event  $n$ .JOINREPLICATION() from  $m$ 
2:   sendto  $succ$ .RETRIEVEITEMS( $pred, n, n$ )
3: end event

4: event  $n$ .LEAVEREPLICATION() from  $m$ 
5:   sendto  $n$ .RETRIEVEITEMS( $pred, n, succ$ )
6: end event
```

```
7: event  $n$ .RETRIEVEITEMS( $start, end, p$ ) from  $m$ 
8:   for  $r := 1$  to  $f$  do
9:      $items[r] := \emptyset$ 
10:     $i := start$ 
11:    while  $i \neq end$  do
12:       $i := i \oplus 1$ 
13:       $items[r][i] := localHashTable[r][i]$ 
14:    end while
15:  end for
16:  sendto  $p$ .REPLICATE( $items, start, end$ )
17: end event

18: event  $n$ .REPLICATE( $items, start, end$ ) from  $m$ 
19:   for  $r := 1$  to  $f$  do
20:      $i := start$ 
21:     while  $i \neq end$  do
22:        $i := i \oplus 1$ 
23:        $localHashTable[r][i] := items[r][i]$ 
24:     end while
25:   end for
26: end event
```

# Replication

- Lookup and Insertion

---

**Algorithm 26** Lookup and item insertion for symmetric replication

---

```
1: event  $n$ .INSERTITEM( $key, value$ ) from  $app$ 
2:   for  $r := 1$  to  $f$  do
3:      $replicaKey := key \oplus (r - 1) \frac{N}{f}$ 
4:      $n$ .LOOKUP( $replicaKey$ , ADDITEM( $replicaKey, value, r$ ))
5:   end for
6: end event

7: procedure  $n$ .ADDITEM( $key, value, r$ )
8:    $localHashTable[key][r] := value$ 
9: end procedure

10: event  $n$ .LOOKUPITEM( $key, r$ ) from  $app$ 
11:    $replicaKey := key \oplus (r - 1) \frac{N}{f}$ 
12:   LOOKUP( $replicaKey$ , GETITEM( $replicaKey, r$ ))
13: end event

14: procedure  $n$ .GETITEM( $key, r$ )
15:   return  $localHashTable[r][key]$ 
16: end procedure
```

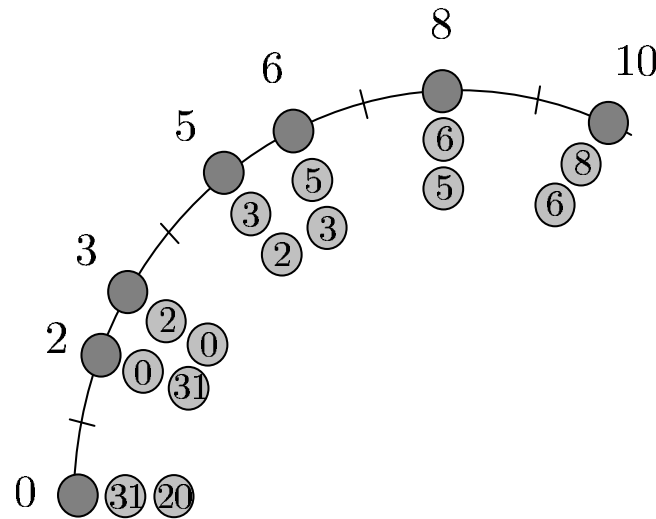
# Replication

- Symmetric replication scheme
  - Symmetric replication enables an application to make parallel lookups to exactly  $k$  replicas of an item, where  $k \leq f$  if the replication degree is  $f$ .
    - to speed up the lookup process
  - a join or a leave only requires the joining or leaving node to exchange data with its successor prior to joining or leaving.



# Replication

- Other replica placement schemas
  - Multiple Hash Functions: hash key with  $f$  hash functions
  - Successor lists: store at the  $f$  closest successors
  - Leaf sets: store on  $\lfloor f/2 \rfloor$  closest successors and  $\lfloor f/2 \rfloor$  closest predecessors



successor-list replication

Applications

# Applications

- Storage system
  - PAST[Storage management and caching in past, a large-scale, persistent peer-to-peer storage utility]
  - CFS[Wide-area cooperative storage with CFS]
- Host discovery and mobility
- Web caching and web servers
- Publish/subscribe systems, e.g. FeedTree
- P2P, e.g. BitTorrent

# Take-home Message

- A Randomized locking mechanism with node's and succ's lock
  - to support atomic ring management (JOIN, LEAVE, LOOKUP)
- Routing maintainence
  - by additional routing pointer augment.
  - With Recursive/Iterative/Transitive/Greedy lookup algo.
- Provide the Algo for Group communication
  - Broadcast, Bulk, Bulk own
- Provide the symmetric replication mechanism
  - to augment the robustness

# References

- ALI GHODSI, Distributed k-ary System: Algorithms for Distributed Hash Tables.
- [https://en.wikipedia.org/wiki/Distributed\\_hash\\_table](https://en.wikipedia.org/wiki/Distributed_hash_table)