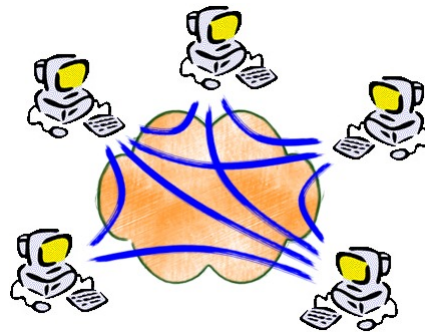**Chapter 6**

# Peer-to-Peer systems

1. Introduction
2. Overlays
3. Distributed Hash Tables

1

---

**Chapter 7**

# Peer-to-Peer systems

1. **Introduction**
   1. Why P2P systems ?
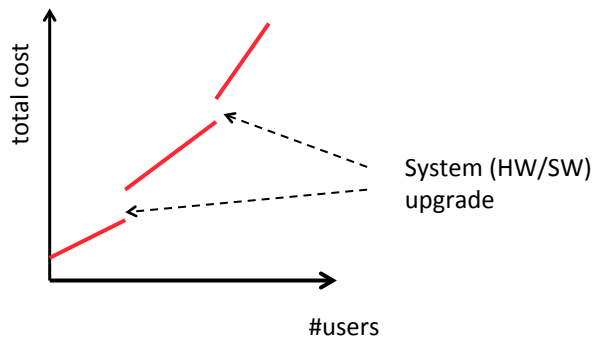   2. P2P generations
2. **Overlays**
3. **Distributed Hash Tables**

2

# The rationale for P2P

**Scaling and cost of client server systems**
popular services need large server infrastructures
-> high investment and operational cost

System (HW/SW)
upgrade

total cost

#users

**Observation**
• performance of edge devices grows (gap client – server closes)
• performance only occasionally fully used

3

---

# The rationale for P2P

"How to unleash the power of Internet's dark matter ?"

**P2P philosophy**
make heavily use of edge resources
            (CPU, storage, I/O devices, Information, ...)
#users ~ #edge devices
        -> infrastructure grows together
            (automatically) with user base

Need for software systems that can survive without central servers.

4

2

## Some numbers

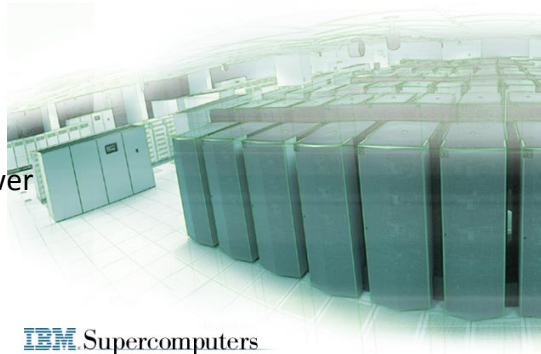**Edge resource estimation**

In total
- #hosts : $10^9$
- processor performance : 2 GFLOPs
- disk space : 500 GB

Available
- 1% connected hosts
- 50% CPU power
- 10% disk space

-> "P2P supercomputer"
- $10^4$ PFLOPs CPU power
- $5 \cdot 10^5$ PB disk space

**A "real" supercomputer**

Nov 2012 :

1.  Sequoia IBM (20 PFLOPs, 8 MW) @DoE

2. K Computer Fujitsu (11 PFLOPs, 12 MW)  @Riken

**IBM** Supercomputers
www.ibm.com/rs6000/supercomputer

5

---

## Why is P2P difficult ?

**Nodes**
- large number of nodes
- unstable
- under control of end users

**Interconnection infrastructure**
- slow
- unreliable

**No central control**
- more complex managment
  (control infrastructure can fail)

=> 
- decentralized control
- self-organization
  - handling failing nodes
  - spreading load dynamically

6

# Characteristics of P2P systems

**Shared characteristics**
- Users contribute to the total pool of available resources (avoid free-riding)
- All nodes in principle equal
- System operation independent of centralized control

**Application areas**
- file sharing
- collaboration tools (Groove)
- communication (VoIP P2P [Skype], chatting [Jabber])
- CPU scavenging (SETI@Home, Folding@Home, ...)

**Important problems**
- data placement and lookup
- routing (use the network bandwidth efficiently)
- provide anonymity
- self-organization (self-management)

7

# The history of P2P

**Generation I : File sharing with centralized control**
- index maintained in a centralized infrastructure
- download purely P2P
        -> poor scalability
        -> vulnarable to attacks

**Generation II : File sharing with decentralized control**
- self-organization through overlay construction
- two variants:
    - pure P2P (all nodes identical)
    - hybrid P2P (hierarchical structuring)
- scalable and robust

**Generation III : P2P middleware**
- middleware services offered:
    - data placement
    - data lookup
    - automatic replication/caching
    - authentication/security
- applications built on top of P2P middleware
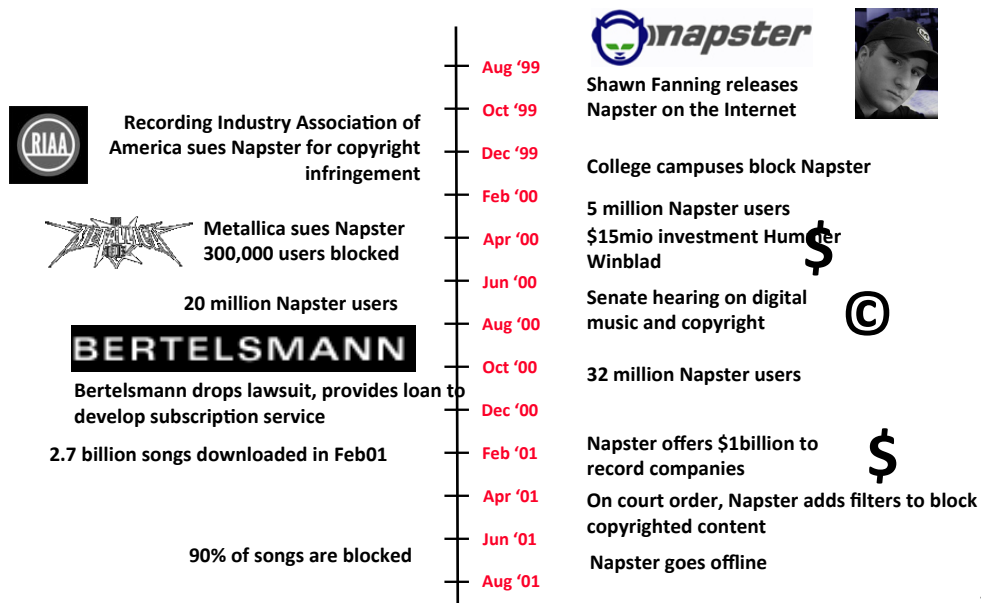
8

## P2P-architectures

|  | mediated | pure | hybrid |
|---|---|---|---|
| data traffic | P2P | P2P | P2P |
| control traffic | client-server | P2P | local : client-server<br>on distance : P2P |
| efficiency | + efficient search<br>+ efficient control | - inefficient search<br>- BW consuming | +/- |
| scalability | - control hot spot<br>(mirrors needed ?) | - BW needed<br>grows rapidly | good compromise |
| robustness | - single point of failure<br>- easy to attack | + graceful degradation<br>+ difficult to attack | ? |
| accountability | easy | difficult | difficult |

Generation I (mediated)      Generation II (pure, hybrid)

9

---

## Generation I : Rise and fall of Napster

**Left side (RIAA / industry):**

Recording Industry Association of America sues Napster for copyright infringement

Metallica sues Napster
300,000 users blocked

20 million Napster users

**BERTELSMANN**

Bertelsmann drops lawsuit, provides loan to develop subscription service

2.7 billion songs downloaded in Feb01

90% of songs are blocked

**Timeline:**

Aug '99
Oct '99
Dec '99
Feb '00
Apr '00
Jun '00
Aug '00
Oct '00
Dec '00
Feb '01
Apr '01
Jun '01
Aug '01

**Right side (Napster):**

Shawn Fanning releases Napster on the Internet

College campuses block Napster

5 million Napster users
$15mio investment Hummer Winblad

Senate hearing on digital music and copyright

32 million Napster users

Napster offers $1billion to record companies

On court order, Napster adds filters to block copyrighted content

Napster goes offline

10

*5*

## Generation I : SETI@Home

**SETI**

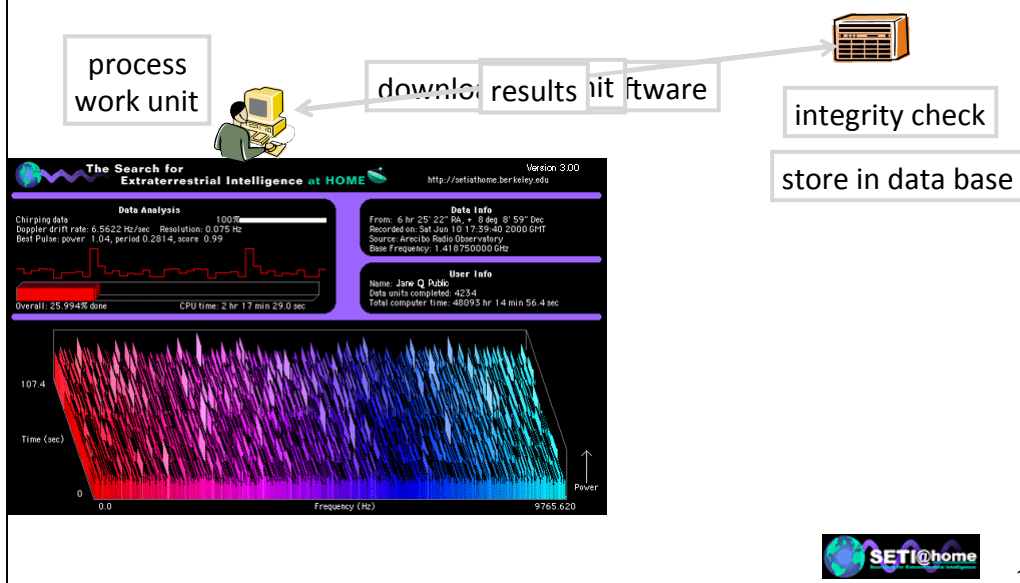="Search for extraterrestrial Intelligence"
- started in 1998 as a 2 year project
- 4 M users signed up
- Radio telescope data sent to clients for digital signal analysis
- Nodes process data when cycles are available
  (works as screen saver)
- Using resources to allow better signal analysis

35 GB/tape
16 hours recorded data

10 tapes/week, 350 GB

≈10 000 0.3 MB work units

11

---

## SETI : How it works

process work unit

download results unit ftware

integrity check

store in data base

**The Search for Extraterrestrial Intelligence at HOME**
Version 3.00
http://setiathome.berkeley.edu

**Data Analysis**  100%
Chirping data
Doppler drift rate: 6.5622 Hz/sec   Resolution: 0.075 Hz
Best Pulse: power 1.04, period 0.2814, score 0.99

Overall: 25.994% done   CPU time: 2 hr 17 min 29.0 sec

**Data Info**
From: 6 hr 25' 22" RA, + 8 deg 8' 59" Dec
Recorded on: Sat Jun 10 17:39:40 2000 GMT
Source: Arecibo Radio Observatory
Base Frequency: 1.418750000 GHz

**User Info**
Name: Jane Q Public
Data units completed: 423.4
Total computer time: 48093 hr 14 min 56.4 sec

107.4

Time (sec)

0

Power

0.0   Frequency (Hz)   9765.620

12

## SETI : Some numbers

computations per work unit                $3.1 \times 10^{12}$ FP-operations
work unit throughput                    700 000/day

➡ $22 \times 10^{17}$ FLOP/day

➡ >25 TFLOPS

|  | SETI@home | ASCI White@DoE |
|---|---|---|
| Processing | 25 TFLOPS | 12.3 TFLOPS |
| Cost | 1 M USD | 110 M USD |

13

---

## Generation III : P2P middleware

**Requirements**
• add/find/remove distributed resources transparently
• globally scalable
• load balancing
• exploitation of locality
• dynamic adaptations (dynamic hosts/resources)
• security of data
• anonymity

**Important platforms**
• Pastry, Tapestry, CAN, Chord, Kademlia

14

**Chapter 7**

# Peer-to-Peer systems

1. Introduction
2. **Overlays**
3. Distributed Hash Tables

---

# Overlay routing

Full mesh overlay topology
scales badly !



**A**

| dest. | next hop |
|-------|----------|
| H | B |
| E | E |

**B**

| dest. | next hop |
|-------|----------|
| H | E |

**E**

| dest. | next hop |
|-------|----------|
| H | H |

**a**

| dest. | next |
|-------|------|
| h | c |
| e | c |

**c**

| dest. | next |
|-------|------|
| h | f |
| e | d |

**f**

| dest. | next |
|-------|------|
| h | g |

**g**

| dest. | next |
|-------|------|
| h | h |

**d**

| dest. | next |
|-------|------|
| e | e |

**b**

| dest. | next |
|-------|------|
| e | c |

## Overlay routing

IP packet

| IP header | IP payload |
|-----------|------------|

overlay packet

IP source address
IP destination address

| overlay header | overlay payload |
|----------------|-----------------|

overlay source address (edge node ID)
overlay destination address (edge node ID)

IP layer and overlay layer have own
• addressing scheme
• routing protocol

17

---

## Content based routing

**Use ID of object to manipulate instead of node ID**
-> overlay can redirect to closest replica
-> overlay can optimize placement and #replica's

**Globally Unique ID (GUID)**
- constructed through hash of object state
    - content itself
    - object description
- overlay = distributed hash table (DHT)

**Responsibilities of DHT**
- find object given GUID (route requests)
    -> map GUID to node ID
- announce GUIDs of new objects
- remove GUIDs of old objects
- manage with nodes becoming unavailable

ff478f2d-e398-449a-93d4-62b0727
adffed61-73f9-42a1-8d54-90f3537
ec359ec3-99be-4c44-8449-be7a1ab
9349e44d-5367-4b20-b6b1-d421c5a
54fbb946-f9e3-498f-b04c-e91584d
c8d1308c-76d2-452e-8d80-6d08483
f1a841a3-b456-4e92-99b5-04f3d86
ba1a4110-01cb-4a4a-967e-b8adbf5
63d06b7a-8cf5-46d5-bcdd-ffd3be7
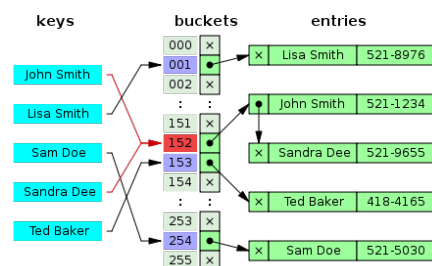4ffa8171-78e0-455a-9cbe-e9fee93

18

# DHT routing vs. IP routing

| | IP routing | DHT routing |
|---|---|---|
| **Size** | IPv4 : $2^{32}$ nodes<br>IPv6 : $2^{128}$ nodes<br>but: hierarchical structured | $> 2^{128}$ GUIDs<br>flat |
| **Load balancing** | routes determined by topology<br>(e.g. OSPF routes) | any algorithm can be used<br>objects can be relocated<br>to optimize routing |
| **Dynamics** | static (time constant 1h)<br>asynchronous w.r.t. application | can be dynamic<br>synchronous or asynchronous |
| **Resilience** | built in | ensured through replication of objects |
| **Target** | IP destination maps to 1 node | GUID can map to different replicas |

19

---

## Chapter 7

# Peer-to-Peer systems

1. Introduction
2. Overlays
3. **Distributed Hash Tables**
   1. Introduction
   2. Circular routing (1D)
   3. Pastry (1D)
   4. Chord (1D)
   5. CAN (nD)

keys    buckets    entries

| John Smith | | 000 × |
| | | 001 ● |
| | | 002 × |
| Lisa Smith | | : : |
| | | 151 × |
| Sam Doe | | 152 |
| | | 153 |
| | | 154 × |
| Sandra Dee | | : : |
| Ted Baker | | 253 × |
| | | 254 ● |
| | | 255 × |

| × | Lisa Smith | 521-8976 |
| ● | John Smith | 521-1234 |
| × | Sandra Dee | 521-9655 |
| × | Ted Baker | 418-4165 |
| × | Sam Doe | 521-5030 |

20

# Data management APIs

**Distributed Hash Table [DHT] API**
>    **put(GUID,data)**
>    **remove(GUID)**
>    **value = get(GUID)**

>    DHT takes care of:
>    - finding good location
>    - number of replicas needed


**Distributed Object Location and Routing [DOLR] API**

>    **publish(GUID)**
>    **unpublish(GUID)**
>    **sendToObj(message,GUID,[#])**

>    node made explicitly responsible for GUID

21

---

# Bootstrapping

The bootstrapping problem
>    how to find a network node ?

Approaches
>    - pre-configured static addresses of stable nodes
>        - should have fixed IP address
>        - should always be on
>    - DNS-service
>        domain name resolves to nodes

Configuration info made available
>    - routing table bootstrap info
>    - GUID space the node is responsible for
>    - protocol information

## DHT : GUID routing

**Concept**

Content Item → hash function → GUID

e.g. Node Public key

Node is responsible for range of GUIDs it belongs to

**Nodes are responsible for part of GUID space**

Node1

Node2

Node3

{GUID}

GUIDs in intersection are replicated

23

## Circular routing

Suppose : GUID consists of 4 hex symbols

0    FFFF

Active node 2

Active node 1

l=4

Active node has Leaf set (size 2l)
= Table (GUID,IP) of neighbours

24

# Circular routing

0    FFFF

k=Destination

l=4

Source

l

j

i

h

a

b

c

d

e

f

g

| GUID | IP |
|------|-----|
| a | IPa |
| b | IPb |
| c | IPc |
| **d** | **IPd** |

| GUID | IP |
|------|-----|
| e | IPe |
| f | IPf |
| g | IPg |
| **h** | **IPh** |

| GUID | IP |
|------|-----|
| i | IPi |
| j | IPj |
| **k** | **IPk** |
| l | IPl |

N active nodes
leaf set size = 2l    =>    worst case
N/2l hops

25

# Prefix routing

Optimize routing efficiency
• allow for longer jumps
• without increasing leaf set

=> need for additional "jump" mechanism

0    FFFF

leaf set

26

# Prefix routing

**Routing table structure @active node**

n-bit GUID         -> n/4 routing table rows
(16 bit            -> 4 rows)

p = longest prefix match
Routing table of node with GUID 53AF

[GUID,IP] of next hop for message destinated for 5Cxx

| p | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | A | B | C | D | E | F | | |
| 0 | x | x | x | x | x | ■ | x | x | x | x | x | x | x | x | x | x | | |
| 1 | x | x | x | ■ | x | x | x | x | x | x | x | x | x | x | x | x | | |
| 2 | x | x | x | x | x | x | x | x | x | x | x | ■ | x | x | x | x | | |
| 3 | x | x | x | x | x | x | x | x | x | x | x | x | x | x | x | ■ | | |

[GUID,IP] of next hop for message
destinated for 53A3 (=closest active node to 53A3)

27

---

# Prefix Routing

Message m from S(ource) to D(estination) arrives in C(urrent) node
Leaf set : size 2l ($L_{-l}$ -> $L_l$)
R : routing matrix

```
if (L₋ₗ ≤ D ≤ Lₗ) {
        forward m
                - to GUIDᵢ found in leaf set
                - to current node C
} else {
        find longest prefix match p of D and C
        c = symbol (p+1) of D
        if(R_pc ≠ null) forward m to R_pc
        else {
                find GUIDᵢ in L or R with |GUIDᵢ – D|<|GUIDᵢ-C|
                forward m to GUIDᵢ
        }

}
```

**Observation:**
        **correct routing possible with empty routing table**
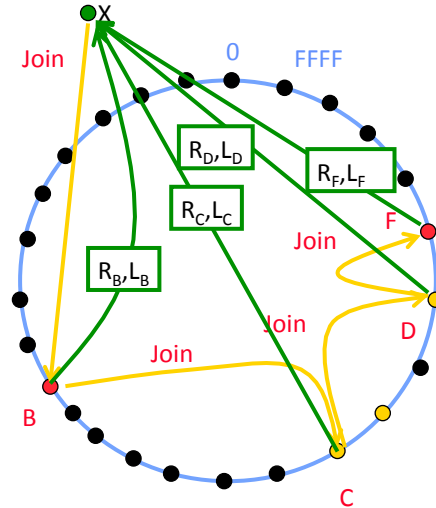
28

*14*

# Node dynamics : Joining

X wants to join
- how to construct $R_X$, $L_X$ ?
- how to adapt R of all other nodes ?

Special "nearest neighbour" discovery
algorithm to find a nearby active node B

1. X sends join(X,$GUID_X$) to B
2. DHT routes join the usual way
   to node F
       $\min |GUID_F - GUID_X|$
3. All nodes on routing path
   send info to X to construct $R_X$,$L_X$



29

---

# Node dynamics : Joining

**Constructing $R_X$**

{$B_0,B_1,B_2, ... ,B_{N-1}$} = set of physical nodes visited by Join message
      $B_0 = B$
      $B_{N-1} = F$

Row 0 :
        - used to route GUID with no common prefix
        - bootstrap node $B_0$ (very) close to X
        => $R_X[0] = R_{B0}[0]$
Row 1 :
        - used to route GUID with common prefix 1
        - prefix($GUID_{B1}$,$GUID_X$)≥1
        => $R_X[1] = R_{B1}[1]$
...
Row i : => $R_X[i] = R_{Bi}[i]$

**Constructing $L_X$**
        X should be neighbour of F
        initial choice : $L_X = L_F$

30

# Node dynamics : Leaving

Only leaf sets are repaired

l=4

$L_D$

| | |
|---|---|
| GUID($D_{-4}$) | IP($D_{-4}$) |
| GUID($D_{-3}$) | IP($D_{-3}$) |
| GUID($D_{-2}$) | IP($D_{-2}$) |
| GUID($D_{-1}$) | IP($D_{-1}$) |
| GUID($D_1$) | IP($D_1$) |
| GUID($D_2$) | IP($D_2$) |
| GUID($D_3$) | IP($D_3$) |
| GUID($D_4$) | IP($D_4$) |

1. D finds node close to $GUID_X$
   -> $D_2$
2. Get $L_{D2}$
3. Adapt $L_D$ based on
   -> Remove $D_3$
   -> Add $D_5$
4. Propagate to neighbours

0    FFFF

$L_{D2}$

Failing node X

unreachable

Detecting node D

$D_4$
$D_3$
$D_2$
$D_1$

$D_{-1}$ $D_{-2}$ $D_{-3}$ $D_{-4}$

$L_D$

31

---

# Chord

**Basic Chord API**

      lookup(key) : maps key -> IP address of node responsible for the key

**Application built on top on Chord**

      - uses lookup(key)

      - gets informed by Chord when key-set of current node changes

      - is responsible for (if desired)

            - authentication

            - caching and replication

**Sample applications**

      - cooperative mirroring (multiple servers cooperate to store content)

      - time shared storage (ensuring data availability, even if server off-line)

      - distributed indexes (content sharing)

      - large-scale combinatorial search (e.g. code breaking)
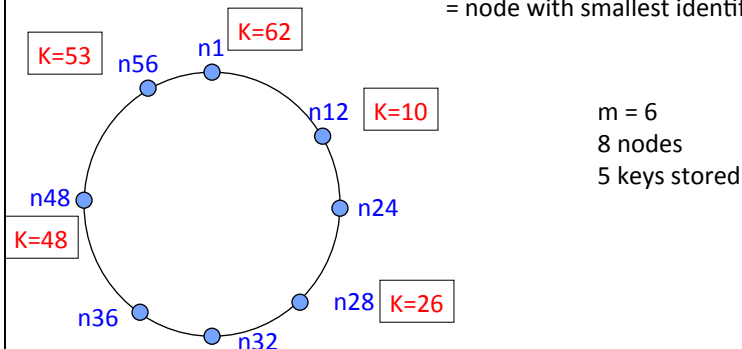
32

# Consistent hashing

Hash function (SHA-1) produces m-bit identifiers

    hash(nodeIP)
    hash(key)

Identifiers mapped to identifier circle modulo $2^m$

Key k -> hash(k) -> node(hash(k)) = successor(k)
                        = node with smallest identifier ≥ hash(k)

K=53    K=62    K=10

m = 6
8 nodes
5 keys stored

n56 n1 n12 n48 n24 n36 n32 n28
K=48 K=26

33

---

# Joining and Leaving

**X Leaving**
        keys associated with X -> reassigned to successor(X)

**X Joining**
        X gets some of the keys assigned to successor(X)

K=53    K=62    K=10

n56 n1 n54 n12 n48 n24 n36 n32 n28
K=48 K=26

34

---

## Simple Routing

**Request lookup(id) arrives at node n**
Node n knows id of its successor

n.findSuccessor(id) {
        if id in (n,successor]  return successor
        else return successor.findSuccessor(id)
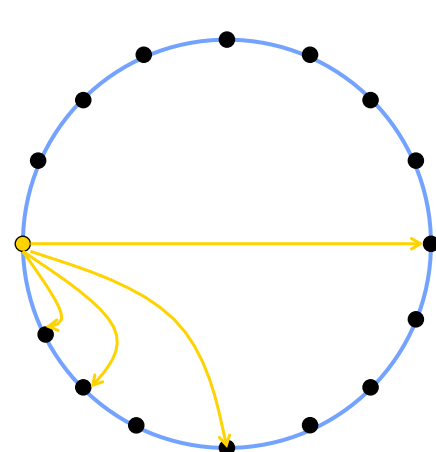}

K=53    K=62    lookup(53)

n56  n1  n12

n48    n24

K=48

n36  n32  n28  K=26

(a,b]
Clockwise interval
Excluding a, including b

## Skiplist routing

Each node has **finger table**
≈ leaf set with non-equidistant node intervals

**Finger table (m entries)**

| Distance | ID | IP |
|---|---|---|
| 1 | | |
| 2 | | |
| 4 | | |
| 8 | | |
| ... | | |

**Routing to GUID**
        send message to successor GUID
        = highest ID in finger table
            with ID≤ GUID
O(log N) hops needed
**Robustness**
        node keeps track of n successors
        = circular routing in case finger table
            not valid

## Skiplist routing

**Finger table**

**DistanceID          IP**

1
2
4

**More formally**

row i in finger table is node at distance larger than $2^{i-1}$
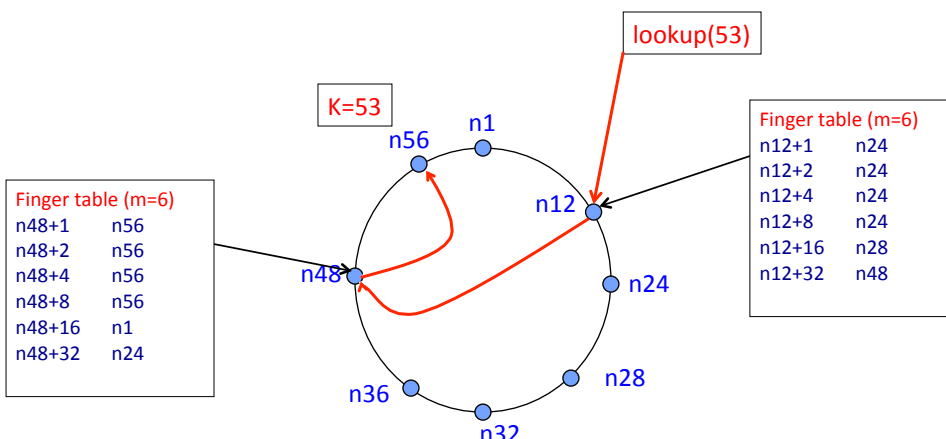
finger[i] =successor($n+2^{i-1}$ )   (i≥1)

Finger table (m=6)

| n48+1 | n56 |
|-------|-----|
| n48+2 | n56 |
| n48+4 | n56 |
| n48+8 | n56 |
| n48+16 | n1 |
| n48+32 | n24 |

Finger table (m=6)

| n12+1 | n24 |
|-------|-----|
| n12+2 | n24 |
| n12+4 | n24 |
| n12+8 | n24 |
| n12+16 | n28 |
| n12+32 | n48 |

n1
n56
n12
n48
n24
n36
n32
n28

37

---

## Skiplist routing

**Routing** : take largest jump possible !

lookup(53)

K=53

Finger table (m=6)

| n48+1 | n56 |
|-------|-----|
| n48+2 | n56 |
| n48+4 | n56 |
| n48+8 | n56 |
| n48+16 | n1 |
| n48+32 | n24 |

Finger table (m=6)

| n12+1 | n24 |
|-------|-----|
| n12+2 | n24 |
| n12+4 | n24 |
| n12+8 | n24 |
| n12+16 | n28 |
| n12+32 | n48 |

n1
n56
n12
n48
n24
n36
n32
n28

38

## Skiplist routing

**Routing** : take largest jump possible !

```
n.findSuccessor(id) {
        if id in (n,successor]  return successor
        else {
                n'=closestPreceedingNode(id)
                return n'.findSuccessor(id)
        }
}

n.closestPreceedingNode(id) {
        for i=m downto 1
                if finger[i] in (n,id] return finger[i]
        return n
}
```

39

## Node dynamics : Joining

**Bootstrapping (create a new ring)**

```
n.create() {
        predecessor=null
        successor=n
}
```

**Joining a ring (containing n')**

```
n.join(n') {
        predecessor=null
        successor=n'.findSuccessor(n)
}
```

40

*20*

**Stabilizing the ring**
-Runs periodically
-Informs nodes about newly joined nodes
-Fix finger tables
-Fix predecessors

n.stabilize() {
        x=successor.predecessor
        if x in (n,successor]
                        successor = x     // better successor found
        successor.notify(n)          // adapt predecessor of successor
}

n.notify(n') {                //n' thinks it is the predecessor of n
        if (predecessor==null) OR (n' in (predecessor,n[)
                      predecessor = n'  // better predecessor found
}

41

---

Steady-State : no changes !

$P(s_P, p_P)$

P.stabilize()
        - adapts successor @P ($s_P$)
        - adapts predecessor @Q ($p_Q$)

$Q(s_Q, p_Q)$

If successor and predecessor @P and Q correct,
        i.e.        $s_P = Q$, $p_Q = P$
        x=successor.predecessor=$(s_P)$.predecessor=$p_{(sP)}$=$p_Q$=P
        -> x in (P,Q] ? -> FALSE
        -> Q.notify(P)

        P in ($p_Q$,Q[ ? -> P in (P,Q[ ? -> FALSE

NO UPDATES

42

21

# Node dynamics : Join

$P(s_P=Q, p_P=?)$

$X(s_X=null, p_X=null)$

$Q(s_Q=?, p_Q=P)$

X joins, P<X<Q
Ring stable

X.join(Q)
-> $p_X=null$
-> $s_X=Q.findSuccessor(X)=Q$

$P(s_P=Q, p_P=?)$

$X(s_X=Q, p_X=null)$

$Q(s_Q=?, p_Q=P)$

43

---

# Node dynamics : Join

$P(s_P=Q, p_P=?)$

$X(s_X=Q, p_X=null)$

$Q(s_Q=?, p_Q=P)$

X.stabilize()
-> $x=successor.predecessor=(s_X).predecessor=p_{(sX)}=p_Q=P$
-> P in (X,Q] -> FALSE
-> Q.notify(X)

Q.notify(X)
-> X in $(p_Q,Q[ = (P,Q[ $ -> TRUE -> $p_Q=X$

$P(s_P=Q, p_P=?)$

$X(s_X=Q, p_X=null)$

$Q(s_Q=?, p_Q=X)$

44

# Node dynamics : Join

$P(s_P=Q, p_P=?)$

$X(s_X=Q, p_X=null)$

$Q(s_Q=?, p_Q=X)$

P.stabilize()

-> x=successor.predecessor=$(s_P)$.predecessor=$p_{(sP)}$=$p_Q$=X

-> X in (P,Q] -> TRUE -> $s_P$=X

-> X.notify(P)

X.notify(P)

-> $p_X$ == null -> $p_X$=P

$P(s_P=X, p_P=?)$

$X(s_X=Q, p_X=P)$

$Q(s_Q=?, p_Q=X)$

45

45

---

# Node dynamics : Join example

N100 {K85,K93,K99}

N200 {K125,K135, K175}

N150 JOINS

N100 {K85,K93,K99}

N150 {}

N200 {K125,K135, K175}

N150 COPIES KEYS FROM N200

46

# Node dynamics : Join example

N100 {K85,K93,K99}

N150 {K125,K135}

N200 {K125,K135, K175}

N100, N150 and N200 stabilize

N100 {K85,K93,K99}

N150 {K125,K135}

N200 {K175}

47

---

# Node dynamics : Stability

```
n.fixFingers() {
        next++                              // f[next] will be fixed
        if (next>m) next = 1
        finger[next]=findSuccessor(n+2next-1)
}


n.checkPredecessors() {          // allows active predecessor to enter
        if (predecessor failed) predecessor = null
}
```

48

# Node dynamics : Failures

**As long as each node knows successor -> correct operation of lookup**
**Consider finger table @P**

| | |
|---|---|
| **P+1** | **Na** |
| **P+2** | **Nb** |
| **P+4** | **Nc** |
| **P+8** | **Nd** |
| **P+16** | **Ne** |
| **P+32** | **Nf** |

**Suppose**

- lookup(id) is launched @P, Nd<id<Ne
- nodes Na, Nb, Nc, Nd and Ne fail simultaneously
-> successor(id)=Nf (instead of Ne)
-> ERROR

**Robustness ?**

- each node has list of r first successors
- to corrupt ring -> r nodes have to fail simultaneously
- changes needed to stabilization and lookup code

49

# Node dynamics : Failures (pseudo code)

```
n.findSuccessor(id) {
        if id in (n,successor]  return successor
        else {
                try {
                        n'=closestPreceedingNode(id)
                        return n'.findSuccessor(id)
                } catch(TimeoutException e) {
                        invalidate n' from finger
                        and successor table
                        return n.findSuccessor(id)
                }
        }
}


n.closestPreceedingNode(id) {
        for i=m downto 1
                if finger[i] in (n,id] return finger[i]
        for i=r downto 1
                if successor[i] in (n,id] return successor[i]
        return n
}
```
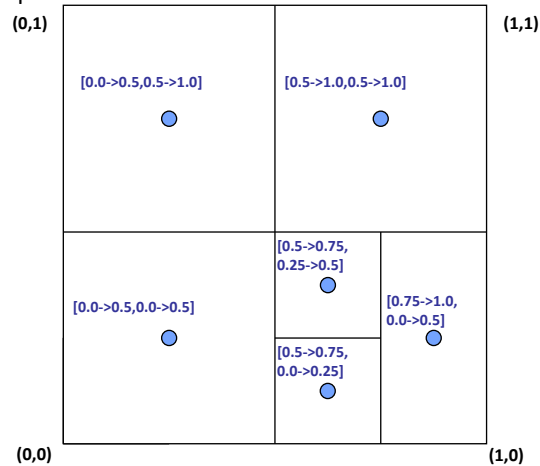
50

## Multi-dimensional routing

CAN : Content addressable network

**Basic idea**

use d-dimensional space to map keys

node is responsible for d-dimensional hypercube
- Cartesian coordinates used
- space is d-torus

**Example 2D**

(0,1)                          (1,1)

[0.0->0.5,0.5->1.0]      [0.5->1.0,0.5->1.0]

[0.5->0.75,
0.25->0.5]

[0.0->0.5,0.0->0.5]      [0.75->1.0,
0.0->0.5]

[0.5->0.75,
0.0->0.25]

(0,0)                          (1,0)

51

---

## Multi-dimensional routing

**Mapping content to node**
- Key-Value pair <K,V> mapped to point P in Cartesian space
    P=hash(K)
- P belongs to region owned by node N
- N stores <K,V>

**Retrieving entry for key K**
- compute P=hash(K)
- if P not owned by requester or neighbours
    -> route request in the CAN network

**Self-organizing routing**

node learns and stores set of neighbours

neighbour : shares hyperplane of dimension (d-1)
-> in (d-1) dimension, intervals overlap
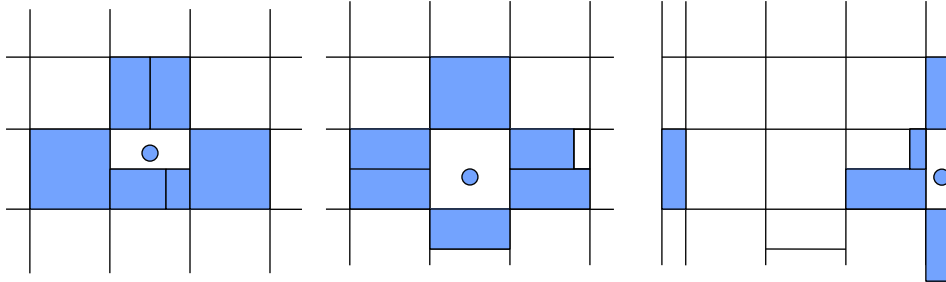info stored per neighbour n
IP(n), zone(n)

A

E

B      IP(A), zone(A)
IP(B), zone(B)
...

D

C         F

## Multi-dimensional routing

**Neighbour nodes : examples**



Suppose equally partitioned zones (regular d-dimensional grid)
-Number of neighbours : 2d
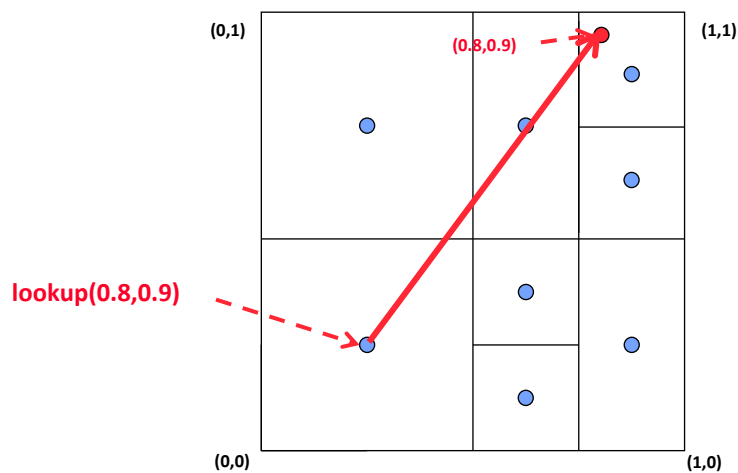-Average routing path length = $(d\, n^{1/d})/4$

---

## Multi-dimensional routing

**Routing in d dimensions**
Follow straight line from source to destination



lookup(0.8,0.9)

(0,1)   (0.8,0.9)   (1,1)

(0,0)   (1,0)

54

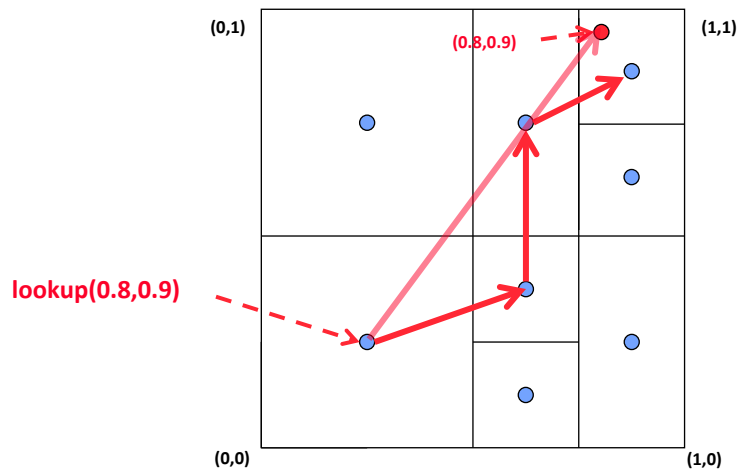## Multi-dimensional routing

**Routing in d dimensions**

Forward to closest neighbour (greedy forwarding)

-> each node keeps track of at least d neighbours

-> space complexity O(d)

-> lookup cost $O(dN^{1/d})$ (N number of nodes)



(0,1)    (0.8,0.9)    (1,1)

lookup(0.8,0.9)

(0,0)    (1,0)    55

---

## Node dynamics : Joining the CAN

**1. Find an active node in the CAN (Bootstrapping)**

- CAN has DNS name
- DNS resolves to IP address of bootstrapping node
- each bootstrap node has list of probably active nodes
  - nodes that joined before
  - nodes can notify before they leave
  - nodes supposed to reply to ping-messages from bootstrap
- bootstrap node replies to "join" by random selection from the active node list
- essential to minimize activity of the bootstrap node

**2. Find a zone to care for**

**3. Publish new node to neighours**
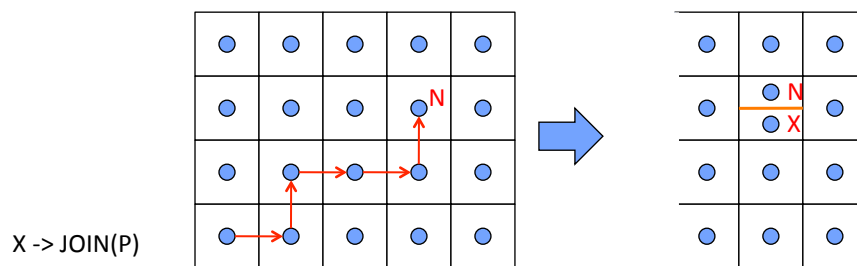
56

## Node dynamics : Joining the CAN

**1. Find an active node in the CAN (Bootstrapping)**
**2. Find a zone to care for**
- X selects random point P
- X sends JOIN(P) -> arrives at node N currently responsible for P
- N splits zone in 2
- N sends to X
- range of keys X will be responsible for
- <Key,Value> pairs X will handle
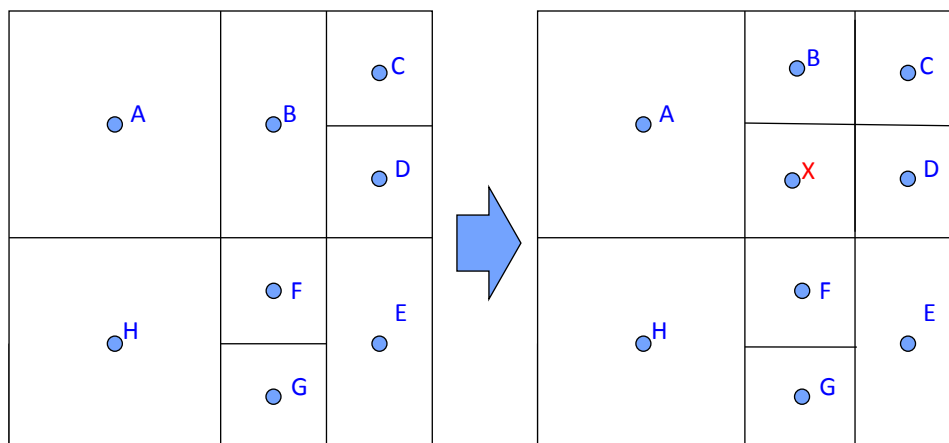3. Publish new node to neighbours



X -> JOIN(P)

57

## Node dynamics : Joining the CAN

**Joining : example**



neighbour(B) = {A,C,D,F,G}

neighbour'(B) = {A,C,G,X}
neighbour'(X) = {A,D,F,B}

58

# Node dynamics : Joining the CAN

**1. Find an active node in the CAN (Bootstrapping)**
**2. Find a zone to care for**
**3. Publish new node to neighbours**

        - X and P update neighbour set

$$neighbour'(X) \subset \big(neighbour(P) \cup \{P\}\big)$$

$$neighbour'(P) \subset \big(neighbour(P) \cup \{X\}\big)$$

        - X and P send immediate update on zone(X) and zone(P) to neighbour(P) and neighbour(P')

        - every node sends periodic refreshes (soft-state updates)
                - n + zone(n)
                - for each neighbour : ID + zone

**Joining is a LOCAL operation**
    -> only neighbours are affected
    -> O(d)

59

# Node dynamics : Leaving the CAN

**Normal exit**
      - X checks if zone(X) can be merged with zone(n)
      - if so, zone(X) is handed to n
      - if not, a neighbour zone will take care of multiple zones
      - which neighbour ?
             minimize maximum zone size of neighbours

**Unexpected departure (failure)**-> Immediate take-over
Failure detected through absence of update messages
When node N detects X has departed (probably multiple N's !)
      - start take over timer
               timeout = C zoneSize(N)
               when time-out {
                    N sends TAKEOVER(zoneSize(N)) to neighbour(X)
                    update zone(N) info
               }
      - on receipt of TAKEOVER(zoneSize(Y)) at node N
               if(zoneSize(Y)<zoneSize(N)) {
                    cancel timer @ N
                    update neighbour info
               }
               else send TAKEOVER(zoneSize(N)) to neighbour(X)

60

## Node dynamics: Leaving the CAN

**Unexpected departure (failure)**

**Immediate take-over algorithm OK to handle single node failures**
If less than 50% of neighbour(X) nodes reachable
   -> locate active neighbours through expanding ring search
    prior to take-over

**Expanded ring search @ node N**
    TTL = 1
    do {
      broadcast request for info to neighbours(N)
      TTL++
    } while neighbourInfo incomplete

**Avoid too much fragmentation**
  background zone-reassignment process

61

## Advanced CAN

1. Increase d
   - reduces hop count and latency
   - increases node state

2. Use multiple "realities"
   - as if separate instances of CAN are running on same node infrastructure
   - in each reality, node is responsible for different zone
       -> improved data availability, robustness
       -> reduced latency

3. RTT based routing metrics
   - use RTT weighted distances when forwarding
       -> favours low latency paths
     -> lower latency

4. Overloading coordinate zones
   - multiple nodes are responsible for same zone ("peers")
   - reduces number of zones
       -> reduced path length
       -> reduced latency
       -> improved fault tolerance

62

# Advanced CAN

5.  Multiple hash functions
    - k different hash functions map same value to k nodes
    - queries are sent to k nodes
        -> improved latency
        -> improved data availability
    - cost : larger node state, increased query traffic
    - instead of launching parallel queries: start with query to closest node

6.  Organize coordinate space based on physical network layout

7.  Uniform coordinate partitioning
    prior to splitting zone, check whether neighbour could be split
    (i.e. has larger zone)
        -> better load balancing

8.  Caching and replication to manage hot spots
    - cache recently accessed data in node, and check cache before forwarding
    - replicate frequently accessed data in neighbouring nodes

63

# Research challenges

1.  Appropriate distance function
2.  Keep system structure simple under frequent joins/leaves
3.  Fault tolerance measures
4.  Concurrent changes
5.  Proximity routing (adapt logical routing to physical topology)
6.  Cope with malicious nodes
7.  Indexing and keyword search (instead of ID's)

64

*32*