# Structured Peer-to-Peer Systems for Distributed Applications

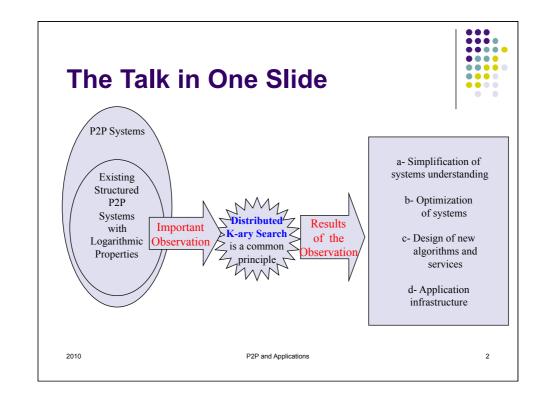
THE ADVENTURES OF



www.ist-selfman.org

Seif Haridi (SICS/KTH)
Sameh El-Ansary (SICS)
Ali Ghodsi (KTH)
Luc Onana Alima (SICS/KTH)
Per Brand (SICS)
Thorsten Schütt (ZIB)
Alexander Reinefeld (ZIB)
Peter Van Roy (UCL)
Boris Mejias (UCL)





#### **Outline**



3

- Overview of P2P systems
  - Three generations of P2P
- Distributed Hash Tables (DHTs)
- DKS: A Realistic DHT
- Broadcast service in DKS
- DHT as application infrastructure
  - Scalaris and Beernet libraries
  - DeTransDraw application on gPhone
  - Distributed Wikipedia application

P2P and Applications

Overview of P2P systems



2010 P2P and Applications

# What is Peer-To-Peer Computing? (1/3)



 A. Oram ("Peer-to-Peer: Harnessing the Power of Disruptive Technologies")

P2P is a class of applications that:

- Takes advantage of resources (storage, CPU, etc,..) available at the edges of the Internet.
- Because accessing these decentralized resources means operating in an environment of unstable connectivity and unpredictable IP addresses, P2P nodes must operate outside the DNS system and have significant or total autonomy from central servers.

2010 P2P and Applications

# What is Peer-To-Peer Computing? (2/3)



- P2P Working Group (A Standardization Effort)
   P2P computing is:
  - The sharing of computer resources and services by direct exchange between systems.
  - Peer-to-peer computing takes advantage of existing computing power and networking connectivity, allowing economical clients to leverage their collective power to benefit the entire enterprise.

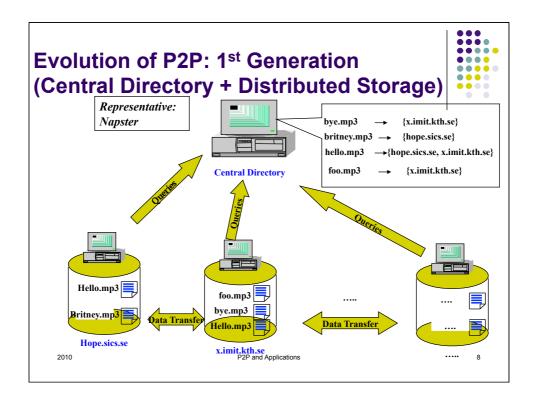
# What is Peer-To-Peer Computing? (3/3)

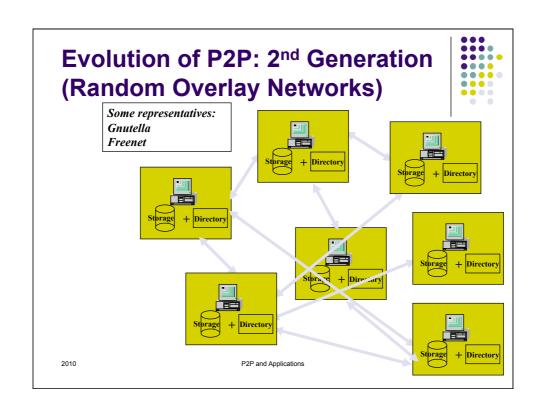


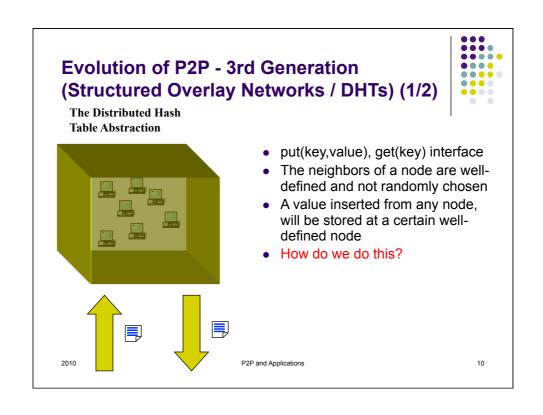
#### Our view

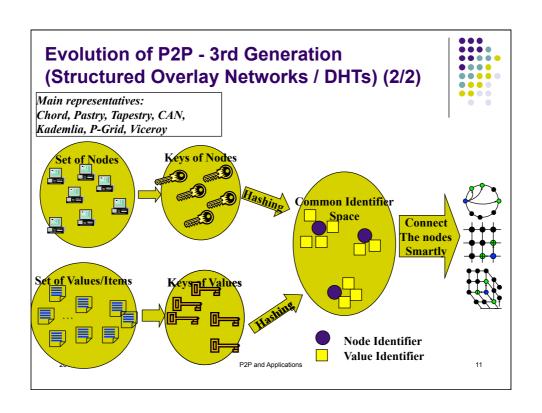
P2P computing is distributed computing with the following desirable properties:

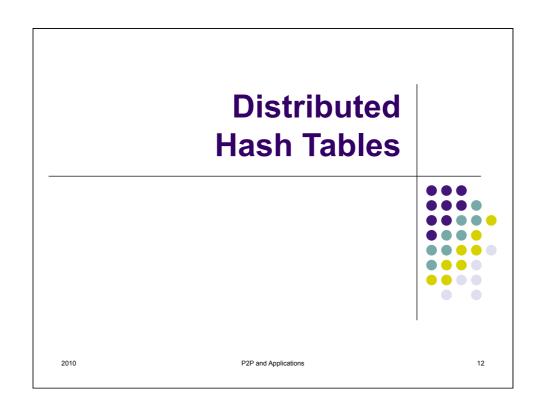
- Resource sharing
- Dual client/server role (same role for all nodes)
- Decentralization/autonomy
- Scalability
- Robustness/self-organization







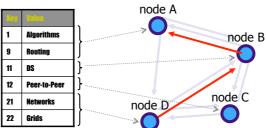




## The Principle of Distributed Hash Tables

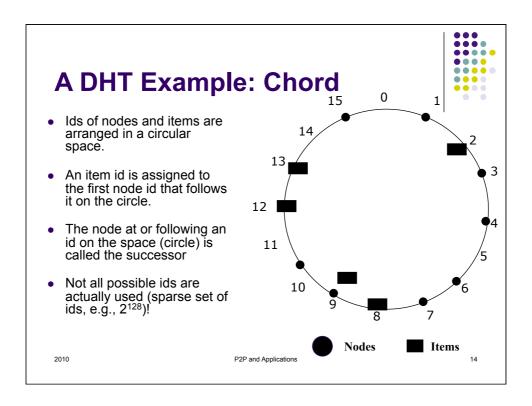


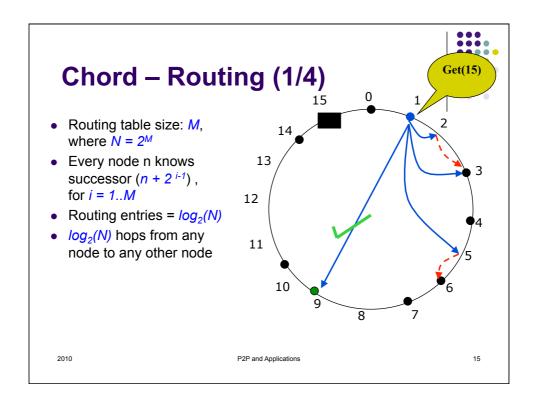
 A dynamic distribution of a hash table onto a set of cooperating nodes

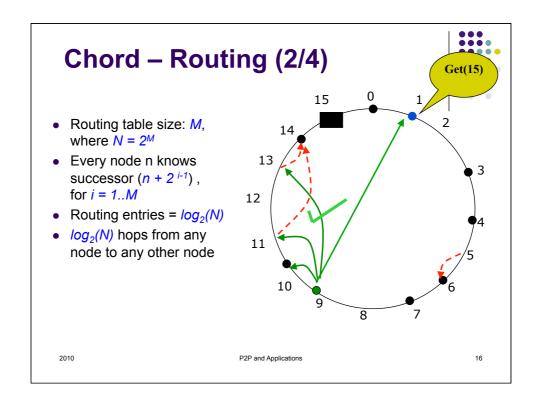


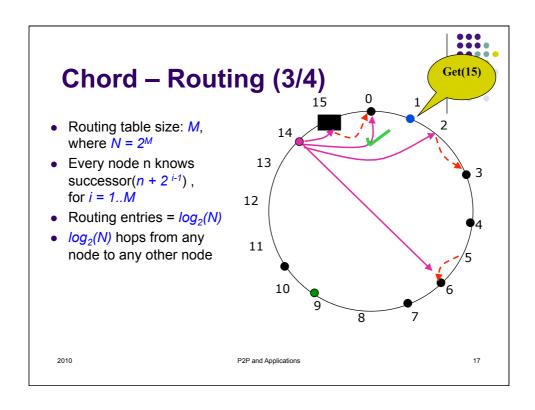
- Basic service: lookup operation
- $\longrightarrow$ Node D: lookup(9)
- Key resolution from any node
- Each node has a routing table
  - Pointers to some other nodes
  - Typically, a constant or a logarithmic number of pointers

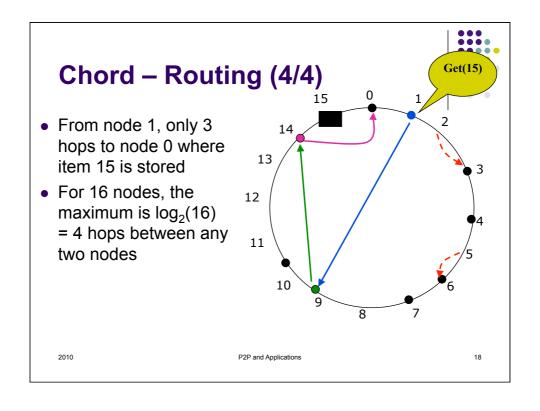
10 P2P and Applications











#### **Taxonomy of P2P Systems**



P2P Systems

Unstructured

Hybrid Decentralized (Napster)

Fully Decentralized (Gnutella)

Partially Decentralized (Kazaa)

Structured (Chord, CAN, Tapestry, Pastry)

2010 P2P and Applications

#### **Comparison of P2P Systems**



	Napster	Gnutella	Pastry	Tapestry	Chord	CAN
Scalable	No	No	Yes	Yes	Yes	Yes
Guarantees	High	Low	High	High	High	High
Locality	Default	Default	Yes	Yes	No	No
Replication	Ad-hoc	Ad-hoc	Well-placed	Well-placed	Well-placed	Well-placed
Routing Table	O(N)	Varies	$O(\log(N))$	$O(\log(N))$	$O(\log(N))$	2d
Item Read	O(1)	O(N)	$O(\log(N))$	$O(\log(N))$	$O(\log(N))$	$O(N^{\frac{1}{d}})$
Item Insert	O(1)	O(1)	$O(\log(N))$	$O(\log(N))$	$O(\log(N))$	$O(N^{\frac{1}{d}})$
Item Delete	O(1)	O(1)	$O(\log(N))$	$O(\log^2(N))$	$O(\log(N))$	$O(N^{\frac{1}{d}})$
Node Insert	O(1)	O(N)	$O(\log(N))$	$O(\log(N))$	$O(\log(N))$	$O(N^{\frac{1}{d}})$
Node Delete	O(1)	Q(1)	$O(\log(N))$	$O(\log(N))$	$O(\log(N))$	O(1)

#### **Research Issues in DHTs**



- Lack of a Common Framework
- Absence of Locality
- Cost of Maintaining the Structure
- Complex Queries
- Heterogeneity
- Group Communication/Higher Level Services
- Cloud/Grid Integration

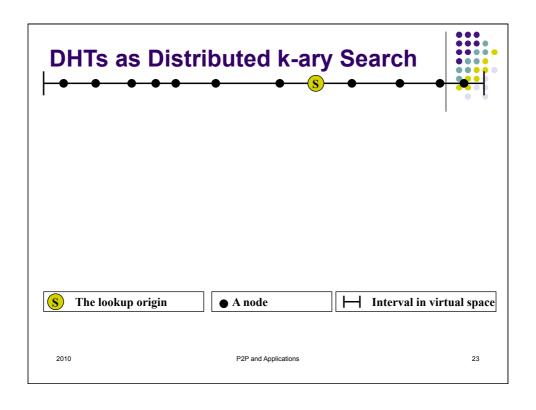


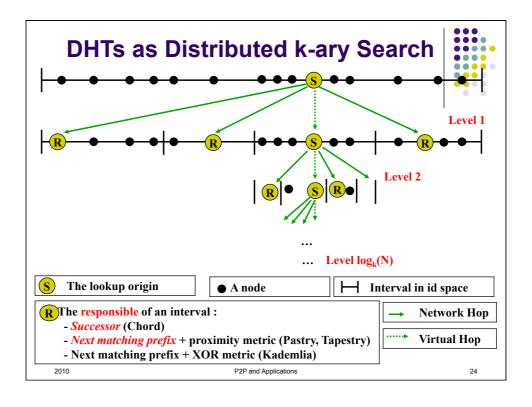
P2P and Applications 21

#### **Framework**



- A Framework for Peer-To-Peer Lookup Services Based On k-ary Search
- Aspects: Understanding, Optimization





#### **The Space-Performance Trade-off**



- We have \( \mathscr{N} \) nodes.
- A node keeps info about a subset of peers
- Lookup length vs. routing table size trade-off
- Extremes:
  - Keep info about all peers
  - Keep info about one peer

2010 P2P and Applications 25

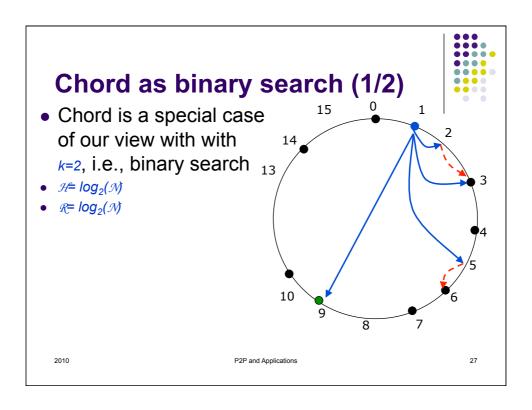
#### Relating $\mathcal{N}$ , $\mathcal{H}$ and $\mathcal{R}$

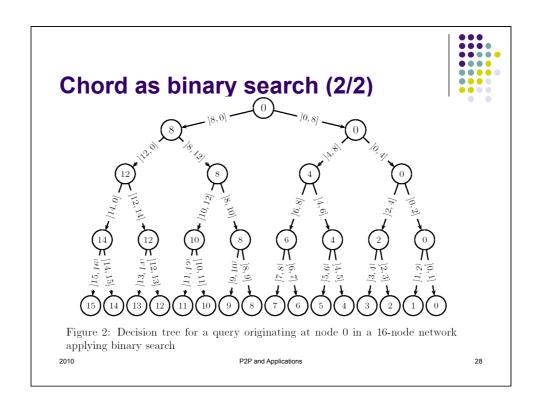


In general, for  $\mathcal{N}$  nodes, the maximum lookup path length  $\mathcal{H}$  and the number of routing entries  $\mathcal{R}$  are as follows:

$$\mathcal{H} = log_k(\mathcal{N})$$
 (Number of levels in the tree)  
 $\mathcal{R} = (k-1)^* log_k(\mathcal{N})$  (k-1 neighbors per level)

$$\mathcal{N} = (\mathcal{R}/\mathcal{H} + 1)^{\mathcal{H}}$$









Suggestion: Increase the search arity k by following the guidelines of our view and put enough info for k-ary search

$$\mathcal{H} = \log_{k}(\mathcal{N})$$

$$\mathcal{R} = (k-1) \log_{k}(\mathcal{N})$$

2010

P2P and Applications

#### Why does routing table size matter?



- Not because of storage capacity
- First reason: because of the number of hops
- Second reason: because of the effort needed to correct an inconsistent routing table after the network changes

2010

P2P and Applications

#### **DKS: A Realistic DHT**



2010 P2P and Applications 31

### **DKS**



- A P2P system that:
  - Realizes the DKS principle
  - Offers strong guarantees because of the local atomic actions
  - Introduces novel technique that avoids unnecessary bandwidth consumption
- Relevance to research issues in state-ofthe-art P2P systems:
  - Common framework
  - Cost of maintaining the structure

#### DKS(N,k,f)



- Title: DKS(N,k,f): Family of Low Communication, Scalable and Fault-Tolerant Infrastructures for P2P Applications
- Authors: Luc Onana Alima, Sameh El-Ansary, Per Brand, and Seif Haridi.
- Place: In The 3rd International Workshop on Global and Peer-To-Peer Computing on Large-scale Distributed Systems - CCGRID2003, Tokyo, Japan, May 2003.

P2P and Applications

Aspects: Understanding, Design

#### Design principles in DKS(N,k,f)



33

- Distributed K-ary Search (DKS) principle
- Local atomic action for joins and leaves
- Correction-on-use technique
- (Replication for fault tolerance) not discussed

#### **Design Principles in DKS**



- Tunability
  - Routing table size vs. lookup length
  - Fault-tolerance degree
- Local atomic join and leave
  - Strong guarantees
- Correction-on-use
  - No unnecessary bandwidth consumption

2010

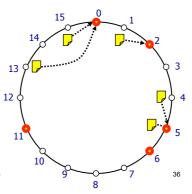
P2P and Applications

35

#### **DKS overlay illustrated (1/3)**



- An *identifier space* of size  $N=k^L$  is used
  - A logical ring of N positions
  - A hash function, H, maps:
    - Nodes onto the logical ring
    - Items onto the logical ring
  - An item (key,value) is stored at successor of H(key), first node that follows H(key) on the ring in the clockwise direction

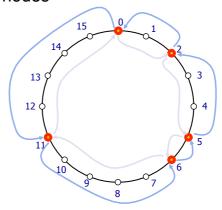


2010

#### DKS overlay illustrated (2/3)



- Basic interconnection
  - Bidirectional linked list of nodes
  - Each node points to its
    - Predecessor
    - Successor
  - Resolving key
    - O(N) hops in an N-node system



2010

P2P and Applications

37

#### Design principle 1: Distributed K-ary Search (DKS) principle



- The DKS is designed from the beginning based on the distributed k-ary search principle
- The system uses the successor of an identifier in a circular space for assigning responsibilities

2010

P2P and Application

#### **DKS** overlay illustrated (3/3)



- Enhanced Interconnection
  - Speeding up key resolution: log<sub>k</sub>(N) hops
  - At each node, a RT of log<sub>k</sub>(N) levels
  - Each level of RT has k intervals
  - For level I and interval i
    - (RT(I))(i) = address of the first node that follows the start of the interval i

(responsible node)

2010

P2P and Applications

30

#### **Notation**



Identifier space I

$$I = \{0,1,...,N-1\}$$

 $N = k^L$  (size of the space)

 $a \oplus b \equiv (a+b) \mod N$ 

**Levels and Views** 

number of levels  $\log_k N = L$ 

at level 1,  $V^1 = I_0^1 \cup I_1^1 \cup ... \cup I_{k-1}^1$ 

$$I_i^1 = [x_i^1, x_{i+1}^1[, x_i^1 = n \oplus i \frac{N}{k}]$$

2010

P2P and Application

#### Levels and views



#### **Levels and Views**

number of levels 
$$\log_k N = L$$
  
at level  $1, V^1 = I_0^1 \cup I_1^1 \cup ... \cup I_{k-1}^1$   
 $I_i^1 = [x_i^1, x_{i+1}^1[, x_i^1 = n \oplus i \frac{N}{k} \text{ for } 0 \le i \le k-1$   
at level  $2 \le l \le L, V^l = I_0^l \cup I_1^l \cup ... \cup I_{k-1}^l$   
 $I_0^l = [x_0^l, x_1^l[, ..., I_{k-1}^l = [x_{k-1}^l, x_0^{l-1}[$   
 $I_i^l = [x_i^l, x_{i+1}^l[, x_i^l(n) = n \oplus i \frac{N}{k^l} \text{ for } 0 \le i \le k-1$ 

2010

2P and Applications

41

#### Responsibility



42

At each level  $V^l$  there is a node  $R(I_i^l)$  that is responsible for each interval  $I_i^l$ 

let x be a node, S(x) is the first node encountered in [x, x]For an arbitrary node n and an arbitrary level l, the responsible for  $I_i^l$  is  $S(x_i^l)$ 

Each node n is itself responsible for  $I_1^l$  for any l Observe that  $S(x_i^l)$  is a function of n can be computed by any node

#### **DKS Overlay illustrated-4**

- Example, *k*=4, *N*=16 (4<sup>2</sup>)
  - At each node an RT of two levels
  - In each level, 4 intervals
  - Let us focus on node 1
  - Level 1, 4 intervals

```
I_0 = [1,5[ : (RT(1))(0) = 1]

I_1 = [5,9[ : (RT(1))(1) = 6]

I_2 = [9,13[ : (RT(1))(2) = 10]

I_3 = [13,1[ : (RT(1))(3) = 15]
```

• level 2, 4 intervals

```
I_0 = [1,2[ : (RT(2))(0) = 1]

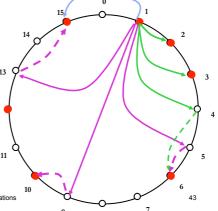
I_1 = [2,3[ : (RT(2))(1) = 2]

I_2 = [3,4[ : (RT(2))(2) = 3]

I_3 = [4,5[ : (RT(2))(3) = 6]
```

 $I_3 = [4,5[ : (RT(2))(3) = 6]$ 





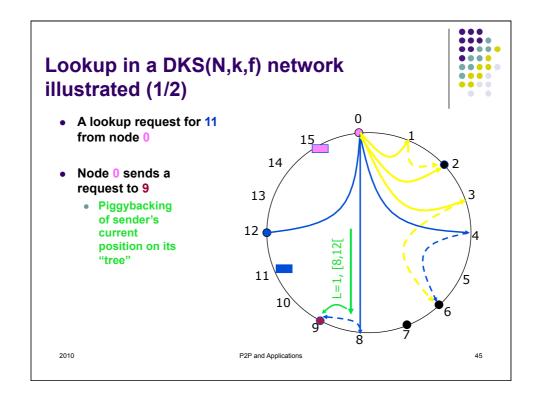
### Lookup in a DKS(N,k,f) network (basic idea)

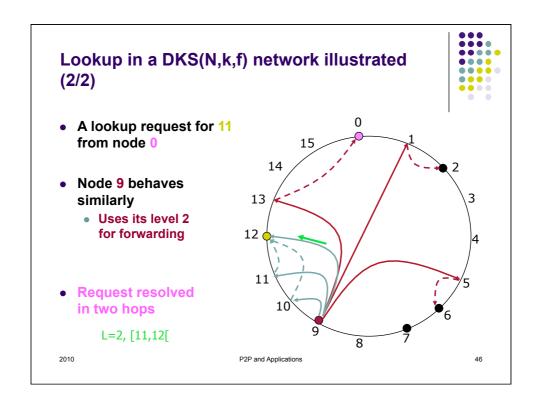


- A predecessor pointer is added at each node
- Interval routing
  - If key between my predecessor and me, done
  - Otherwise, systematic forwarding level by level

2010

P2P and Application





#### Design principle 2: Local atomic action for guarantees



- To ensure that any key-value pair previously inserted is found despite concurrent joins and leaves
  - We use local atomic operation for
    - Node join
    - Node leave
- Stabilization-based systems do not ensure this

2010 P2P and Applications 47

#### **DKS(N,k,f)** network construction

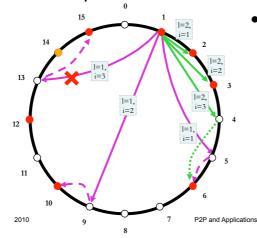


- A joining node is atomically inserted by its current successor on the virtual space
- The atomic insertion involves only three nodes in fault-free scenarios
- The new node receives approximate routing information from its current successor
- Concurrent joins on the same segment are serialized by means of local atomic action

#### **DKS** routing table maintenance



- Node 14 joins the system
  - Example: node 1 in DKS(*N*=16, *k*=4, *f*)



- Node 1's pointer on level=1, interval=3 becomes invalid
  - Will be corrected by Correction-on-use

49

# **Design principle 3: Correction-on-use**



- A node always talks to a responsible node
  - Knowledge of responsible may be erroneous
- "If you tell me from where (in your "tree",) you are contacting me, then I can tell you whether you know the correct responsible"
  - Help others to correct themselves
- "If I heard from you, I learn about your existence"
  - Help to correct myself

2010

P2P and Applications

#### **Correction on use**

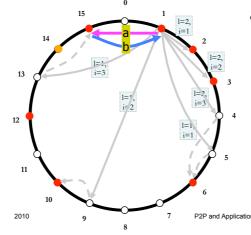
- Look-up or insert messages from node n to node n'
- Add the following to the message
  - *i* (interval) and *l* (level)
- Node n' can compute :  $x_i^l(n)$
- Node *n'* maintains a list of predecessors *BL*

2010 P2P and Applications 51

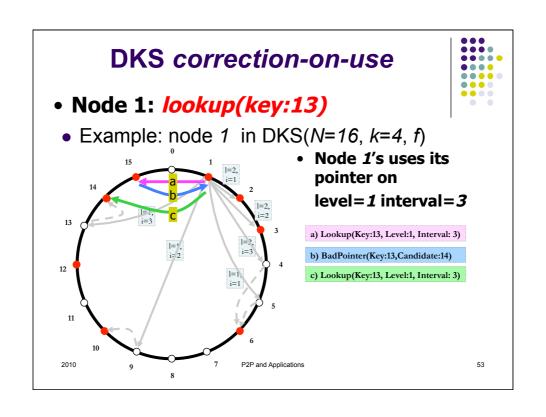
#### **DKS** correction-on-use

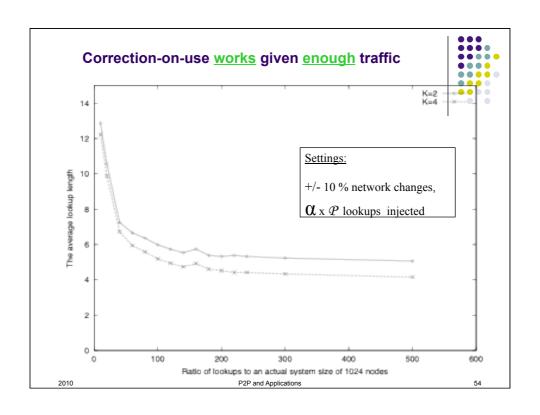


- Node 1: lookup(key:13)
  - Example: node 1 in DKS(N=16, k=4, f)



- Node 1's uses its pointer on level=1 interval=3
  - a) Lookup(Key:13, Level:1, Interval: 3)
  - b) BadPointer(Key:13,Candidate:14)





#### **Broadcast in DKS**



2010 P2P and Applications 55

#### **Efficient Broadcast**



- Title: Efficient Broadcast in Structure P2P Systems
- Authors: Sameh El-Ansary, Luc Onana Alima, Per Brand, and Seif Haridi.
- Place: In The 2nd International Workshop on Peer-to-Peer Systems (IPTPS '03), February 2003.
- Related aspects: Design

#### Motivation: Why broadcast is needed for DHTs?



- In general, support for **global dissemination/ collection** of info in DHTs.
- In particular, the ability to perform arbitrary queries in DHTs.

2010 P2P and Applications 57

#### The Broadcast Problem in DHTs



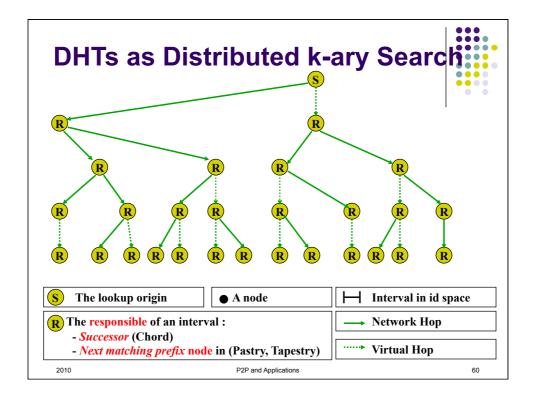
**Problem:** Given an overlay network constructed by a P2P DHT system, find an efficient algorithm for broadcasting messages. The algorithm should not depend on global knowledge of membership and should be of equal cost for any member in the system.

#### **The Efficient Broadcast Solution**



Construct a spanning tree derived from the decision tree of the distributed k-ary search after removal of the virtual hops.

This is equivalent to a Best-Effort Broadcast: it relies on the correctness of the initial node.



#### **Other Solutions for Broadcast**



- Gnutella-like Flooding in DHT
  - (Pro) Known diameter → Correct TTL → High Guarantees
  - (Con) The traffic is high with redundant messages
- Traversing the ring in Chord or Pastry
  - (Pro) No redundant messages
  - (Con) Sequential execution time
  - (Con) Highly sensitive to failure

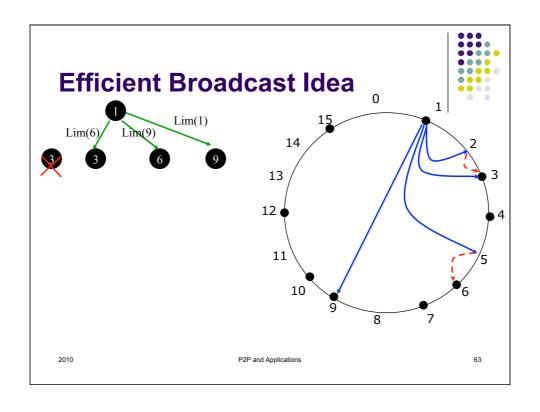
2010 P2P and Applications

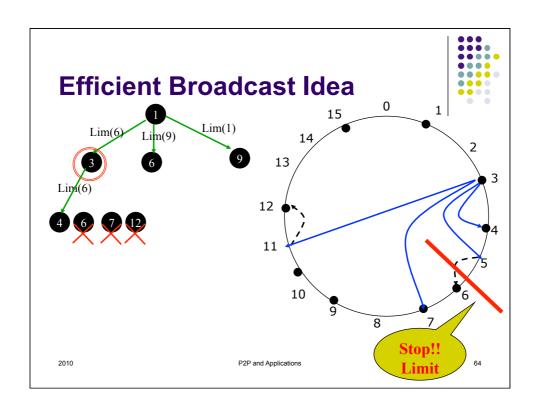
# Efficient Broadcast Algorithm Invariants

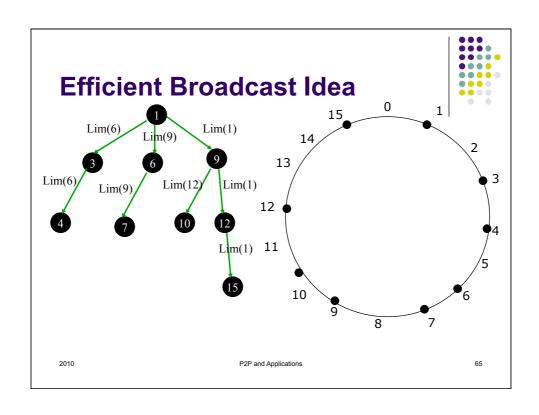


- Any node sends to distinct routing entries.
- Any sender informs a receiver about a forwarding limit, that should not be crossed by the receiver or the neighbors of the receiver.

Forwarding within disjoint intervals where every node receives a message exactly once.







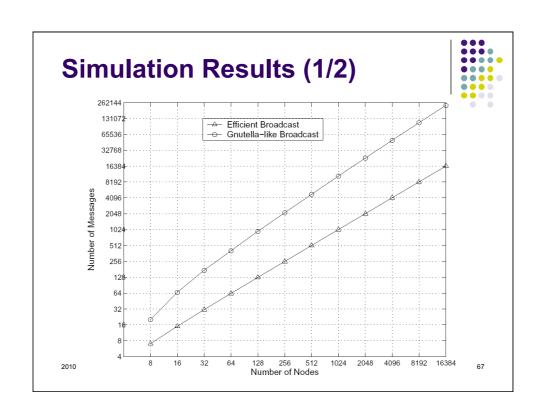
#### **Cost Versus Guarantees**

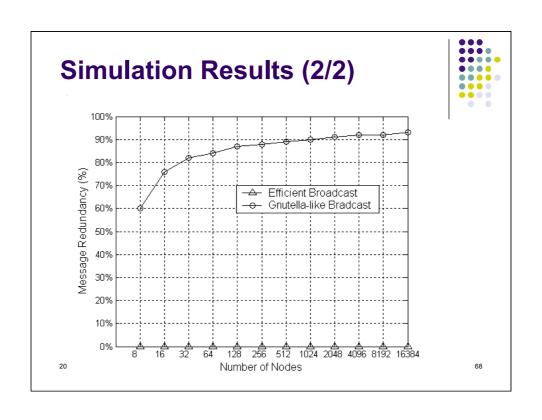


- Q: Is N-1 messages tolerable for any application?
- A1: Broadcast is a costly basic service, if necessary, broadcast wisely.
- A2: If less guarantees are desirable, prune or traverse the spanning tree differently.

2010

P2P and Applications





#### **Broadcast Contributions**



- Presents an optimal algorithm for broadcasting in DHTs
- Relevance to research issues in state-ofthe-art P2P systems:
  - Group communication
  - Complex queries

2010 P2P and Applications 6

#### **Conclusion**



- By using the distributed k-ary search framework for the understanding, optimization and design of existing structured P2P systems with logarithmic performance properties, we were able to provide solutions to current research issues in state-of-the-art systems namely:
  - Lack of a common framework
  - Group communication
  - Complex queries
  - Cost of maintaining the structure

# DHT as Application Infrastructure



2010 P2P and Applications

# DHT as Application Infrastructure



- A DHT can be used as an infrastructure for building large-scale distributed applications
  - Basic services: storage, replication, transactions, broadcast, ...
  - It is easy to build a scalable decentralized application on top
  - In the SELFMAN project we built a Distributed Wikipedia
- We start with a component model
  - Services and applications are concurrent components
  - We create a layered structure (previous project for SINF2345!)
  - Let us show how it is done by building our DHT from reusable components
- We built several applications using this architecture
  - Collaborative drawing (DeTransDraw), Distributed Wikipedia

.

P2P and Applications



## DHT (monolithic)

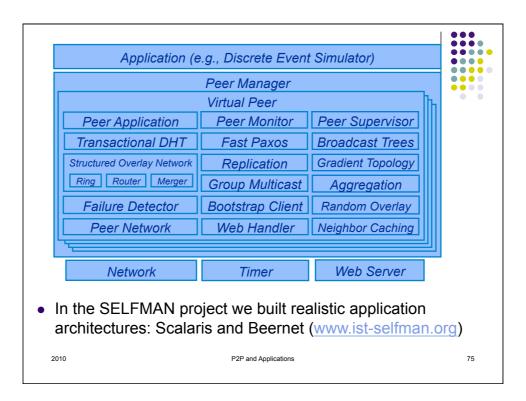
- Now let us look inside our DHT's implementation
- We want to build the DHT from components

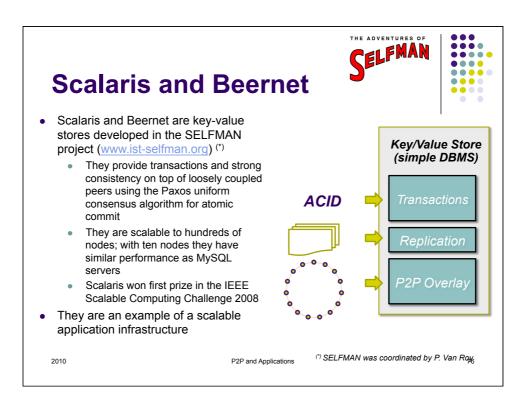
2010 P2P and Applications 73

# Application Structured Overlay Network Ring Periodic Stabilization Router Topology Maintenance Unification Failure Detector Network Timer



- DHT = Structured overlay network
- Reconstruct DHT with building blocks





# DeTransDraw on Beernet for gPhone

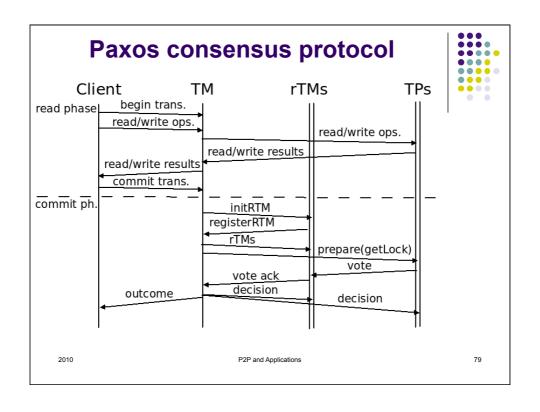


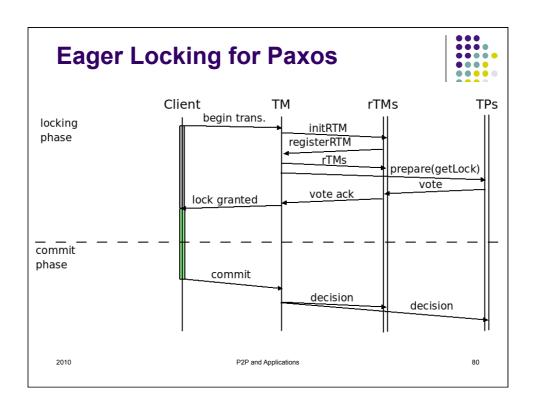
2010 P2P and Applications 77

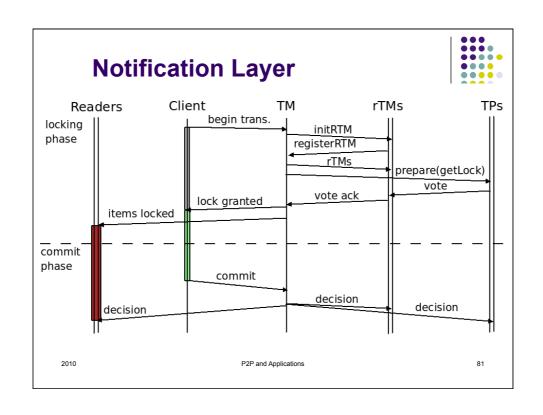
#### **DeTransDraw**

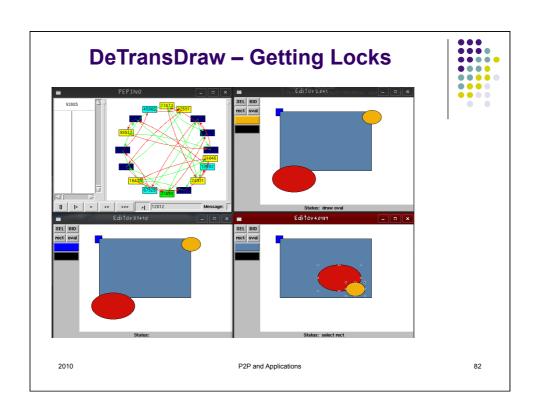


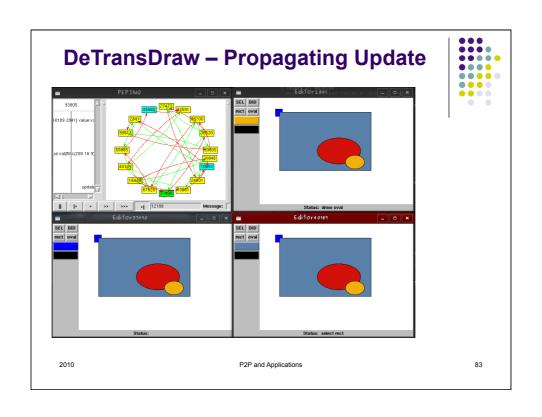
- DeTransDraw is a collaborative drawing application
  - Each user sees exactly the same drawing space
  - Users update the drawing space using transactions
  - For quick response time, the transaction is initiated concurrently with the display update
- Prototype application implemented on top of Beernet
  - Beernet implemented in Mozart, ported to gPhone with Android operating system

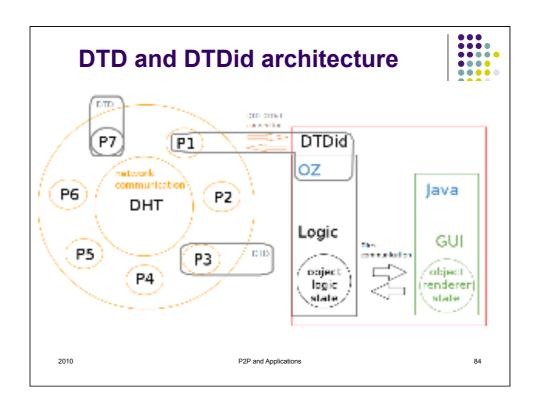






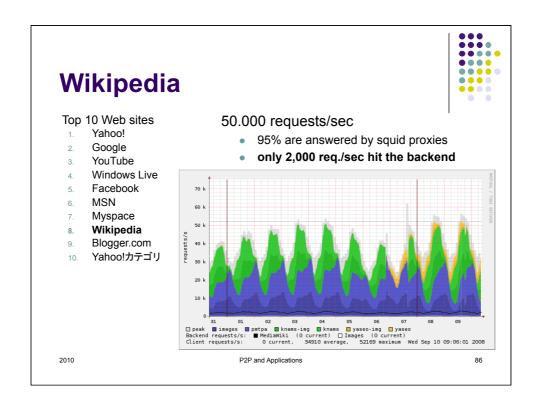


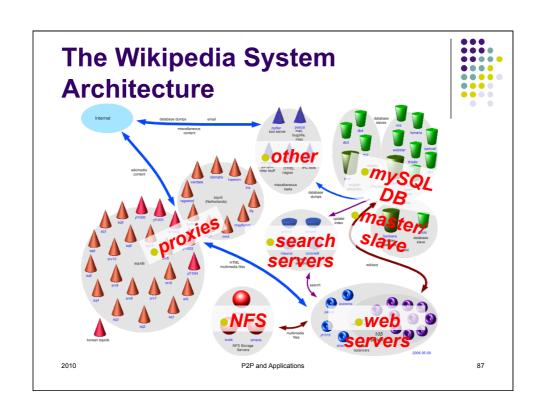


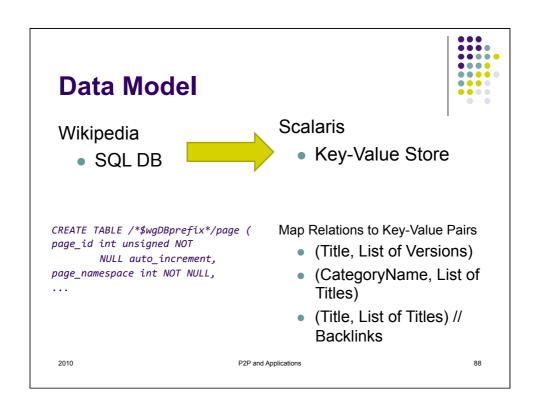


## Distributed Wikipedia on Scalaris









## Data Model (Simple Query Layer)

```
void updatePage(string title, int oldVersion, string newText)
    //new transaction
   Transaction t = new Transaction();
    //read old version
   Page p = t.read(title);
    //check for concurrent update
   if(p.currentVersion != oldVersion)
       t.abort();
   else{
       //write new text
       t.write(p.add(newText));
        //update categories
       foreach(Category c in p)
           t.write(t.read(c.name).add(title));
       t.commit();
   }
}
2010
                                  P2P and Applications
```

89

#### **Self-\* Architecture**

#### Database:

- Chord#
- Mapping
  - Wiki -> Key-Value Store

#### Renderer:

- Java
  - Tomcat
  - Plog4u
- Jinterface
  - Interface to Erlang

Chord#

Load Balancer

Request for page A

Client

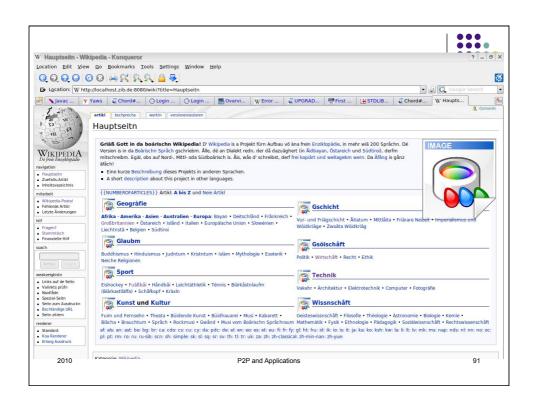
Replica of page A

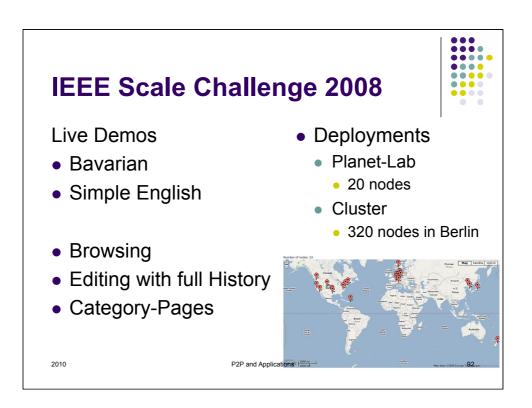
Webserver

2010

P2P and Applications

90









2010



93

- DHT is a practical base to build scalable decentralized applications
  - DHT provides efficient communication and storage
  - Important services, including broadcast and transactions can be built on top of a DHT
  - We built two implementations, Scalaris and Beernet
    - Both are available as open-source systems
- We built several applications in the SELFMAN project
  - DeTransDraw with Beernet on gPhone
  - Distributed Wikipedia with Scalaris (won a prize!)

### **Stub Slides**



2010 P2P and Applications 95

#### **Notation**



Identifier space I

$$I = \{0,1,...,N-1\}$$

 $N = k^L$  (size of the space)

 $a \oplus b \equiv (a+b) \mod N$ 

#### **Levels and Views**

number of levels  $\log_k N = L$ 

at level 1,  $V^1 = I_0^1 \cup I_1^1 \cup ... \cup I_{k-1}^1$ 

$$I_i^1 = [x_i^1, x_{i+1}^1[, x_i^1 = n \oplus i \frac{N}{k}]$$

#### Levels and views



#### **Levels and Views**

number of levels 
$$\log_k N = L$$
  
at level  $1, V^1 = I_0^1 \cup I_1^1 \cup ... \cup I_{k-1}^1$   
 $I_i^1 = [x_i^1, x_{i+1}^1[, x_i^1 = n \oplus i \frac{N}{k} \text{ for } 0 \le i \le k-1$   
at level  $2 \le l \le L, V^l = I_0^l \cup I_1^l \cup ... \cup I_{k-1}^l$   
 $I_0^l = [x_0^l, x_1^l[, ..., I_{k-1}^l = [x_{k-1}^l, x_0^{l-1}[$   
 $I_i^l = [x_i^l, x_{i+1}^l[, x_i^l(n) = n \oplus i \frac{N}{k^l} \text{ for } 0 \le i \le k-1$ 

2010

P2P and Application

9

#### Responsibility



At each level  $V^l$  there is a node  $R(I_i^l)$  that is responsible for each interval  $I_i^l$ 

let x be a node, S(x) is the first node encountered in [x, x]For an arbitrary node n and an arbitrary level l, the responsible for  $I_i^l$  is  $S(x_i^l)$ 

Each node n is itself responsible for  $I_i^l$  for any l Observe that  $S(x_i^l)$  is a function of n can be computed by any node

2010

P2P and Application

## **Routing table**



At each level  $V^l$  there is a node  $R(I_i^l)$  that is responsible for

$$RT_n: \Lambda \to K \to I$$

+ a predecessor pointer P(n)

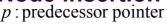
 $\Lambda$ : set of levels

 $K:\{1,..,k\}$ 

I:set of nodes

2010 P2P and Applications

## **Node insertion I** *p*: predecessor pointer



*Insert*: nodes to be inserted

 $n_j$ : node to be inserted

**Empty Network**  $(n_j)$ 

for each  $1 \le l \le L, i \in K$  do

$$RT(l)(i) := n_j$$

$$p := n_j$$

*Insert* := *nil* 

2010 100





#### **Node insertion II**

**NonEmpty Network with a single node**(*n*)

Node *n* computes the RT of  $n_i$ :

for  $l \in L$  do

$$RT_{n_i}(l)(0) := n_j$$

for  $l \in L, i \in K \setminus \{0\}$  do

if 
$$x_i^l(n_i) \in [n_i, n[$$
 then

$$RT_{n_i}(l)(i) := n$$

else 
$$RT_{n_i}(l)(i) := n_j$$

$$p_{n_i} := n$$

n adapts its  $RT_n$  in a similar way

101

**Node insertion III NonEmpty Network with multiple nodes** 

Node *n* computes the RT of  $n_i$ :

for 
$$l \in L$$
 do  $RT_{n_i}(l)(0) := n_j$ 

for  $l \in L, i \in K \setminus \{0\}$  do

if 
$$x_i^l(n_j) \in [n_j, n[$$
 then  $RT_{n_j}(l)(i) := n$ 

elseif 
$$x_i^l(n_j) \in [p, n_j[ \text{ then } RT_{n_i}(l)(i) := n_j]$$

elseif 
$$x_i^l(n_j) \in [n, p[$$
 then  $RT_{n_j}(l)(i) := RT_n(l)(i)$  (approximation)

$$p_{n_j} := p_n$$

n adapts its  $RT_n$  in a similar way

$$p_n := n_j$$

P2P and Applications

#### **Node insertion IV**

- Node insertion is an atomic operation
- Coordinated and serialized by n
- p is informed of  $n_i$
- Other insertion requests to n wait
- *n* is the coordinator of 2PC
- Clients p and n<sub>i</sub>

2010 P2P and Applications 103



#### **Correction on use**

- Look-up or insert messages from node n to node n'
- Add the following to the message
  - *i* (interval) and *l* (level)
- Node n' can compute :  $x_i^l(n)$
- Node n' maintains a list of predecessors BL