

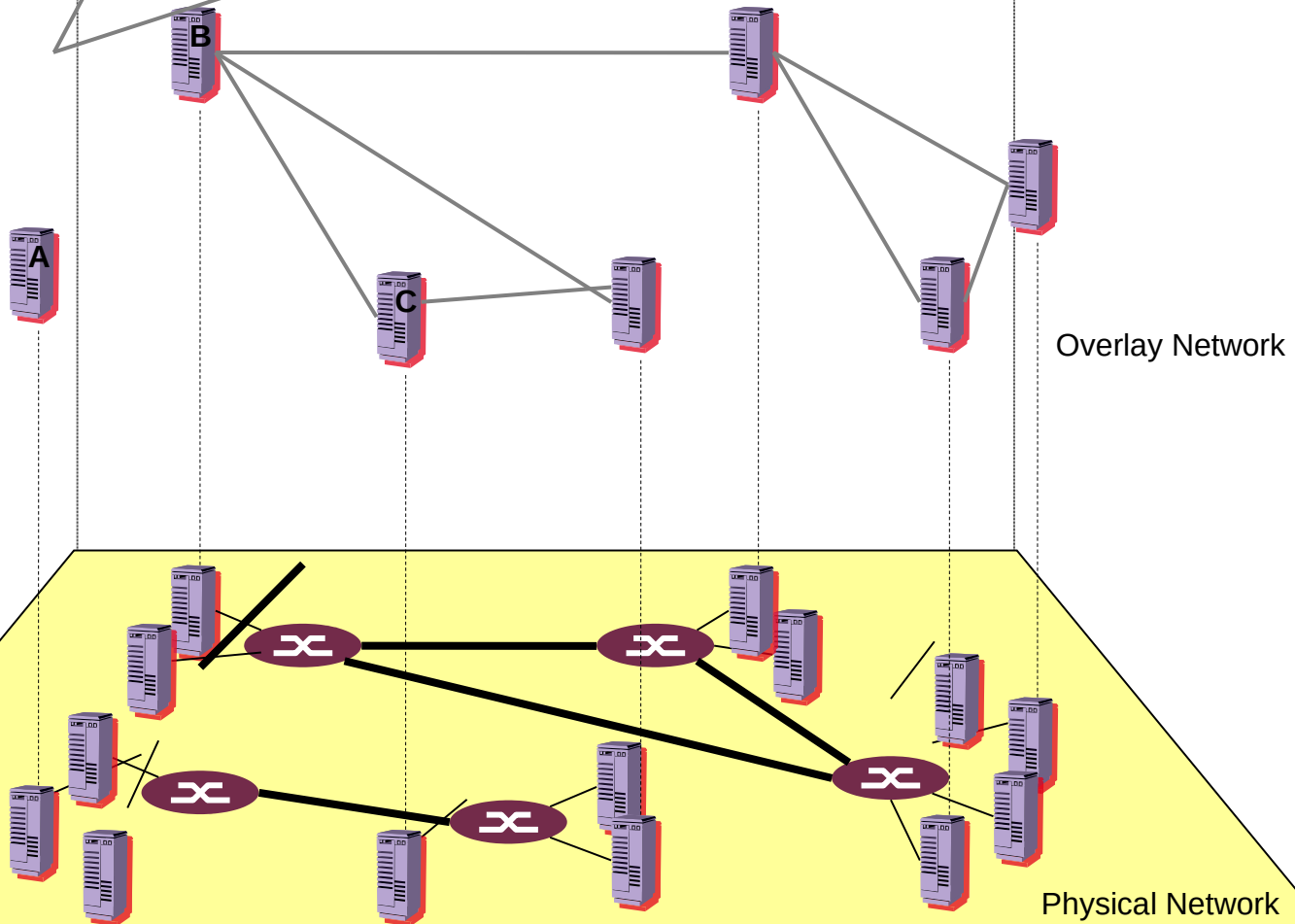


TLC Pastry

Quentin Dufour
Credits : Davide Frey



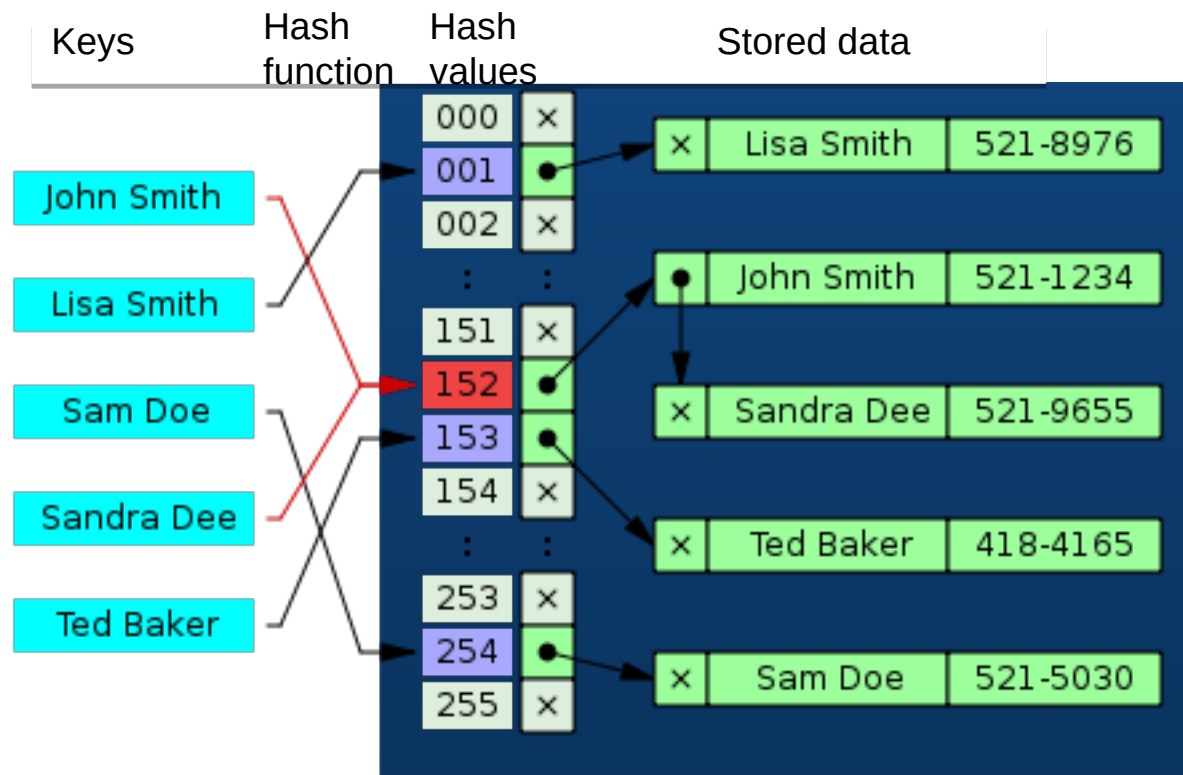
Key Concept: Overlay Network



Overlay types

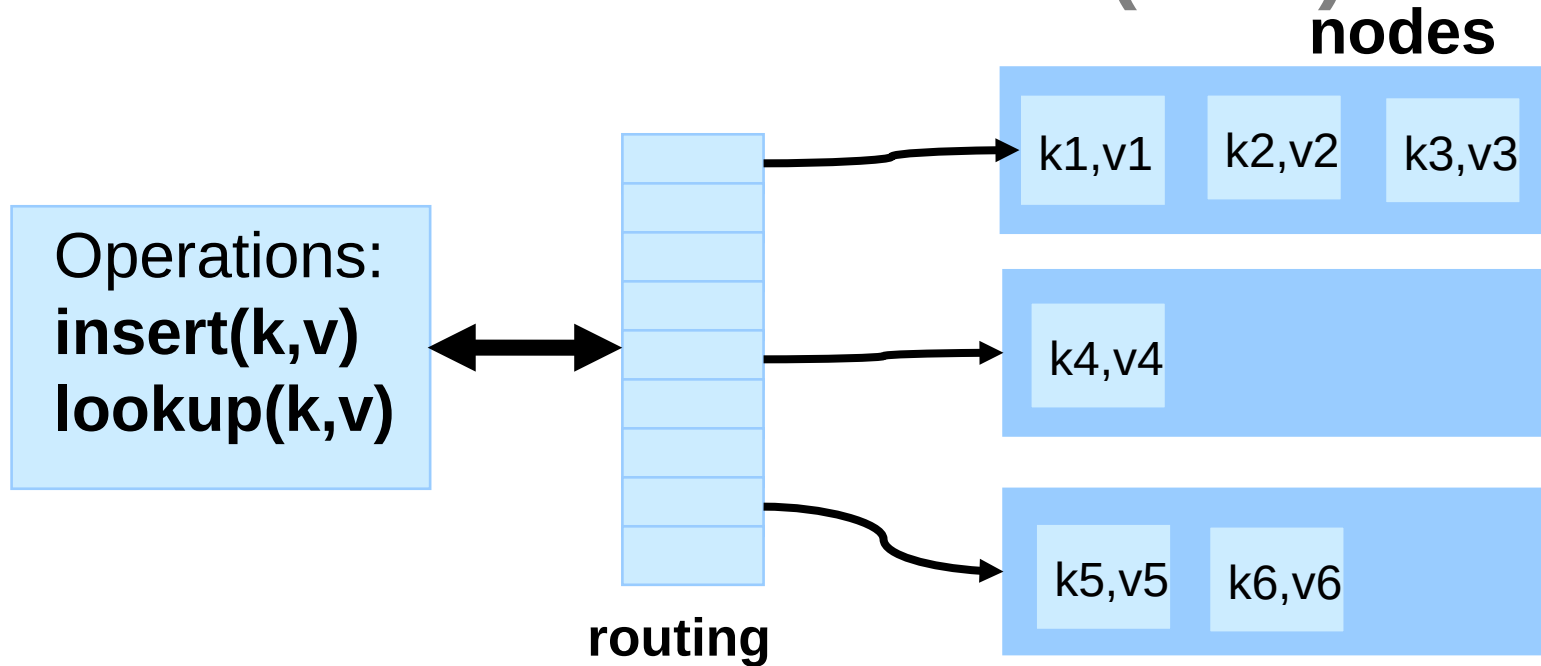
Unstructured P2P	Structured P2P
<ul style="list-style-type: none">○ Any two nodes can establish a link<ul style="list-style-type: none">○ Topology evolves at random○ Topology reflects desired properties of linked nodes	<ul style="list-style-type: none">○ Topology strictly determined by node IDs

Hash Table



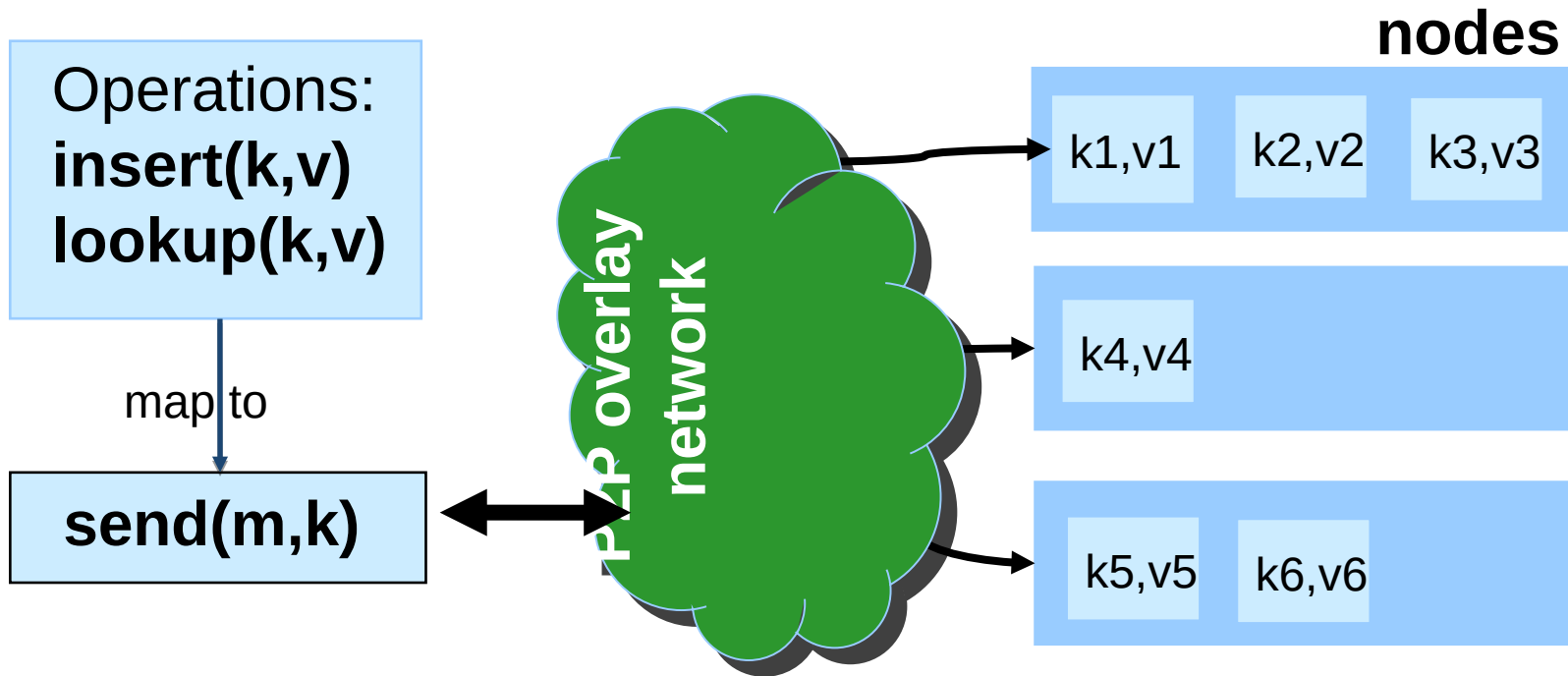
Efficient information lookup

Distributed Hash Table (DHT)



- Store `<key,value>` pairs
- Efficient access to a value given a key
- Must route hash keys to nodes.

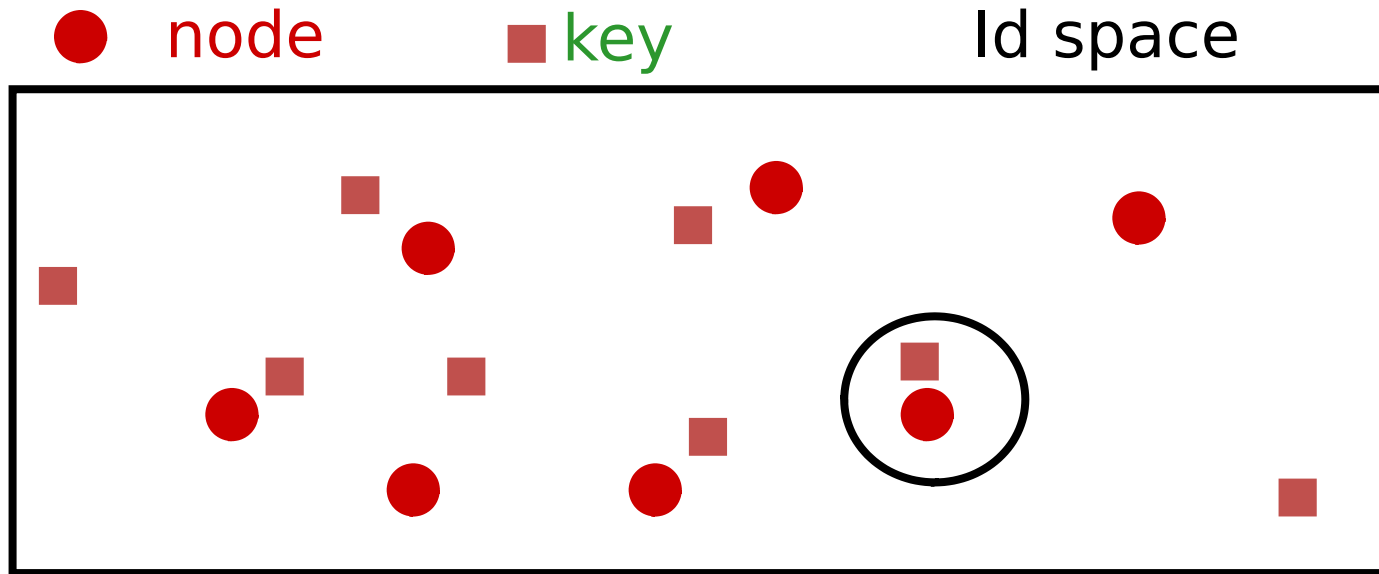
Distributed Hash Table



- Insert and Lookup send messages keys
- P2P Overlay defines mapping between keys and physical nodes

DHT Examples

Pastry (MSR/RICE)



NodeId = 128 bits

Nodes and key place in a linear space (ring)

Mapping : a key is associated to the node with the numerically closest nodeId to the key

Pastry (MSR/Rice)

Naming space :

- Ring of 128 bit integers
- *nodeIds* chosen at random
- Identifiers are a set of digits in base 16

Key/node mapping

- key associated with the node with the numerically closest node id

Routing table:

- Matrix of 128/4 lines et 16 columns
- $routeTable(i,j)$:
 - nodeId* matching the current node identifier up to level i
 - with the next digit is j

Leaf set

- 8 or 16 closest numerical neighbors in the naming space

Proximity Metric

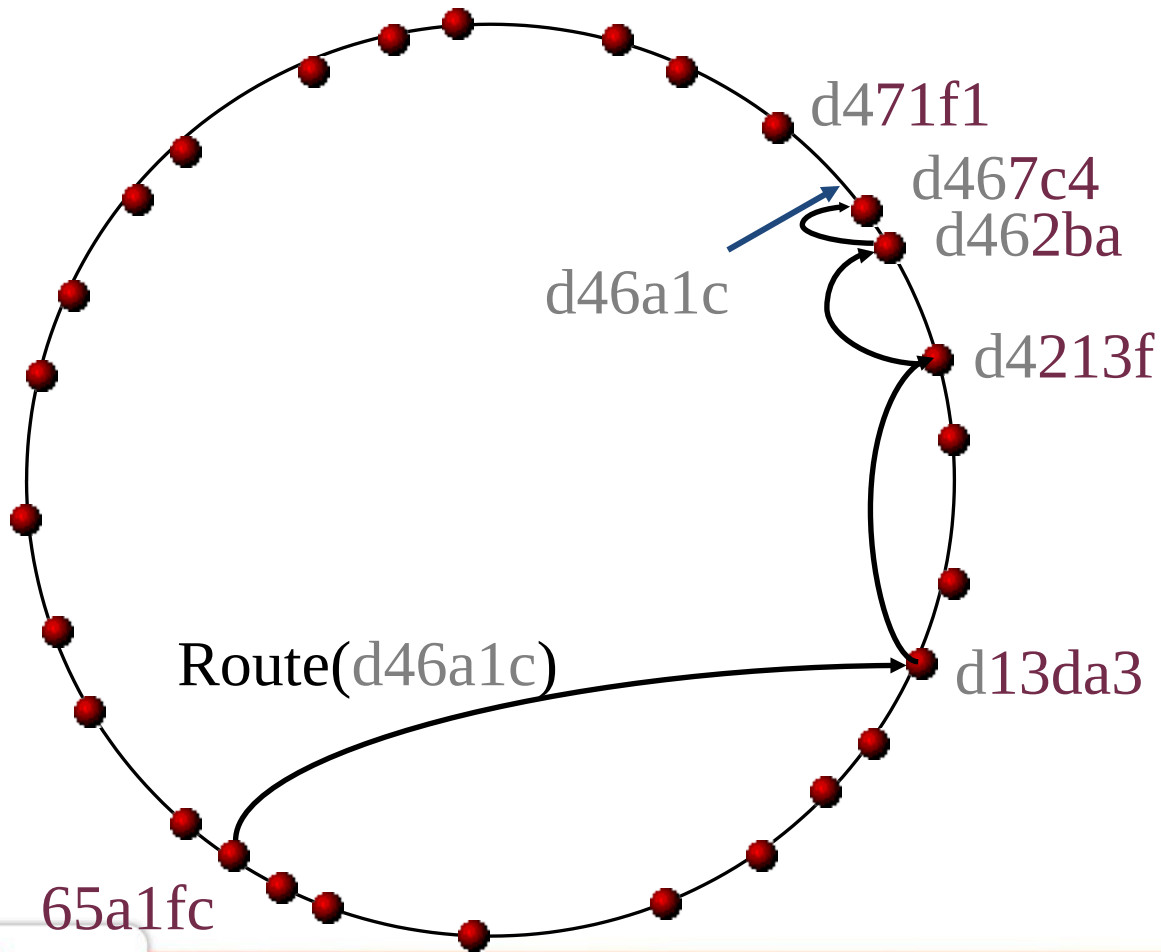
- *Bias selection of nodes*

Pastry: Routing table(#65a1fcx)

Line 0	0 x	1 x	2 x	3 x	4 x	5 x		7 x	8 x	9 x	a x	b x	c x	d x	e x	f x
Line 1	6 0 x	6 1 x	6 2 x	6 3 x	6 4 x		6 6 x	6 7 x	6 8 x	6 9 x	6 a x	6 b x	6 c x	6 d x	6 e x	6 f x
Line 2	6 5 0 x	6 5 1 x	6 5 2 x	6 5 3 x	6 5 4 x	6 5 5 x	6 5 6 x	6 5 7 x	6 5 8 x	6 5 9 x		6 5 b x	6 5 c x	6 5 d x	6 5 e x	6 5 f x
Line 3	6 5 a 0 x		6 5 a 2 x	6 5 a 3 x	6 5 a 4 x	6 5 a 5 x	6 5 a 6 x	6 5 a 7 x	6 5 a 8 x	6 5 a 9 x	6 5 a a x	6 5 a b x	6 5 a c x	6 5 a d x	6 5 a e x	6 5 a f x

$\log_{16} N$
lines

Pastry: Routing



Properties

✂ $\log_{16} N$ hops

✂ Size of the state maintained
(routing table): $O(\log N)$

Routing algorithm (on node A)

- (1) if ($L_{-\lfloor |L|/2 \rfloor} \leq D \leq L_{\lfloor |L|/2 \rfloor}$) {
- (2) // D is within range of our *leaf set*
- (3) forward to L_i , s.th. $|D - L_i|$ is minimal;
- (4) } else {
- (5) // use the routing table
- (6) Let $l = shl(D, A)$;
- (7) if ($R_l^{D_l} \neq null$) {
- (8) forward to $R_l^{D_l}$;
- (9) }
- (10) else {
- (11) // rare case
- (12) forward to $T \in L \cup R \cup M$, s.th.
- (13) $shl(T, D) \geq l$,
- (14) $|T - D| < |A - D|$
- (15) }
- (16) }

R_i^i : entry of the routing table R , $0 \leq i \leq 2^b$,

line l , $0 \leq l \leq \lfloor 128/b \rfloor$

L_i : i th closest nodeId in the leafset

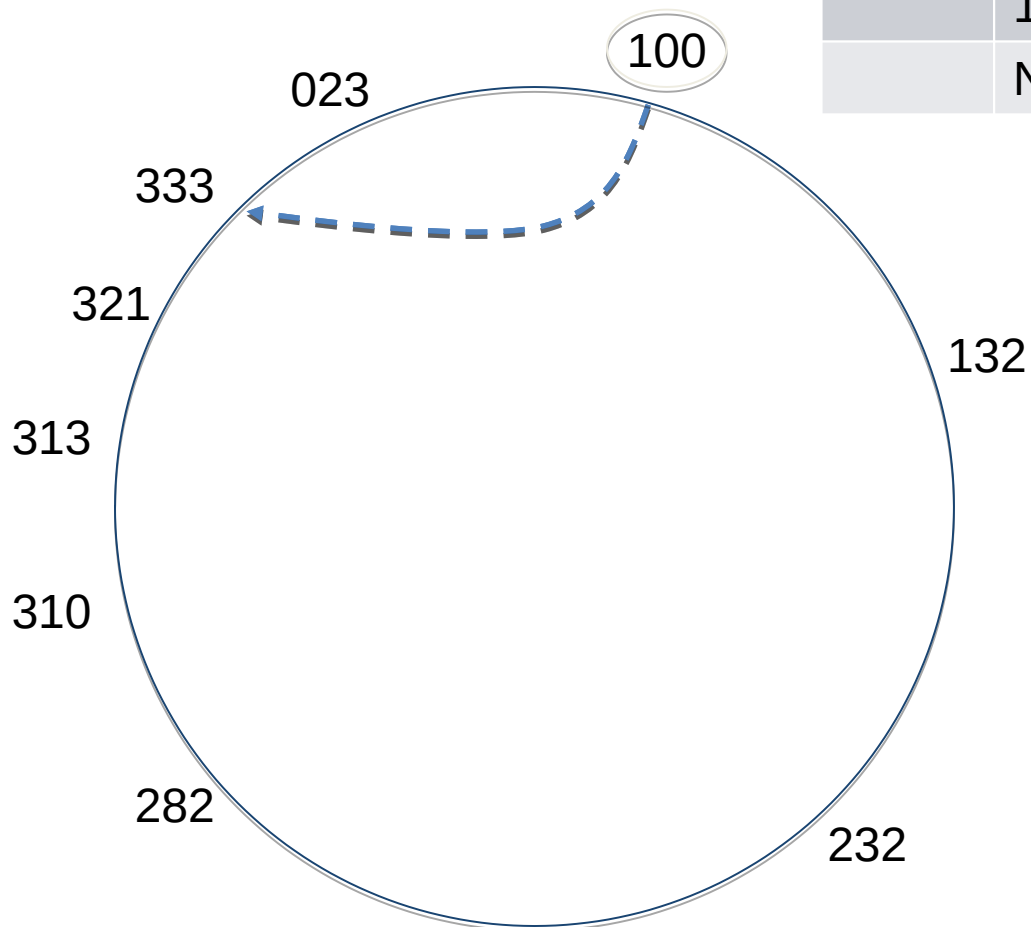
D_l : value of the l digits of key D

$SHL(A, B)$: length of the shared prefix between A and B

Pastry Example

b=4

Route to 311



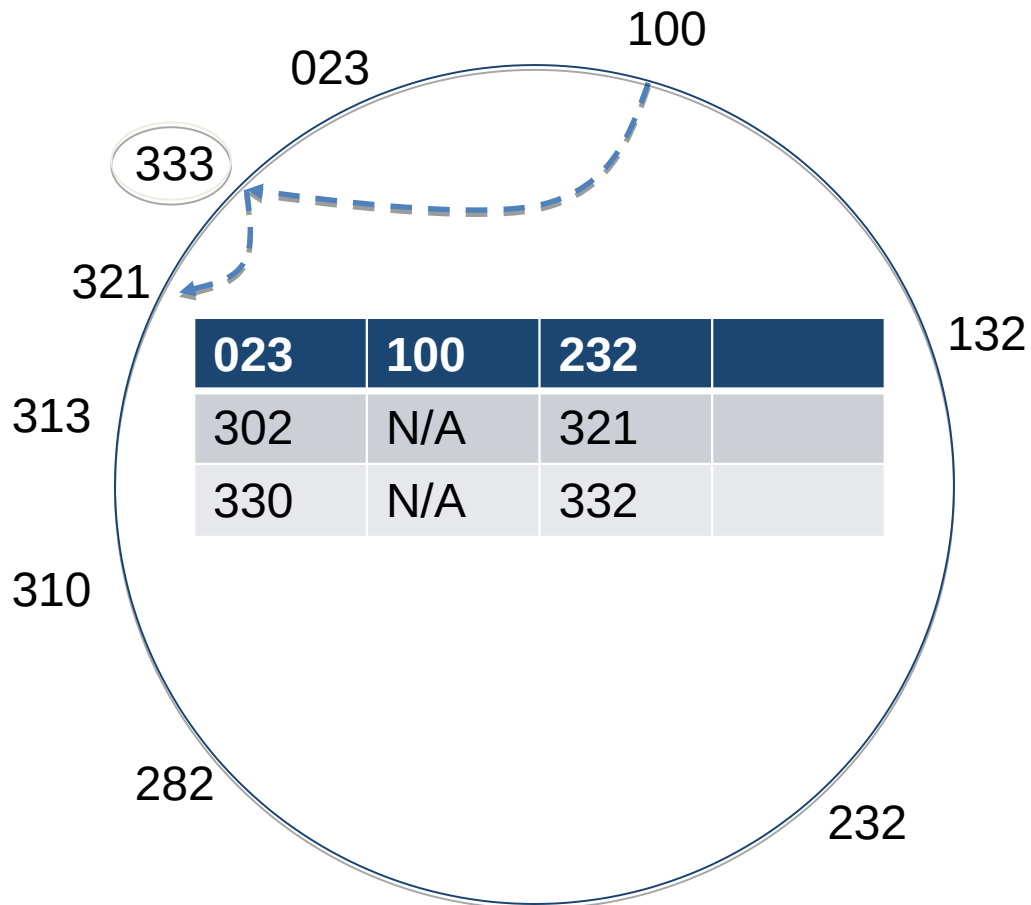
023		232	333
	113	122	132
	N/A	102	103

For space reasons some nodes are not shown on the ring

Pastry Example

b=4

Route to 311

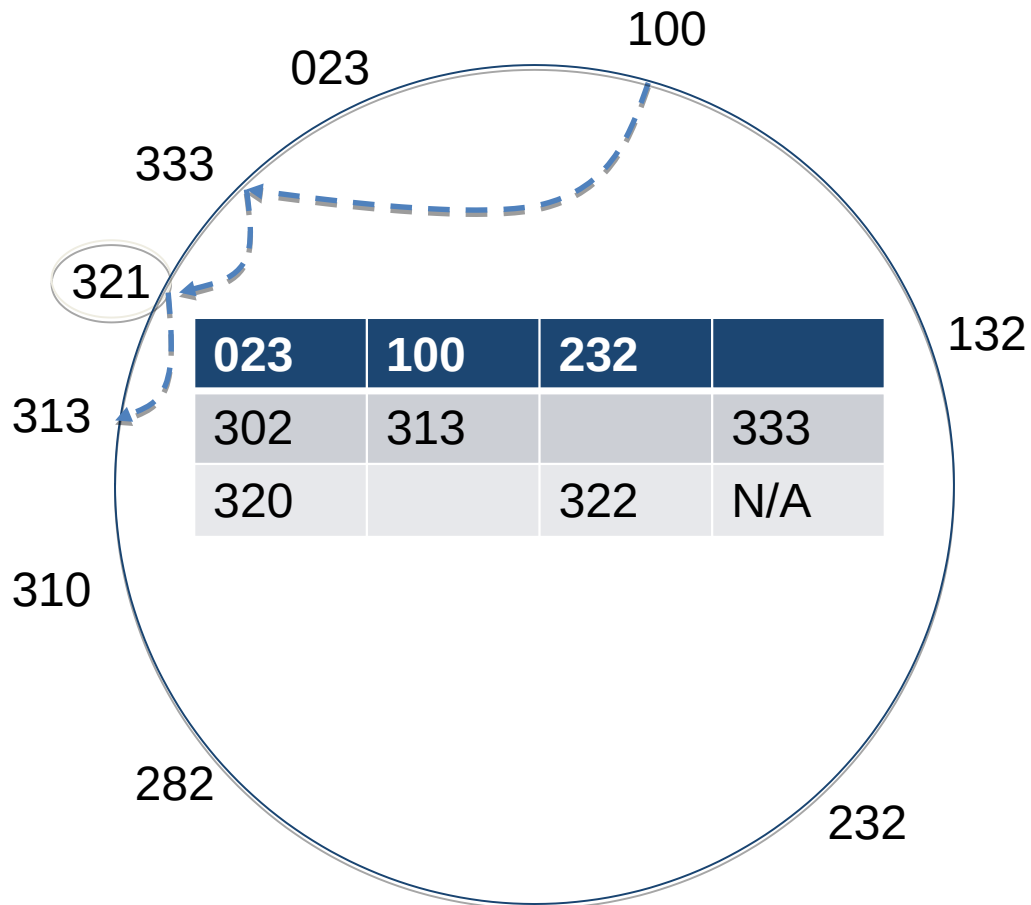


For space reasons some nodes are not shown on the ring

Pastry Example

b=4

Route to 311

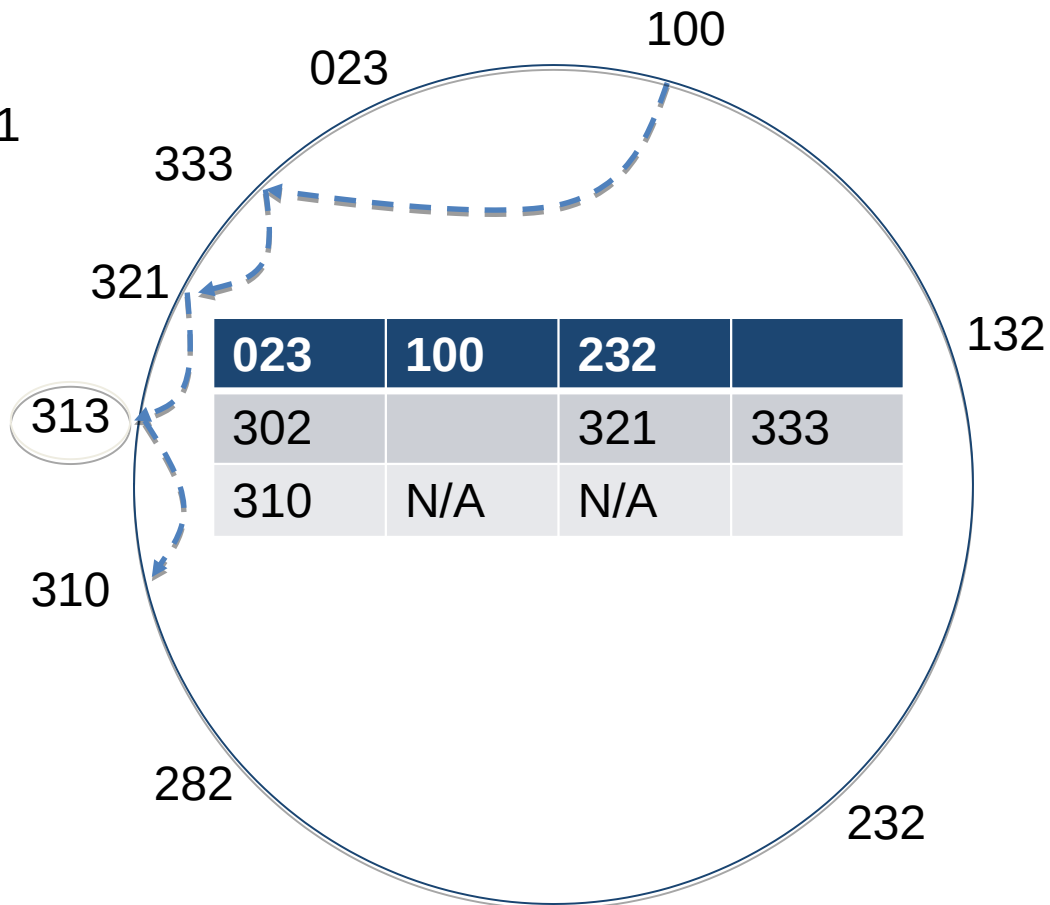


For space reasons some nodes are not shown on the ring

Pastry Example

b=4

Route to 311

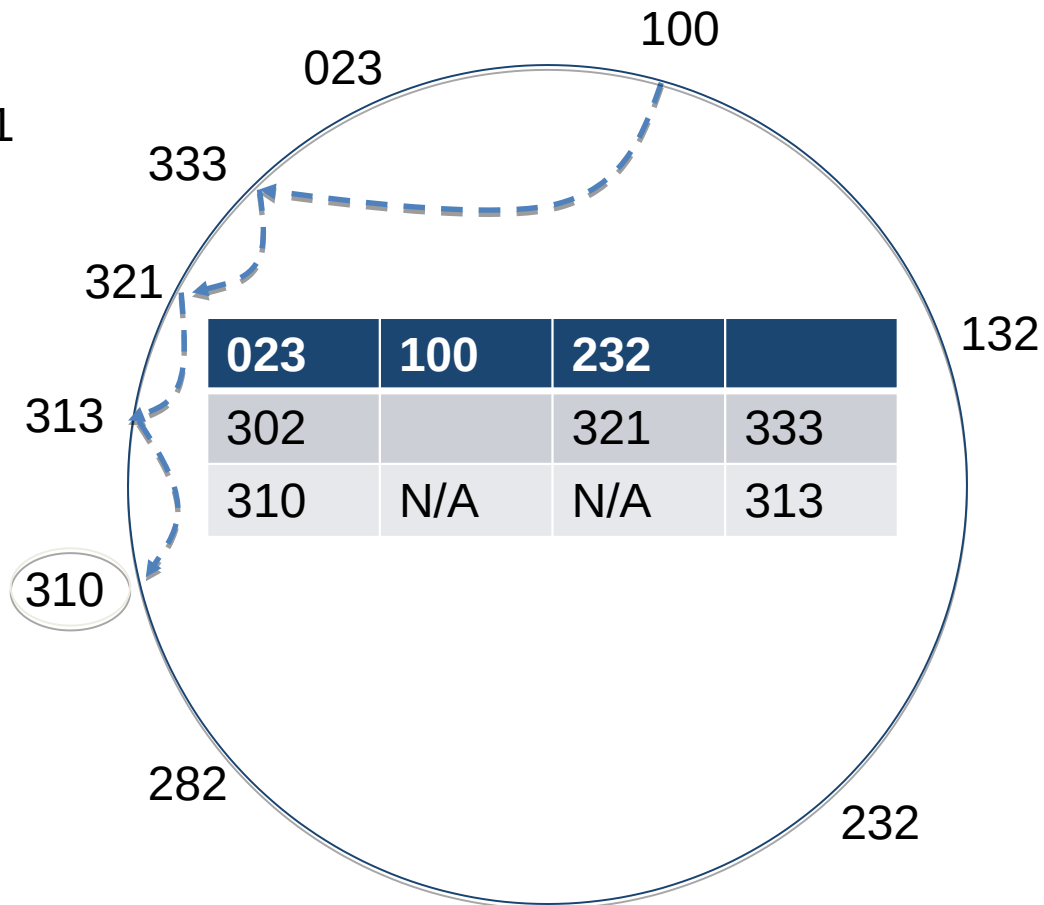


For space reasons some nodes are not shown on the ring

Pastry Example

b=4

Route to 311



For space reasons some nodes are not shown on the ring

Node departure

Explicit departure or failure

Replacement of a node

The leafset of the closest node in the leafset contains the
closest new node, not yet in the leafset

Update from the leafset information

Update the application

Failure detection

Detected when immediate neighbours in the name space
(leafset) can no longer communicate

Detected when a contact fails during the routing

Routing uses an alternative route

Fixing the routing table of A

- Repair

R_l^d : entry of the routing table of A to repair

A contacts another entry (at random) R_l^i from the same line

so that ($i \neq d$) and asks for entry R_l^d , otherwise another entry

from R_{l+1}^i ($i \neq d$) if no node in line l answers the request.

State maintenance

Leaf set

- is aggressively monitored and fixed

Routing table

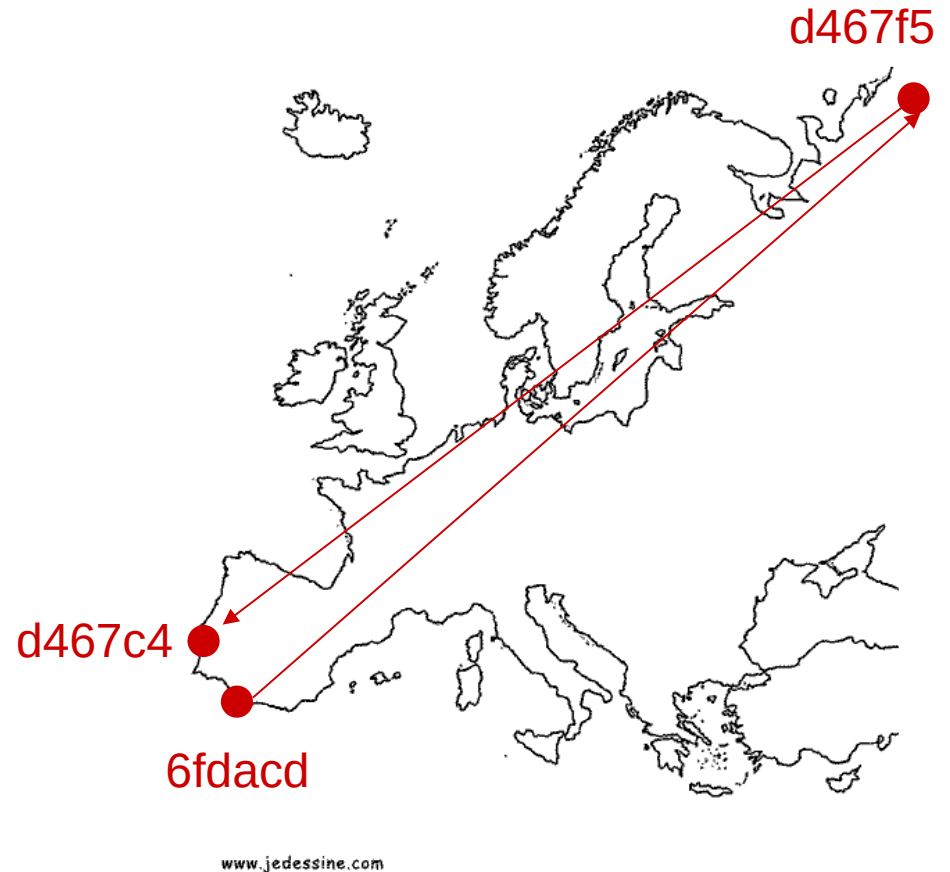
- are lazily repaired

When a hole is detected during the routing

- Periodic gossip-based maintenance

Reducing latency

- **Random assignment of nodeId:** Nodes numerically close are geographically (topologically) distant
- **Objective:** fill the routing table with nodes so that routing hops are as short (latency wise)



Exploiting locality in Pastry

Neighbour selected based of a network proximity

metric:

- Closest topological node
- Satisfying the constraints of the routing table

routeTable(i,j):

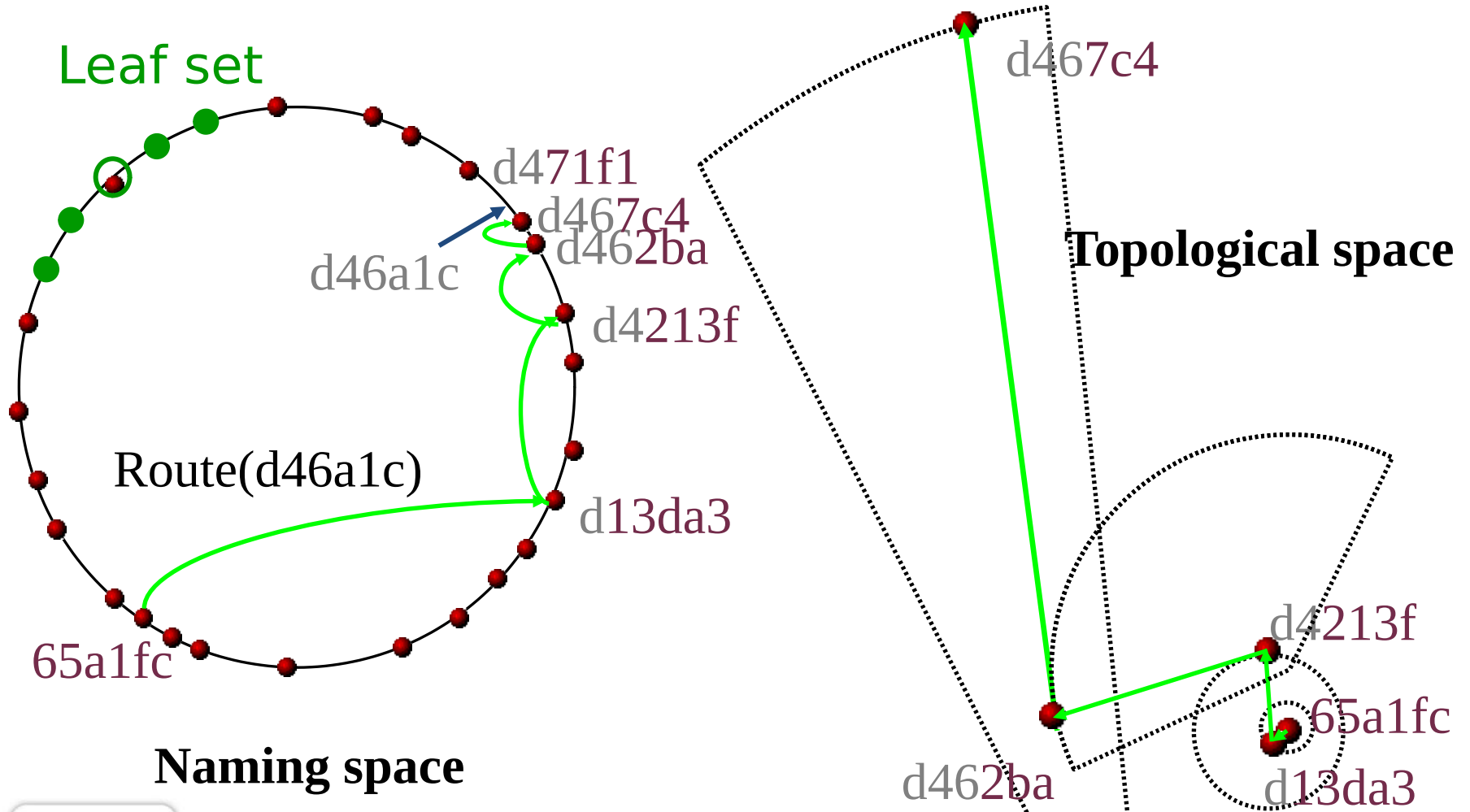
- *nodeId* corresponding to the current *nodeId* up to level i

next digit = j

- nodes are close at the top level of the routing table

- Farther nodes at the bottom levels of the routing tables

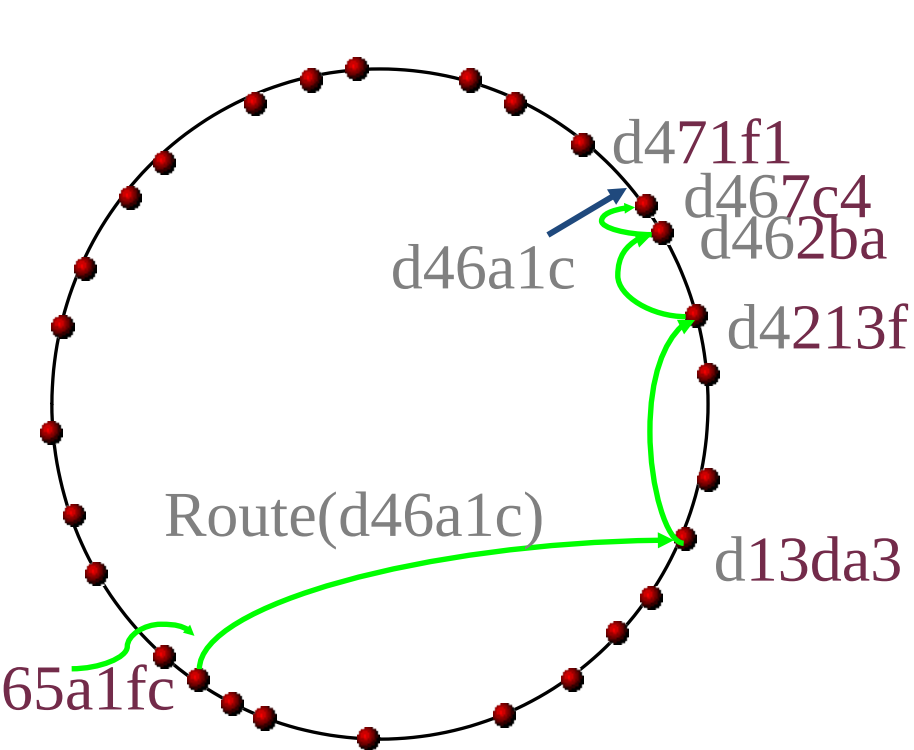
Proximity routing in Pastry



Locality

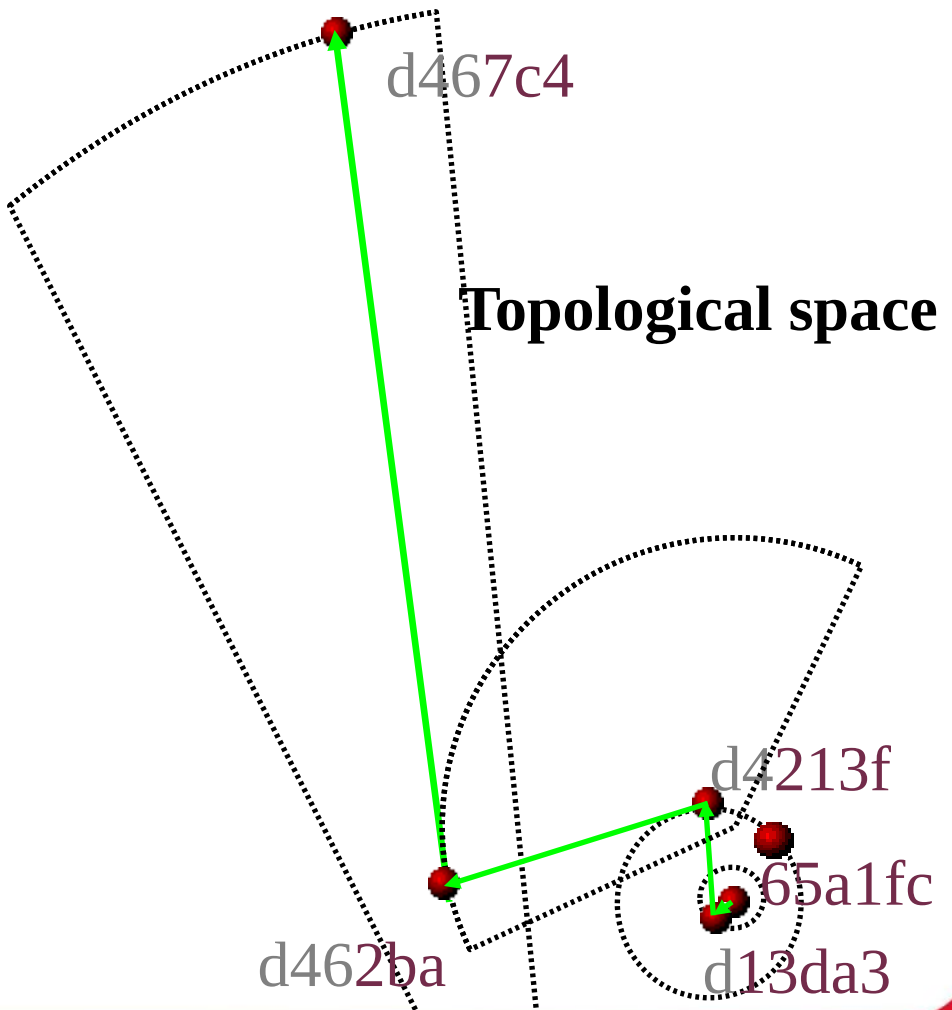
1. Joining node X routes asks A to route to X
 - Path A,B,... -> Z
 - Z numerically closest to X
 - X initializes line i of its routing table with the contents of line i of the routing table of the *ith* node encountered on the path
2. Improving the quality of the routing table
 - X asks to each node of its routing table its own routing state and compare distances
 - Gossip-based update for each line (20mn)
 - Periodically, an entry is chosen at random in the routing table
 - Corresponding line of this entry sent
 - Evaluation of potential candidates
 - Replacement of better candidates
 - New nodes gradually integrated

Node insertion in Pastry



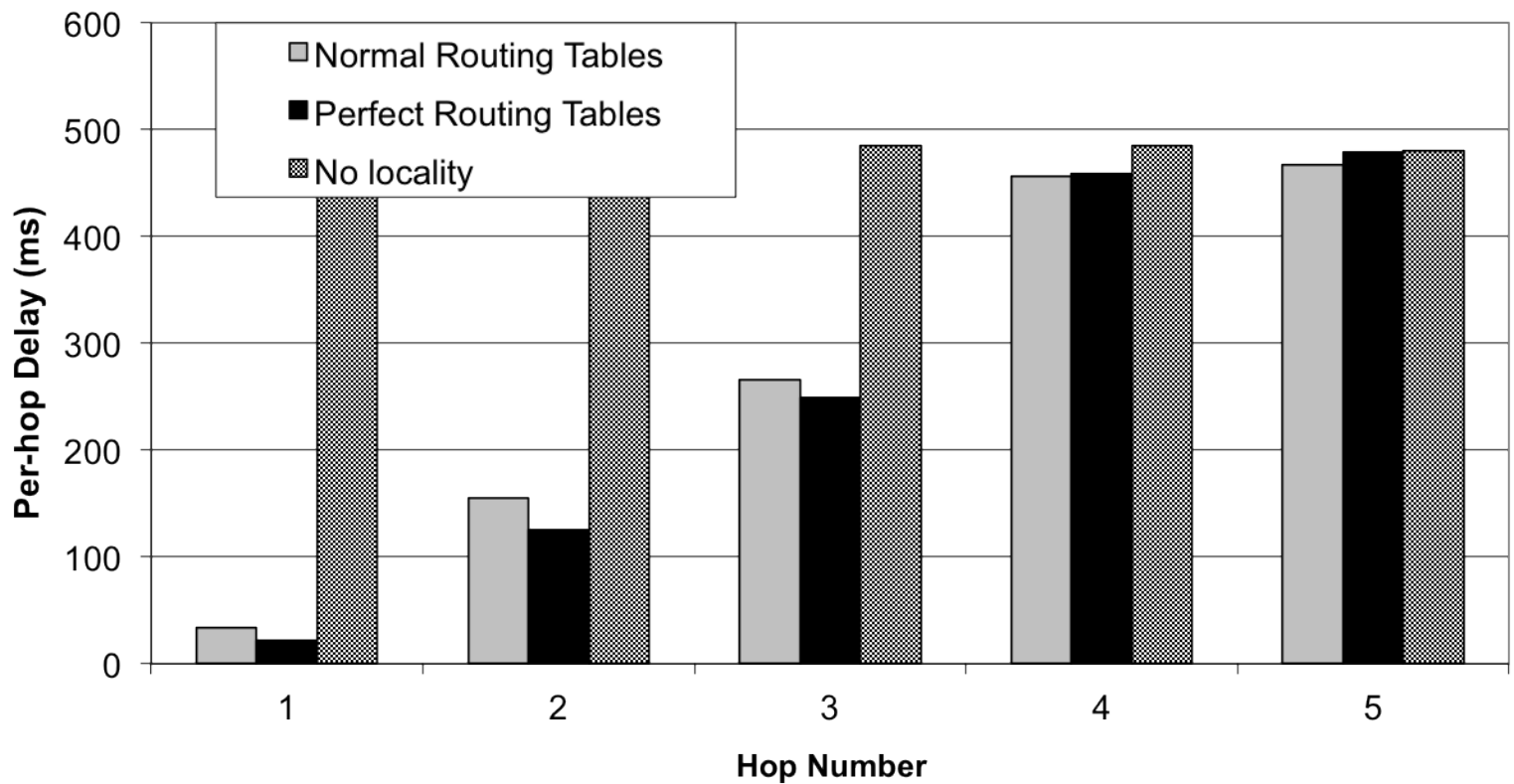
● New node: d46a1c

Naming space



Topological space

Performance 1.59 slower than IP on average



References

- Rowstron and P. Druschel, "**Pastry: Scalable, distributed object location and routing for large-scale peer-to-peer systems**", *Middleware'2001*, Germany, November 2001.