

Peer-to-Peer Networks

05 Pastry

Christian Schindelhauer

Technical Faculty

Computer-Networks and Telematics

University of Freiburg

- Peter Druschel
 - Rice University, Houston, Texas
 - now head of Max-Planck-Institute for Computer Science, Saarbrücken/
Kaiserslautern
- Antony Rowstron
 - Microsoft Research, Cambridge, GB
- Developed in Cambridge (Microsoft Research)
- Pastry
 - Scalable, decentralized object location and routing for large scale peer-to-peer-network
- PAST
 - A large-scale, persistent peer-to-peer storage utility
- Two names one P2P network
 - PAST is an application for Pastry enabling the full P2P data storage functionality
 - We concentrate on Pastry

↳ "Flaxton - Routing"

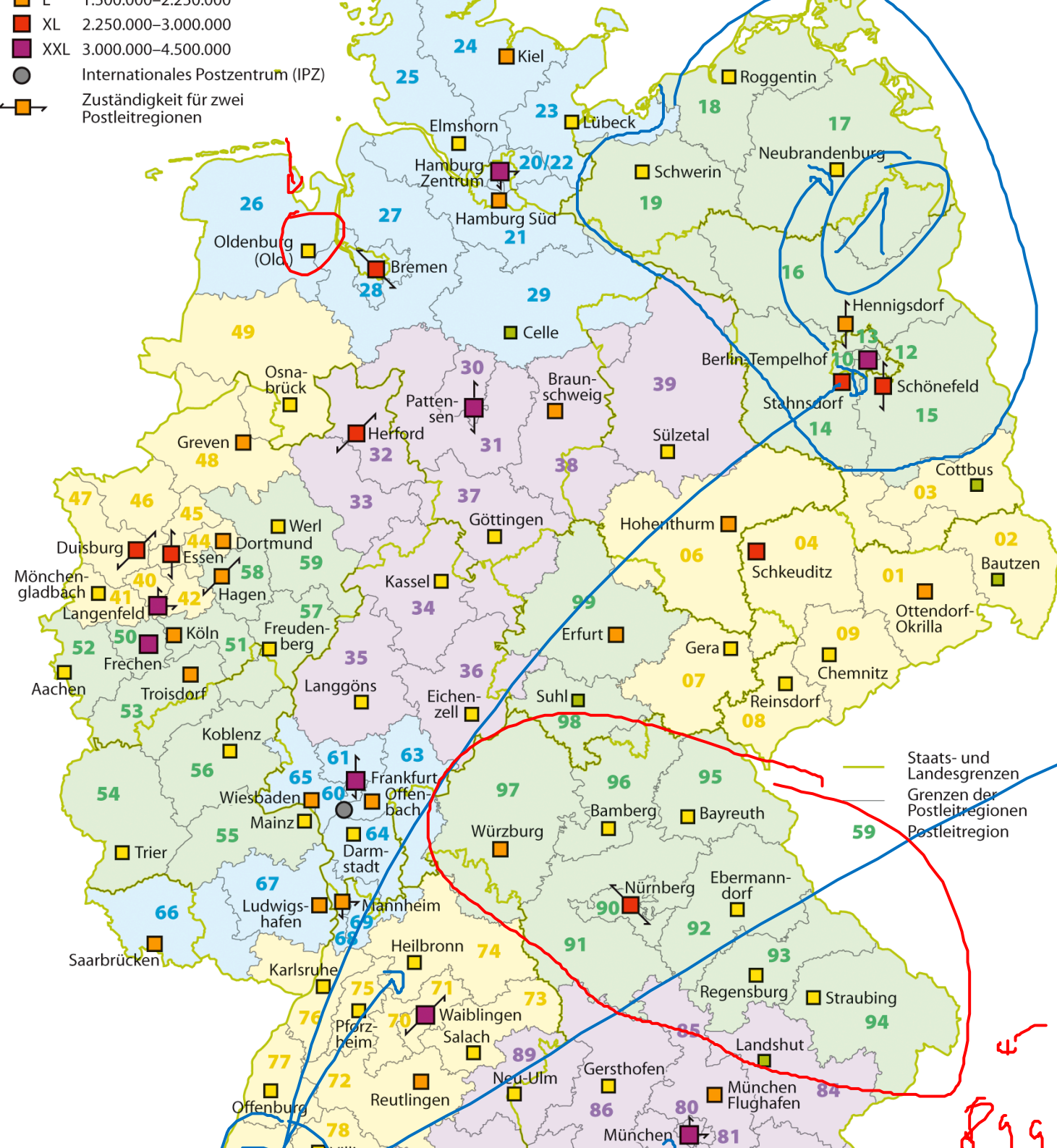
Pastry, Tapestry

ZIP-codes

PLZ

iburg

- L 1.500.000–2.250.000
- XL 2.250.000–3.000.000
- XXL 3.000.000–4.500.000
- Internationales Postzentrum (IPZ)
- ↔ Zuständigkeit für zwei Postleitregionen

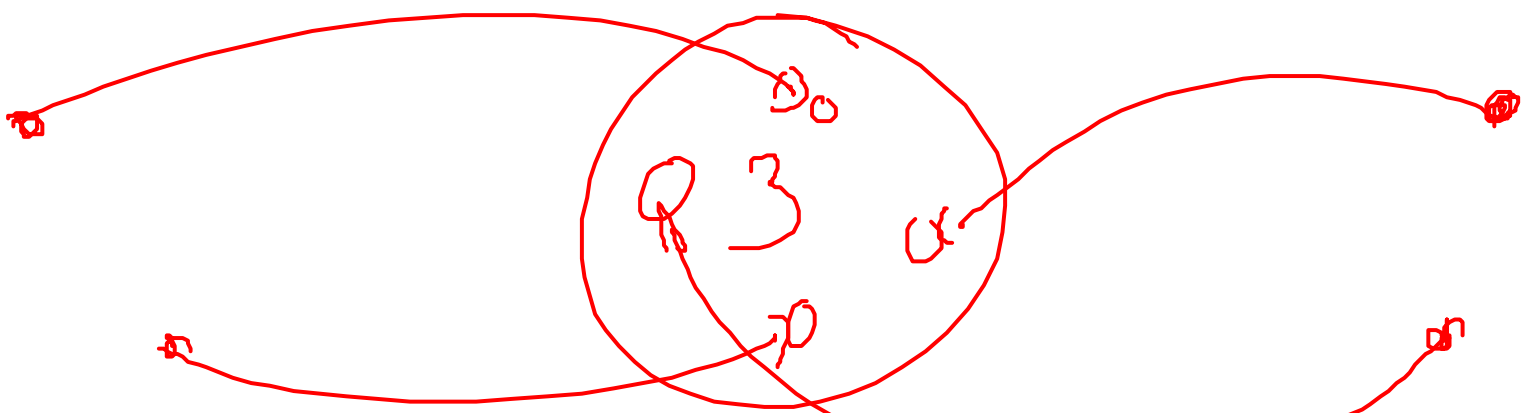
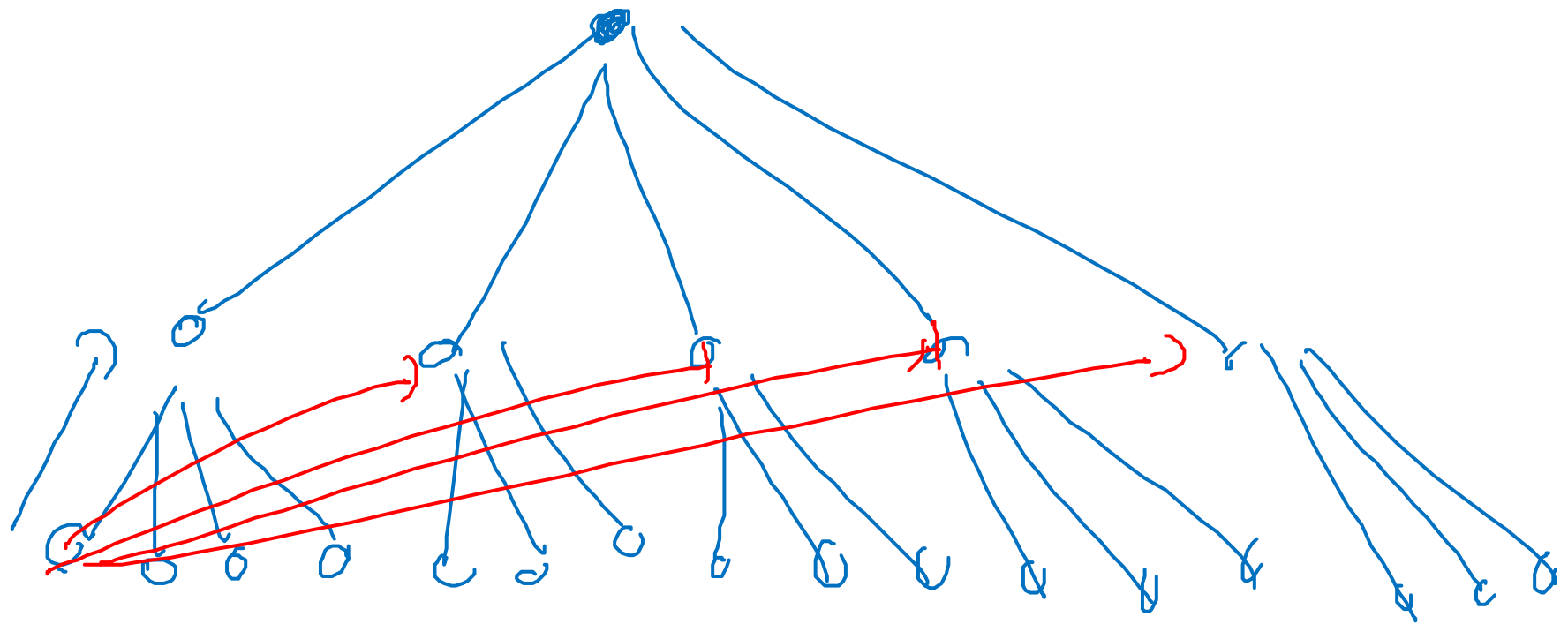


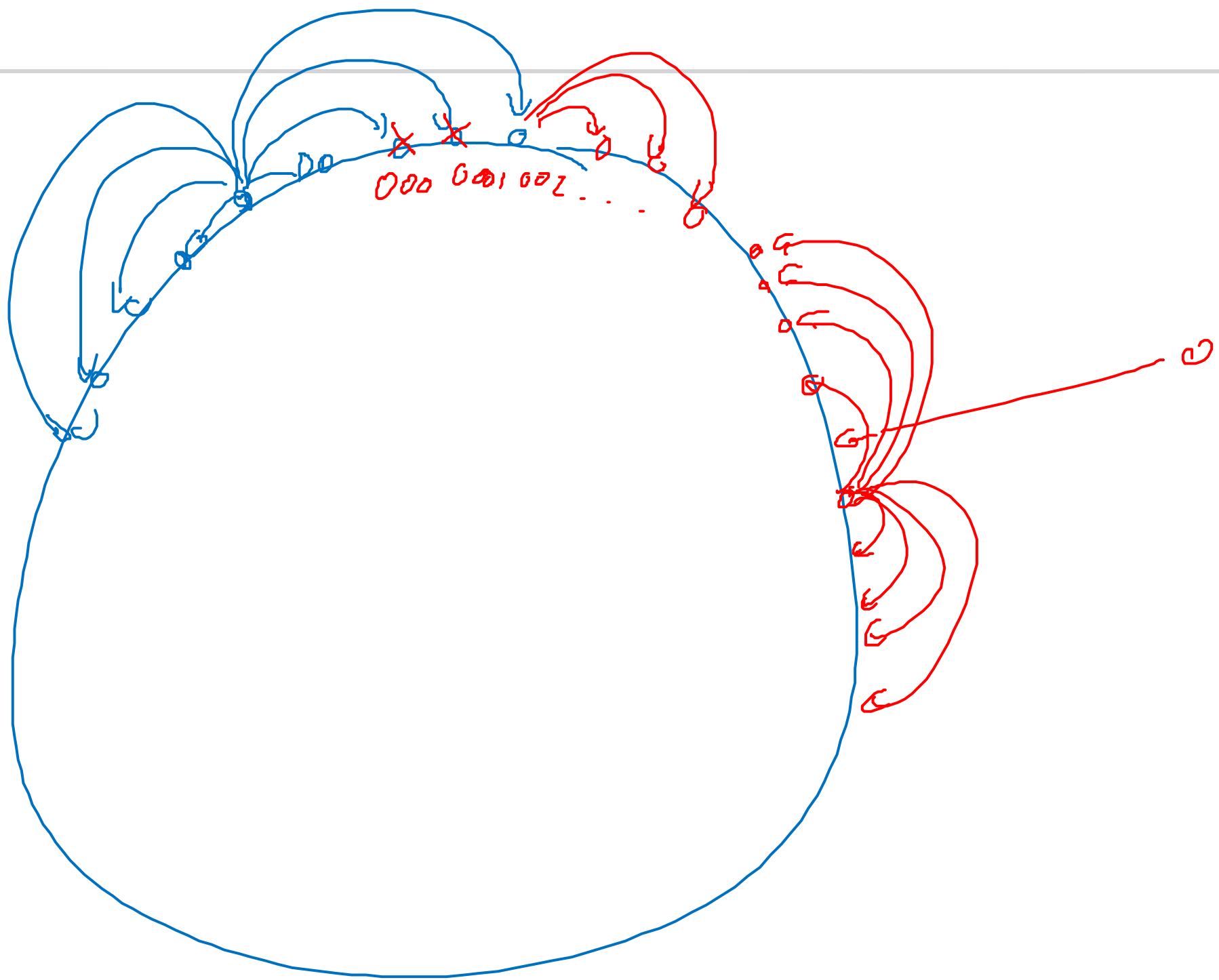
(79)
799
795
756

97532 → 0
9912
9.....

— Staats- und Landesgrenzen
— Grenzen der Postleitregionen
— Postleitregion

WRONG





- Each peer has a 128-bit ID: nodeID
- unique and uniformly distributed
- e.g. use cryptographic function applied to IP-address

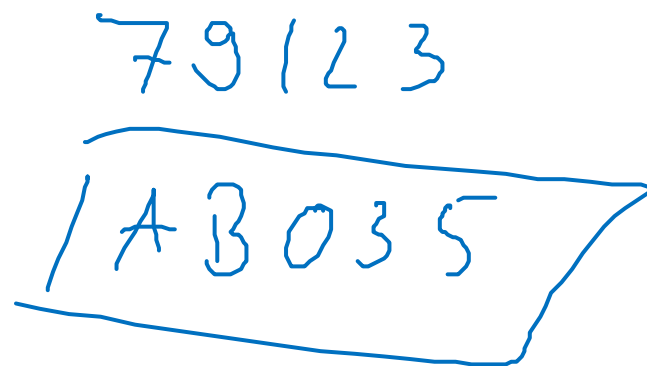
$$\left[\frac{1}{2^{128}} \cdot n \right] \quad \frac{1}{2^{128}}$$

- Routing
- Keys are matched to $\{0,1\}^{128}$
- According to a metric messages are distributed to the neighbor next to the target

Routing table has $O(2^b(\log n)/b) + \ell$ entries

- n: number of peers
- ℓ : configuration parameter

- b: word length
 - typical: $b=4$ (base 16), $\ell=16$



- message delivery is guaranteed as long as less than $\ell/2$ neighbored peers fail

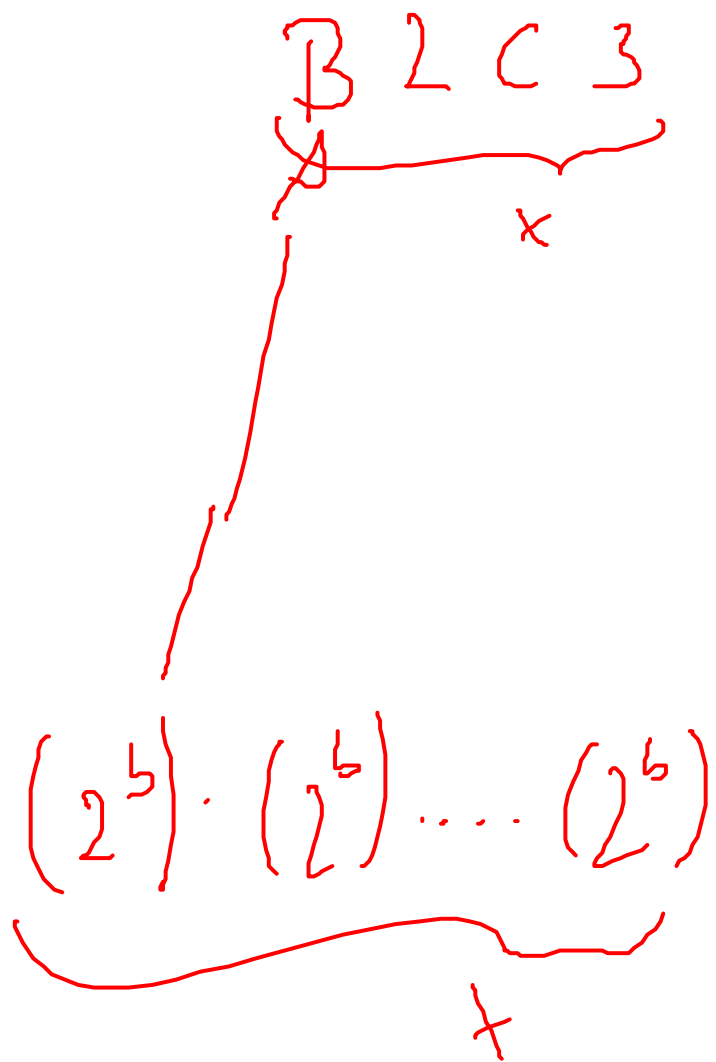
Inserting a peer and finding a key needs $O((\log n)/b)$ messages

~~65A0BA14~~
70000000

- NodeID presented in base 2^b
 - e.g. NodeID: 65A0BA13
- For each prefix p and letter $x \in \{0, \dots, 2^b - 1\}$ add an peer of form px^* to the routing table of NodeID, e.g.
 - $b=4, 2^b=16$
 - 15 entries for $0^*, 1^*, \dots, F^*$
 - 15 entries for $60^*, 61^*, \dots, 6F^*$
 - ...
 - if no peer of the form exists, then the entry remains empty
- Choose next neighbor according to a distance metric
 - metric results from the RTT (round trip time)
- In addition choose ℓ neighbors
 - $\ell/2$ with next higher ID
 - $\ell/2$ with next lower ID

<u>0</u>	<u>1</u>	<u>2</u>	3	4	5	6	7	8	9	a	b	c	d	e
x	x	x	x	x	x	x	x	x	x	x	x	x	x	x
6	6	6	6	6		6	6	6	6	6	6	6	6	6
0	1	2	3	4		6	7	8	9	a	b	c	d	e
x	x	x	x	x		x	x	x	x	x	x	x	x	x
6	6	6	6	6	6	6	6	6	6		6	6	6	6
5	5	5	5	5	5	5	5	5	5		5	5	5	5
0	1	2	3	4	5	6	7	8	9		b	c	d	e
x	x	x	x	x	x	x	x	x	x		x	x	x	x
6		6	6	6	6	6	6	6	6	6	6	6	6	6
5		5	5	5	5	5	5	5	5	5	5	5	5	5
a		a	a	a	a	a	a	a	a	a	a	a	a	a
0		2	3	4	5	6	7	8	9	a	b	c	d	e
x		x	x	x	x	x	x	x	x	x	x	x	x	x

Union Bound



~~Union Bound~~

$$\frac{n}{(2^b)^x} \leq \frac{1}{n^c}$$

$$2^{b \cdot x} \geq n^{c+1}$$

$$b \cdot x \geq (\log_2 n) \cdot (c+1)$$

$$x \geq \frac{1}{b} \cdot (\log_2 n) \cdot (c+1)$$

Example b=2

Routing Table

- For each prefix p and letter $x \in \{0, \dots, 2^b - 1\}$ add an peer of form px^* to the routing table of NodeID

In addition choose ℓ neighbors

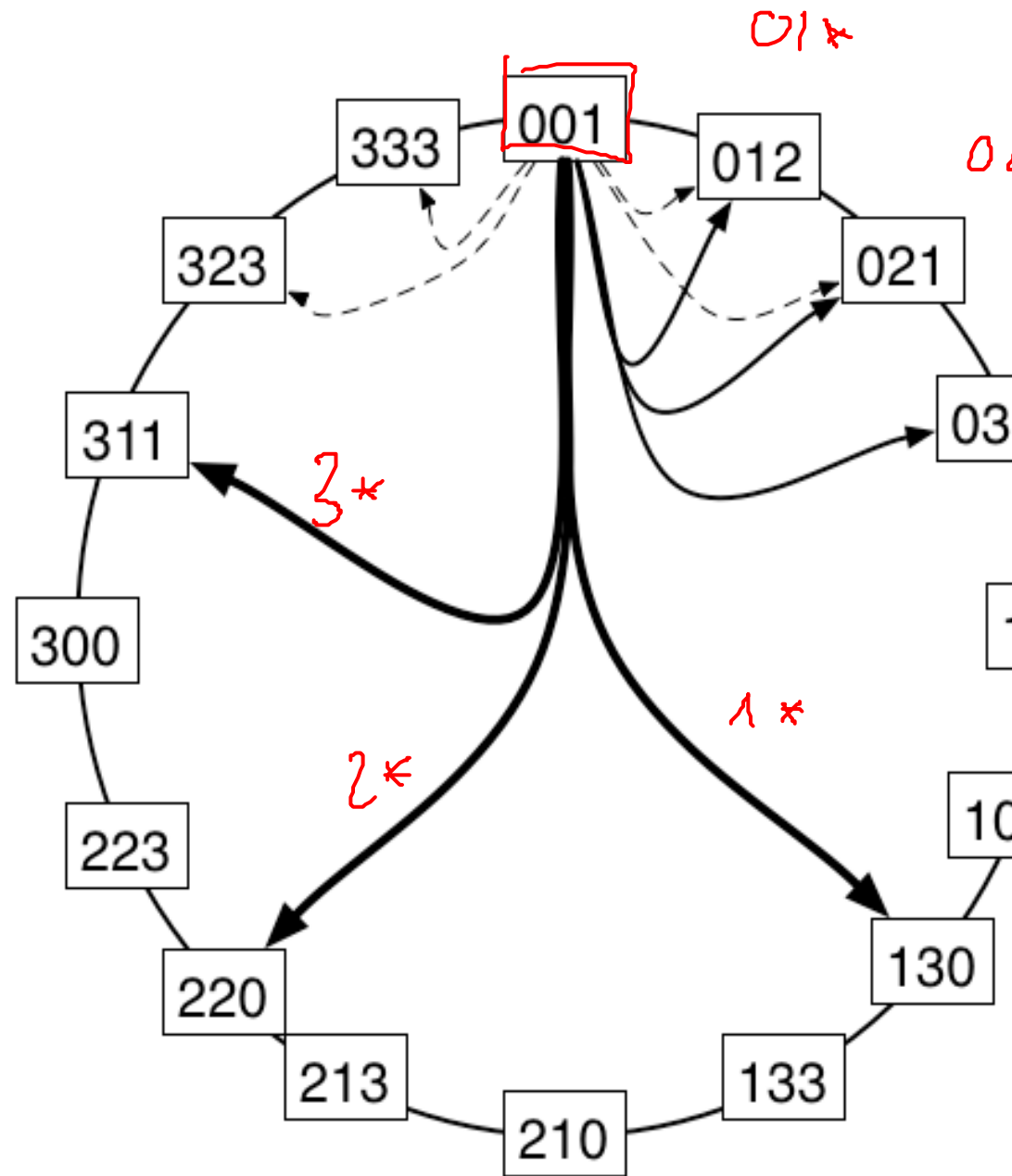
- $\ell/2$ with next higher ID
- $\ell/2$ with next lower ID

Observation

- The leaf-set alone can be used to find a target

Theorem

- With high probability there are at most $O(2^b (\log n)/b)$ entries in



Theorem

With high probability there are at most $O(2^b (\log n)/b)$ entries in each routing table

Proof

The probability that a peer gets the same m-digit prefix is 2^{-bm}

The probability that a m-digit prefix is unused is $(1 - 2^{-bm})^n \leq e^{-n/2^{bm}}$

For $m=c (\log n)/b$ we get

$$e^{-n/2^{bm}} \leq e^{-n/2^{c \log n}} \leq e^{-n/n^c} \leq e^{-n^{c-1}}$$

With (extremely) high probability there is no peer with the same prefix of length $(1+\epsilon)(\log n)/b$

Hence we have $(1+\epsilon)(\log n)/b$ rows with $2^b - 1$ entries each

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

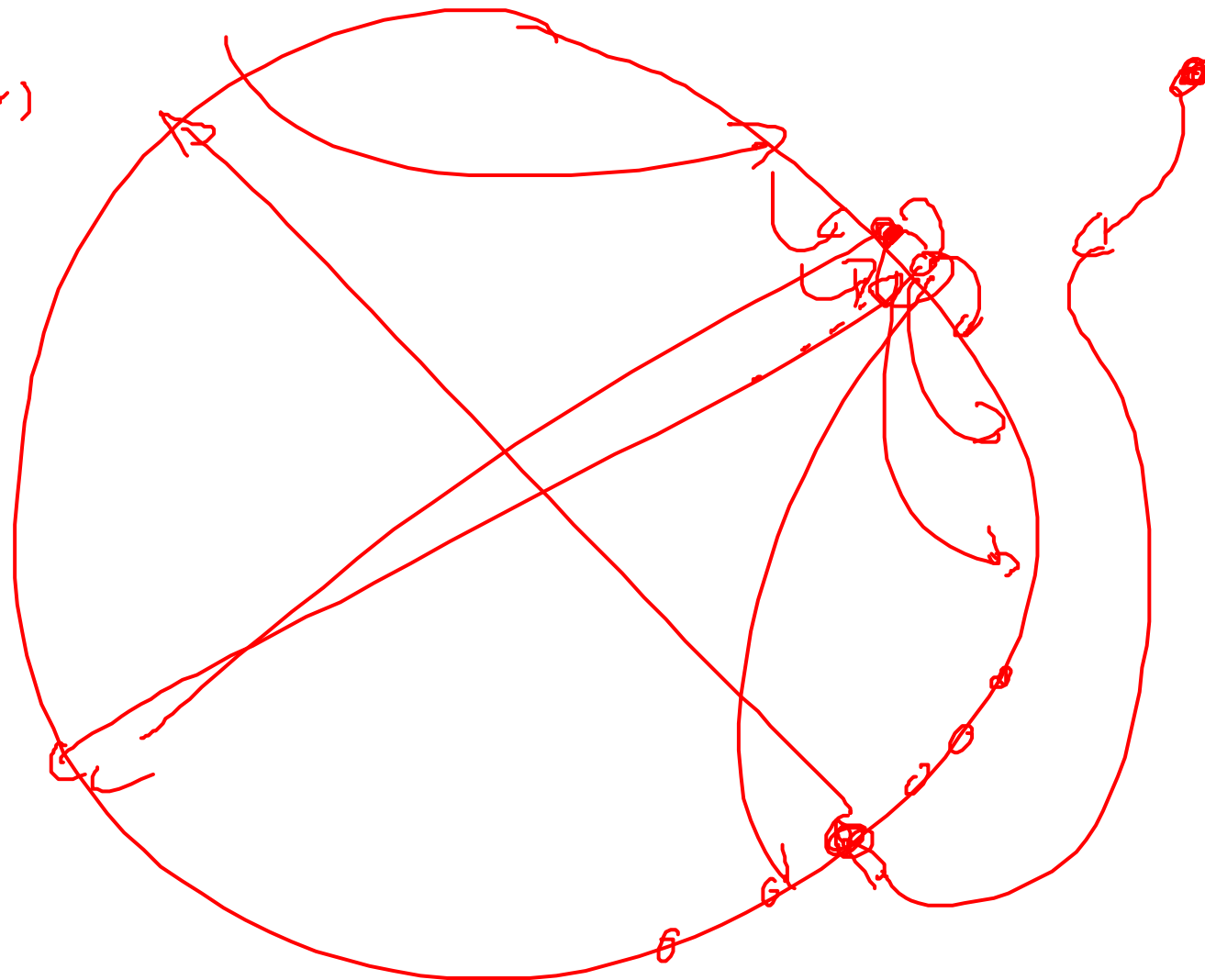
okup
ing neighbors

Fingers

coming
Finger

You

17



New node x sends message to the node z with the longest common prefix p

x receives

- routing table of z
- leaf set of z

x updates leaf-set

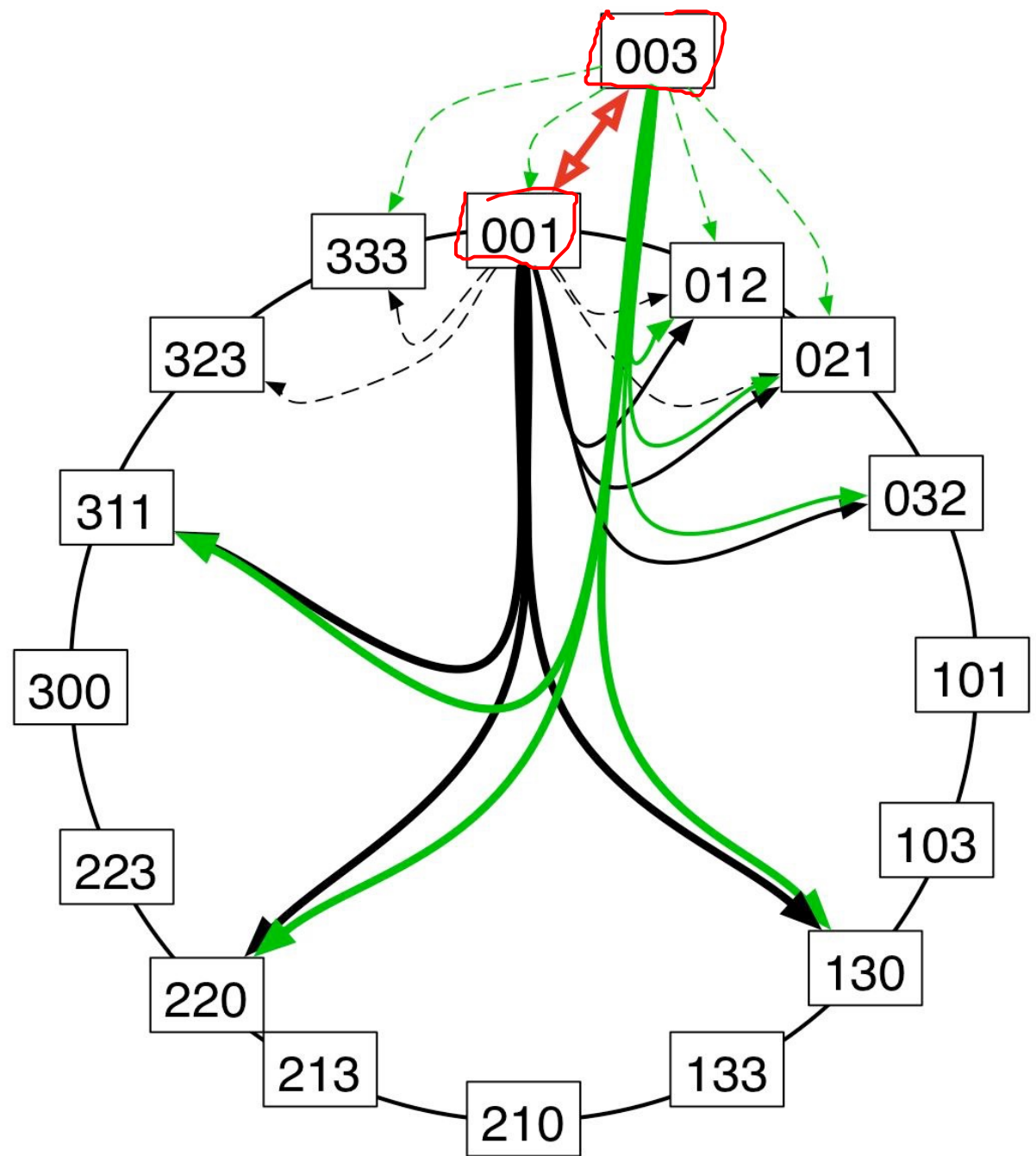
x informs ~~informant~~ ℓ -leaf set

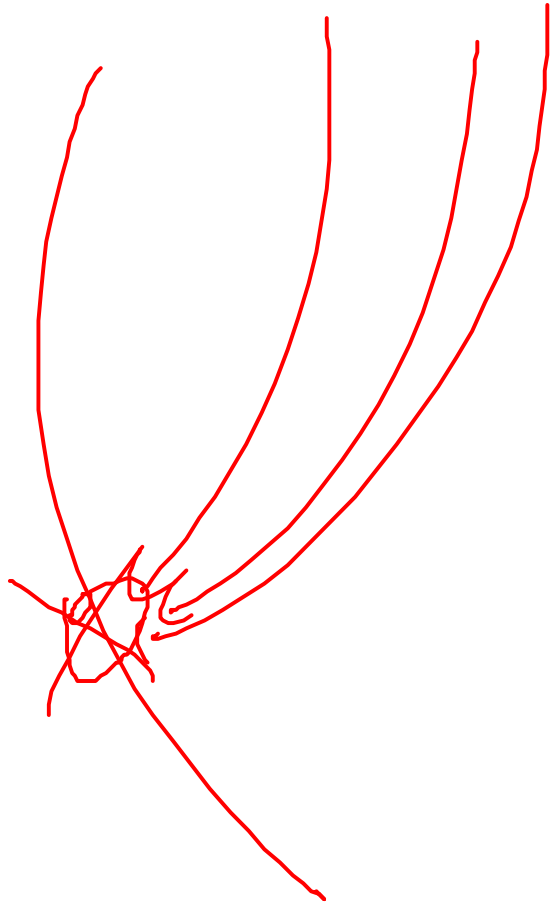
x informs peers in routing table

- with same prefix p (if $\ell/2 < 2^b$)

Number of messages for adding a peer

- ℓ messages to the leaf-set
- expected $(2^b - \ell/2)$ messages to nodes with common prefix
- one message to z with answer





■ Inheriting the next neighbor routing table does not allow work perfectly

■ Example

- If no peer with 1* exists then all other peers have to point to the new node

- Inserting 11

- 03 knows from its routing table

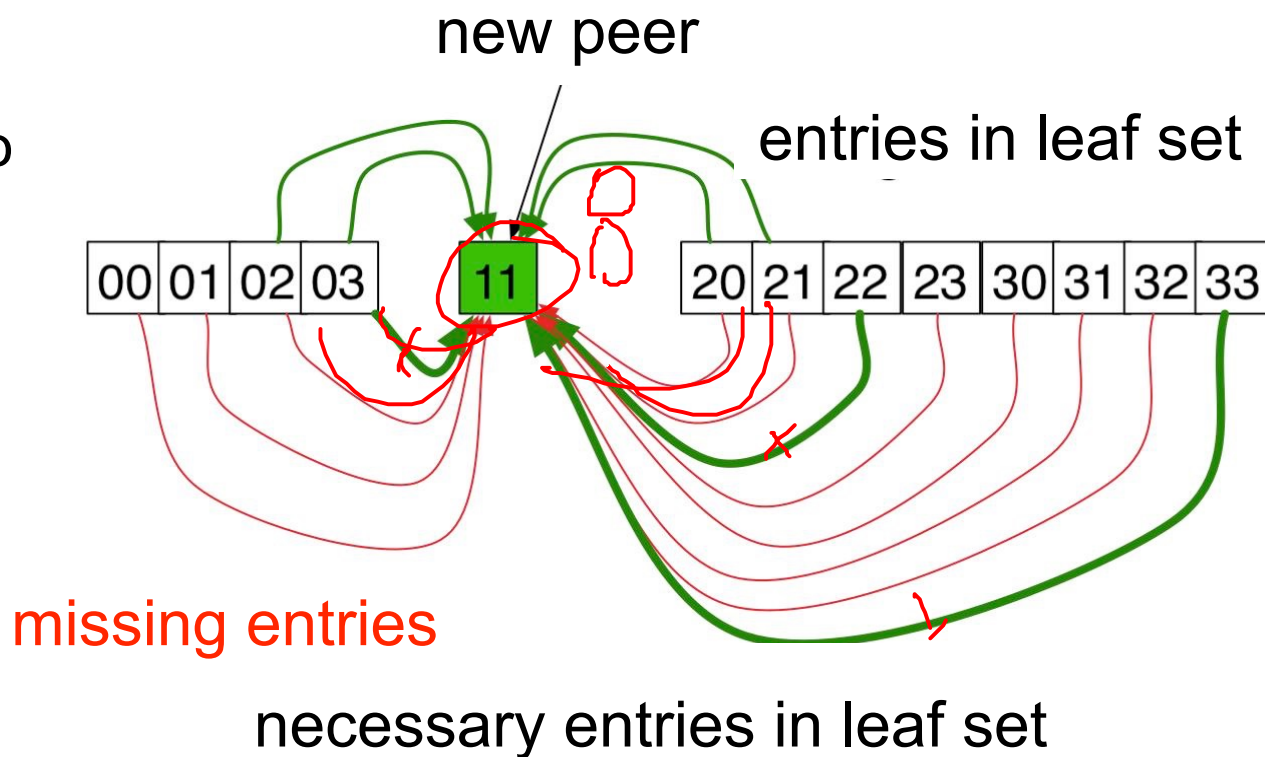
- 22,33

- 00,01,02

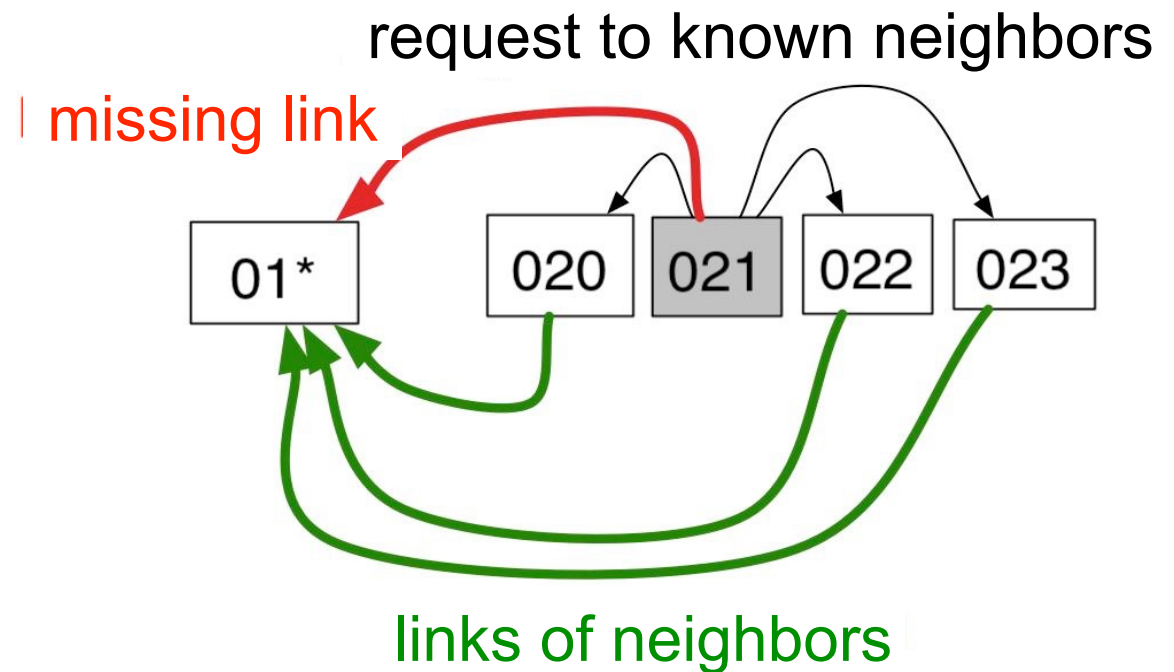
- 02 knows from the leaf-set

- 01,02,20,21

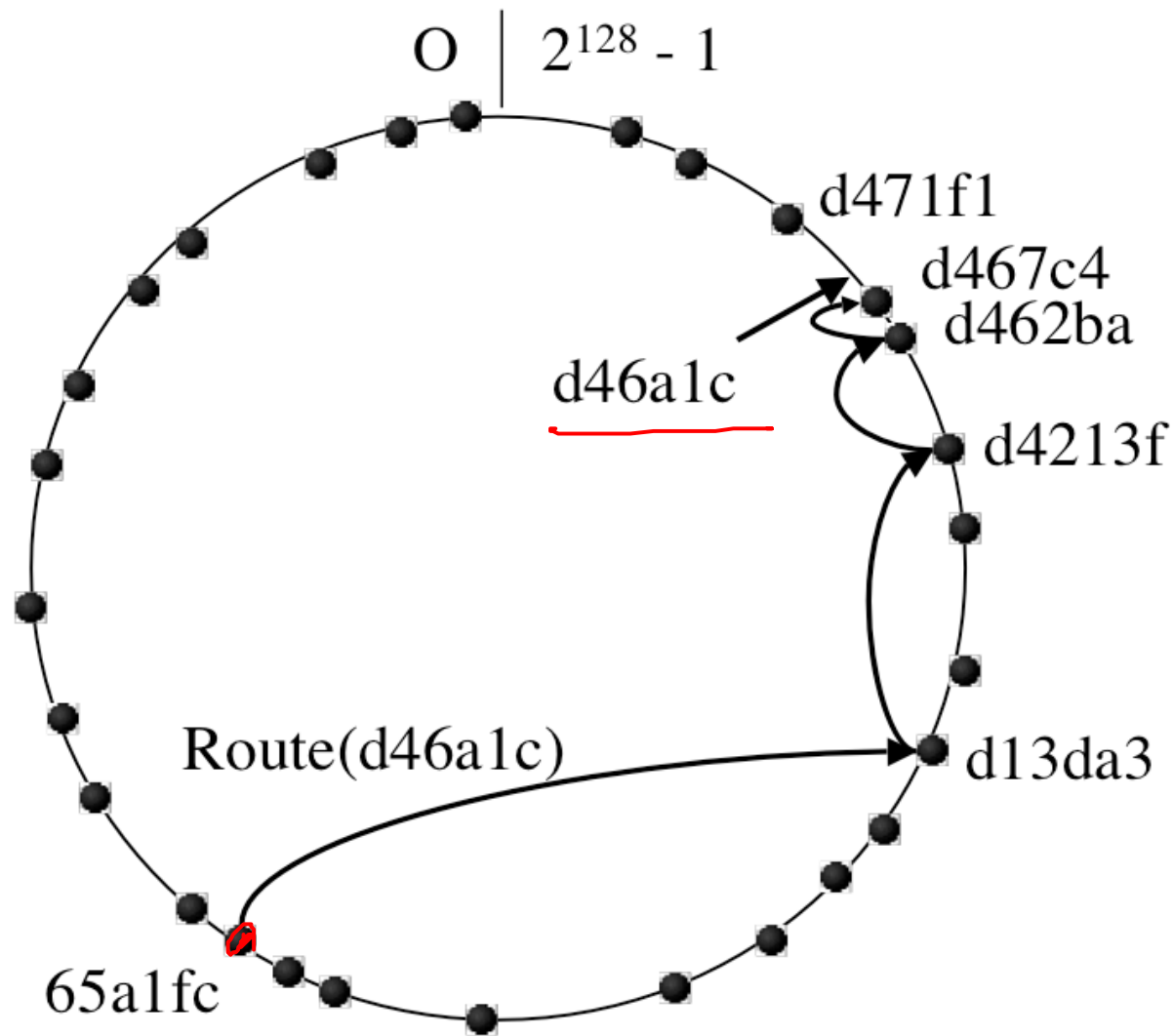
■ 11 cannot add all necessary links to the routing tables



- Assume the entry R_i^j is missing at peer D
 - j-th row and i-th column of the routing table
- This is noticed if message of a peer with such a prefix is received
- This may also happen if a peer leaves the network
- Contact peers in the same row
 - if they know a peer this address is copied
- If this fails then perform routing to the missing link



- Compute the target ID using the hash function
- If the address is within the ℓ -leaf set
 - the message is sent directly
 - or it discovers that the target is missing
- Else use the address in the routing table to forward the message
- If this fails take best fit from all addresses



l-leafset

~~routing table~~

nodes in the vicinity of D
(according to RTT)

key

nodeID of current peer

j-th row and i-th column of
the routing table

numbering of the leaf set

i-th digit of key D

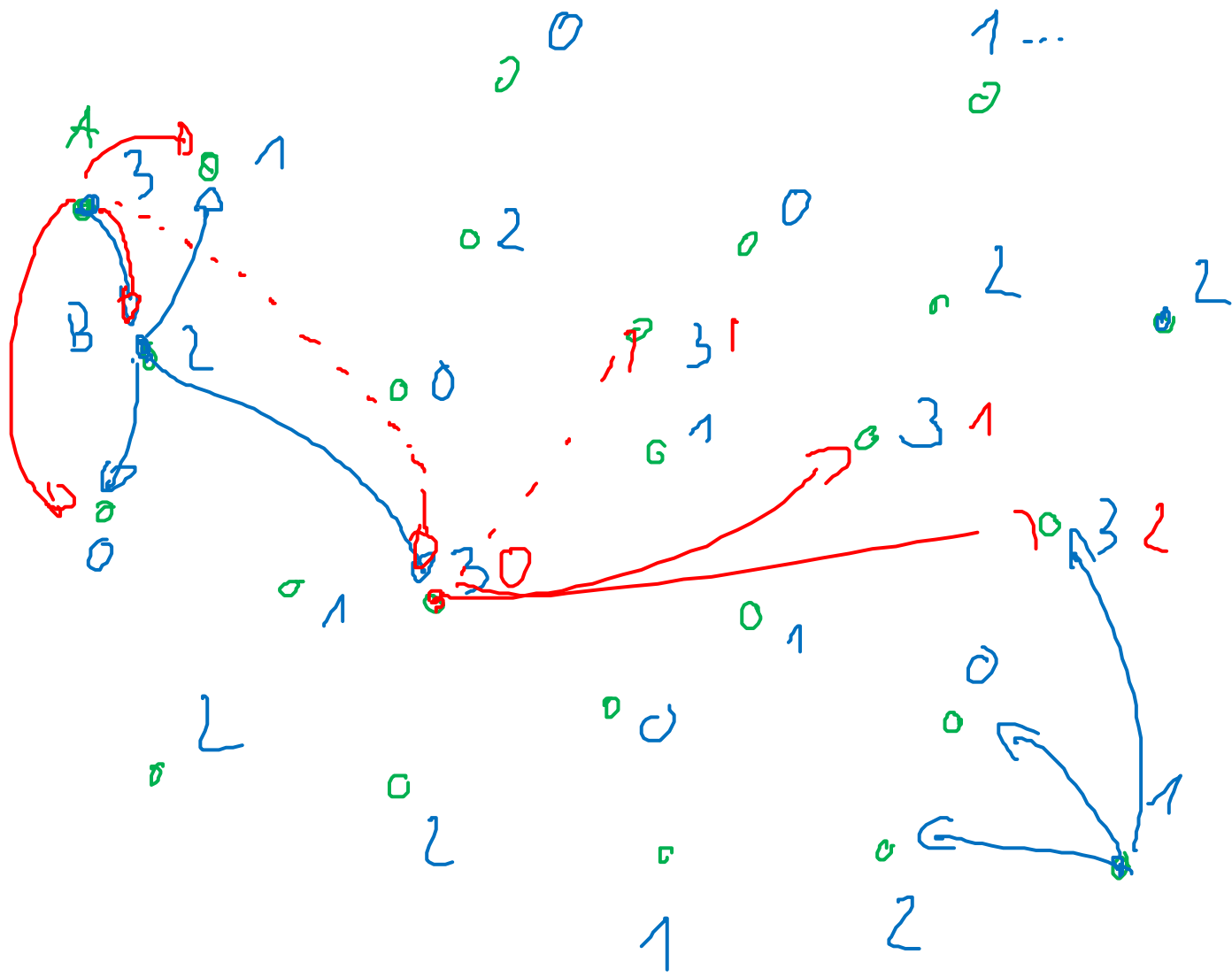
l(A): length of the largest

common

prefix of A and D
(shared header length)

- (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) }

„magic set“



- If the Routing-Table is correct
 - routing needs ~~$O((\log n)/b)$~~ ^{$+ O(1)$} messages
- As long as the leaf-set is correct
 - routing needs $O(n/l)$ messages
 - unrealistic worst case since even damaged routing tables allow dramatic speedup
- Routing does not use the real distances
 - M is used only if errors in the routing table occur
 - using locality improvements are possible
- Thus, Pastry uses heuristics for improving the lookup time
 - these are applied to the last, most expensive, hops

(x, y, z)

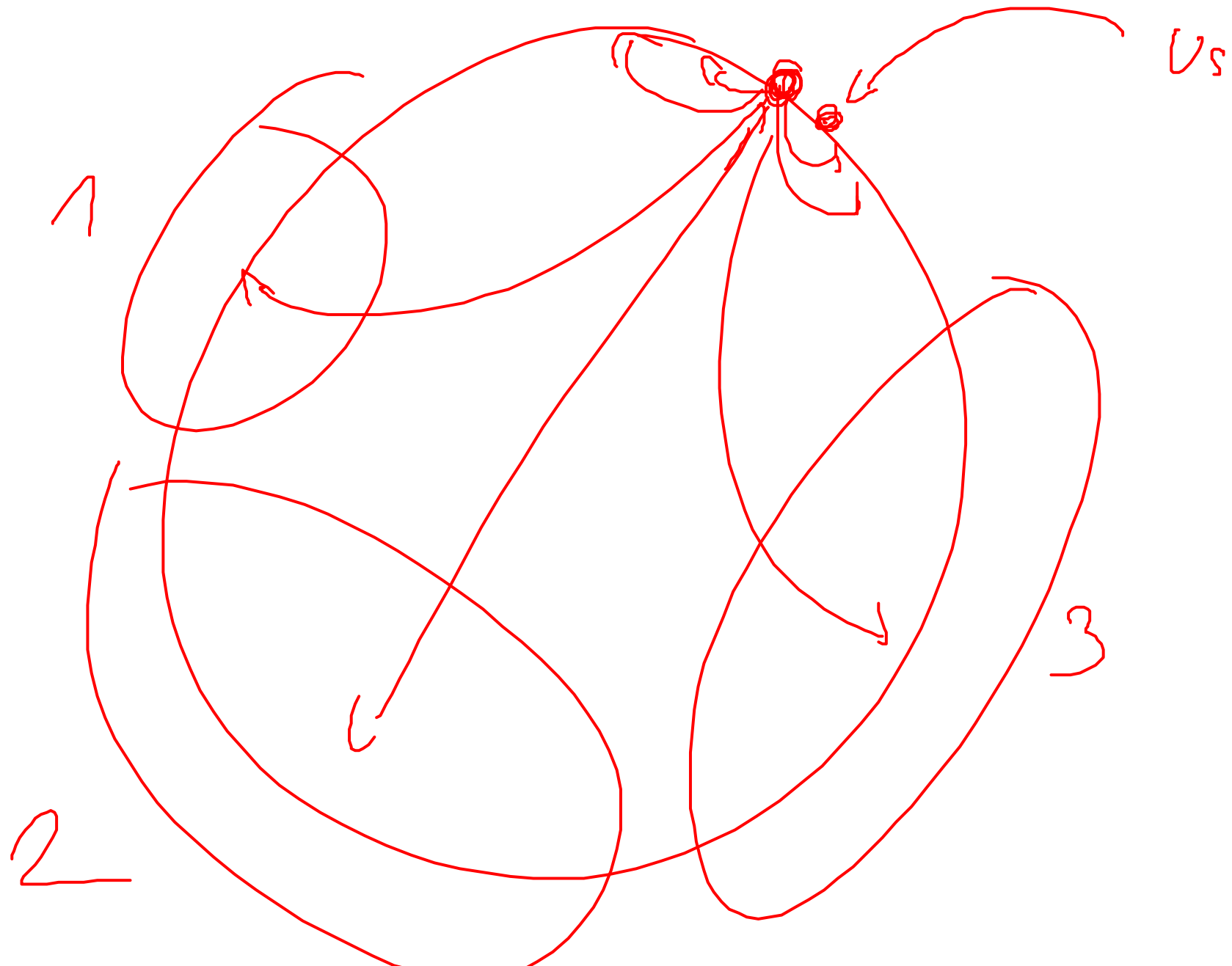
- Leaf-set peers are not near, e.g.
 - New Zealand, California, India, ...
- TCP protocol measures latency
 - latencies (RTT) can define a metric
 - this forms the foundation for finding the nearest peers
- All methods of Pastry are based on heuristics
 - i.e. no rigorous (mathematical) proof of efficiency
- Assumption: metric is Euclidean

metric $d(A, B)$ ≥ 0

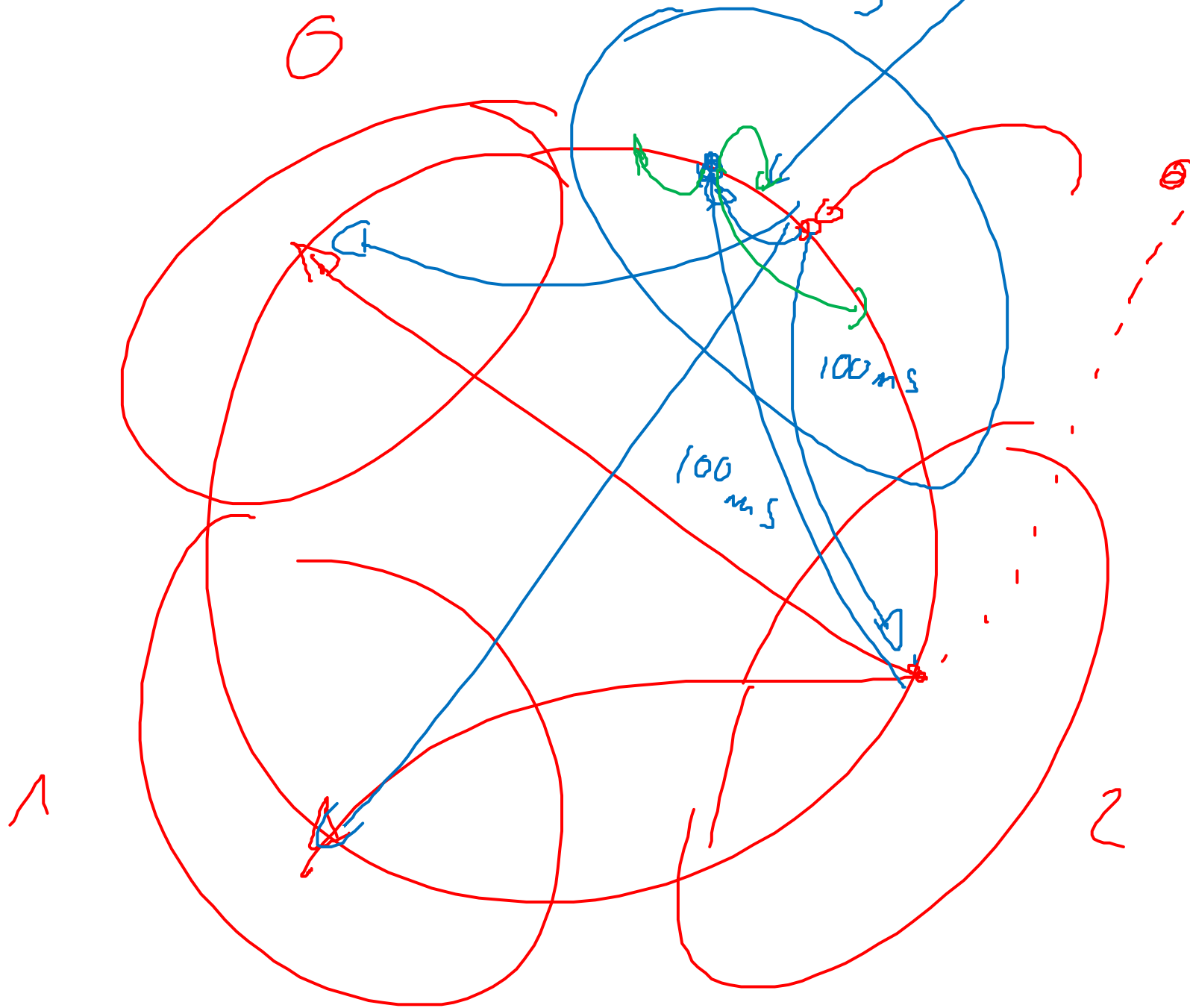
1. $d(A, A) = 0$
2. $d(A, B) = d(B, A)$

$d(A, C) \leq d(A, B) + d(B, C)$

3. Triangle Property



$\leq 200 \text{ms}$



Assumption

- When a peer is inserted the peers contacts a near peer
- All peers have optimized routing tables

But:

- The first contact is not necessary near according to the node-ID

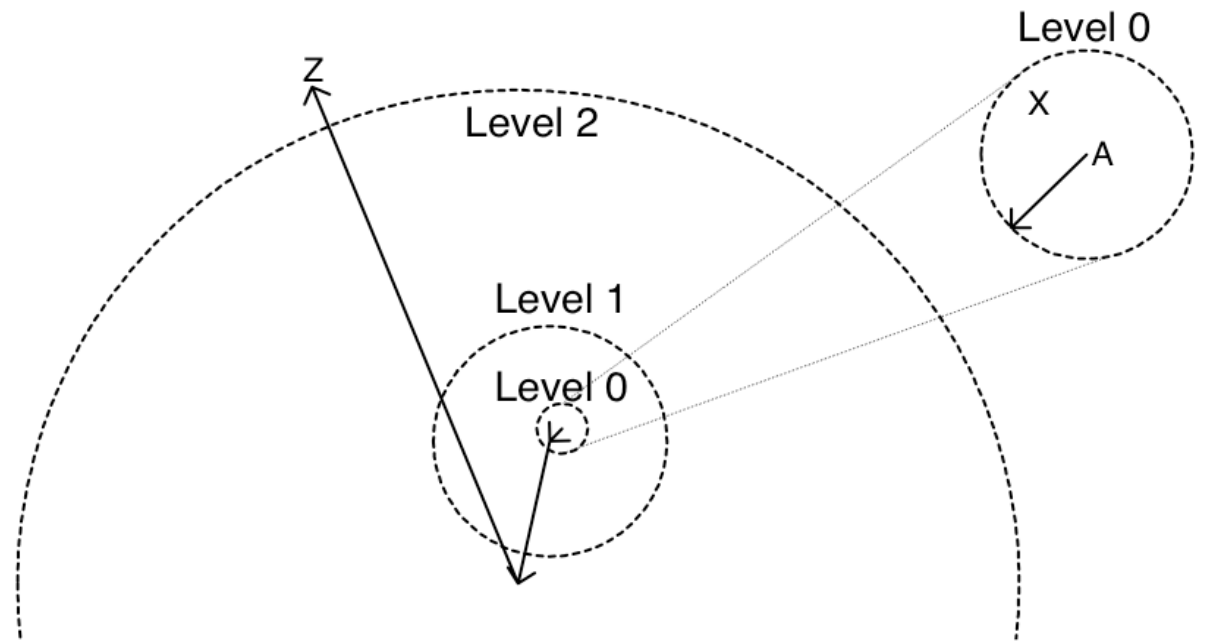
1st step

- Copy entries of the first row of the routing table of P
 - good approximation because of the triangle inequality (metric)

2nd step

- Contact fitting peer p' of p with the same first letter
- Again the entries are relatively close

Repeat these steps until all entries

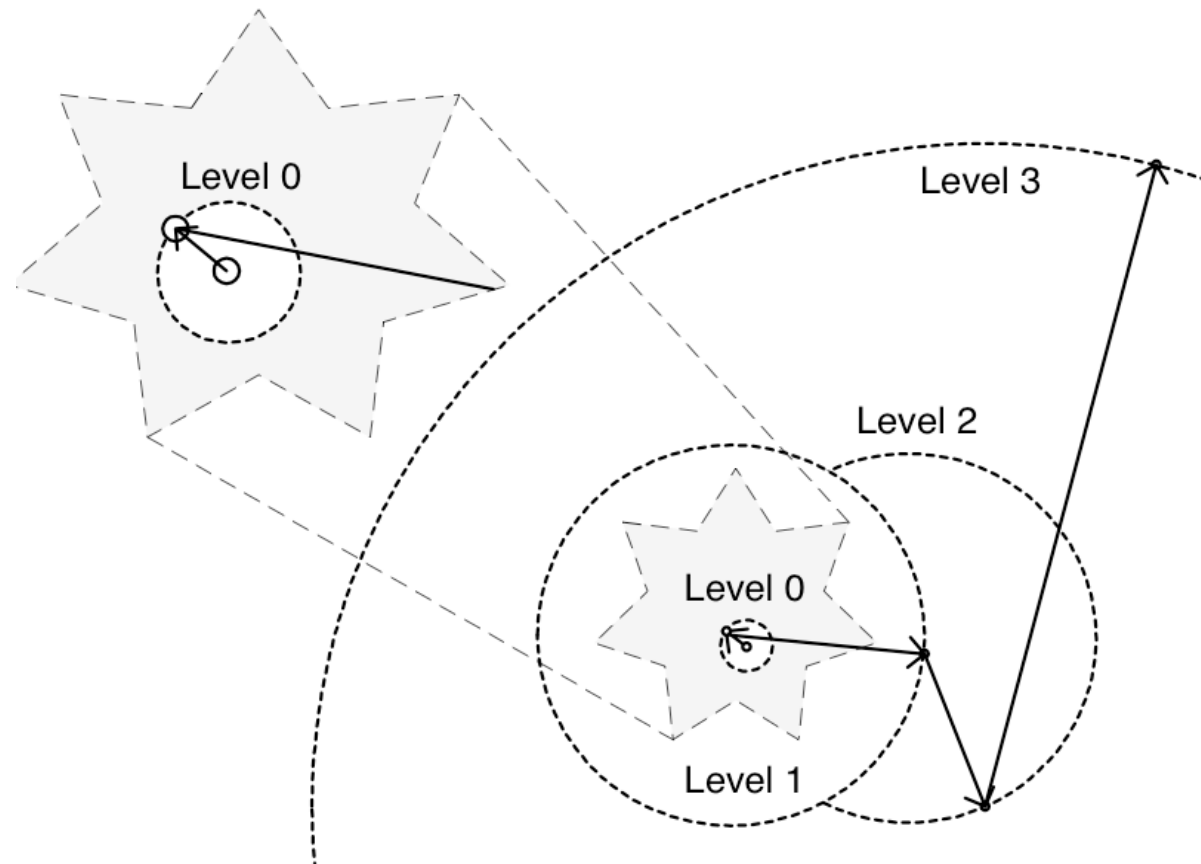


in the best case

- each entry in the routing table is optimal w.r.t. distance metric
- this does not lead to the shortest path

there is hope for short lookup times

- with the length of the common prefix the latency metric grows exponentially
- the last hops are the most expensive ones
- here the leaf-set entries help



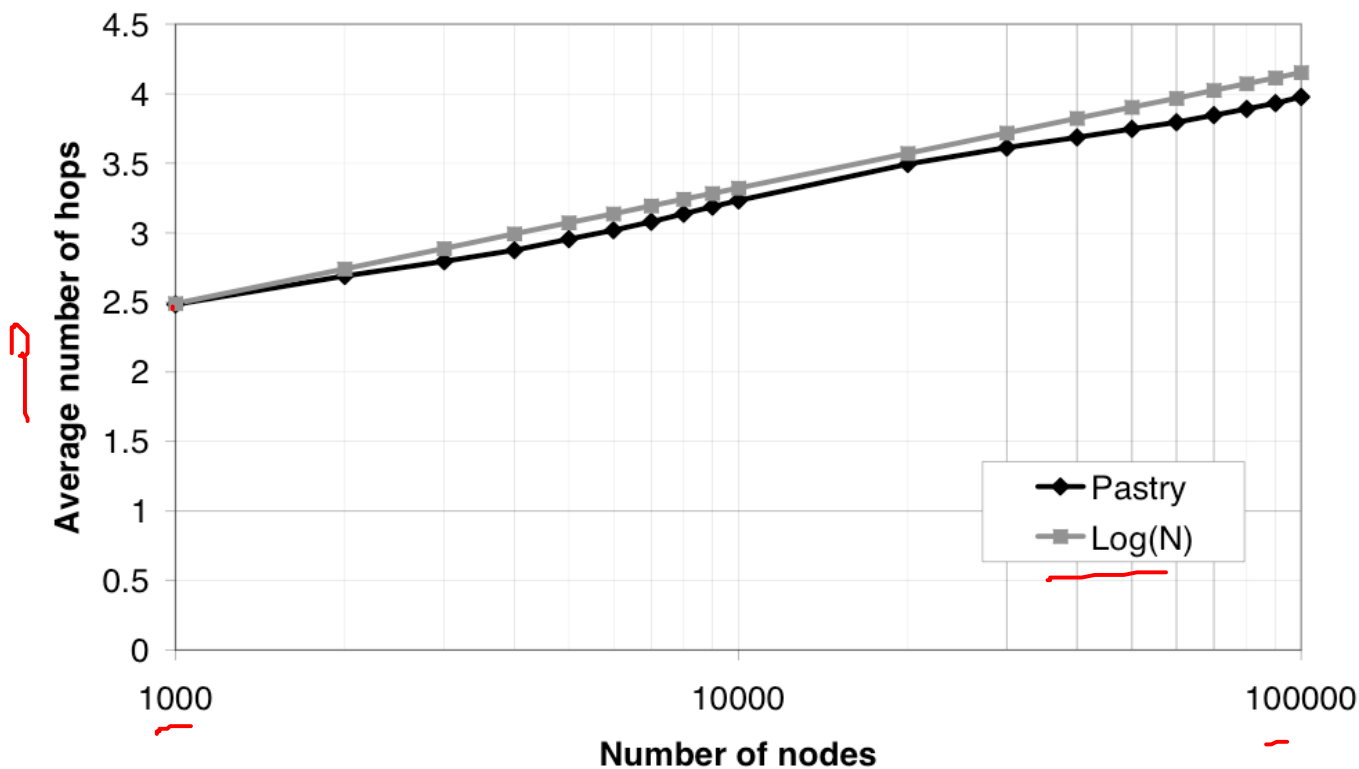
- Node-ID metric and latency metric are not compatible
- If data is replicated on k peers then peers with similar Node-ID might be missed
- Here, a heuristic is used
- Experiments validate this approach

Parameter $b=4$,
 $l=16$, $M=32$

In this experiment
 the hop distance
 grows
 logarithmically with
 the number of
 nodes

The analysis
 predicts $O(\log n)$

Fits well

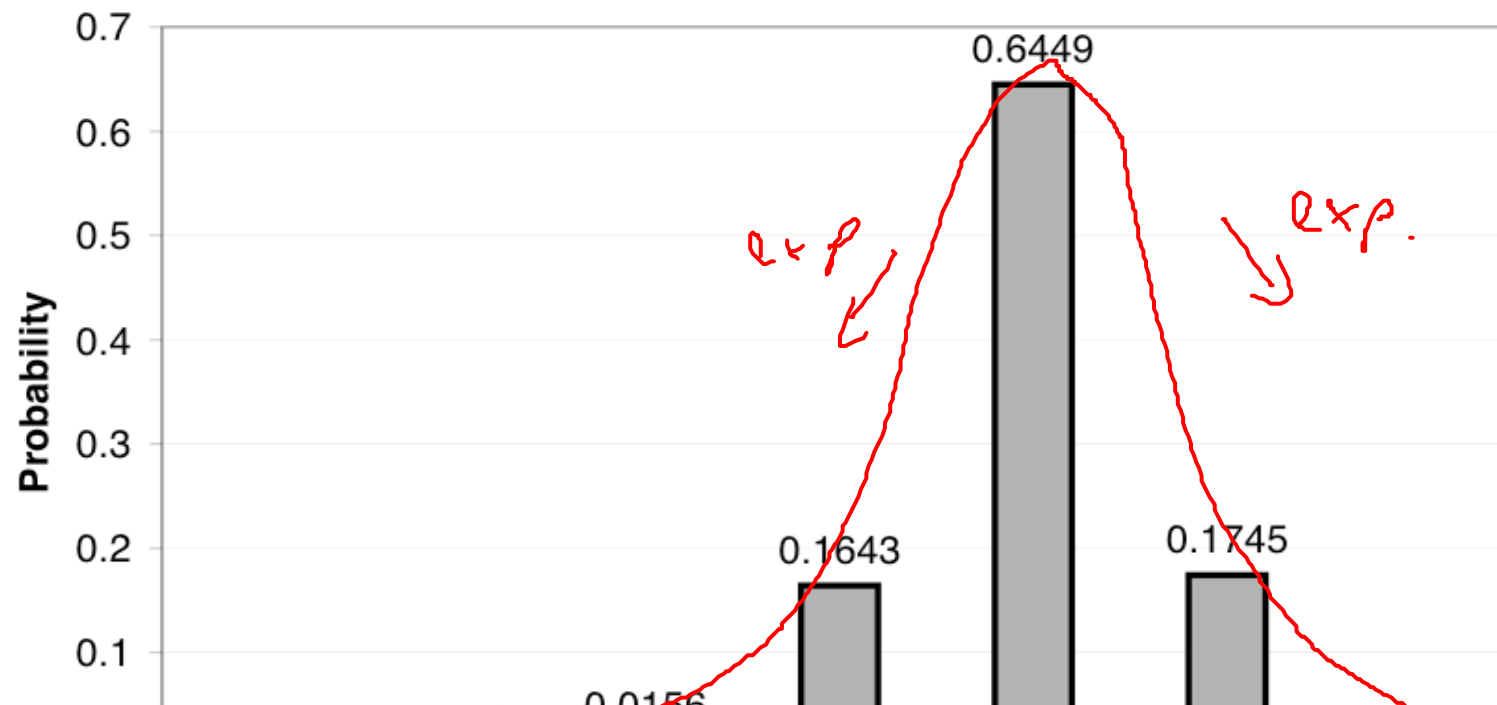


iburg Distribution of Hops

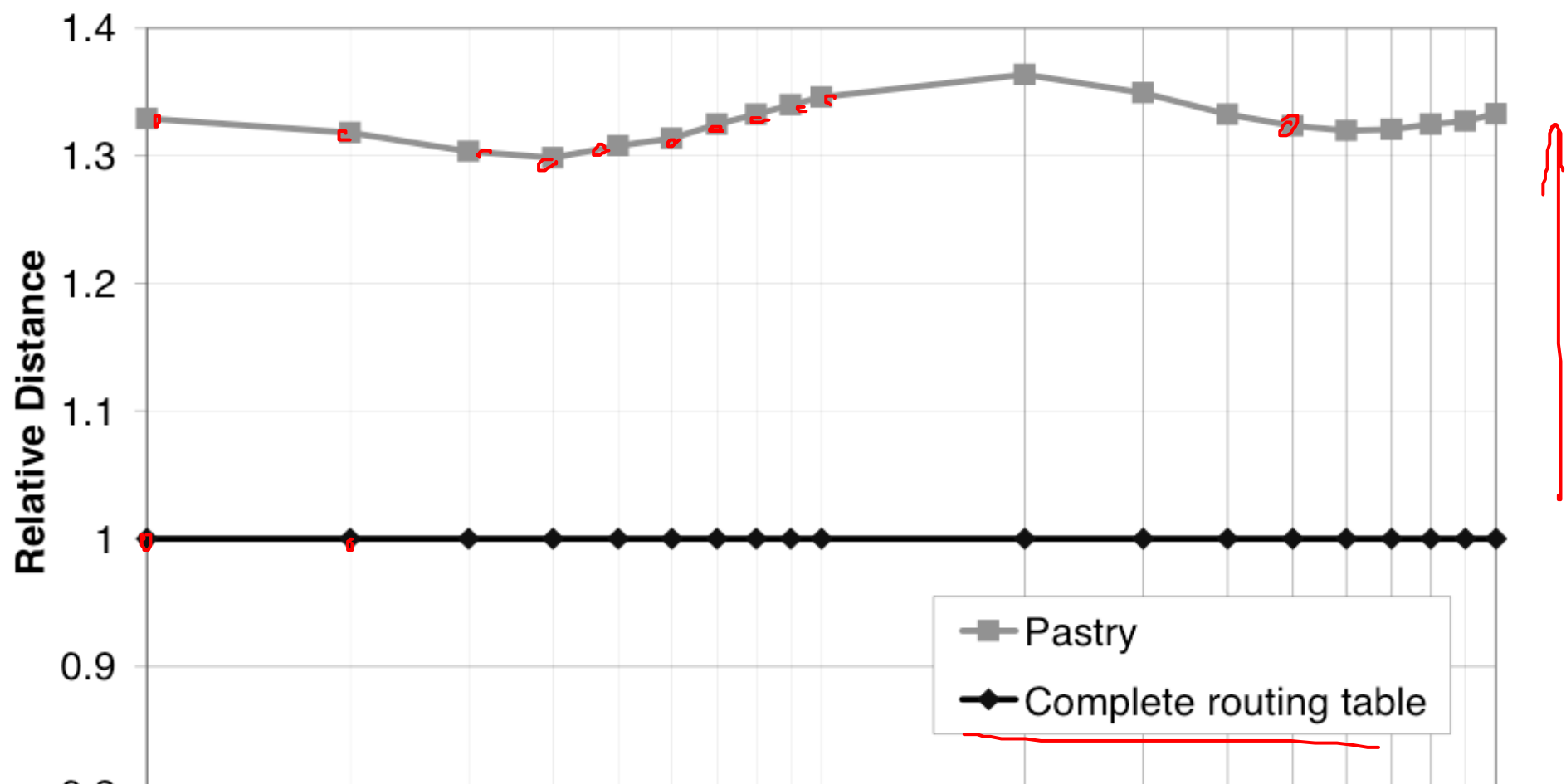
Parameter $b=4$, $l=16$, $M=32$, $n = \underline{100,000}$

Result

- deviation from the expected hop distance is extremely small
- Analysis predicts difference with extremely small probability
- fits well



- Parameter $b=4$, $l=16$, $M=3$
- Compared to the shortest path astonishingly small
 - seems to be constant



Peer-to-Peer Networks

05 Pastry

Christian Schindelhauer

Technical Faculty

Computer-Networks and Telematics

University of Freiburg