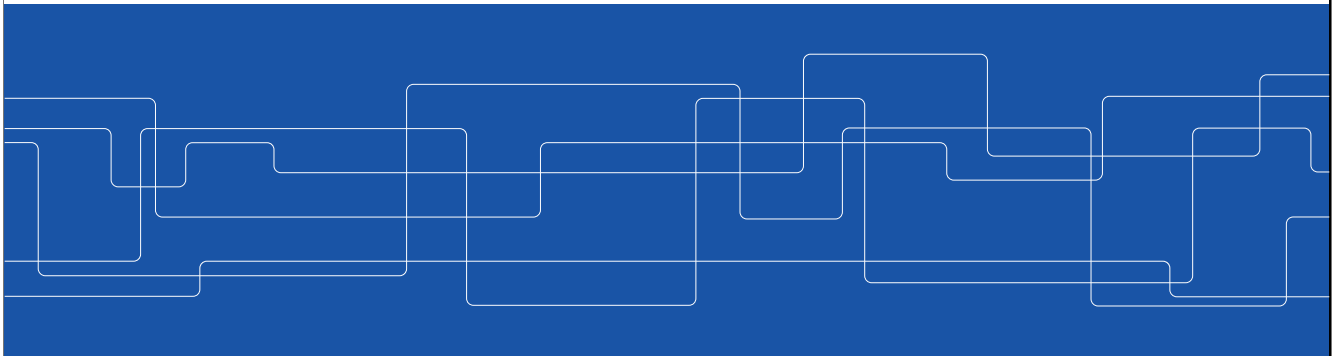# Distributed Hash Tables

Johan Montelius and Vladimir Vlassov

# Distributed Hash Tables

- Large-scale databases (key-value stores)
  - hundreds of servers
- High churn rate
  - servers will come and go
- Benefits
  - fault tolerant
  - high performance
  - self administrating

# A key-value store

Associative array to store *key-value pairs*, a data structure known as a *hash table* (array of buckets) that maps keys to values.

Operations:
**put (key, object)** – store a given object with a given key
**object: = get (key)** – read an object given key.

Design issues:
- Identify: how to uniquely identify an object
- Store: how to distribute objects among servers
- Route: how to find an object

ID2201 DISTRIBUTED SYSTEMS / DISTRIBUTED HASH TABLES

# Unique identifiers

We need *unique identifiers* to identify objects, i.e., to find a bucket to get/put an object with a given key

$$identifier = f(key, size\_of\_hash\_table)$$
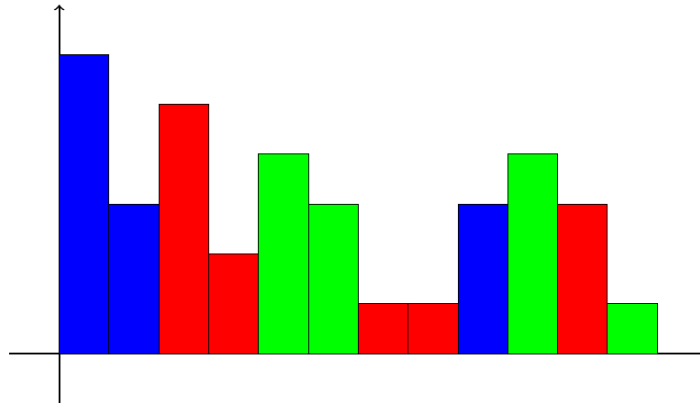
How to select identifiers:
- use a key (a name)
- a cryptographic hash of the key
- a cryptographic hash of the object

*Why hash?*

# Key distribution – direct map

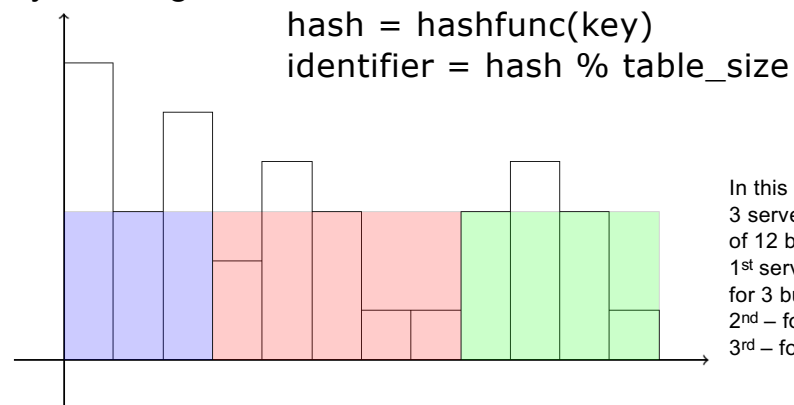A direct map of keys to identifiers (buckets) might give a non-uniform (uneven) distribution of keys among buckets.



ID2201 DISTRIBUTED SYSTEMS / DISTRIBUTED HASH TABLES

3 servers (RGB) to store buckets

# Key distribution – hashing keys

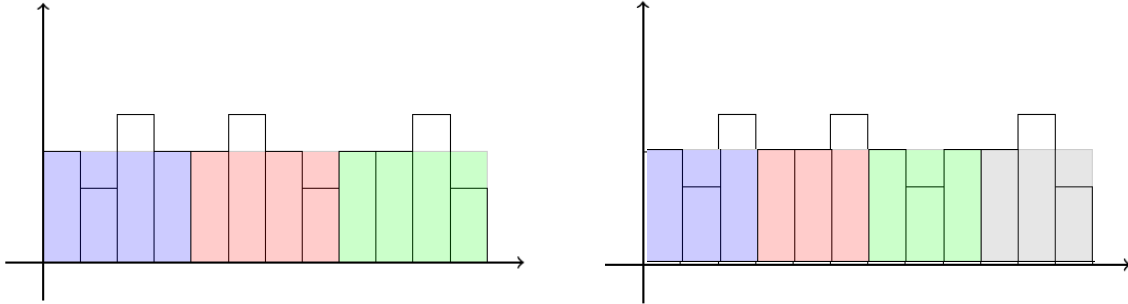*A cryptographic hash function* gives a uniform (even) distribution of the keys among buckets

hash = hashfunc(key)
identifier = hash % table_size

In this example:
3 servers to store a DHT of 12 buckets.
1st server is responsible for 3 buckets;
2nd – for 5 buckets,
3rd – for 4 buckets.

3 servers to store 12 buckets: 1st server store 3 buckets; 2nd – 5, etc.

# Add a server

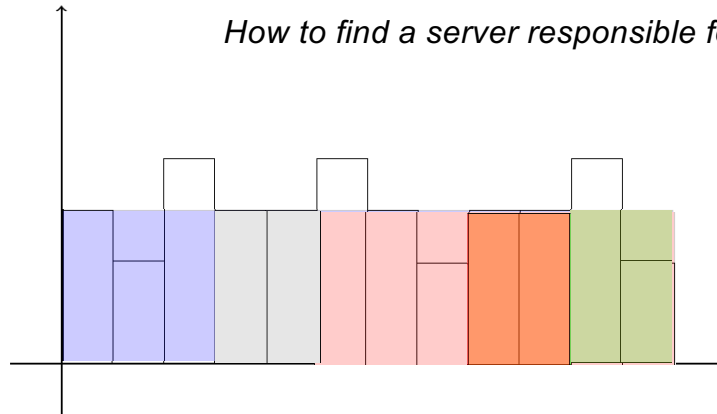*At three o'clock in the morning, do:*

3 servers + 1 = 4 servers; need to redistribute buckets, i.e. move data.

# Random distribution
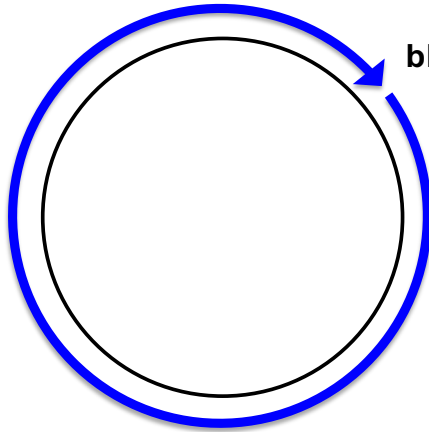
Random distribution of key ranges among servers

*How to find a server responsible for a given key?*
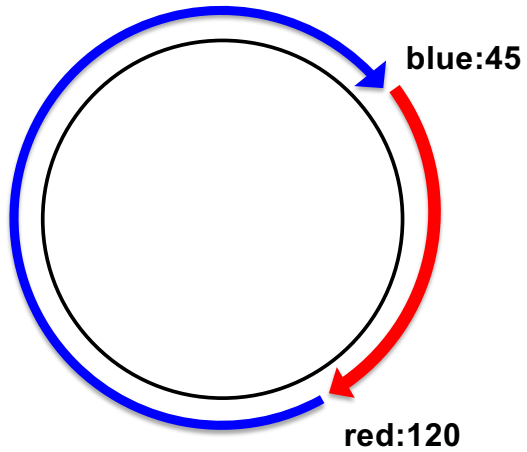
# Circular domain

**blue:45**

- ID domain: 0,1,2,..., size-1
- clockwise step along the ring
  $$i = (i + 1)\% \text{ size}$$
- *responsibility*: from your predecessor to your number
- when inserted: take over responsibility

ID domain of keys (and servers)
Blue is its own predecessor

# Circular domain

**blue:45**

**red:120**
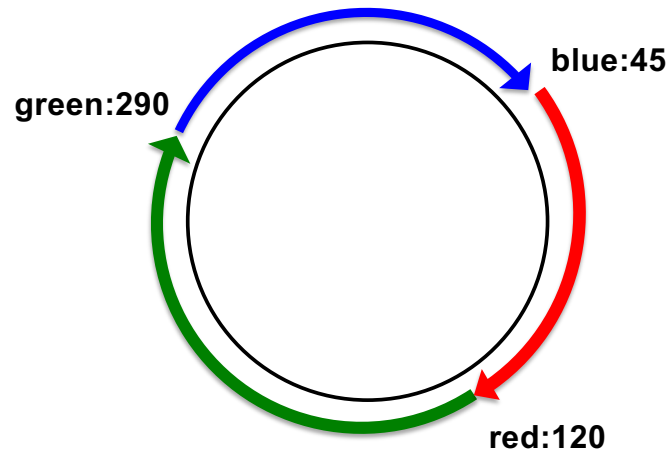
- *responsibility*: from your predecessor to your number
- when inserted: take over responsibility
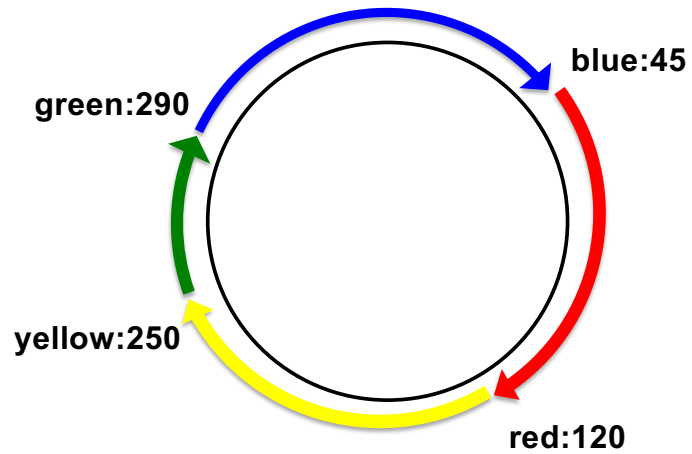- e.g., red:120 is responsible for keys in the range (45, 120]

# Circular domain



green:290

blue:45

red:120

- *responsibility*: from your predecessor to your number
- when inserted: take over responsibility

# Circular domain

blue:45

green:290

yellow:250

red:120
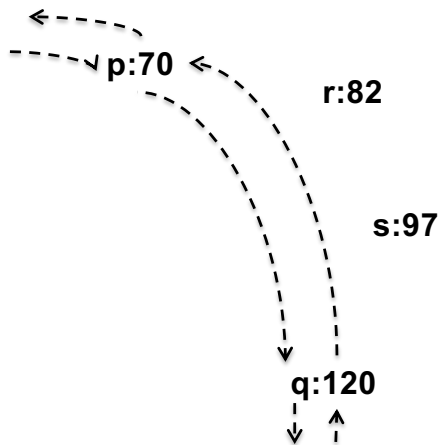
- *responsibility*: from your predecessor to your number
- when inserted: take over responsibility
- talk to the node in front of you

ID2201 DISTRIBUTED SYSTEMS / DISTRIBUTED HASH TABLES

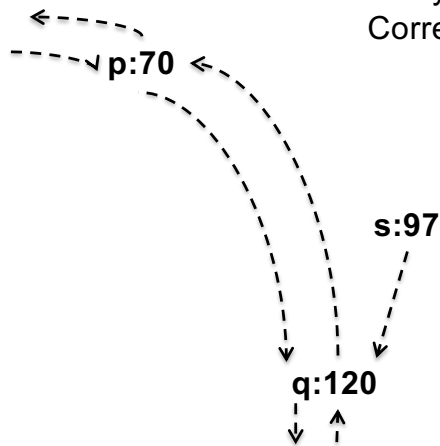# Double linked circle

p:70

r:82

s:97

q:120

- predecessor
- successor
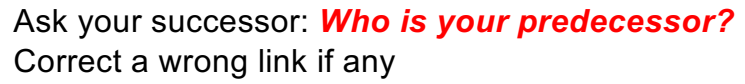- how do we insert a new node
- concurrently

# Stabilization

Ask your successor: ***Who is your predecessor?***
Correct a wrong link if any

s:97:   - Who is your predecessor?
q:120: - It's p at 70(p).
s:97:   - Why don't you point to me!

**p:70**

**s:97**

**q:120**

# Stabilization

Ask your successor: ***Who is your predecessor?***
Correct a wrong link if any

**p:70**

**s:97**

**q:120**

s:97:   - Who is your predecessor?
q:120: - It's p at 70(p).
s:97:   - Why don't you point to me!
p:70:   - Who is your predecessor?
q:120: - It's s at 97(s).
p:70:   - Hmmm, that's a better successor.

ID2201 DISTRIBUTED SYSTEMS / DISTRIBUTED HASH TABLES

# Stabilization

**Ask your successor: *Who is your predecessor?***
Correct a wrong link if any

**p:70**

**s:97**

**q:120**

*Let's play a game!*

s:97:  - Who is your predecessor?
q:120: - It's p at 70(p).
s:97:  - Why don't you point to me!
p:70:  - Who is your predecessor?
q:120: - It's s at 97(s).
p:70:  - Hmmm, that's a better successor.
p:70:  - Who is your predecessor?
s:97:  - I don't have one.
p:70:  - Why don't you point to me!

Rule of thumb: new node should find its successor (the node responsible for its id) and contact it to become its predecessor, i.e., force the successor to correct its predecessor;
If a node receives an answer from its successor with an id different from its own, it should correct its successor and become its new successor's predecessor.
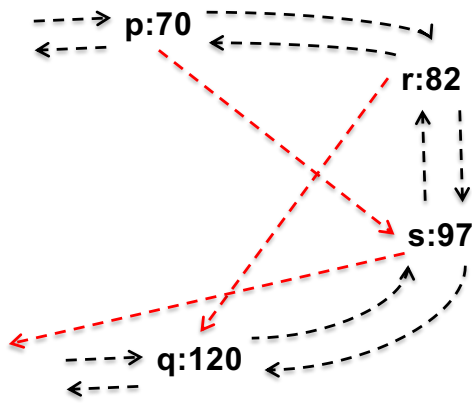
# Stabilization

Stabilization is run periodically: allow nodes to be inserted concurrently.

The inserted node will take over responsibility for part of a segment.
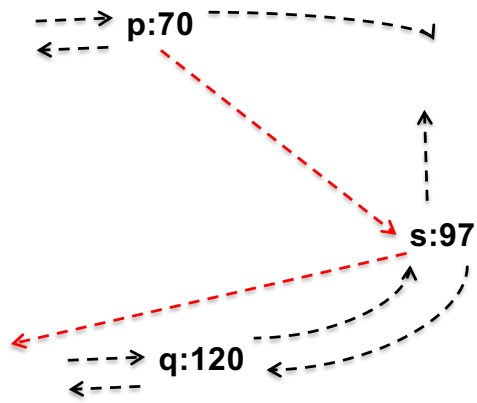
# Crashing nodes

p:70

r:82

s:97

q:120

- monitor neighbors
- safety pointer

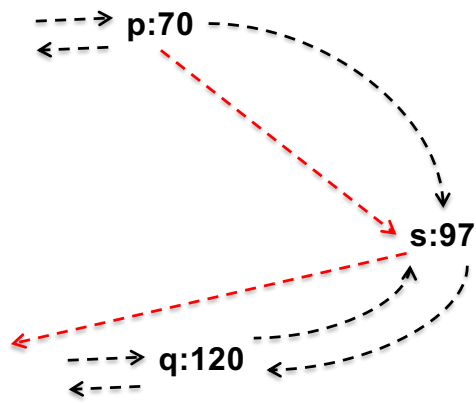Safety pointers is successor pointers

# Crashing nodes

- monitor neighbors
- safety pointer
- detect crash

# Crashing nodes
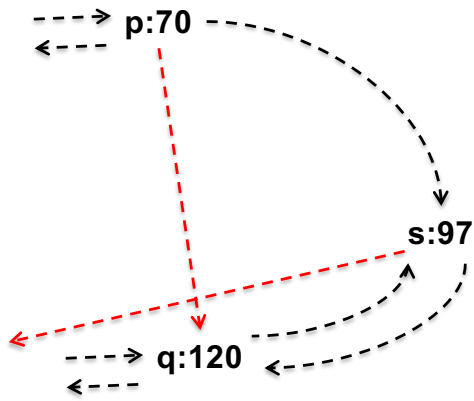
p:70

s:97

q:120

- monitor neighbors
- safety pointer
- detect crash
- update forward pointer

# Crashing nodes

p:70

s:97

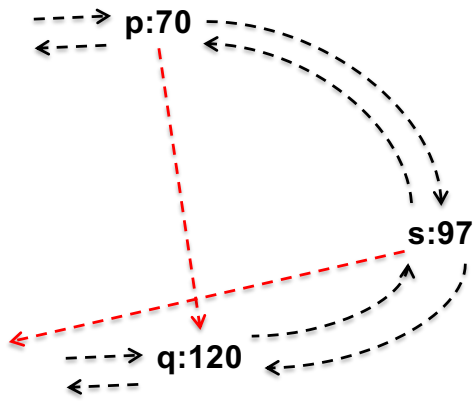q:120

- monitor neighbors
- safety pointer
- detect crash
- update forward pointer
- update safety pointer

# Crashing nodes

p:70

s:97

q:120

- monitor neighbors
- safety pointer
- detect crash
- update forward pointer
- update safety pointer
- **stabilize**

# Russian roulette

How many safety pointers do we need?

The more bullets you have – the bigger chance of being killed and vice versa.

# Replication

Where should we store a replica of our data?

Successor replication. In this example, 120 has a replica of data stored at 97.
Symmetric replication

# Routing overlay

- The problem of finding an object in our distributed table:
  - Nodes can join and crash
  - A trade-off between routing overhead and updating overhead

*In the worst case, we can always forward a request to our successor.*

# Leaf set

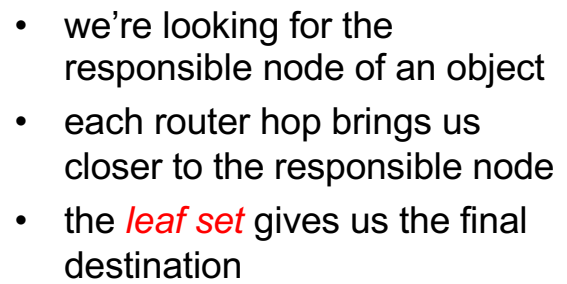Assume that each node holds a leaf set of its closest (±$l$ ) neighbors (a.k.a. a finger table).

We can jump $l$ nodes in each routing step, but we still have O(n) complexity.

The leaf set is updated in *O*($l$).

*The leaf set could be as small as only the immediate neighbors but is often chosen to be a handful.*

# Improvement

350  10
337      20
310           40
             50   get (222)
280
267               70
250               85
238               112
224               120
210               130
195  170  158  145

- we're looking for the responsible node of an object
- each router hop brings us closer to the responsible node
- the *leaf set* gives us the final destination

# Pastry

In a routing table, each row represents one level of routing.

- 32 rows
- 16 entries per row
- any node found in 32 hops
- maximal number of nodes $16^{32}$ or $2^{128}$ (more than enough)
- Search is $O(log(n))$, where $n$ is the number of nodes

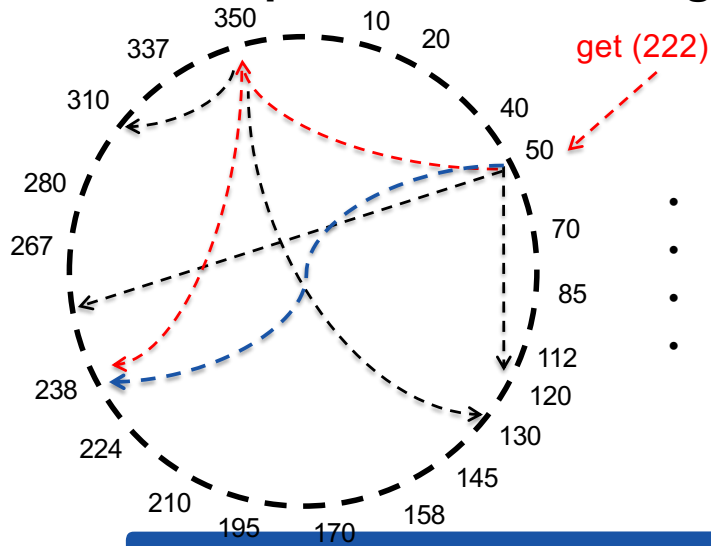Binary search – two regions; k-nary search – k regions. Works only in ordered sets.
Pastry - like K-nary search; k =16
32 rows; start with last row 16 large regions send to one of them; upon receive read 15th row smaller regions, etc.
Each row (level) has 16 regions of diminishing size. Each further level narrows the search region.
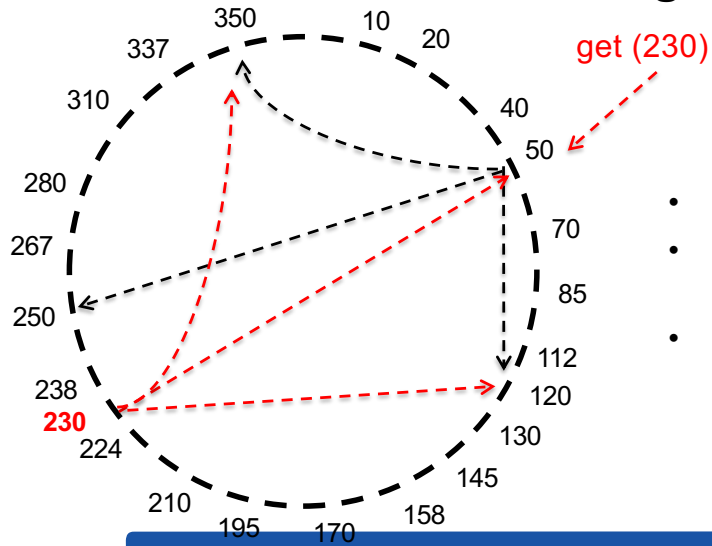
# The price of fast routing



get (222)

- be lazy
- detect failed nodes when used
- route in the alternative direction
- ask neighbors of alternative node

In this example: a node that the "middle" finger points to has crushed - ask other node (350) about an alternative node. (238) – correction on use.

# Network aware routing

350  10  20
337
310
280
267
250
238
**230**
224
210
195  170  158
145
130
120
112
85
70
50
40

get (230)

- when inserting a new node
- attach to the network-wise closest node
- adopt the routing entries on the way down

# Overlay networks

Structured
- a well-defined structure
- takes time to add or delete nodes
- takes time to add objects
- easy to find objects

Unstructured
- a random structure
- easy to add or delete nodes
- easy to add objects
- takes time to find objects

**ID2201 DISTRIBUTED SYSTEMS / DISTRIBUTED HASH TABLES**

# DHT usage

Large scale key-value store.

- fault tolerant system in the high churn rate environment
- high availability, low maintenance

## The Pirate Bay

- replaces the tracker with a DHT
- clients connect as part of the DHT
- DHT keeps track of peers that share content

# Riak

**riak**

- large scale key-value store
- inspired by Amazon Dynamo
- implemented in Erlang

# Summary DHT

- Why hashing?
- Distribute storage in the ring
- Replication
- Routing