

RestFS

Fabrizio Manfredi Furuholmen
Federico Mosca

Beolink.org



RestFS

- Introduction
 - Goals
 - Principals
- RestFS
 - Architecture
 - Internals
 - Sub project
- Conclusion
 - Developments

Zetabyte

1000^7 bytes

10^{21} bytes

1,000,000,000,000,000,000 bytes

All of the data on Earth today 150GB of data per person

2% of the data on Earth in 2020

Unsolved problem

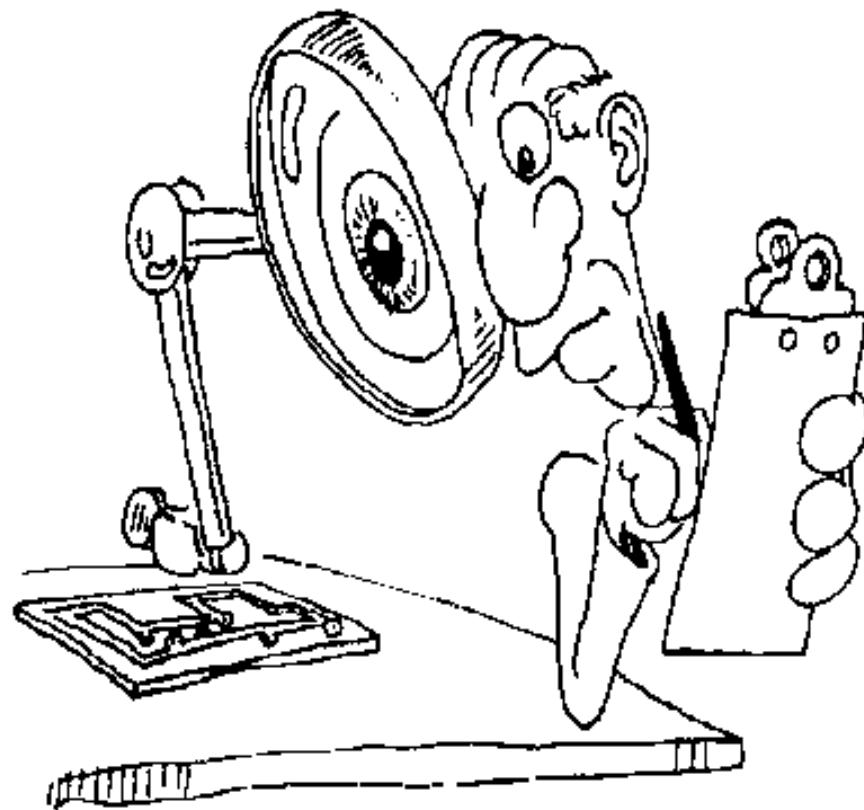
Beolink.org





**RestFS
is
High scalable, high available
network object storage**

Create a framework for
testing a new
technologies and
paradigm



**“Moving Computation
is
Cheaper than Moving Data”**

“There is always a failure waiting around the corner”

*Werner Vogel

“Decompose into small loosely coupled, stateless building blocks”

*' Leaving a Legacy System Revisited' Chad Fowler



Objects

- Separation btw data and metadata
- Each element is marked with a revision
- Each element is marked with an hash.



Cache

- Client side
- Callback/ Notify
- Persistent



Transmission

- Parallel operation
- Http like protocol
- Compression
- Transfer by difference



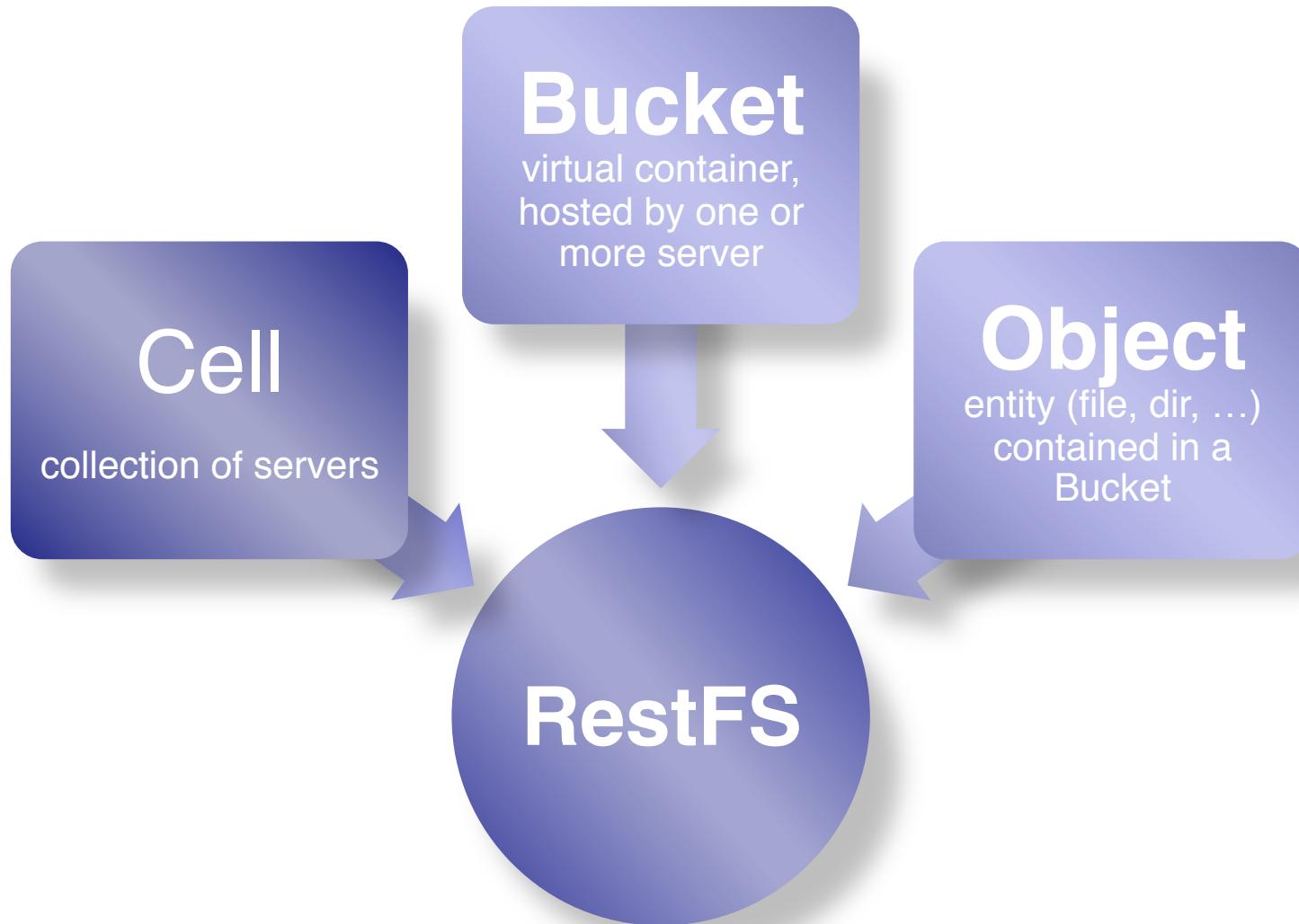
Distribution

- Resource discovery by DNS
- Data spread on multi node cluster
- Decentralize
- Independents cluster
- Data Replication



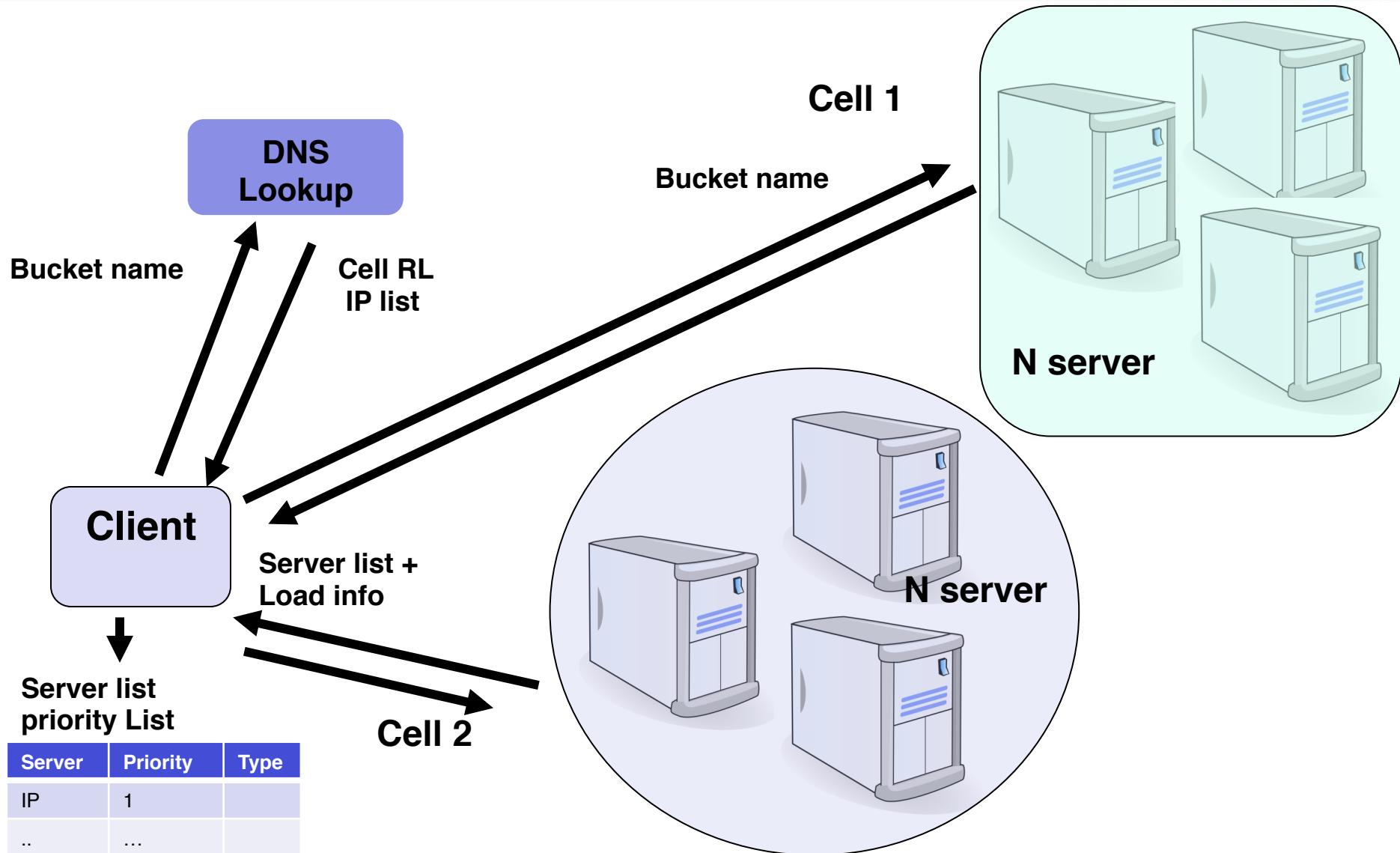
Security

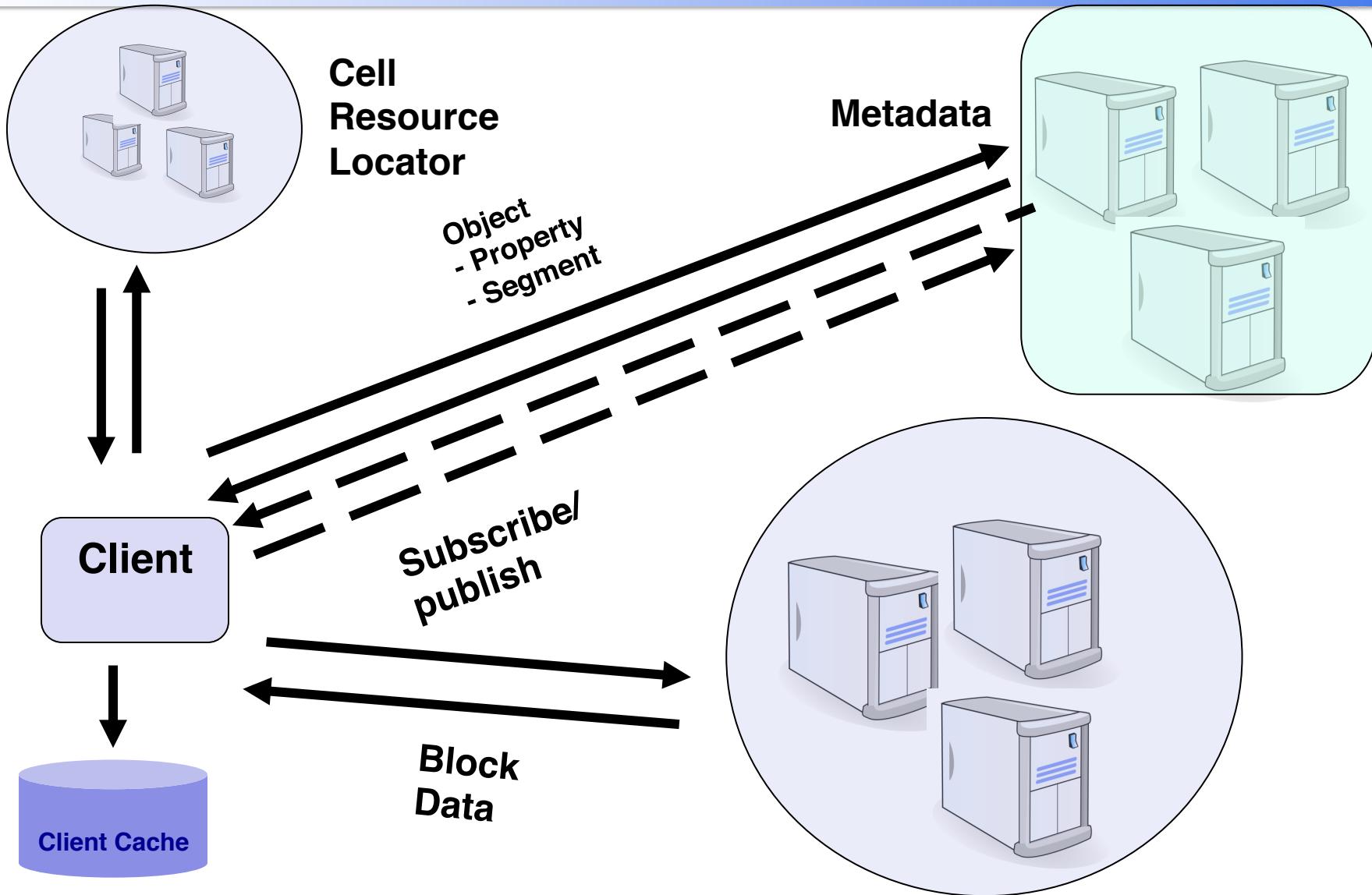
- Secure connection
- Encryption client side,
- Extend ACL
- Delegation/ Federation
- Admin Delegation

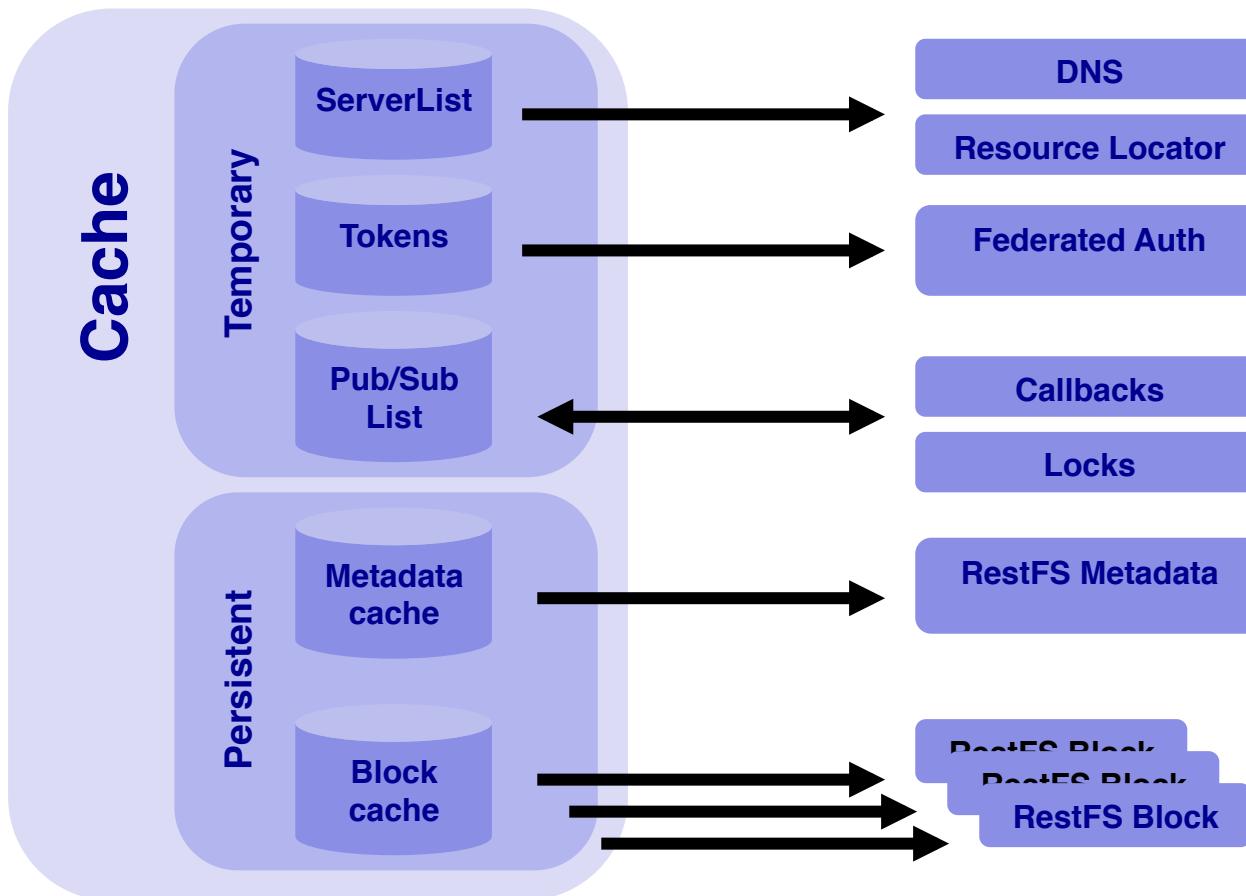


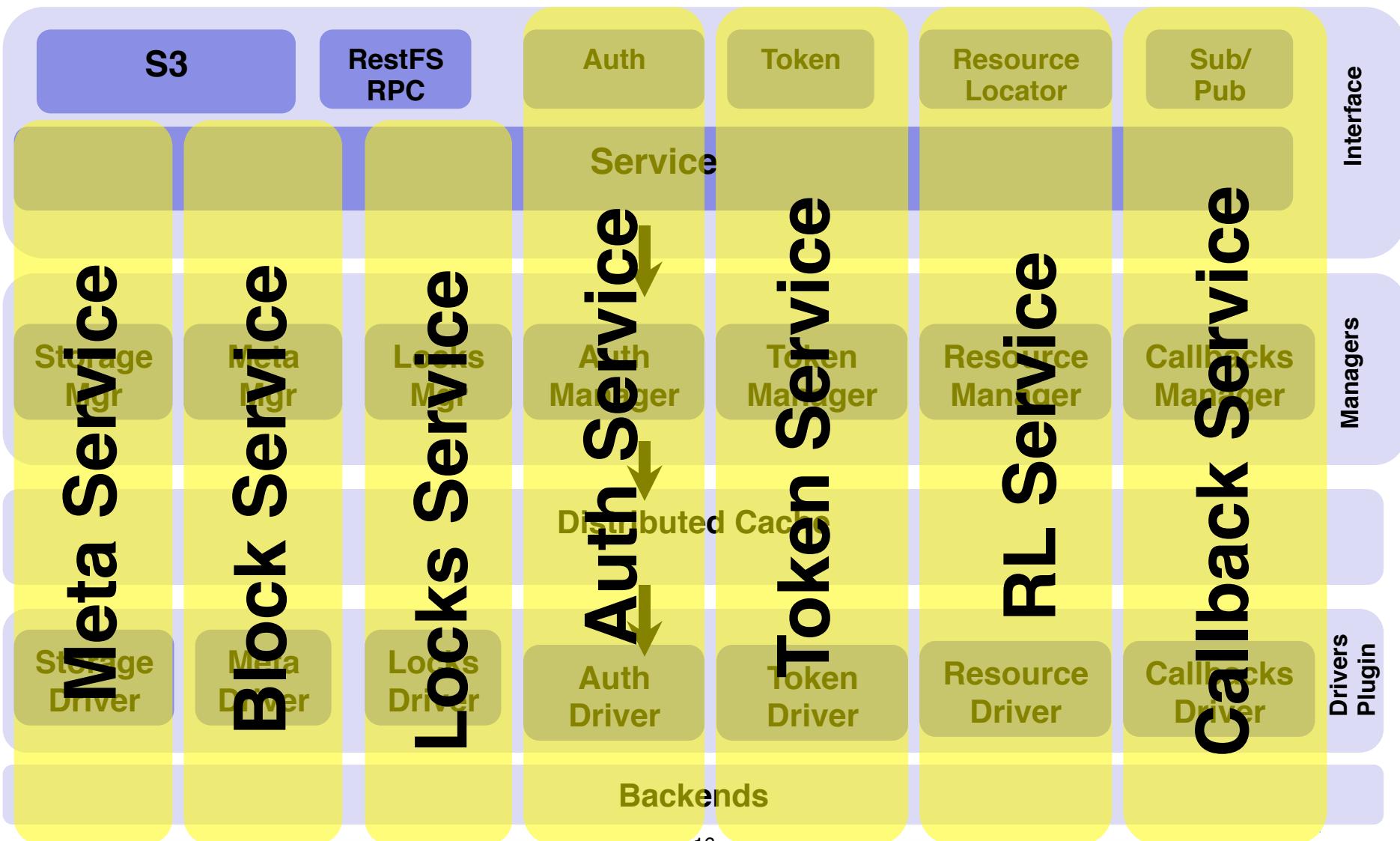
Bucket Discovery

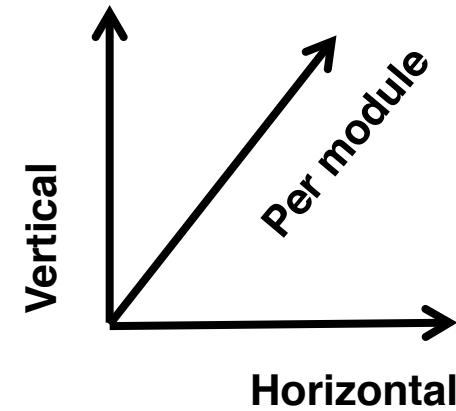
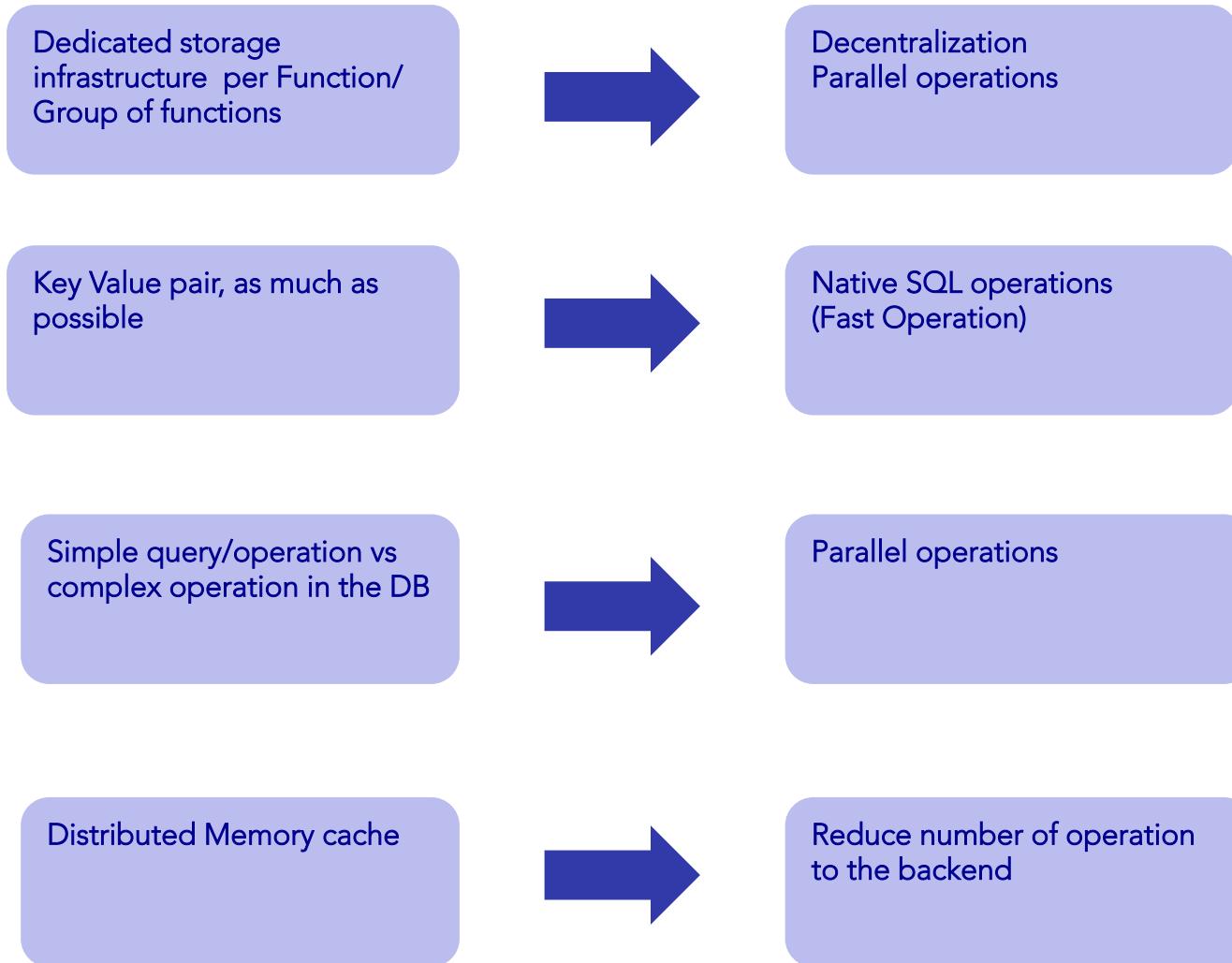
Beolink.org





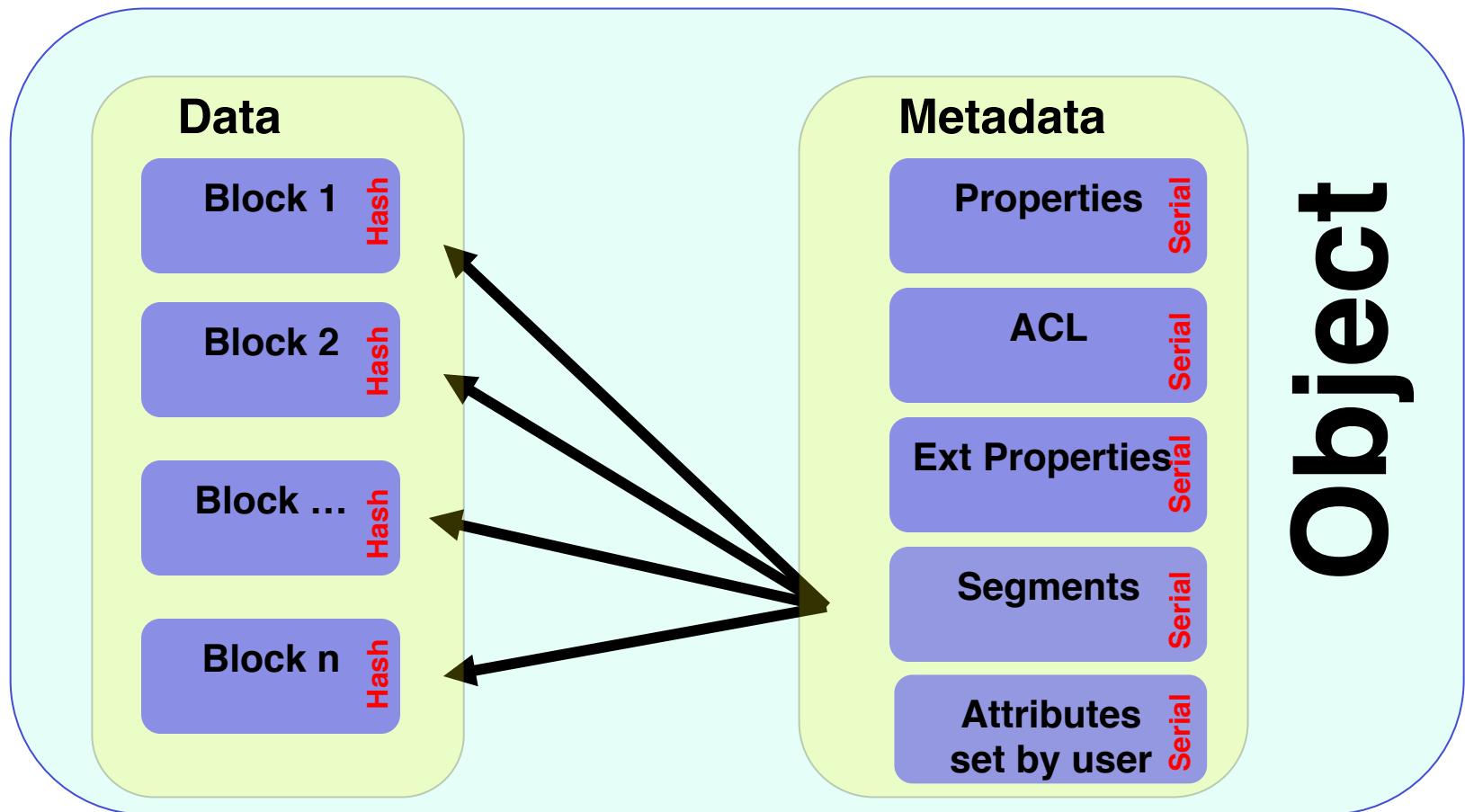


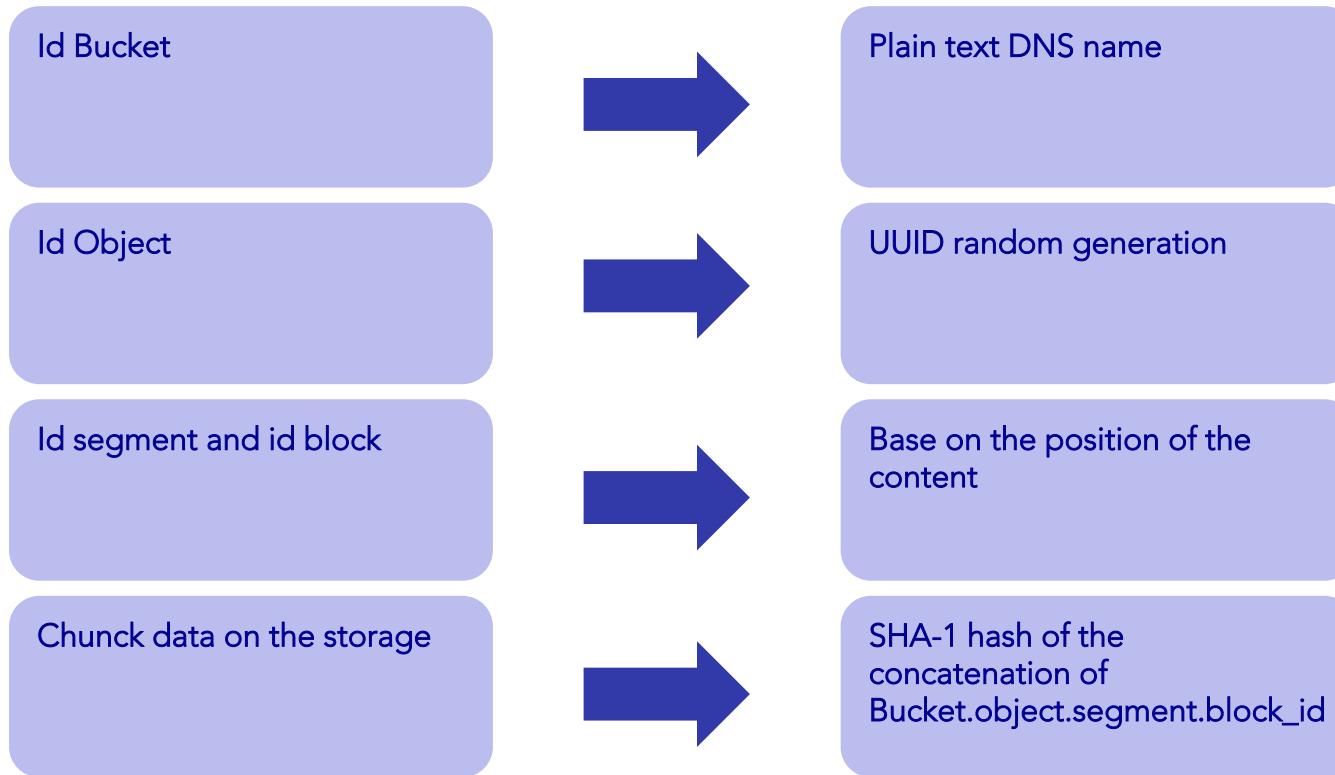


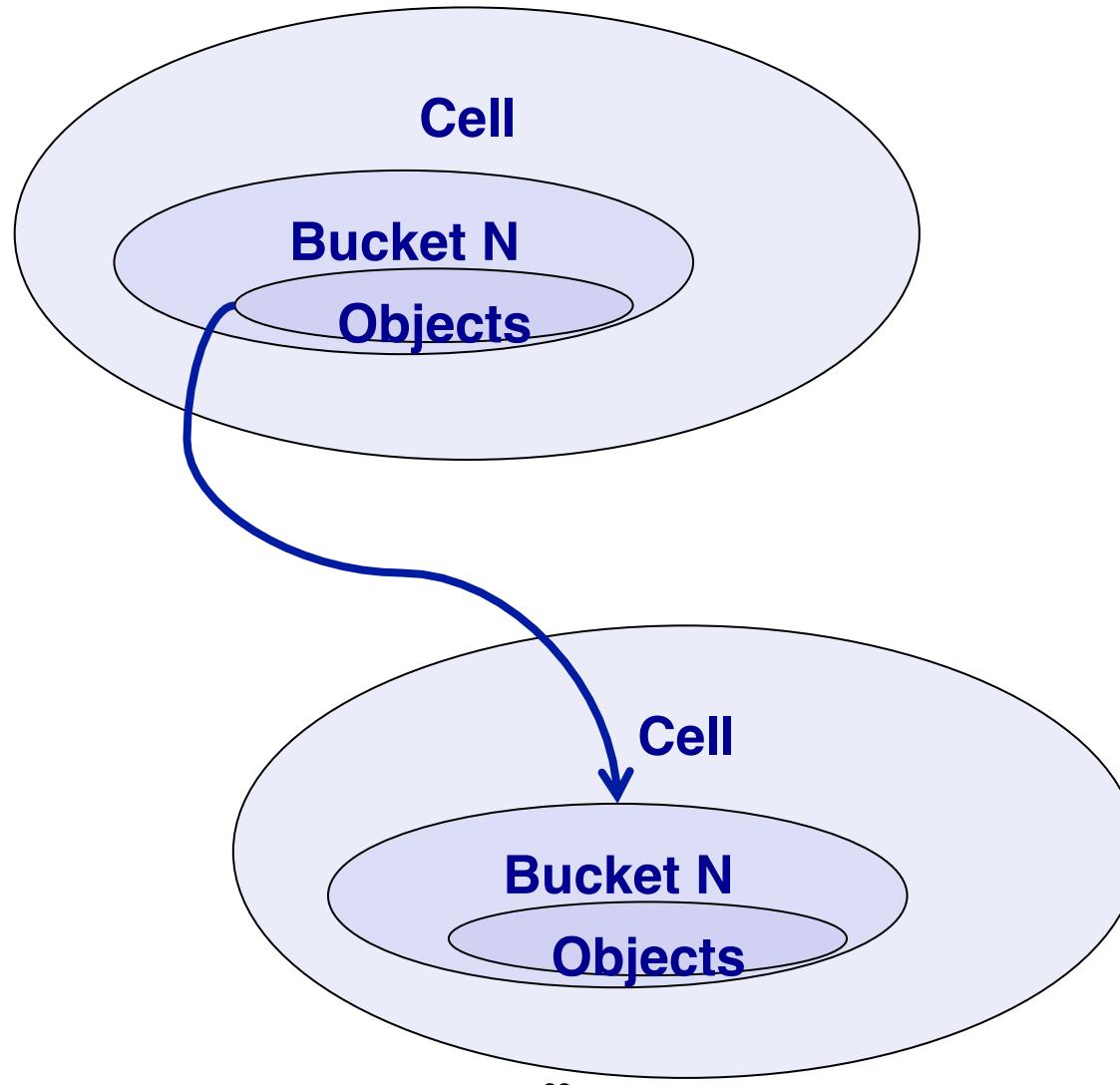


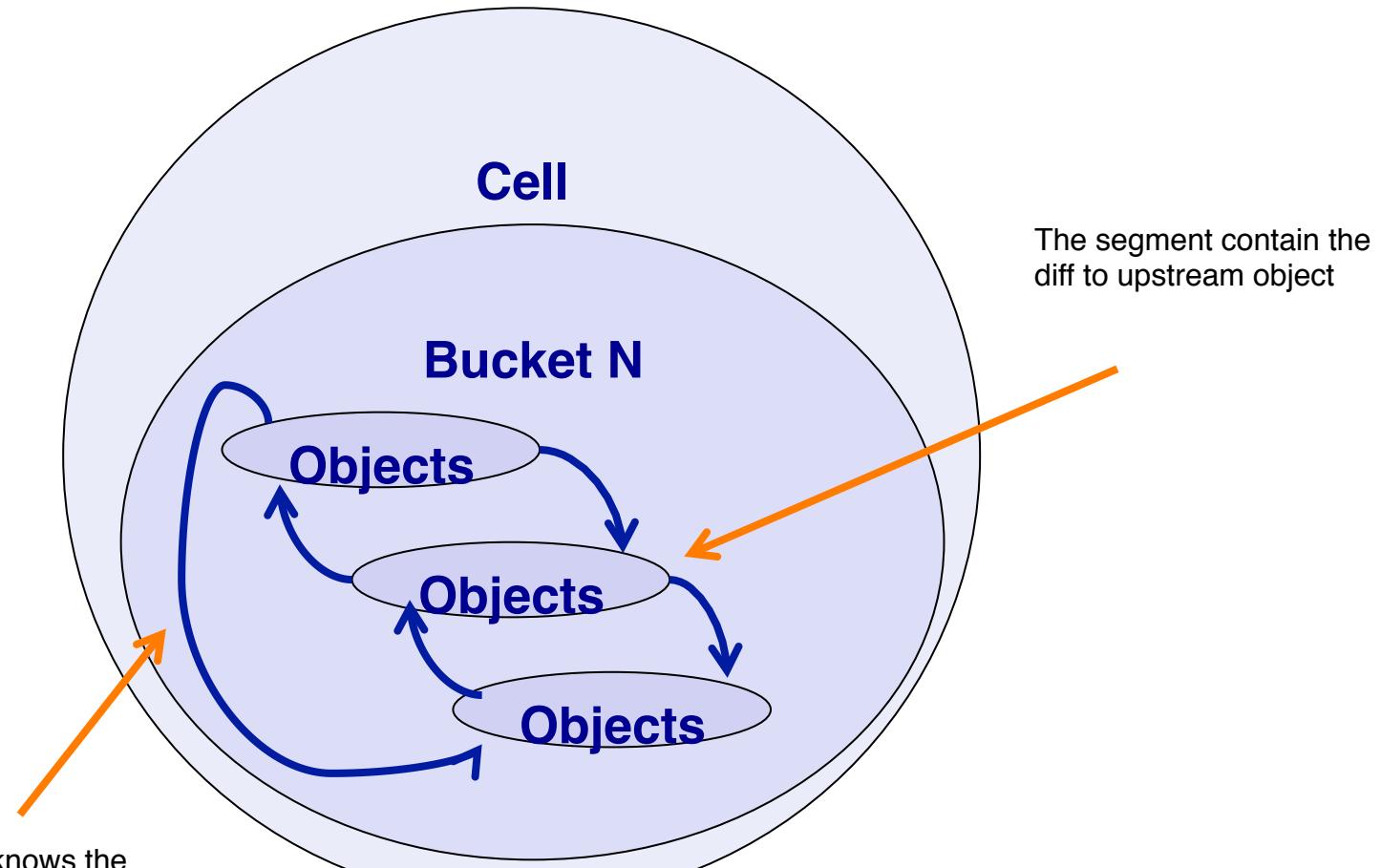
Bucket

Objects







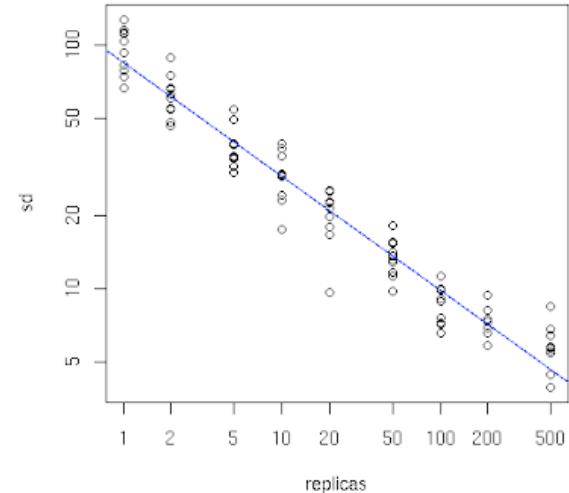
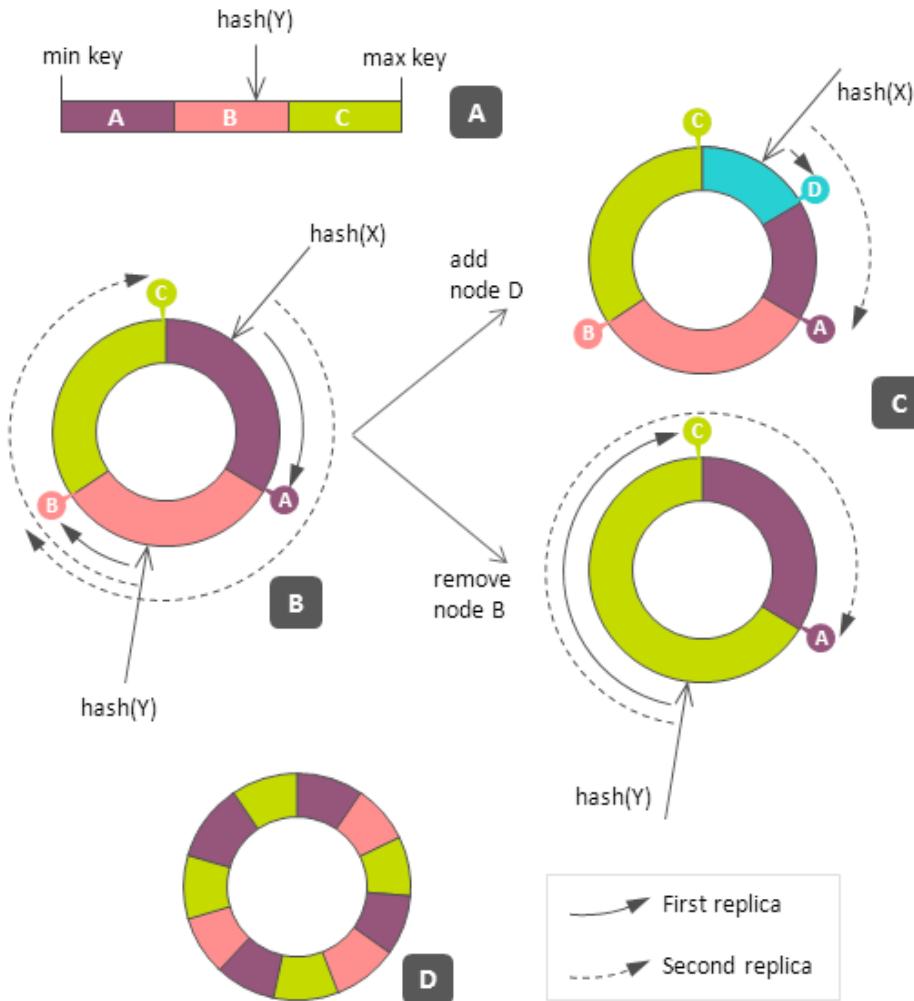


Each object knows the previous and the next.
The current object knows the previous and the last

Block Storage

Backend: Consistent Hashing

Beolink.org



Number of key to move for add/
remove a node :

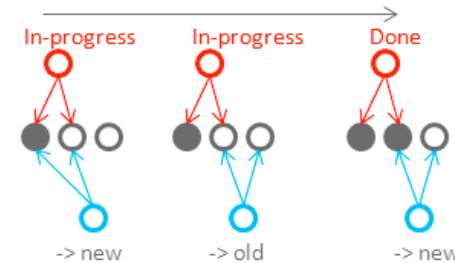
Keys/Node= keys to relocate

Blocks are collected in shards

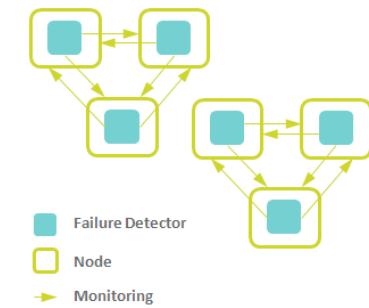
3 Copies

Configurable read and write consistent level and security:

- 2W1R
- 2W2R
- 1W1R
- ...

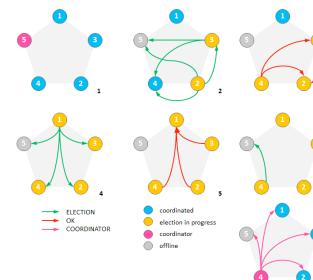


Monitor of neighbored small cluster of 3 nodes (GOSSIP)

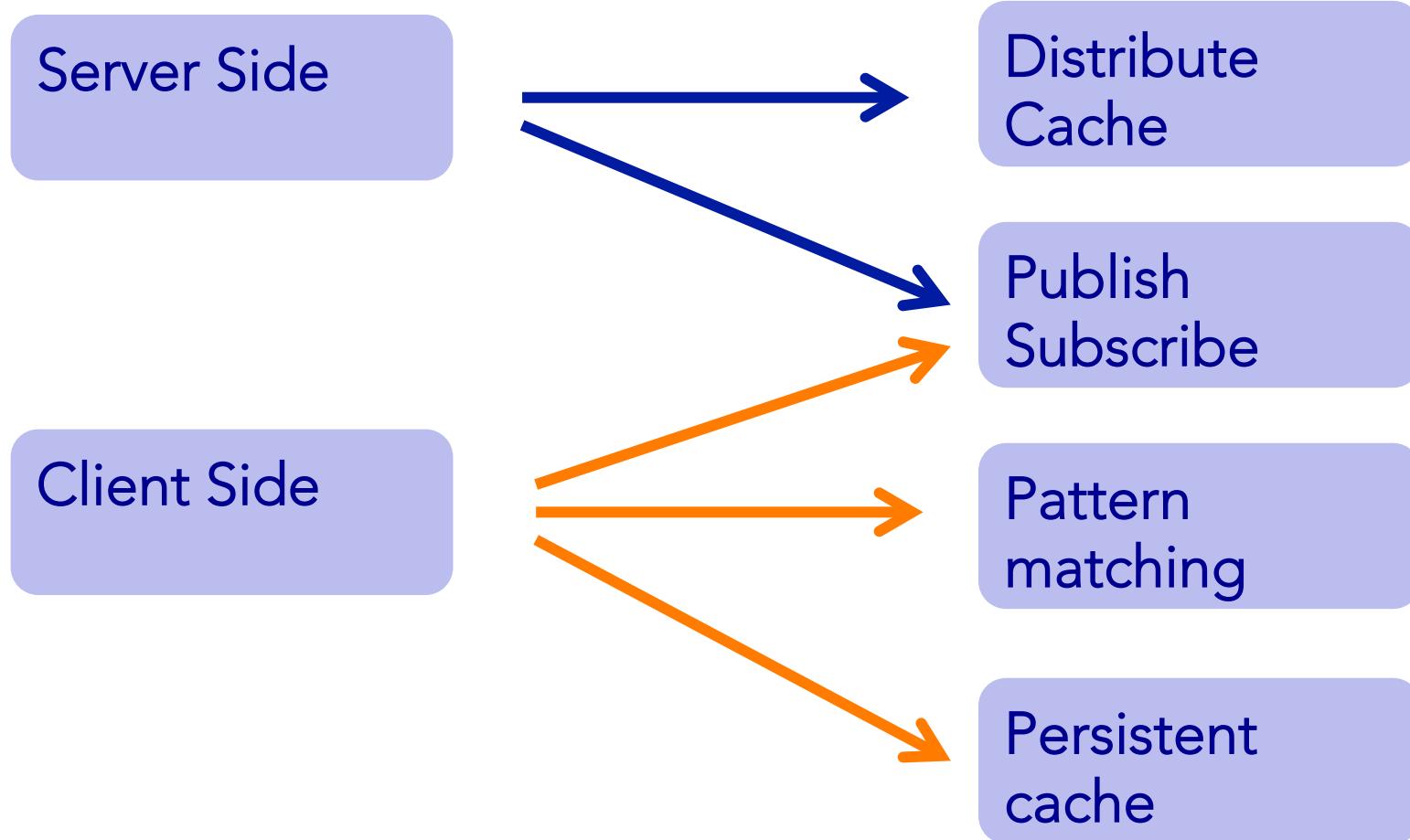


Mini cluster election

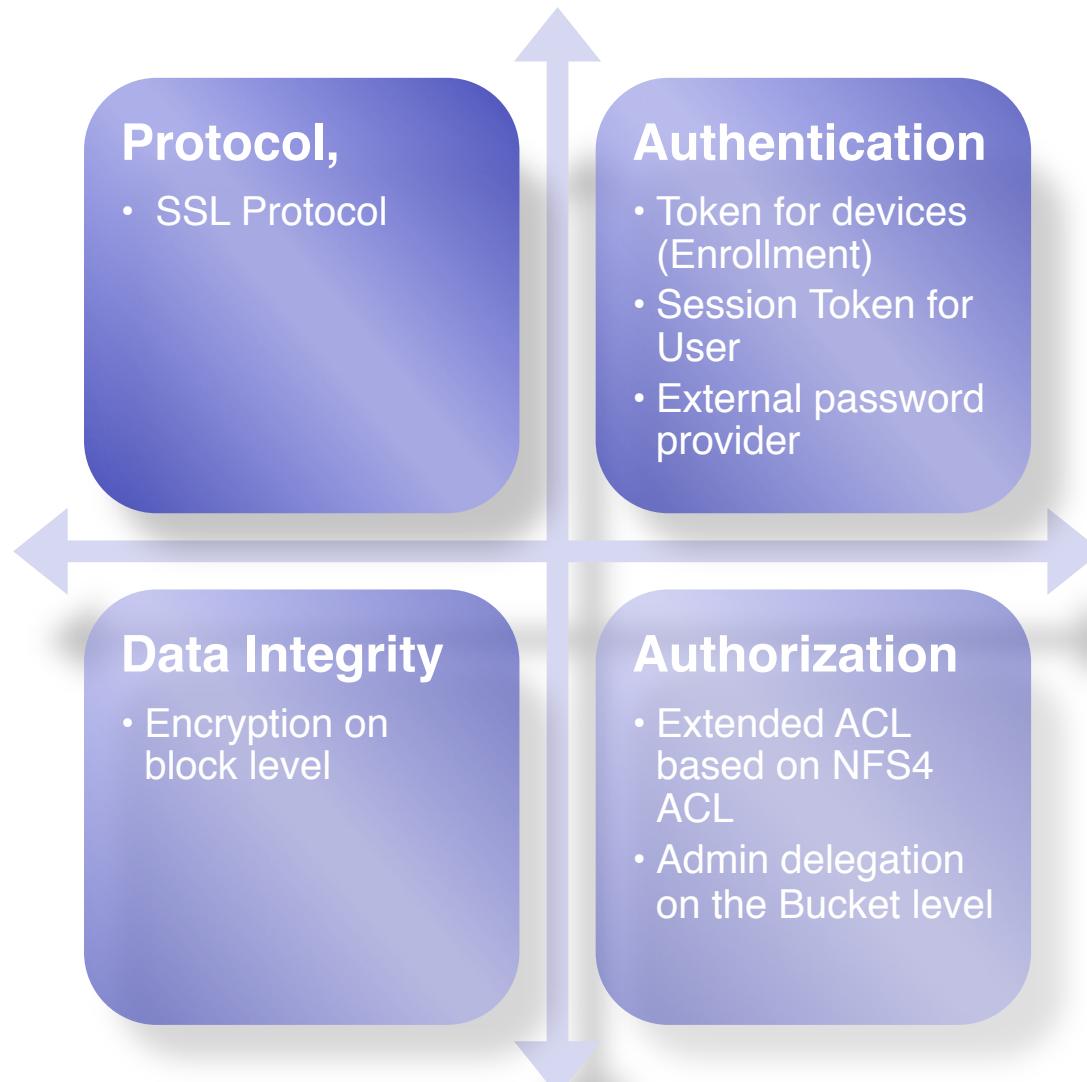
key space reclaim for replica coordination,
leave join cluster



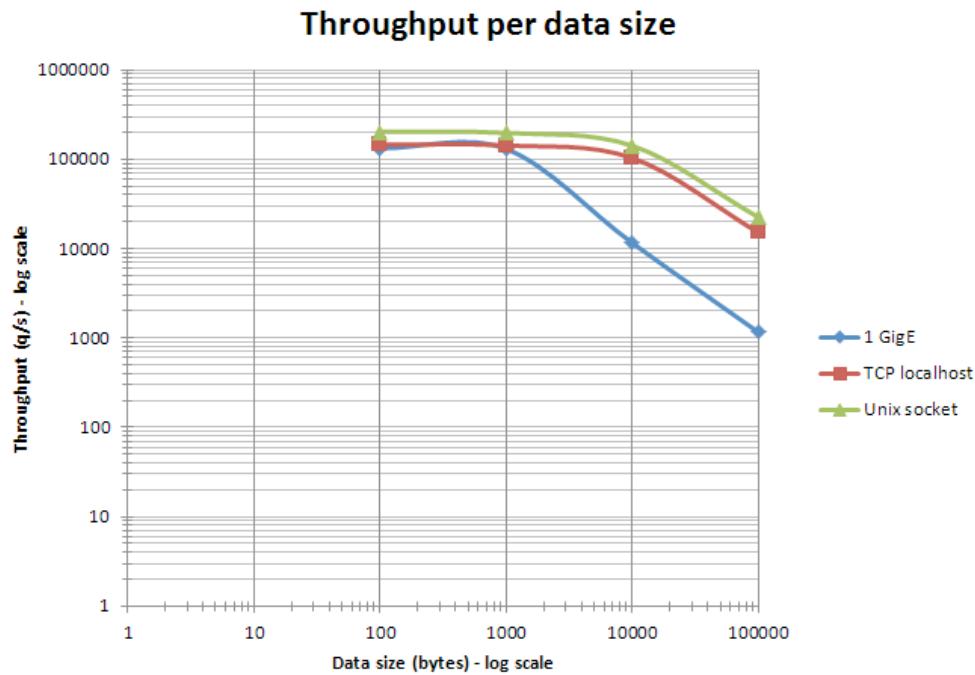
Cache



Security



NOSQL DB

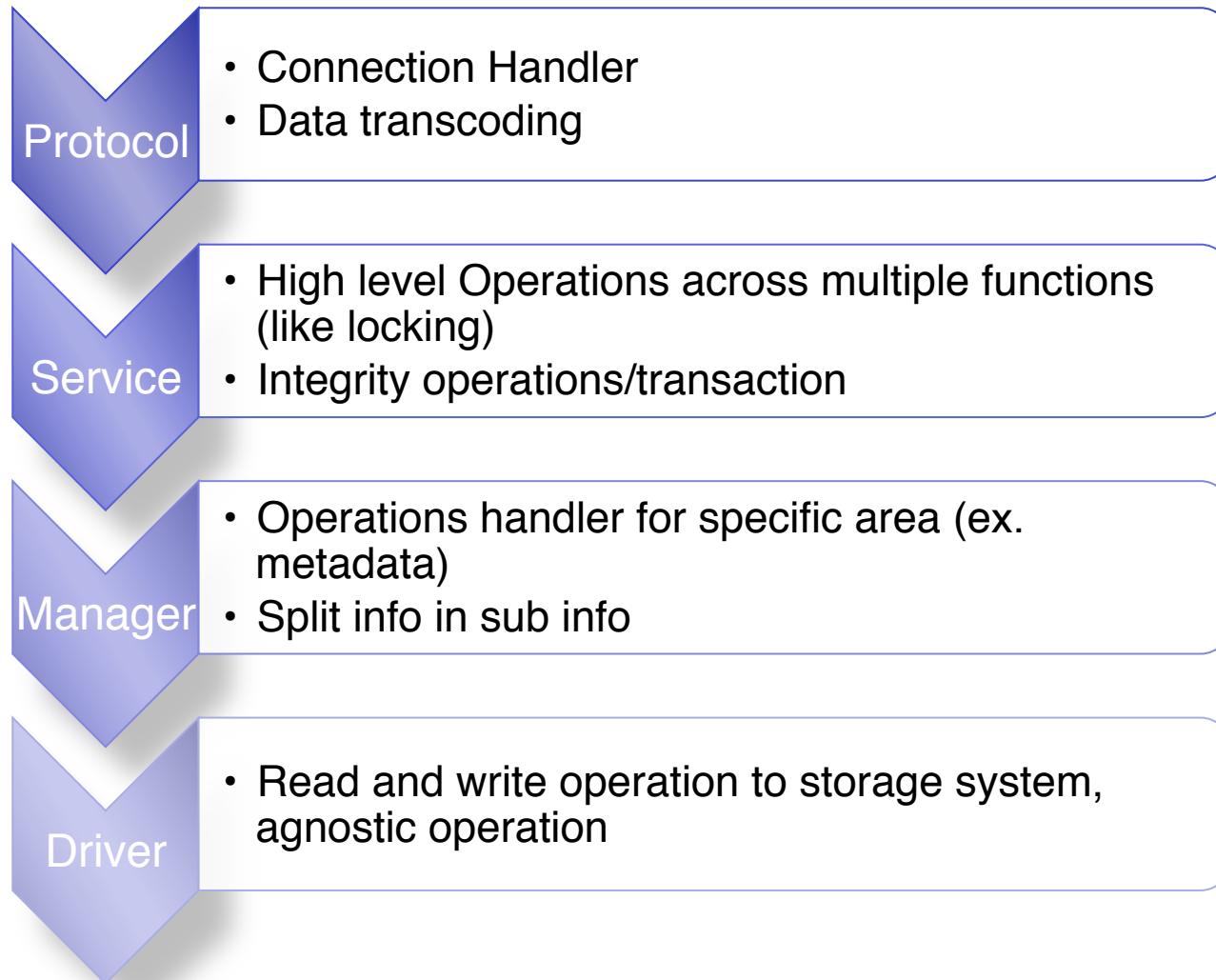


```
$ ./redis-benchmark -r 1000000 -n 2000000 -t get,set,lpush,lpop -P 16 -q
```

SET: 552028.75 requests per second
GET: 707463.75 requests per second
LPUSH: 767459.75 requests per second
LPOP: 770119.38 requests per second

Code

Interface, dynamic load



What we are using

Beolink.org



Module	Software
Storage	Filesystem, DHT (kademlia, Pastry*)
Metadata	SQL(mysql,sqlite), Nosql (Redis)
Auth	Oauth(google, twitter, facebook), kerberos*, internal
Protocol	Websocket
Message Format	JSON-RPC 2.0, Amazon S3
Encoding	Plain, bson
CallBack	Subscribe/Publish Websocket/Redis, Async I/O TornadoWeb, ZeroMQ*
HASH	Sha-XXX, MD5-XXX, AES
Encryption	SSL, ciphers supported by crypto++
Discovery	DNS, file base

* are planned

What is it good for ?

Beolink.org

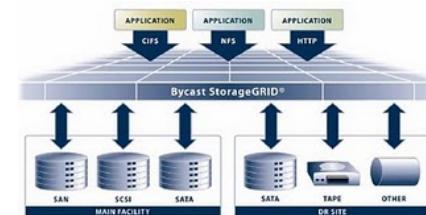
User

- Home directory
- Remote/Internet disks



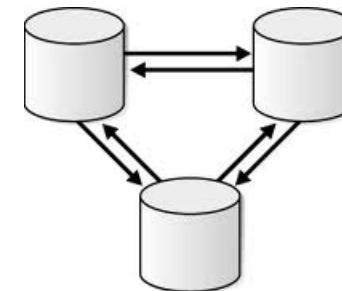
Application

- Object storage
- Shared space
- Virtual Machine



Distribution

- CDN (Multimedia)
- Data replication
- Disaster Recovery



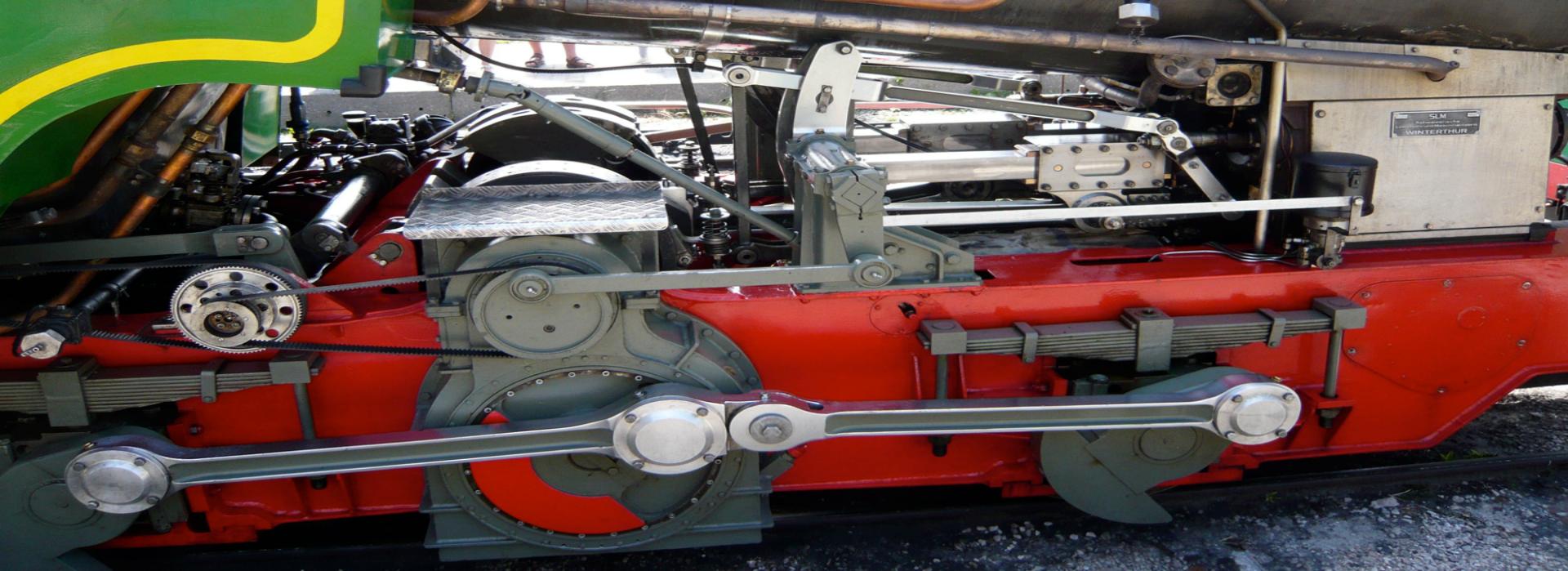


European AFS & Kerberos Conference 2014

26-28 March 2014
CERN Genève

<http://indico.cern.ch/conferenceDisplay.py?confId=271400>





Thank you

<http://restfs.beolink.org>

manfred.furuholmen@gmail.com
fege85@gmail.com

Beolink.org

Bucket

The bucket has many properties, the property element is a collection of object information, with this element you can retrieve the default value for the bucket (logging level, security level, ect).

Bucket Name

Bucket Name
zebra

Property

```
segment_size= 512
block_size = 16k
max_read'=1000
Bucket_size=0
Bucket_quota=10000
storage_class=STANDARD
compression= none
logging=enable
bucket_type=fs
...
```

Default parameters

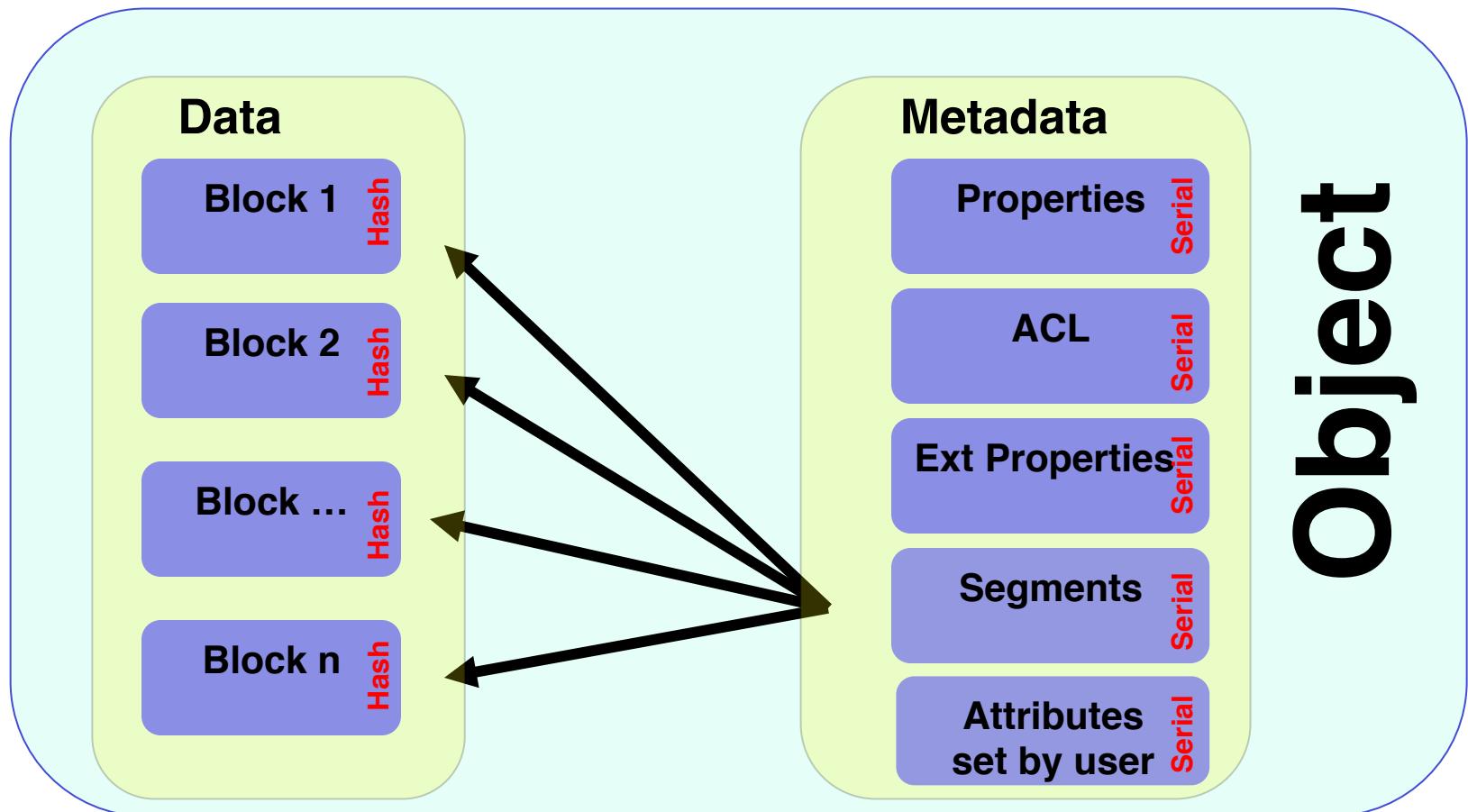
Python Dict

Properties objects:

- Property
- Property Ext
- Property ACL
- Property Stats

- **Filesystem**, The bucket is used as a filesystem
- **Logging**, Logging operation done on the specific Bucket
- **Replica RO**, Bucket shadow replication
- ...Custom definition

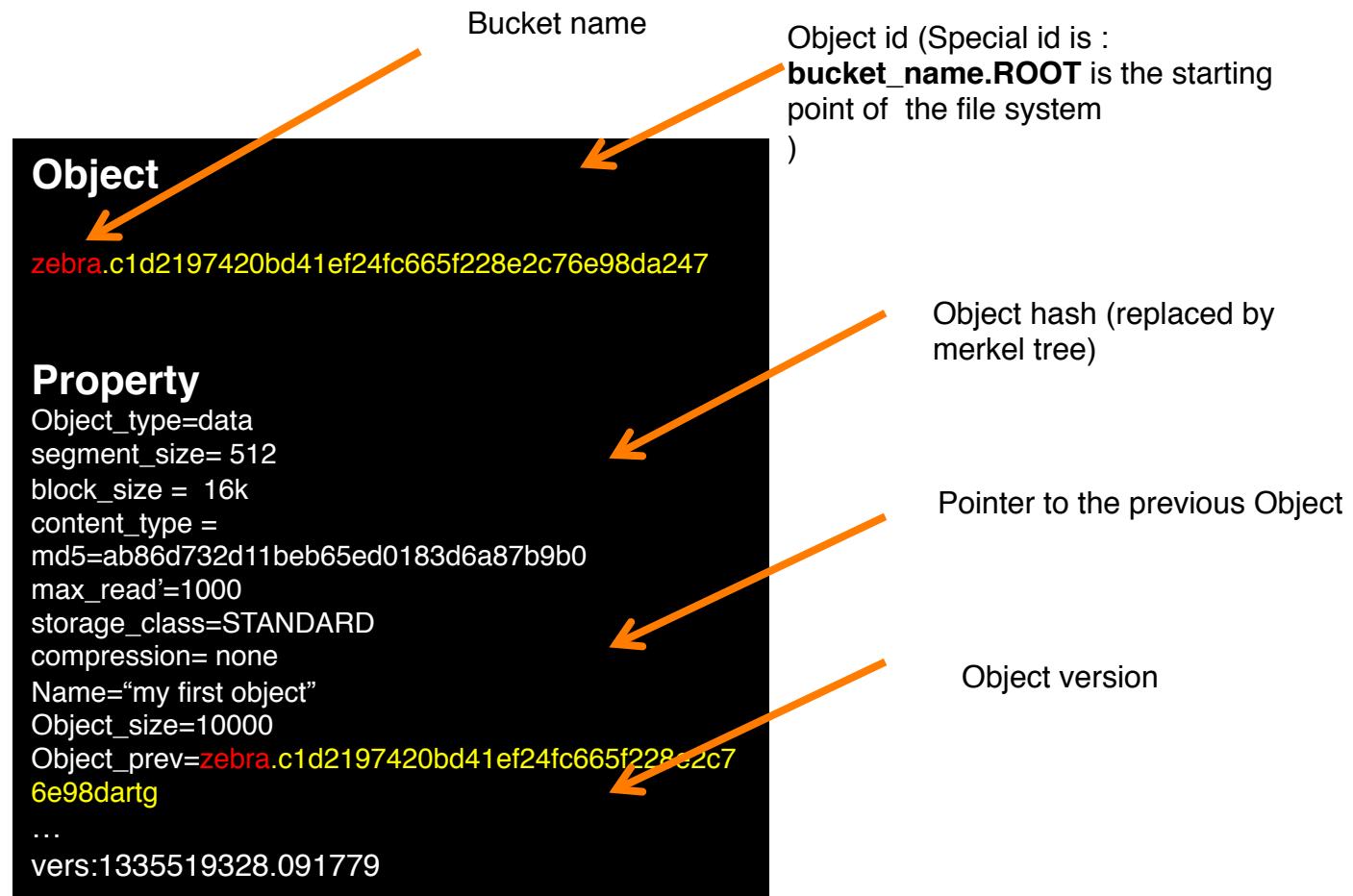
Objects



Object type:

- **Data**, Contains files
- **Folder**, Special object that contain others objects
- **Mount point**, Contains the name of the buckets
- **Link**, Contains the name of the objects
- **Immutable**, Gold image
- Custom**, Defined by the users

Object default

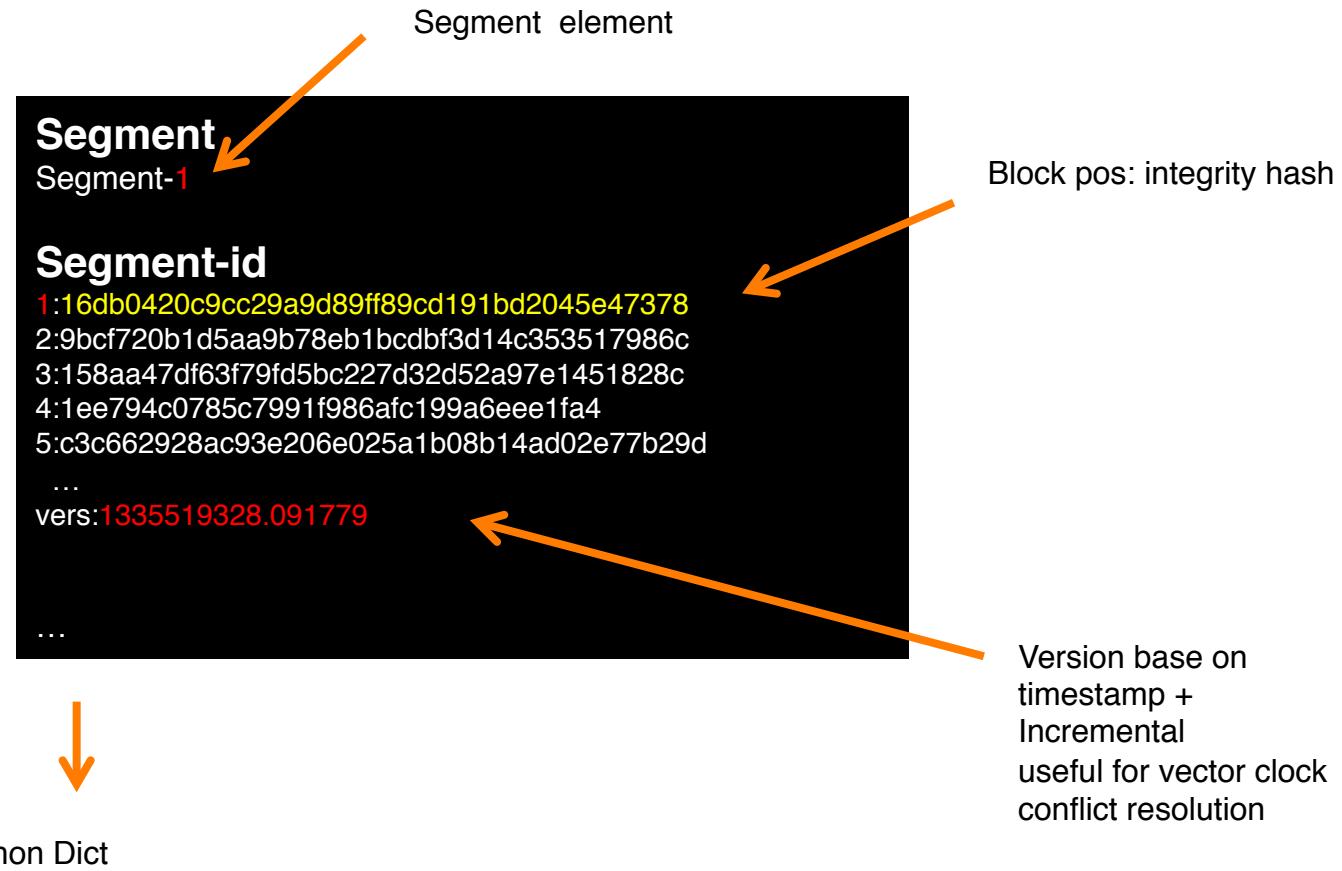


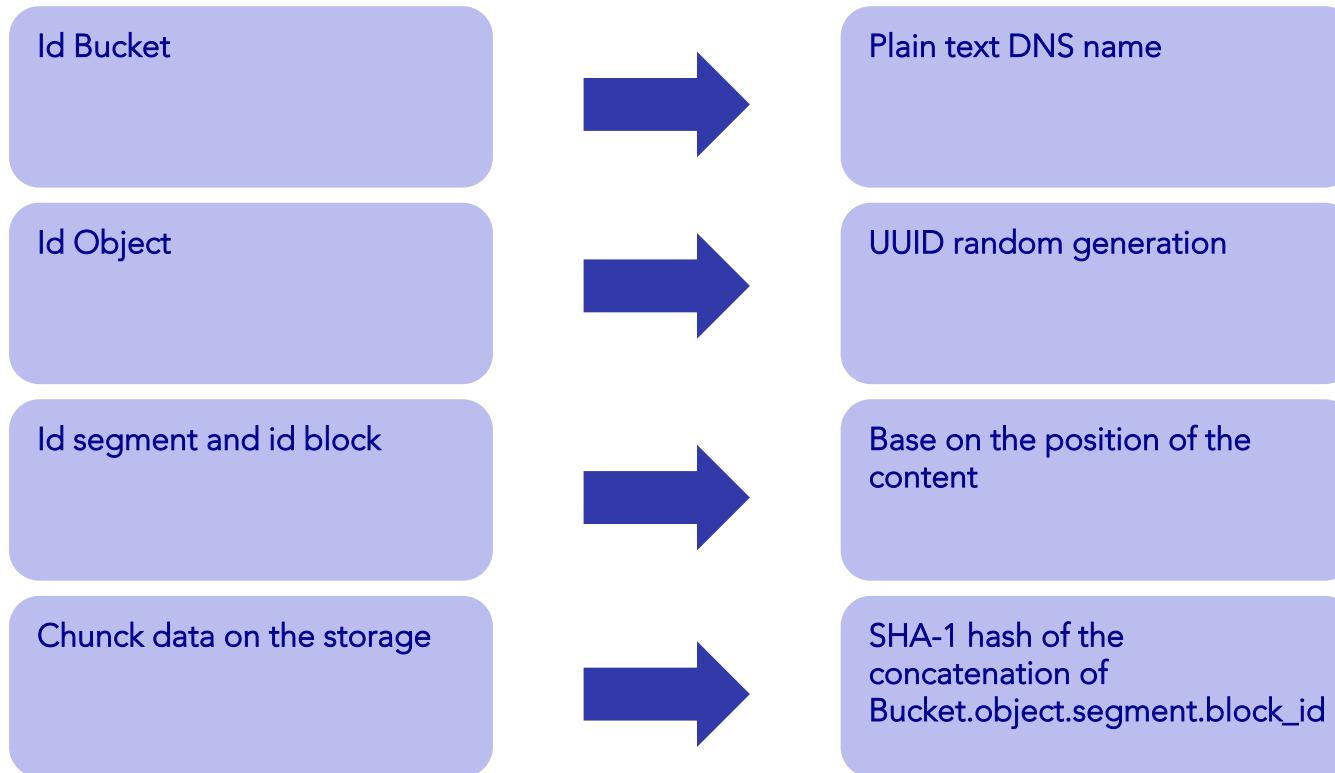
Metadata Segment

Beolink.org

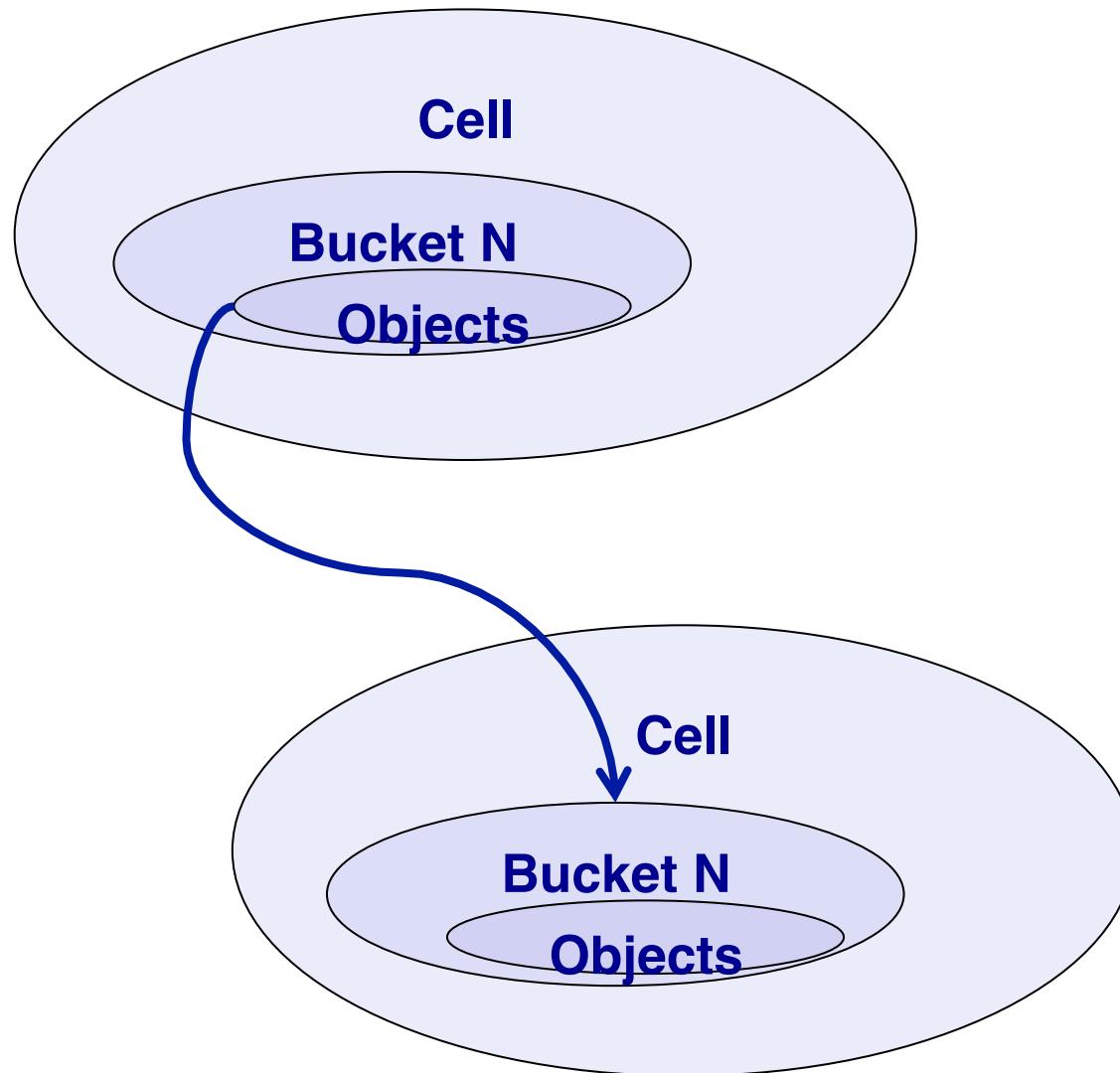
Data_size

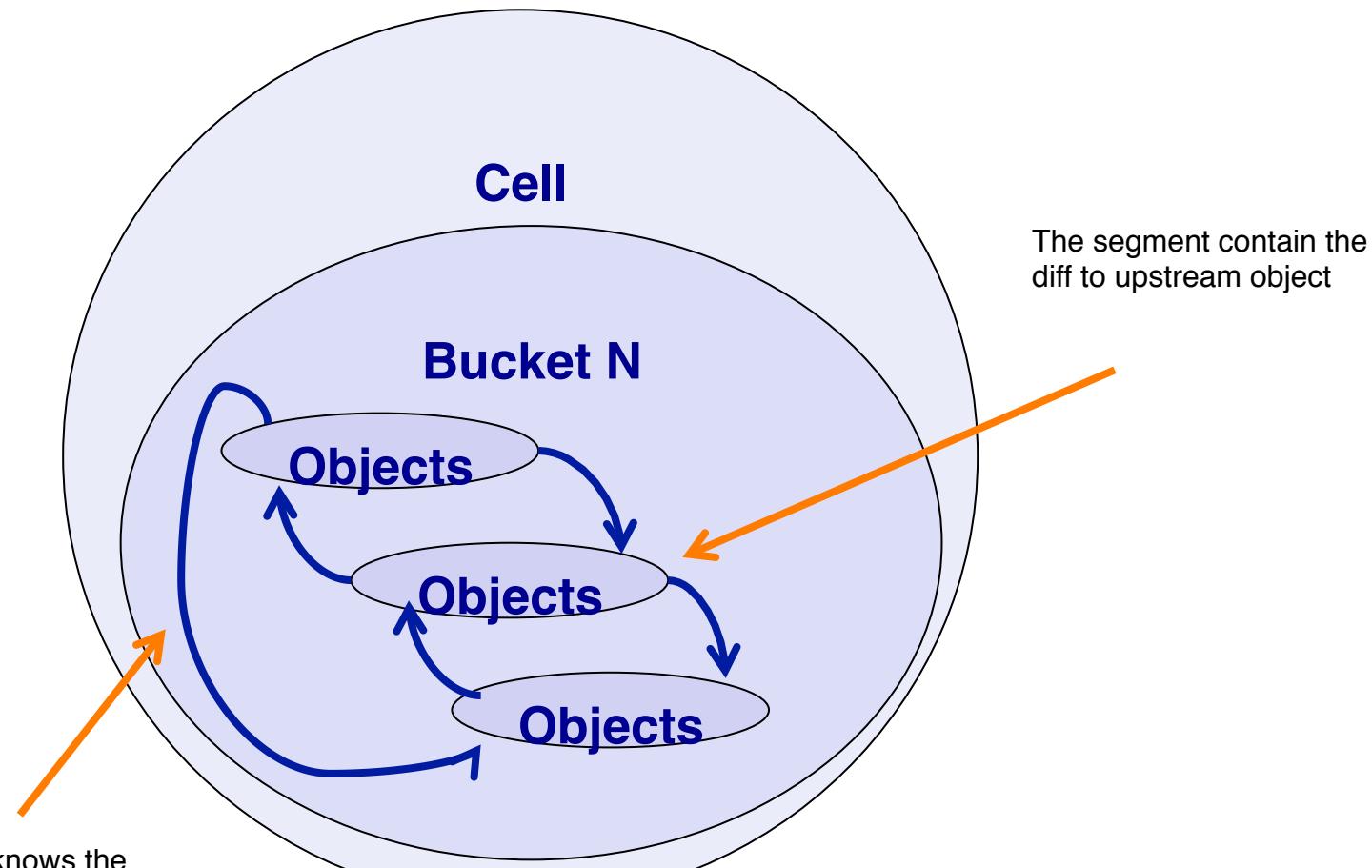
----- = Total Segment
block_size*segment_size





Id Object is unique inside of the Bucket, with bucket name the id is a UUID





Each object knows the previous and the next.
The current object knows the previous and the last

Protocols

Standard HTTP/
HTTPS port

```
GET /mychat HTTP/1.1
Host: server.example.com
Upgrade: websocket
Connection: Upgrade
Sec-WebSocket-Key: x3JJHMbDL1EzLkh9GBhXDw==
Sec-WebSocket-Protocol: chat
Sec-WebSocket-Version: 13
Origin: http://example.com
```

Simple to convert in
python dict

```
--> { "method": "readBlock", "params": [...], "id": 1 }
<-- { "result": [...], "error": null, "id": 1 }
```

WebSocket

is a web technology for multiplexing bi-directional, full-duplex communications channels over a single TCP connection.

```
{"hello": "world"}
→
"\x16\x00\x00\x00\x02hello\x00
\x06\x00\x00\x00world\x00\x00"
```

JSON-RPC

is lightweight remote procedure call protocol similar to XML-RPC. It's designed to be simple

BSON

short for Binary JSON, is a binary-encoded serialization of JSON-like documents..

BSON can be compared to binary interchange formats

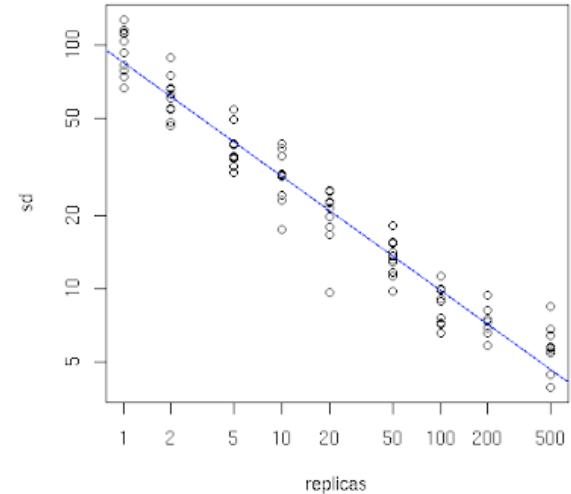
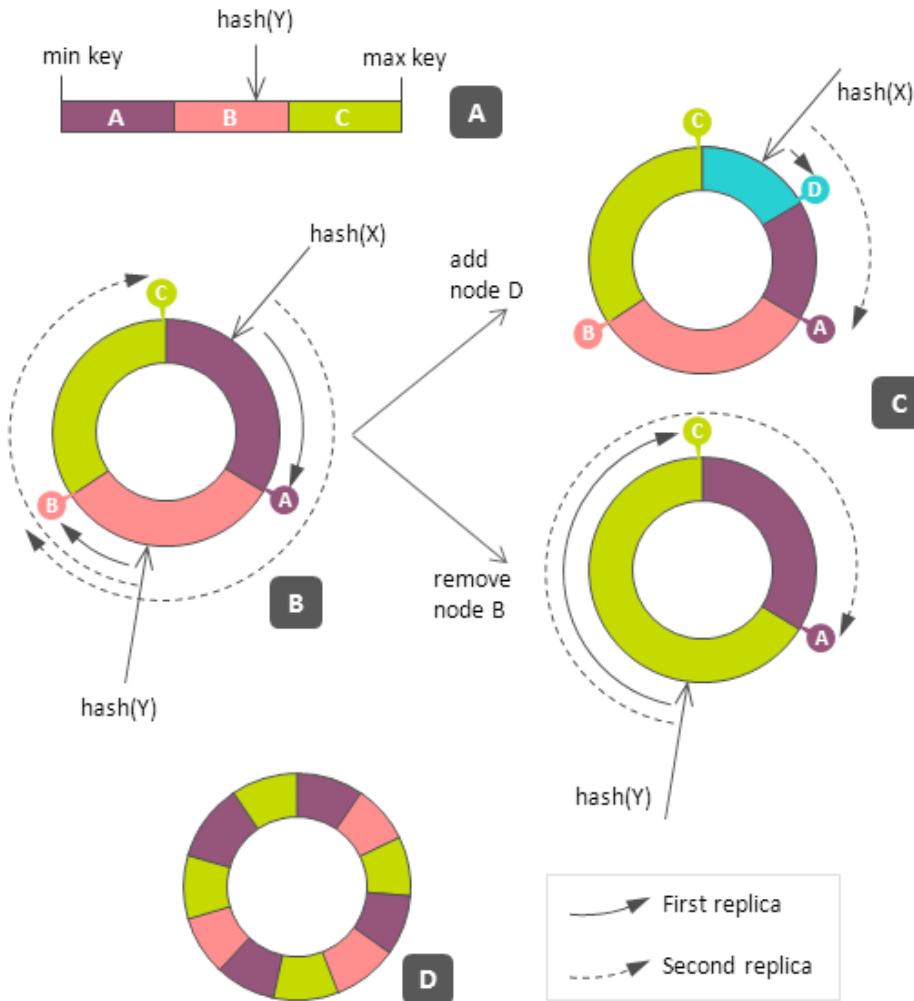
*Compression is a long story...



Block Storage

Backend: Consistent Hashing

Beolink.org



Number of key to move for add/
remove a node :

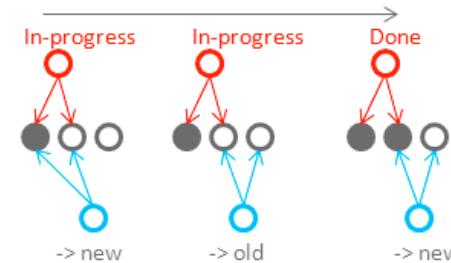
Keys/Node= keys to relocate

Blocks are collected in shards

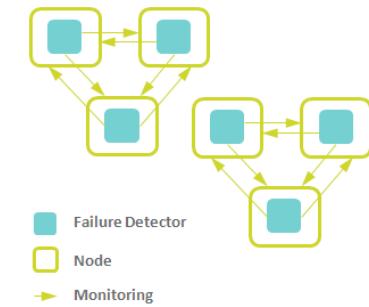
3 Copies

Configurable read and write consistent level and security:

- 2W1R
- 2W2R
- 1W1R
- ...

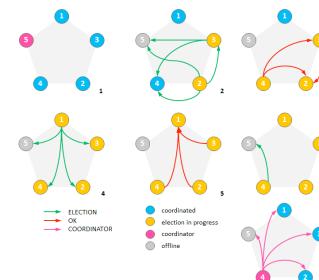


Monitor of neighbored small cluster of 3 nodes (GOSSIP)

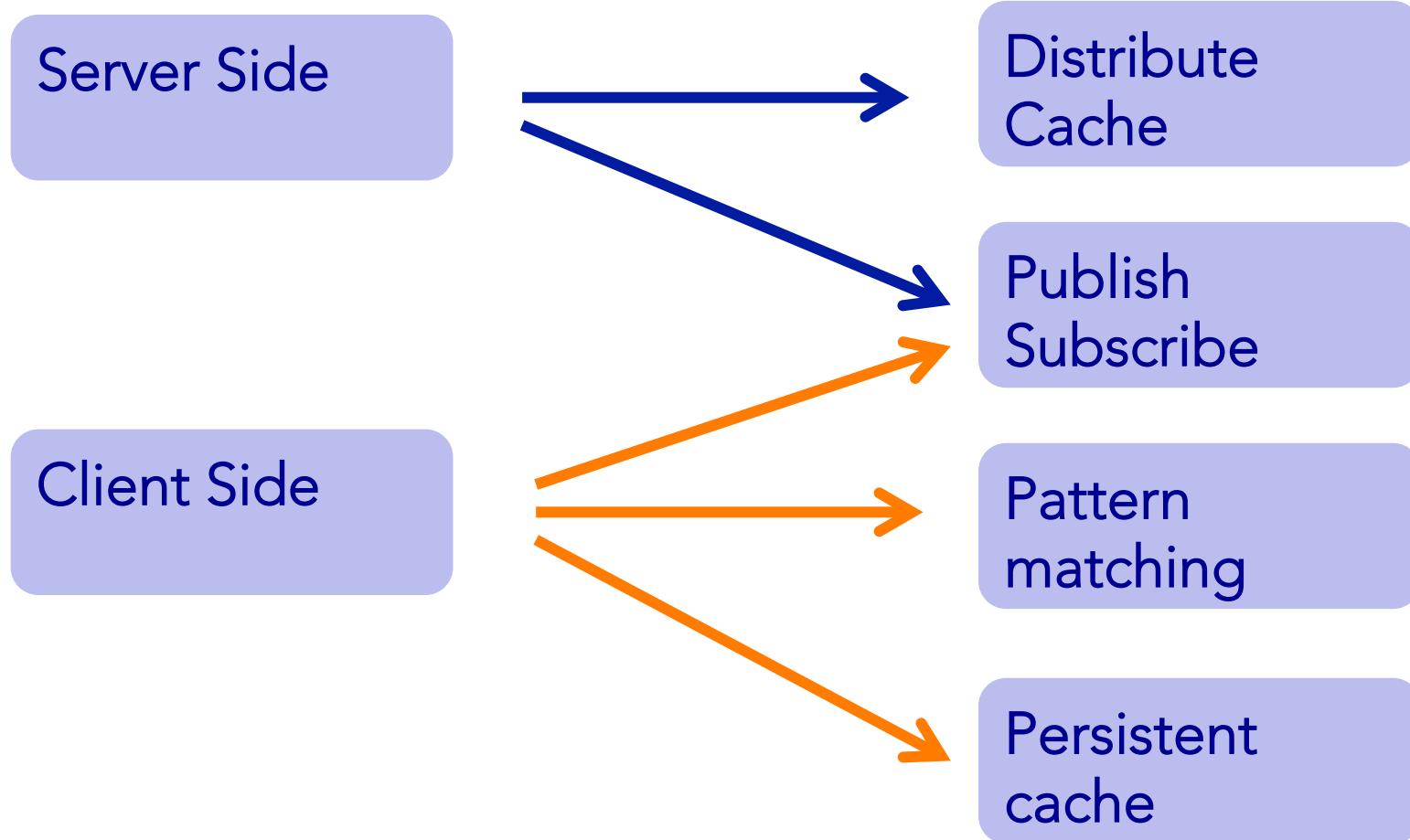


Mini cluster election

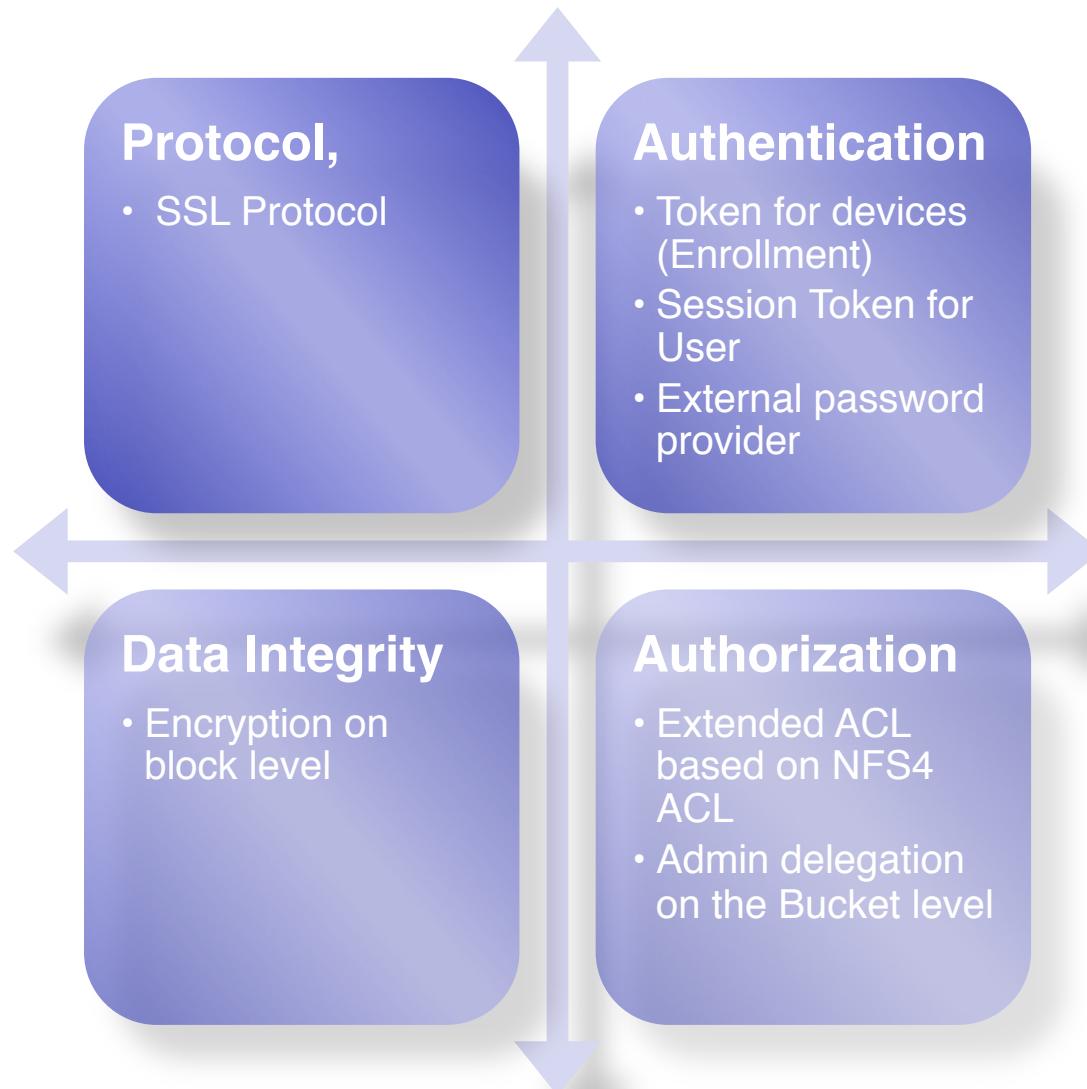
key space reclaim for replica coordination,
leave join cluster



Cache



Security

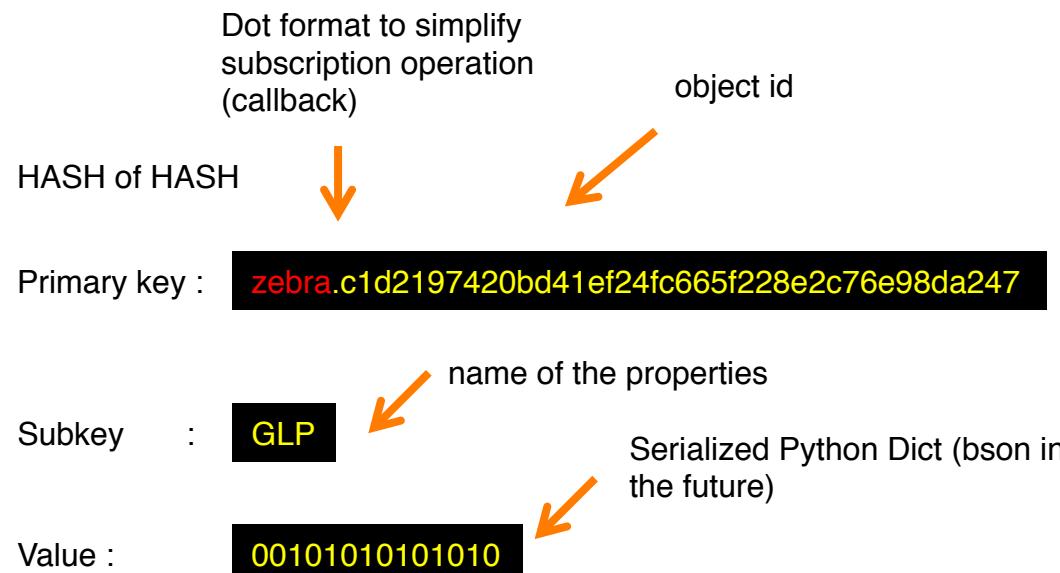


NOSQL DB

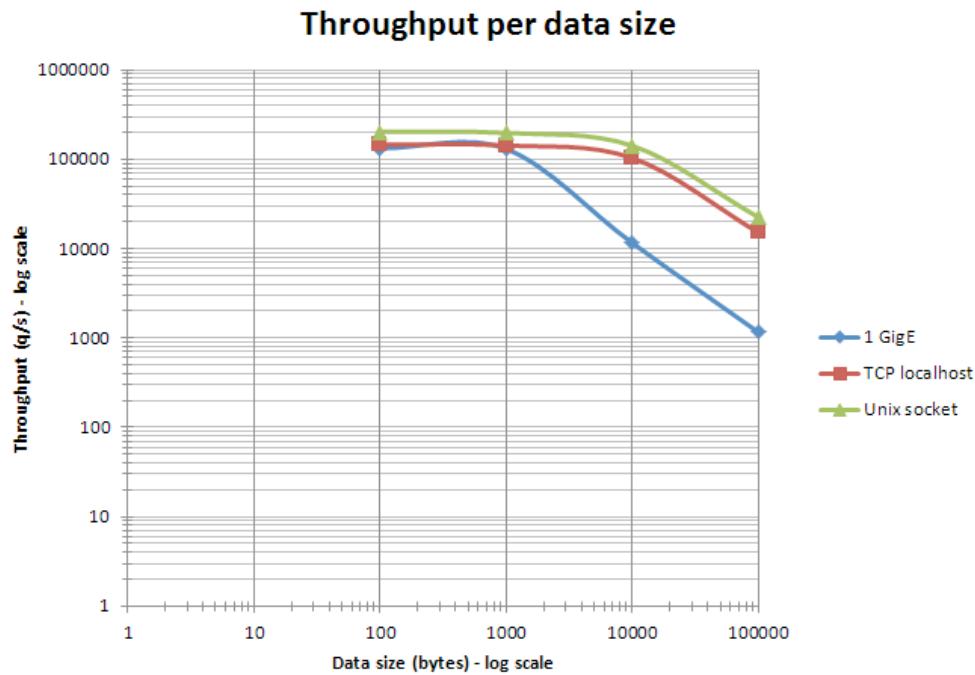


Main characteristics

- Fast
- Store Hash of HASH
- Atomic operation
- Sub/Pub primitives



* Version and Hash of the objects has a dedicated subkey, no serialization

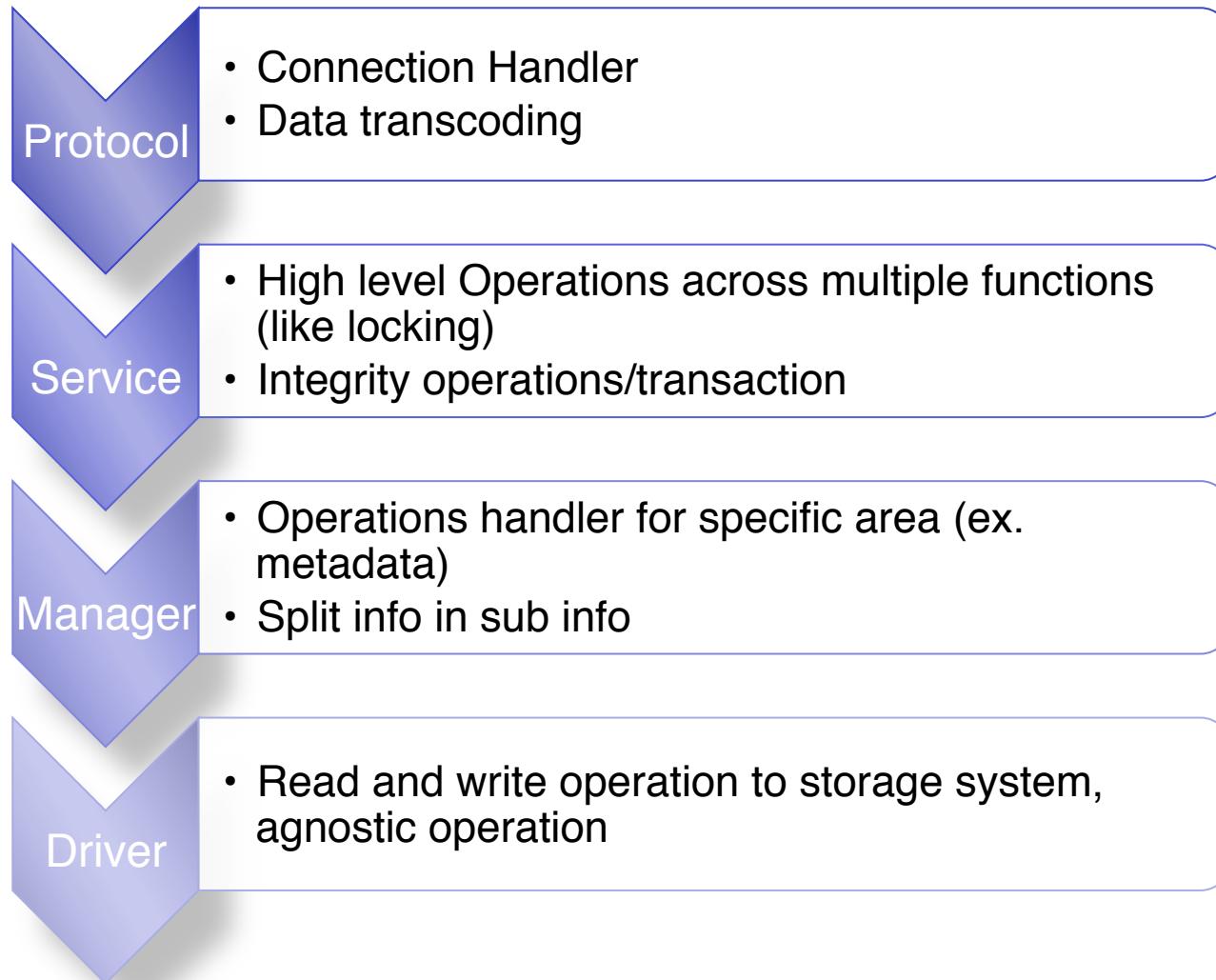


```
$ ./redis-benchmark -r 1000000 -n 2000000 -t get,set,lpush,lpop -P 16 -q
```

SET: 552028.75 requests per second
GET: 707463.75 requests per second
LPUSH: 767459.75 requests per second
LPOP: 770119.38 requests per second

Code

Interface, dynamic load



What we are using

Beolink.org



Module	Software
Storage	Filesystem, DHT (kademlia, Pastry*)
Metadata	SQL(mysql,sqlite), Nosql (Redis)
Auth	Oauth(google, twitter, facebook), kerberos*, internal
Protocol	Websocket
Message Format	JSON-RPC 2.0, Amazon S3
Encoding	Plain, bson
CallBack	Subscribe/Publish Websocket/Redis, Async I/O TornadoWeb, ZeroMQ*
HASH	Sha-XXX, MD5-XXX, AES
Encryption	SSL, ciphers supported by crypto++
Discovery	DNS, file base

* are planned

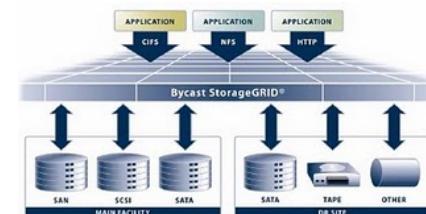
User

- Home directory
- Remote/Internet disks



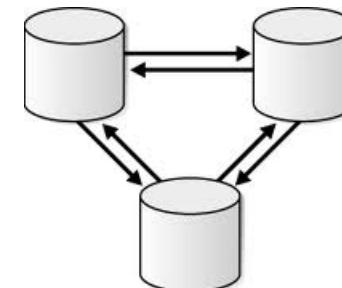
Application

- Object storage
- Shared space
- Virtual Machine



Distribution

- CDN (Multimedia)
- Data replication
- Disaster Recovery



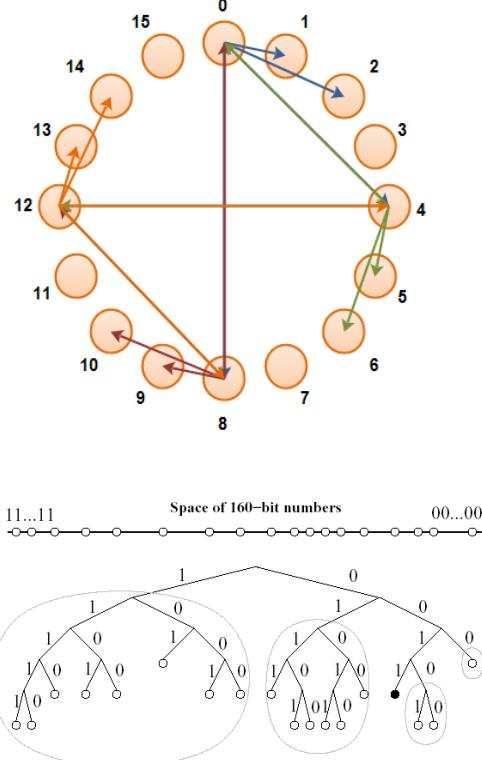


Fig. 1: Kademlia binary tree. The black dot shows the location of node 0011... in the tree. Grey ovals show subtrees in which node 0011... must have a contact.

Transport Layer ZeroMQ

Storage Compressed DAta