



Alex Auvolat, Deuxfleurs Association

`https://garagehq.deuxfleurs.fr/`
Matrix channel: `#garage:deuxfleurs.fr`

Who I am



Alex Auvolat

PhD; co-founder of Deuxfleurs



Deuxfleurs

A non-profit self-hosting collective,
member of the CHATONS network



Our objective at Deuxfleurs

**Promote self-hosting and small-scale hosting
as an alternative to large cloud providers**

Our objective at Deuxfleurs

**Promote self-hosting and small-scale hosting
as an alternative to large cloud providers**

Why is it hard?

Our objective at Deuxfleurs

**Promote self-hosting and small-scale hosting
as an alternative to large cloud providers**

Why is it hard?

Resilience

we want good uptime/availability with low supervision

Building a resilient system with cheap stuff

- ▶ Commodity hardware (e.g. old desktop PCs)

Building a resilient system with cheap stuff



Building a resilient system with cheap stuff



Building a resilient system with cheap stuff

- ▶ Commodity hardware (e.g. old desktop PCs)
(can die at any time)

Building a resilient system with cheap stuff

- ▶ Commodity hardware (e.g. old desktop PCs)
(can die at any time)

- ▶ Regular Internet (e.g. FTTB, FTTH) and power grid connections

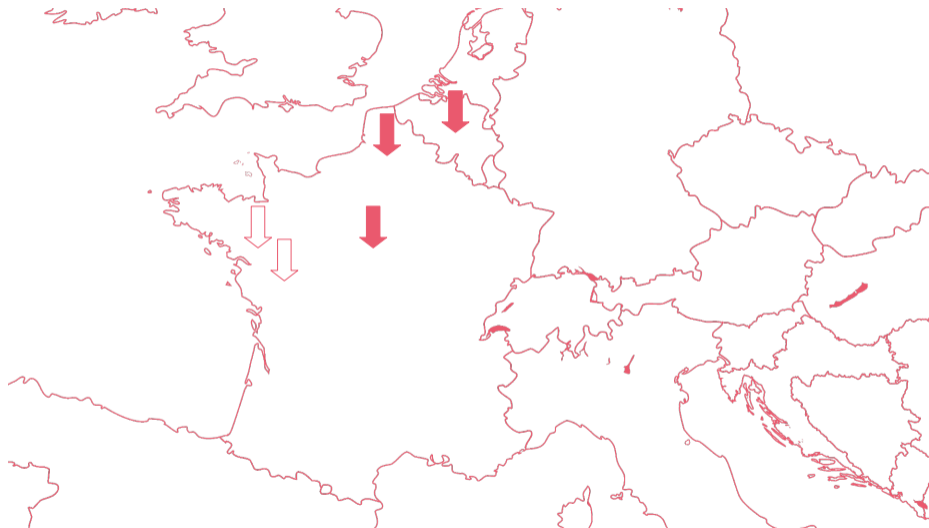
Building a resilient system with cheap stuff

- ▶ Commodity hardware (e.g. old desktop PCs)
(can die at any time)
- ▶ Regular Internet (e.g. FTTB, FTTH) and power grid connections
(can be unavailable randomly)

Building a resilient system with cheap stuff

- ▶ Commodity hardware (e.g. old desktop PCs)
(can die at any time)
- ▶ Regular Internet (e.g. FTTB, FTTH) and power grid connections
(can be unavailable randomly)
- ▶ **Geographical redundancy** (multi-site replication)

Building a resilient system with cheap stuff



Object storage: a crucial component



S3: a de-facto standard, many compatible applications

Object storage: a crucial component



S3: a de-facto standard, many compatible applications

MinIO is self-hostable but not suited for geo-distributed deployments

Object storage: a crucial component



S3: a de-facto standard, many compatible applications

MinIO is self-hostable but not suited for geo-distributed deployments

Garage is a self-hosted drop-in replacement for the Amazon S3 object store

CRDTs / weak consistency instead of consensus

Consensus can be implemented reasonably well in practice, so why avoid it?

CRDTs / weak consistency instead of consensus

Consensus can be implemented reasonably well in practice, so why avoid it?

- ▶ **Software complexity**

CRDTs / weak consistency instead of consensus

Consensus can be implemented reasonably well in practice, so why avoid it?

- ▶ **Software complexity**
- ▶ **Performance issues:**

CRDTs / weak consistency instead of consensus

Consensus can be implemented reasonably well in practice, so why avoid it?

- ▶ **Software complexity**
- ▶ **Performance issues:**
 - ▶ The leader is a **bottleneck** for all requests

CRDTs / weak consistency instead of consensus

Consensus can be implemented reasonably well in practice, so why avoid it?

- ▶ **Software complexity**
- ▶ **Performance issues:**
 - ▶ The leader is a **bottleneck** for all requests
 - ▶ **Sensitive to higher latency** between nodes

CRDTs / weak consistency instead of consensus

Consensus can be implemented reasonably well in practice, so why avoid it?

- ▶ **Software complexity**

- ▶ **Performance issues:**

- ▶ The leader is a **bottleneck** for all requests
- ▶ **Sensitive to higher latency** between nodes
- ▶ **Takes time to reconverge** when disrupted (e.g. node going down)

CRDTs / weak consistency instead of consensus

Consensus can be implemented reasonably well in practice, so why avoid it?

- ▶ **Software complexity**
- ▶ **Performance issues:**
 - ▶ The leader is a **bottleneck** for all requests
 - ▶ **Sensitive to higher latency** between nodes
 - ▶ **Takes time to reconverge** when disrupted (e.g. node going down)

Internally, Garage uses only CRDTs (conflict-free replicated data types)

The data model of object storage

Object storage is basically a **key-value store**:

Key: file path + name	Value: file data + metadata
index.html	Content-Type: text/html; charset=utf-8 Content-Length: 24929 <binary blob>
img/logo.svg	Content-Type: text/svg+xml Content-Length: 13429 <binary blob>
download/index.html	Content-Type: text/html; charset=utf-8 Content-Length: 26563 <binary blob>

The data model of object storage

Object storage is basically a **key-value store**:

Key: file path + name	Value: file data + metadata
index.html	Content-Type: text/html; charset=utf-8 Content-Length: 24929 <binary blob>
img/logo.svg	Content-Type: text/svg+xml Content-Length: 13429 <binary blob>
download/index.html	Content-Type: text/html; charset=utf-8 Content-Length: 26563 <binary blob>

- ▶ Maps well to CRDT data types

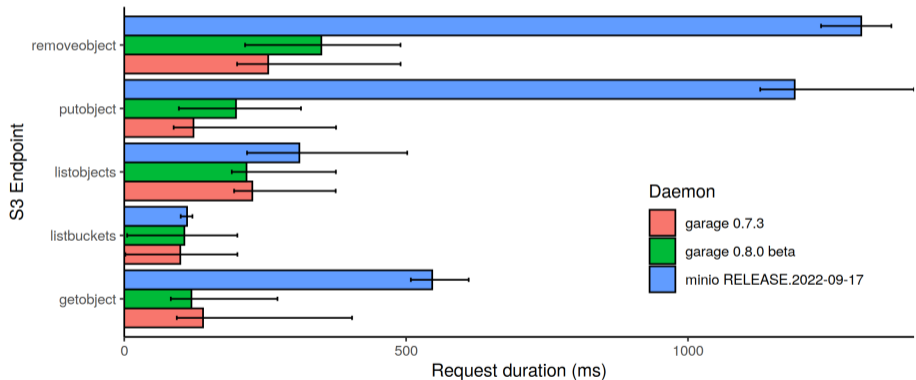
Performance gains in practice

S3 endpoint latency in a simulated geo-distributed cluster

100 measurements, 5 nodes, 50ms RTT + 10ms jitter between nodes

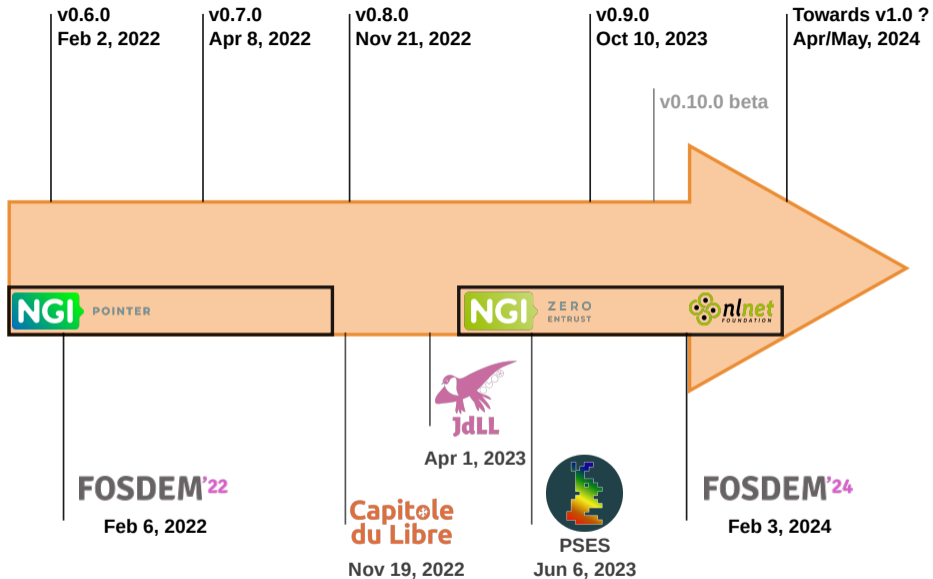
no contention: latency is due to intra-cluster communications

colored bar = mean latency, error bar = min and max latency



Get the code to reproduce this graph at <https://git.deuxfleurs.fr/Deuxfleurs/mknet>

Recent developments

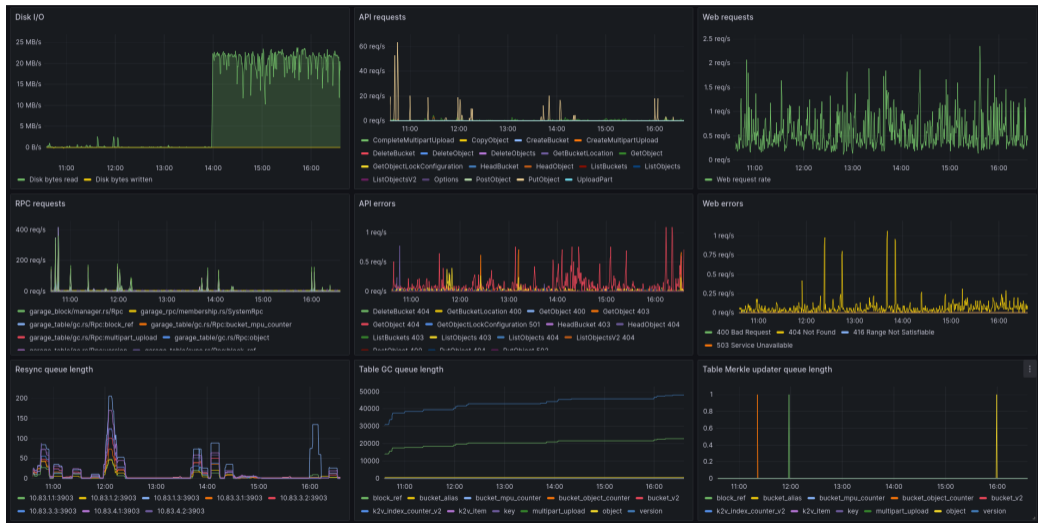


April 2022 - Garage v0.7.0

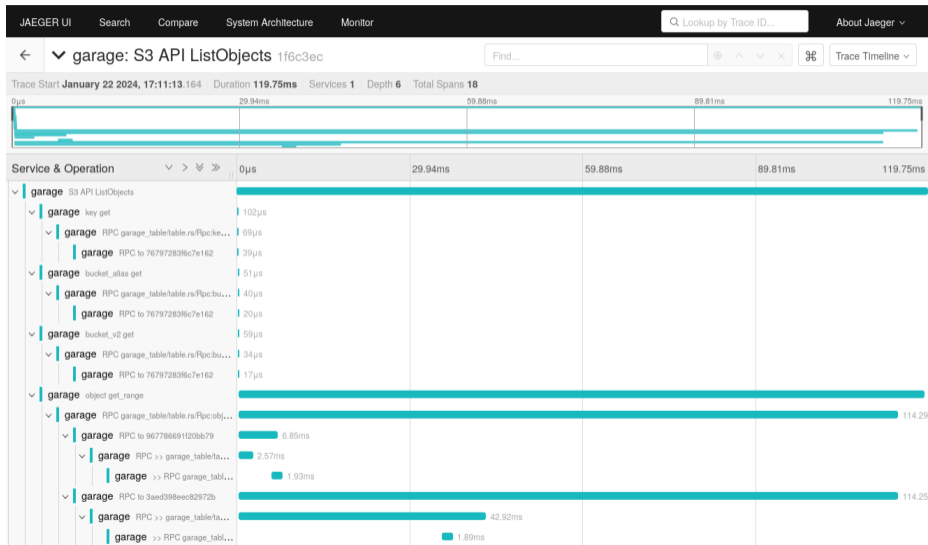
Focus on observability and ecosystem integration

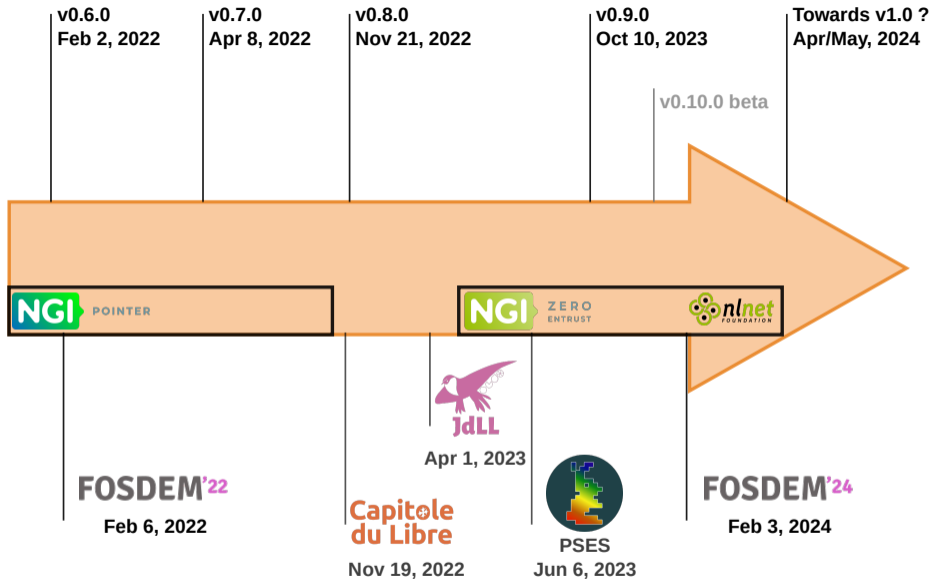
- ▶ **Monitoring:** metrics and traces, using OpenTelemetry
- ▶ Replication modes with 1 or 2 copies / weaker consistency
- ▶ Kubernetes integration
- ▶ Admin API (v0.7.2)
- ▶ Experimental K2V API (v0.7.2)

Metrics (Prometheus + Grafana)



Traces (Jaeger)





November 2022 - Garage v0.8.0

Focus on performance

- ▶ **Alternative metadata DB engines** (LMDB, Sqlite)
- ▶ **Performance improvements:** block streaming, various optimizations...
- ▶ Bucket quotas (max size, max #objects)
- ▶ Quality of life improvements, observability, etc.

About metadata DB engines

Issues with Sled:

- ▶ Huge files on disk
- ▶ Unpredictable performance, especially on HDD
- ▶ API limitations
- ▶ Not actively maintained

LMDB: very stable, good performance, reasonably small files on disk

Sled will be removed in Garage v1.0

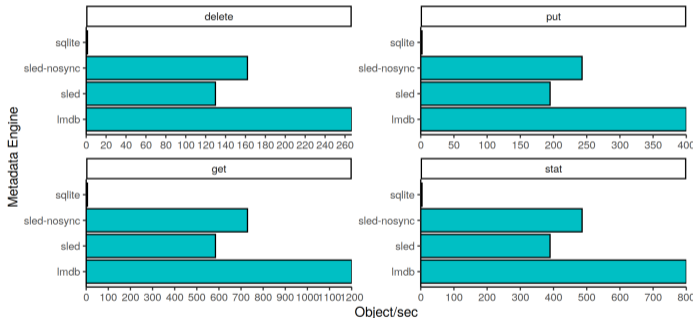
DB engine performance comparison

Comparison of Garage's metadata engines with "minio/warp"

Daemon: Garage v0.8 no-fsync to avoid being impacted by block manager

Benchmark: warp, mixed mode, 5min bench, 256B objects, initialized with 200 objects.

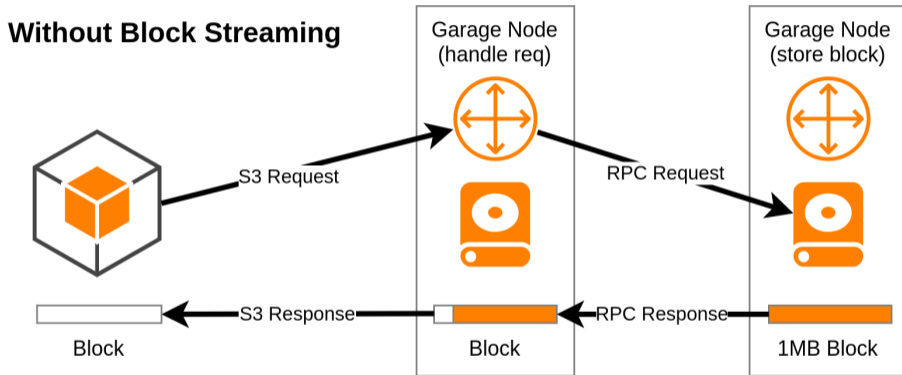
Environment: mknet (Ryzen 5 1400, 16GB RAM, SSD). DC topo (3 nodes, 1Gb/s, 1ms latency).



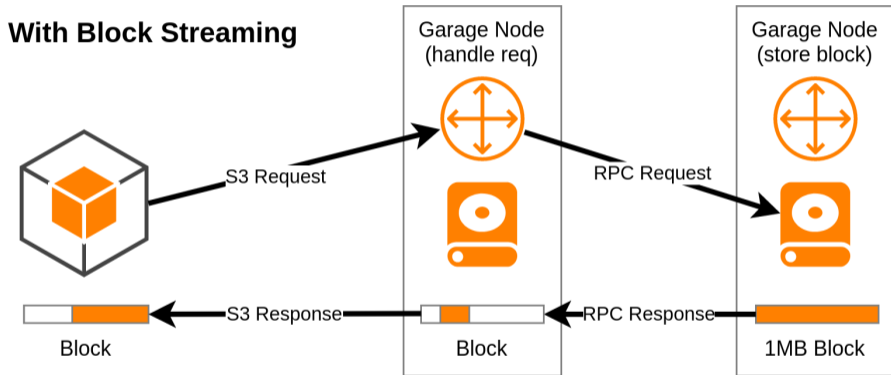
Get the code to reproduce this graph at <https://git.deuxfleurs.fr/Deuxfleurs/mknet>

NB: Sqlite was slow due to synchronous journaling mode, now configurable

Block streaming



Block streaming

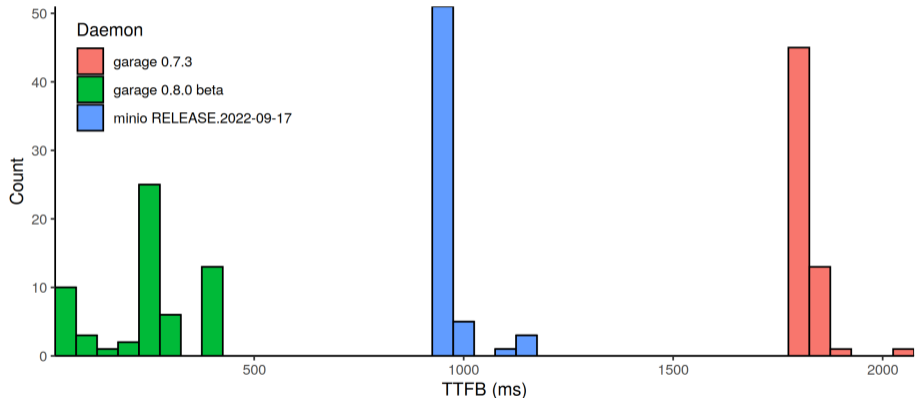


TTFB benchmark

TTFB (Time To First Byte) on GetObject over a slow network (5 Mbps, 500 μ s)

A 1MB file is uploaded and then fetched 60 times.

Except for Minio, the queried node does not store any data (gateway) to force net. communications.



Get the code to reproduce this graph at <https://git.deuxfleurs.fr/Deuxfleurs/mknet>

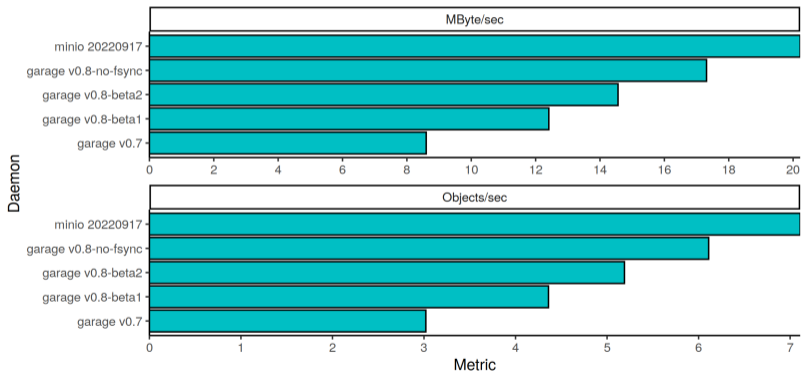
Throughput benchmark

"minio/warp" benchmark, "cluster total" result

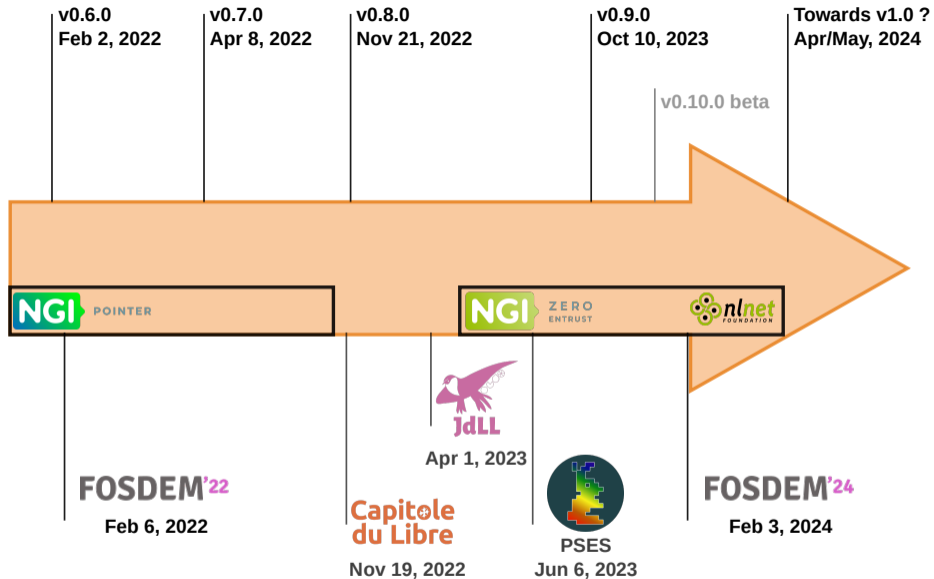
Ran on a local machine (Ryzen 5 1400, 16GB RAM, SSD) with mknet

DC topology (3 nodes, 1GB/s, 1ms lat)

warp in mixed mode, 5min bench, 5MB objects, initialized with 200 objects



Get the code to reproduce this graph at <https://git.deuxfleurs.fr/Deuxfleurs/mknet>



October 2023 - Garage v0.9.0

Focus on streamlining & usability

- ▶ Support multiple HDDs per node
- ▶ S3 compatibility:
 - ▶ support basic lifecycle configurations
 - ▶ allow for multipart upload part retries
- ▶ LMDB by default, deprecation of Sled
- ▶ New layout computation algorithm

Layout computation

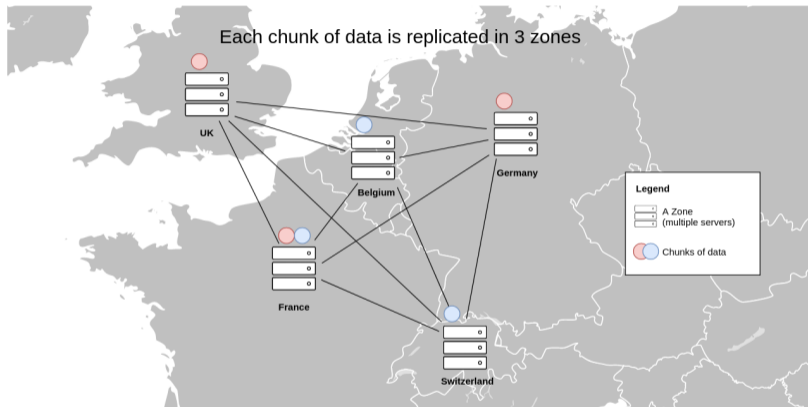
```
[root@celeri:/home/lx]# docker exec -ti e338 /garage status
```

```
==== HEALTHY NODES ====
```

ID	Hostname	Address	Tags	Zone	Capacity
5fcb3b6e39db3dcb	concombre	[2001:470:ca43::31]:3901	[concombre,neptune,france,alex]	neptune	500.0 GB
942dd71ea95f4904	df-ymf	[2a02:a03f:6510:5102:6e4b:90ff:fe3a:6174]:3901	[df-ymf,bespin,belgium,max]	bespin	500.0 GB
fdfaf7832d8359e0	df-ymk	[2a02:a03f:6510:5102:6e4b:90ff:fe3b:e939]:3901	[df-ymk,bespin,belgium,max]	bespin	500.0 GB
0a03ab7c082ad929	ananas	[2a01:e0a:e4:2dd0::42]:3901	[ananas,scorpio,france,adrien]	scorpio	2.0 TB
a717e5b618267806	courgette	[2001:470:ca43::32]:3901	[courgette,neptune,france,alex]	neptune	500.0 GB
2032d0a37f249c4a	abricot	[2a01:e0a:e4:2dd0::41]:3901	[abricot,scopio,france,adrien]	scorpio	2.0 TB
8cf284e7df17d0fd	celeri	[2001:470:ca43::33]:3901	[celeri,neptune,france,alex]	neptune	2.0 TB
17ee03c6b81d9235	df-ykl	[2a02:a03f:6510:5102:6e4b:90ff:fe3b:e86c]:3901	[df-ykl,bespin,belgium,max]	bespin	500.0 GB

Garage stores replicas on different zones when possible

Layout computation



Garage stores replicas on different zones when possible

What a "layout" is

A layout is a precomputed index table:

Partition	Node 1	Node 2	Node 3
Partition 0	lo (jupiter)	Drosera (atuin)	Courgette (neptune)
Partition 1	Datura (atuin)	Courgette (neptune)	lo (jupiter)
Partition 2	lo(jupiter)	Celeri (neptune)	Drosera (atuin)
⋮	⋮	⋮	⋮
Partition 255	Concombre (neptune)	lo (jupiter)	Drosera (atuin)

What a "layout" is

A layout is a precomputed index table:

Partition	Node 1	Node 2	Node 3
Partition 0	lo (jupiter)	Drosera (atuin)	Courgette (neptune)
Partition 1	Datura (atuin)	Courgette (neptune)	lo (jupiter)
Partition 2	lo(jupiter)	Celeri (neptune)	Drosera (atuin)
⋮	⋮	⋮	⋮
Partition 255	Concombre (neptune)	lo (jupiter)	Drosera (atuin)

The index table is built centrally using an optimal algorithm, then propagated to all nodes

An algorithm for geo-distributed and redundant storage in Garage

Mendes Oulamara^{*} and Alex Auvolat[†]

Deuxfleurs

Abstract

This paper presents an optimal algorithm to compute the assignment of data to storage nodes in the Garage geo-distributed storage system. We discuss the complexity of the different steps of the algorithm and metrics that can be displayed to the user.

1 Introduction

Garage^[1] is an open-source distributed object storage service tailored for self-hosting. It was designed by the Deuxfleurs association^[2] to enable small structures (associations, collectives, small companies) to share storage resources to reliably self-host their data, possibly with old and non-reliable machines. To achieve these reliability and availability goals, the data is broken into *partitions* and every partition is replicated over 3 different machines (that we call *nodes*). When the data is queried, it is fetched from one of the nodes. A replication factor of 3 ensures good guarantees regarding node failure^[5]. But this parameter can be another (preferably larger and odd) number.

Moreover, if the nodes are spread over different *zones* (different houses, offices, cities...), we can require the data to be replicated over nodes belonging to different zones. This improves the storage robustness against

```
*mendes@deuxfleurs.fr
†alex@deux.fr
1https://garagehq.deuxfleurs.fr/
2https://deuxfleurs.fr/
```

1

zone failures (such as power outages). To do so, we define a *scattering factor*, that is no more than the replication factor, and we require that the replicas of any partition are spread over this number of zones at least.

In this work, we propose an assignment algorithm that, given the nodes specifications and the replication and scattering factors, computes an optimal assignment of partitions to nodes. We say that the assignment is optimal in the sense that it maximizes the size of the partitions, and hence the effective storage capacity of the system.

Moreover, when a former assignment exists, which is not optimal anymore due to node or zone changes, our algorithm computes a new optimal assignment that minimizes the amount of data to be transferred during the assignment update (the *transfer load*).

We call the set of nodes cooperating to store the data a *cluster*, and a description of the nodes, zones and the assignment of partitions to nodes a *cluster layout*.

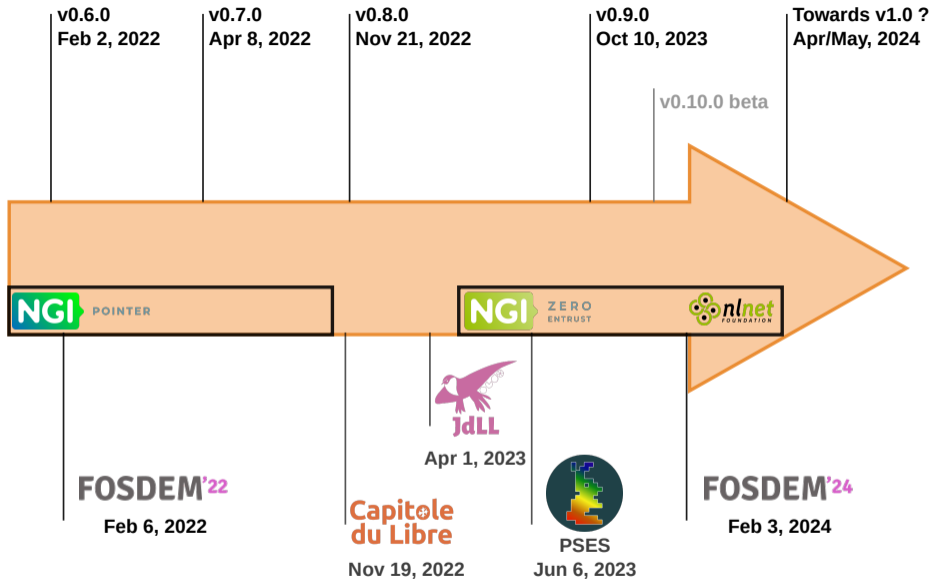
1.1 Notations

Let k be some fixed parameter value, typically 8, that we call the "partition bits". Every object to be stored in the system is split into data blocks of fixed size. We compute a hash $h(b)$ of every such block b , and we define the k first bits of this hash to be the partition number $p(b)$ of the block. This label can take $P = 2^k$ different values, and hence there are P different partitions. We denote \mathbf{P} the set of partition labels (i.e. $\mathbf{P} = [1, P]$).

We are given a set N of N nodes and a set Z of Z zones. Every node n has a non-negative storage capacity $c_n \geq 0$ and belongs to a zone $z_n \in Z$. We are also given a replication factor ρ_N and a scattering factor ρ_Z such that $1 \leq \rho_Z \leq \rho_N$ (typical values would be $\rho_N = \rho_Z = 3$).

Our goal is to compute an assignment $\alpha = (\alpha_p^1, \dots, \alpha_p^N)_{p \in \mathbf{P}}$ such that every partition p is associated to ρ_N distinct nodes $\alpha_p^1, \dots, \alpha_p^N \in N$ and these nodes belong to at least ρ_Z distinct zones. Among the possible assignments, we choose one that *maximizes* the effective storage capacity of the cluster. If the layout contained a previous assignment α' , we *minimize* the amount of data to transfer during the layout update by making α as close as possible to α' . These maximization and minimization are described more formally in the following section.

2



October 2023 - Garage v0.10.0 beta

Focus on consistency

- ▶ Fix consistency issues when reshuffling data

Working with weak consistency

Not using consensus limits us to the following:

Working with weak consistency

Not using consensus limits us to the following:

- ▶ **Conflict-free replicated data types (CRDT)**

Non-transactional key-value stores such as S3 are equivalent to a simple CRDT:
a map of **last-writer-wins registers** (each key is its own CRDT)

Working with weak consistency

Not using consensus limits us to the following:

- ▶ **Conflict-free replicated data types (CRDT)**

Non-transactional key-value stores such as S3 are equivalent to a simple CRDT:
a map of **last-writer-wins registers** (each key is its own CRDT)

- ▶ **Read-after-write consistency**

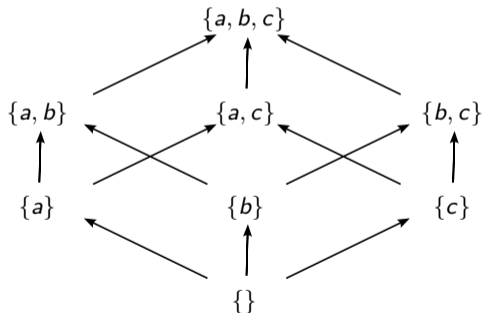
Can be implemented using quorums on read and write operations

CRDT read-after-write consistency using quorums

Property: If node A did an operation $write(x)$ and received an OK response, and node B starts an operation $read()$ after A received OK, then B will read a value $x' \supseteq x$.

CRDT read-after-write consistency using quorums

Property: If node A did an operation $write(x)$ and received an OK response, and node B starts an operation $read()$ after A received OK, then B will read a value $x' \supseteq x$.

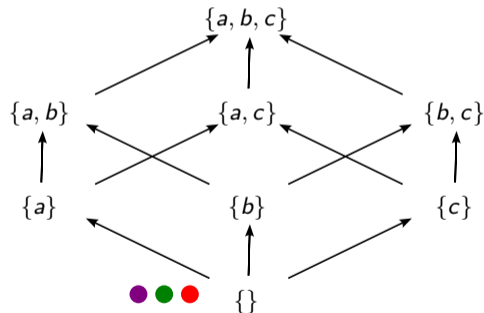


CRDT read-after-write consistency using quorums

Property: If node A did an operation $write(x)$ and received an OK response, and node B starts an operation $read()$ after A received OK, then B will read a value $x' \supseteq x$.

$write(\{a\})$:

- $\not\supseteq \{a\}$
- $\supseteq \{a\}$
- $\not\supseteq \{a\}$

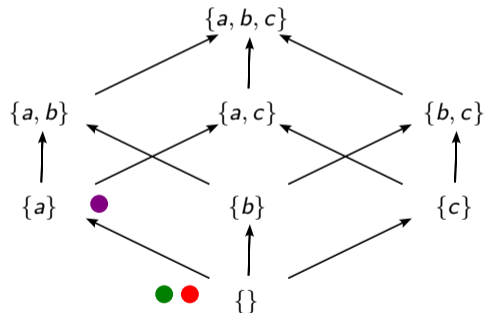


CRDT read-after-write consistency using quorums

Property: If node A did an operation $write(x)$ and received an OK response, and node B starts an operation $read()$ after A received OK, then B will read a value $x' \supseteq x$.

$write(\{a\})$:

- $\supseteq \{a\} \rightarrow \text{OK}$
- $\not\supseteq \{a\}$
- $\not\supseteq \{a\}$



CRDT read-after-write consistency using quorums

Property: If node A did an operation $write(x)$ and received an OK response, and node B starts an operation $read()$ after A received OK, then B will read a value $x' \supseteq x$.

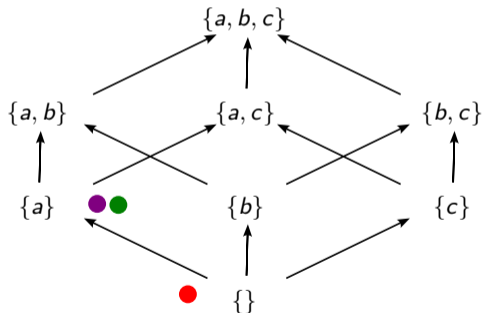
$write(\{a\})$:

● $\supseteq \{a\} \rightarrow$ OK

● $\supseteq \{a\} \rightarrow$ OK

● $\not\supseteq \{a\}$

return OK



CRDT read-after-write consistency using quorums

Property: If node A did an operation $write(x)$ and received an OK response, and node B starts an operation $read()$ after A received OK, then B will read a value $x' \supseteq x$.

$write(\{a\})$:

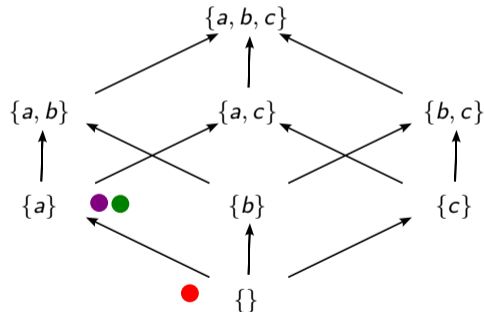
● $\supseteq \{a\} \rightarrow \text{OK}$

● $\supseteq \{a\} \rightarrow \text{OK}$

● $\not\supseteq \{a\}$

return OK

$read()$:



CRDT read-after-write consistency using quorums

Property: If node A did an operation $write(x)$ and received an OK response, and node B starts an operation $read()$ after A received OK, then B will read a value $x' \supseteq x$.

$write(\{a\})$:

● $\supseteq \{a\} \rightarrow \text{OK}$

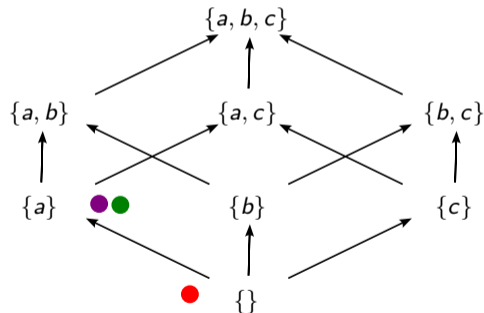
● $\supseteq \{a\} \rightarrow \text{OK}$

● $\not\supseteq \{a\}$

return OK

$read()$:

● $\rightarrow \{\}$



CRDT read-after-write consistency using quorums

Property: If node A did an operation $write(x)$ and received an OK response, and node B starts an operation $read()$ after A received OK, then B will read a value $x' \supseteq x$.

$write(\{a\})$:

● $\supseteq \{a\} \rightarrow \text{OK}$

● $\supseteq \{a\} \rightarrow \text{OK}$

● $\not\supseteq \{a\}$

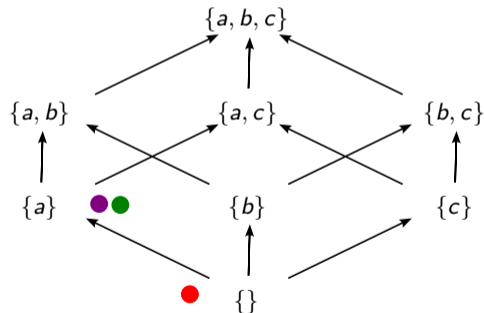
return OK

$read()$:

● $\rightarrow \{\}$

● $\rightarrow \{a\}$

return $\{\} \sqcup \{a\} = \{a\}$



CRDT read-after-write consistency using quorums

Property: If node A did an operation $write(x)$ and received an OK response, and node B starts an operation $read()$ after A received OK, then B will read a value $x' \supseteq x$.

$write(\{a\})$:

● $\supseteq \{a\} \rightarrow \text{OK}$

● $\supseteq \{a\} \rightarrow \text{OK}$

● $\supseteq \{a\}$

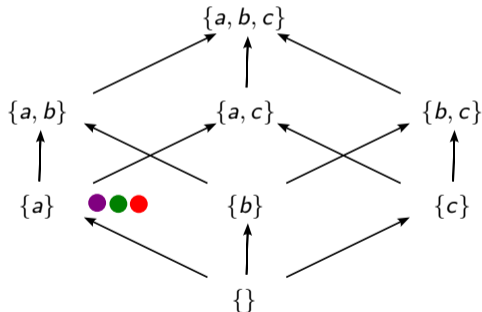
return OK

$read()$:

● $\rightarrow \{\}$

● $\rightarrow \{a\}$

return $\{\} \sqcup \{a\} = \{a\}$



CRDT read-after-write consistency using quorums

Property: If node A did an operation $write(x)$ and received an OK response, and node B starts an operation $read()$ after A received OK, then B will read a value $x' \supseteq x$.

Algorithm $write(x)$:

1. Broadcast $write(x)$ to all nodes
2. Wait for $k > n/2$ nodes to reply OK
3. Return OK

Algorithm $read()$:

1. Broadcast $read()$ to all nodes
2. Wait for $k > n/2$ nodes to reply with values x_1, \dots, x_k
3. Return $x_1 \sqcup \dots \sqcup x_k$

A hard problem: layout changes

- ▶ We rely on quorums $k > n/2$ within each partition:

$$n = 3, \quad k \geq 2$$

A hard problem: layout changes

- ▶ We rely on quorums $k > n/2$ within each partition:

$$n = 3, \quad k \geq 2$$

- ▶ When rebalancing, the set of nodes responsible for a partition can change:

$$\{A, B, C\} \rightarrow \{A, D, E\}$$

A hard problem: layout changes

- ▶ We rely on quorums $k > n/2$ within each partition:

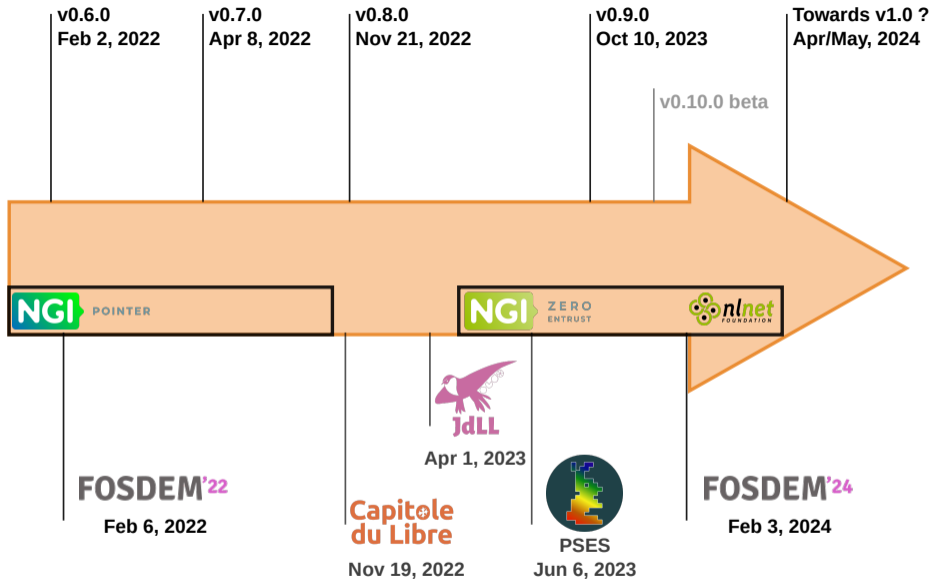
$$n = 3, \quad k \geq 2$$

- ▶ When rebalancing, the set of nodes responsible for a partition can change:

$$\{A, B, C\} \rightarrow \{A, D, E\}$$

- ▶ During the rebalancing, D and E don't yet have the data,
and B and C want to get rid of the data to free up space

→ risk of inconsistency, **how to coordinate?**



Towards v1.0

Focus on security & stability

- ▶ **Security audit** in progress by Radically Open Security
- ▶ Misc. S3 features (SSE-C, ...) and compatibility fixes
- ▶ Improve UX
- ▶ Fix bugs

Operating big Garage clusters

Operating Garage

```
$ garage status
==== HEALTHY NODES ====
ID                Hostname  Address                               Tags                Zone    Capacity  DataAvail
ec5753c546756825 df-pw5    [2a02:a03f:6510:5102:223:24ff:feb0:e8a7]:3991 [df-pw5]  bespin  500.0 GB  429.1 GB (89.0%)
76797283f6c7e162 carcajou  [2001:470:ca43::22]:3991                [carcajou] neptune 200.0 GB  166.3 GB (73.5%)
8073f25ffb7d6944 piranha   [2a01:cb05:911e:ec00:223:24ff:feb0:ea82]:3991 [piranha]  corrin  500.0 GB  457.3 GB (94.0%)
3aed398eec82972b origan    [2a01:e0a:5e4:1d0:223:24ff:feaf:fdec]:3991 [origan]   jupiter 500.0 GB  457.1 GB (93.1%)
967786691f20bb79 caribou   [2001:470:ca43::23]:3991                [caribou]  neptune 500.0 GB  453.1 GB (92.3%)
```

Operating Garage

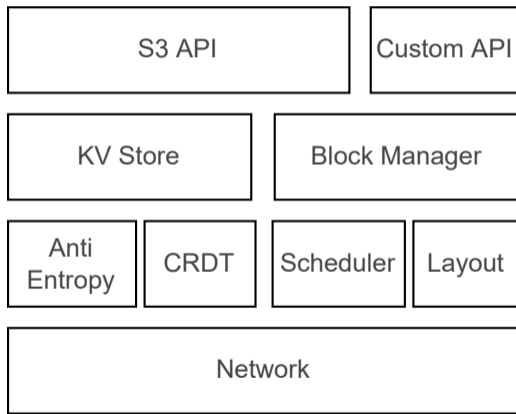
```
$ garage status
==== HEALTHY NODES ====
ID           Hostname  Address                               Tags           Zone    Capacity  DataAvail
ec5753c546756825 df-pw5   [2a02:a03f:6510:5102:223:24ff:feb0:e8a7]:3991 [df-pw5]     bespin  500.0 GB  429.1 GB (89.0%)
76797283f6c7e162 carcajou [2001:470:ca43::22]:3991                [carcajou]  neptune 200.0 GB  166.3 GB (73.5%)
8073f25ffb7d6944 piranha  [2a01:cb05:911e:ec00:223:24ff:feb0:ea82]:3991 [piranha]   corrin  500.0 GB  457.3 GB (94.0%)
3aed398eec82972b origan   [2a01:e0a:5e4:1d0:223:24ff:feaf:fdec]:3991  [origan]    jupiter 500.0 GB  457.1 GB (93.1%)
967786691f20bb79 caribou  [2001:470:ca43::23]:3991                [caribou]   neptune 500.0 GB  453.1 GB (92.3%)
```

```
$ garage status
==== HEALTHY NODES ====
ID           Hostname  Address                               Tags           Zone    Capacity  DataAvail
76797283f6c7e162 carcajou [2001:470:ca43::22]:3991                [carcajou]  neptune 200.0 GB  166.3 GB (73.5%)
8073f25ffb7d6944 piranha  [2a01:cb05:911e:ec00:223:24ff:feb0:ea82]:3991 [piranha]   corrin  500.0 GB  457.3 GB (94.0%)
3aed398eec82972b origan   [2a01:e0a:5e4:1d0:223:24ff:feaf:fdec]:3991  [origan]    jupiter 500.0 GB  457.1 GB (93.1%)
967786691f20bb79 caribou  [2001:470:ca43::23]:3991                [caribou]   neptune 500.0 GB  453.1 GB (92.3%)

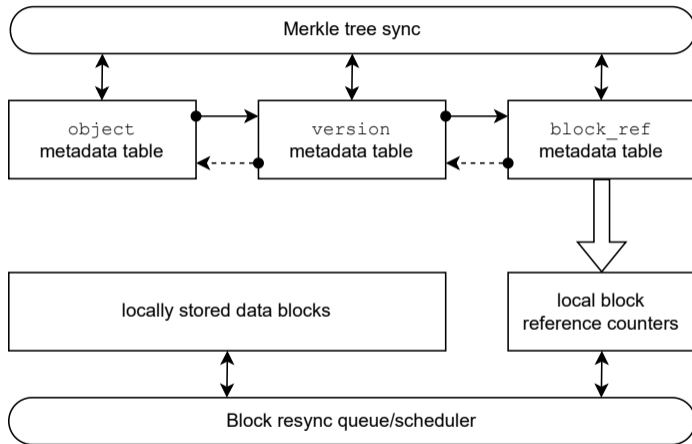
==== FAILED NODES ====
ID           Hostname  Address                               Tags           Zone    Capacity  Last seen
ec5753c546756825 df-pw5   [2a02:a03f:6510:5102:223:24ff:feb0:e8a7]:3991 [df-pw5]     bespin  500.0 GB  5 minutes ago
```

Garage's architecture

Garage as a set of components



Garage's architecture



Digging deeper

```
$ garage stats

Garage version: 20240116133343 [features: k2v, sled, lmbd, sqlite, consul-discovery, kubernetes-discovery, metrics, telemetry-otlp, bundled-libs]
Rust compiler version: 1.68.0

Database engine: LMDB (using Heed crate)

Table stats:
Table      Items  MklItems  MklTodo  GcTodo
bucket_v2  19     20        0        0
key        12     14        0        0
object     67391  80964    0        0
version    33909  42045    0        0
block_ref  334735 370927    0        0

Block manager stats:
number of RC entries (~= number of blocks): 42376
resync queue length: 0
blocks with resync errors: 0

If values are missing above (marked as NC), consider adding the --detailed flag (this will be slow).

Storage nodes:
ID          Hostname  Zone      Capacity  Part.  DataAvail          MetaAvail
ec5753c546756825  df-pw5   bespin    500.0 GB  175    429.1 GB/482.1 GB (89.0%)  429.1 GB/482.1 GB (89.0%)
76797283f6c7e162  carcajou neptune   200.0 GB  70     166.3 GB/226.2 GB (73.5%)  166.3 GB/226.2 GB (73.5%)
8073f25fffb7d6944  piranha  corrin    500.0 GB  173    457.3 GB/486.4 GB (94.0%)  457.3 GB/486.4 GB (94.0%)
3aed398eecd82972b  origan   jupiter   500.0 GB  175    457.1 GB/490.7 GB (93.1%)  457.1 GB/490.7 GB (93.1%)
967786691f20bb79  caribou  neptune   500.0 GB  175    453.1 GB/490.8 GB (92.3%)  453.1 GB/490.8 GB (92.3%)

Estimated available storage space cluster-wide (might be lower in practice):
data: 608.3 GB
metadata: 608.3 GB
```

Digging deeper

```
$ garage worker list
```

TID	State	Name	Tranq	Done	Queue	Errors	Consec	Last
1	Idle	Block resync worker #1	0	-	0	-	-	
2	Idle	Block resync worker #2	0	-	0	-	-	
3	Idle	Block resync worker #3	0	-	0	-	-	
4	Idle	Block resync worker #4	0	-	0	-	-	
5	Idle	Block resync worker #5	-	-	-	-	-	
6	Idle	Block resync worker #6	-	-	-	-	-	
7	Idle	Block resync worker #7	-	-	-	-	-	
8	Idle	Block resync worker #8	-	-	-	-	-	
9	Idle	Block scrub worker	4	-	-	-	-	
10	Idle	bucket_v2 Merkle	-	-	0	-	-	
11	Idle	bucket_v2 sync	-	-	0	1	0	17 hours ago
12	Idle	bucket_v2 GC	-	-	0	-	-	
13	Idle	bucket_v2 queue	-	-	0	-	-	
14	Idle	bucket_alias Merkle	-	-	0	-	-	
15	Idle	bucket_alias sync	-	-	0	1	0	17 hours ago
16	Idle	bucket_alias GC	-	-	0	-	-	
17	Idle	bucket_alias queue	-	-	0	-	-	
18	Idle	key Merkle	-	-	0	-	-	
19	Idle	key sync	-	-	0	1	0	17 hours ago
20	Idle	key GC	-	-	0	-	-	
21	Idle	key queue	-	-	0	-	-	
22	Idle	object Merkle	-	-	0	-	-	
23	Idle	object sync	-	-	0	4	0	17 hours ago
24	Idle	object GC	-	-	0	-	-	
25	Idle	object queue	-	-	0	-	-	
26	Idle	bucket_object_counter Merkle	-	-	0	-	-	
27	Idle	bucket_object_counter sync	-	-	0	4	0	17 hours ago
28	Idle	bucket_object_counter GC	-	-	0	-	-	
29	Idle	bucket_object_counter queue	-	-	0	-	-	
30	Idle	multipart upload Merkle	-	-	0	-	-	
31	Idle	multipart upload sync	-	-	0	5	0	17 hours ago
32	Idle	multipart upload GC	-	-	0	-	-	
33	Idle	multipart upload queue	-	-	0	-	-	
34	Idle	bucket_mpu_counter Merkle	-	-	0	-	-	
35	Idle	bucket_mpu_counter sync	-	-	0	-	-	
36	Idle	bucket_mpu_counter GC	-	-	0	-	-	
37	Idle	bucket_mpu_counter queue	-	-	0	-	-	
38	Idle	version Merkle	-	-	0	-	-	
39	Idle	version sync	-	-	0	50	0	17 hours ago
40	Idle	version GC	-	-	0	-	-	
41	Idle	version queue	-	-	0	-	-	
42	Idle	block_ref Merkle	-	-	0	-	-	
43	Idle	block_ref sync	-	-	0	45	0	17 hours ago
44	Idle	block_ref GC	-	-	0	-	-	
45	Idle	block_ref queue	-	-	0	-	-	
46	Idle	object lifecycle worker	-	-	-	-	-	

Digging deeper

```
$ garage worker get
8073f25ffb7d6944 lifecycle-last-completed 2024-01-23
8073f25ffb7d6944 resync-tranquility 1
8073f25ffb7d6944 resync-worker-count 4
8073f25ffb7d6944 scrub-corruptions_detected 0
8073f25ffb7d6944 scrub-last-completed 2023-12-27T13:49:33.234Z
8073f25ffb7d6944 scrub-next-run 2024-01-31T03:23:02.234Z
8073f25ffb7d6944 scrub-tranquility 4

$ garage worker get -a resync-tranquility
3aed398eec82972b resync-tranquility 1
76797283f6c7e162 resync-tranquility 1
8073f25ffb7d6944 resync-tranquility 1
967786691f20bb79 resync-tranquility 1
ec5753c546756825 resync-tranquility 1
```

Potential limitations and bottlenecks

- ▶ Global:
 - ▶ Max. ~ 100 nodes per cluster (excluding gateways)
- ▶ Metadata:
 - ▶ One big bucket = bottleneck, object list on 3 nodes only
- ▶ Block manager:
 - ▶ Lots of small files on disk
 - ▶ Processing the resync queue can be slow

Deployment advice for very large clusters

- ▶ Metadata storage:
 - ▶ ZFS mirror (x2) on fast NVMe
 - ▶ Use LMDB storage engine
- ▶ Data block storage:
 - ▶ Use Garage's native multi-HDD support
 - ▶ XFS on individual drives
 - ▶ Increase block size (1MB → 10MB, requires more RAM and good networking)
 - ▶ Tune `resync-tranquility` and `resync-worker-count` dynamically
- ▶ Other :
 - ▶ Split data over several buckets
 - ▶ Use less than 100 storage nodes
 - ▶ Use gateway nodes

Current deployments: < 10 TB, we don't have much experience with more

Where to find us



Garage

`https://garagehq.deuxfleurs.fr/`
`mailto:garagehq@deuxfleurs.fr`
`#garage:deuxfleurs.fr` on Matrix

