

Design information

Design documentation outline

- Storage service
 - Redundancy and performance
 - S3 API data path
 - Storage cluster
 - Scaling the Storage service
- Compute service
 - Managed Compute
 - Bring Your Own Compute
 - Scaling the Compute service
- Backend services
 - Network
 - Management
 - Control plane
- Site operations
 - Bringing up a site
 - Maintenance operations

Storage service

Redundancy and performance

S3 API data path

blocked URL

The storage service uses BGP per-flow ECMP routing to achieve redundancy and high throughput for frontend services. A pair of BGP routers will serve in front of the L3-L7 Træfik <https://traefik.io/> load balancers. Using ECMP, we will route each TCP flow to one of the 3 frontend Træfik instances. Normally, each Træfik process will route to a radosgw endpoint on the same host, as drawn above. However, Træfik can also route traffic to a remote radosgw on another node, if needed.

In the simplified example below, Træfik health checks has detected that the backend radosgw on node2 is down (e.g for maintenance) and rerouted traffic to radosgw3.

blocked URL

The last step in the data path is the Ceph storage cluster. All physical drives in the cluster are controlled by a corresponding ceph-osd process. Each radosgw instance write directly to the osd daemons, being coordinated by Ceph cluster monitors.

Storage cluster

Redundancy levels in the storage cluster is controlled by Ceph configuration parameters. The planned start configuration for the system is using Erasure encoding with a 4+2 k/m (data/parity) value. This st

Scaling the Storage service

Compute service

<Picture of software-level services/processes>

Backend services

Network

Operational process

SUNET operates the TOR network switches. Read access will be given to Safespring, e.g SNMP level monitoring. Changes are initiated by both parties, but executed by SUNET.

Safespring operates the management switches. If SUNET wants read access to the management switches it will be provided.

High level design goals

1. Verifiable, repeatable performance for the service layer defined by testable SLO values (throughput with thread counts, quality expectancy)
2. Use BGP routing for server link redundancy, as opposed to L2 protocol mechanisms.

- a. BFD for failure detection ([example docs](#))
3. LLDP to verify topology
4. ECMP based failover for service IPs. The same service IP will be announced from multiple hosts to achieve multipath, each of the hosts are both able to terminate and load balance traffic towards the same service IP. The frontend routers will need to be configured for per-flow based ECMP so that each flow reaches the same load balancer frontend.
5. Fabric that enable easy expansion to a spine+leaf structure if needed

Below is a simple conceptual drawing showing BGP routers and connectivity for what Safespring considers a solution to the above. It is a simple L2 fabric with multiple distinct L2 channels, with no use of L2 redundancy protocols. A BGP router on each server node provides redundancy by having two available paths to all other server nodes. There's also two paths to the customer core through the TOR border leaf routers.

Two route reflectors per L2 switching plane share their routes using iBGP, thus we avoid having to configure a full mesh.

blocked URL

This conceptual illustration shows a full leaf/spine expanding on the same, where BGP is used on the border leafs and on all server nodes.

[blocked URL](#)

These concepts are built on Safespring operational experience where we have seen the benefits of keeping the network as simple as possible.

- <https://code.facebook.com/posts/360346274145943/introducing-data-center-fabric-the-next-generation-facebook-data-center-network/>
- <https://vincent.bernat.im/en/blog/2018-l3-routing-hypervisor>
- <https://docs.projectcalico.org/v3.1/reference/private-cloud/l2-interconnect-fabric>

Operational design goals

1. Vendor qualified binary drivers available for server NICs
 - Greatly simplifies operations of physical nodes
 - Keeps complexity down

Switch hardware

The TOR switches will be acquired on a joint spec with SUNET. The management switches will be acquired by Safespring.

Control plane

The control plane will be based on a four node 2U cluster with redundant frontend services for the storage cluster. It will also be used as a starting point for serving the Compute service

Site operations

Bringing up a site

Maintenance operations