



From VMware, SAN, OpenAFS/NFS to Proxmox and Ceph/CephFS

A Practical Migration Path



- Introduction & context
- Starting point: VMware + SAN
- Why Consider Change
- Proxmox + Ceph
- Migration Path + Lessons learned
- Backup
- Any mishaps?
- Some final thoughts

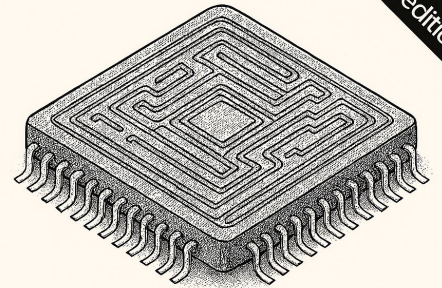
Introduction & Context

A large, light gray silhouette of the London skyline is visible in the background. It includes the Tower Bridge, the Gherkin, Big Ben, and St. Paul's Cathedral. A large, semi-transparent version of the red 'ceph' logo is overlaid on the skyline, centered behind the main title.

```
$ whoami
I'm Wannes Smet
Senior ICT Systems Engineer
at ICsense
Strong interest in
  * Linux
  * Ceph & Proxmox
  * Data center infrastructure
  * Open-Source Software
$
```

Introduction & Context

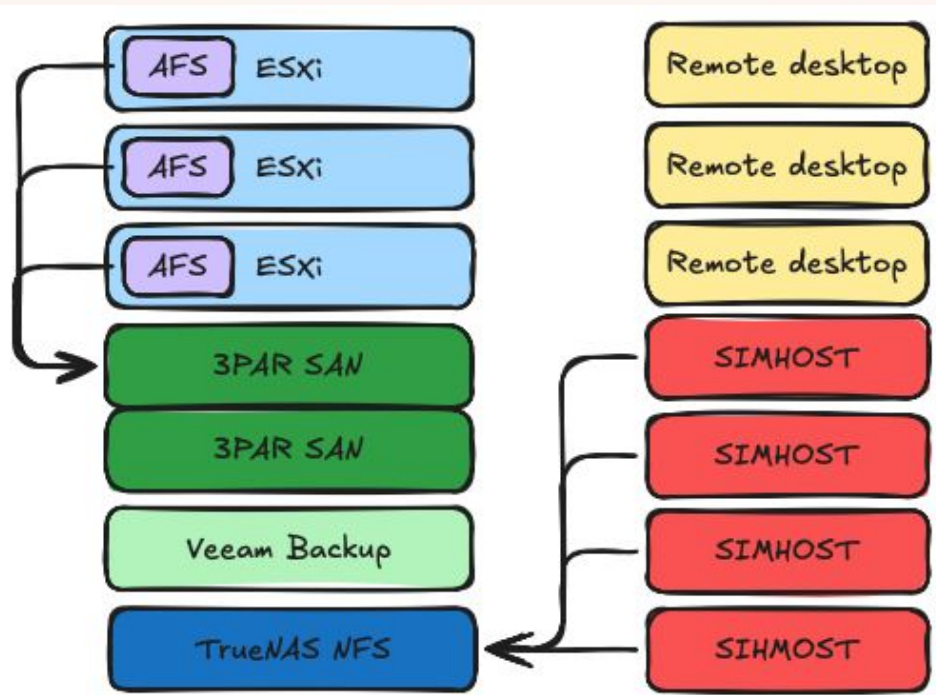
- IC / ASIC design company
- Extensive chip simulation for design verification
- Goal: eliminate design errors before fabrication
- Path from idea to tape-out is costly.
- ASICs don't run code → no post-silicon fixes possible
- We're in medical, automotive, aerospace.
- Errors are not acceptable
- Fixes \approx 6 months delay
- First-Time-Right and automation is key for us.



ICsense IN A NUTSHELL

Understanding IC design
in less than 60 seconds

Starting point: VMware + SAN



- 3 ESXi servers, 512GB RAM, FC HBA, ~120VMs
- OpenAFS VMs for network file storage
- 3PAR SAN, 36 HDD 60TB net capacity (1200IOPS/20MBps throughput avg)
- Veeam Backup server (D2D2T)
- TrueNAS NFS for simulation data (gentle IO workload as well)
- ~10 Remote desktop machines
- ~20 simulation servers

How does that translate into “workload”?

- We want simulations to finish ASAP:
 - Many and fast cores (bare metal)
 - Lots of RAM
 - Lots of disk space
- Relatively frequent hardware refreshes for simulation workloads leave us with many still very capable decommissioned servers

How does that translate into storage workload?

7-day average measurement on our 3PAR SAN:

```
15:45:09 09/14/2023 r/w I/O per second KBytes per sec Svt ms IOSz KB
      client: 24 MiB/s rd, 12 MiB/s wr, 3.04k op/s rd, 497 op/s wr Domain
Cur Avg  Max  Cur  Avg  Max  Cur  Avg  Cur  Avg Qlen
-----
      -   t 1059 626  1081 17731 11545 64977 0.65 0.55 16.7 18.4  1
-----
      1   t 1059 626          17731 11545          0.65 0.55 16.7 18.4  1
Press the enter key to stop...
```

For reference today:

```
io:
      client: 24 MiB/s rd, 12 MiB/s wr, 3.04k op/s rd, 497 op/s wr
```

Why change the “winning” team?

A faint, light grey silhouette of the London skyline is visible in the background. It includes the Tower Bridge, the Gherkin, Big Ben, and St. Paul's Cathedral. A large, semi-transparent version of the red circular icon from the logo is overlaid on the skyline, centered behind the main text.

2022: “VMware by Broadcom”

- Broadcom announced its intentions to acquire VMware
- We weren't comfortable with that idea
- Renewed for 3 years at a ~20% price increase to buy us some time
- Look for alternatives.

An aging SAN

- Performance: more than sufficient
- 80% capacity
- Almost EOL: no upgrades/expansions possible
- Entry level SAN is ~affordable, but lacks features (Veeam integration)
- Mid range SAN: expensive and overkill

- OpenAFS does scale, but not in a way that suits our workload.
- Linux, macOS and Windows all require custom kernel modules/extensions, drivers.
- Small install base, no “vivid community”
- Lots of problems on Linux and macOS:
 - Shell freezes (Kernel ring buffer: `waiting for busy volume`)
 - Memory leaks (resulting in crashing simulation servers)
 - macOS requires a kernel extension

NFS simulation data

File sharing NFS:

- One major flaw: **SPOF**



Enter Proxmox

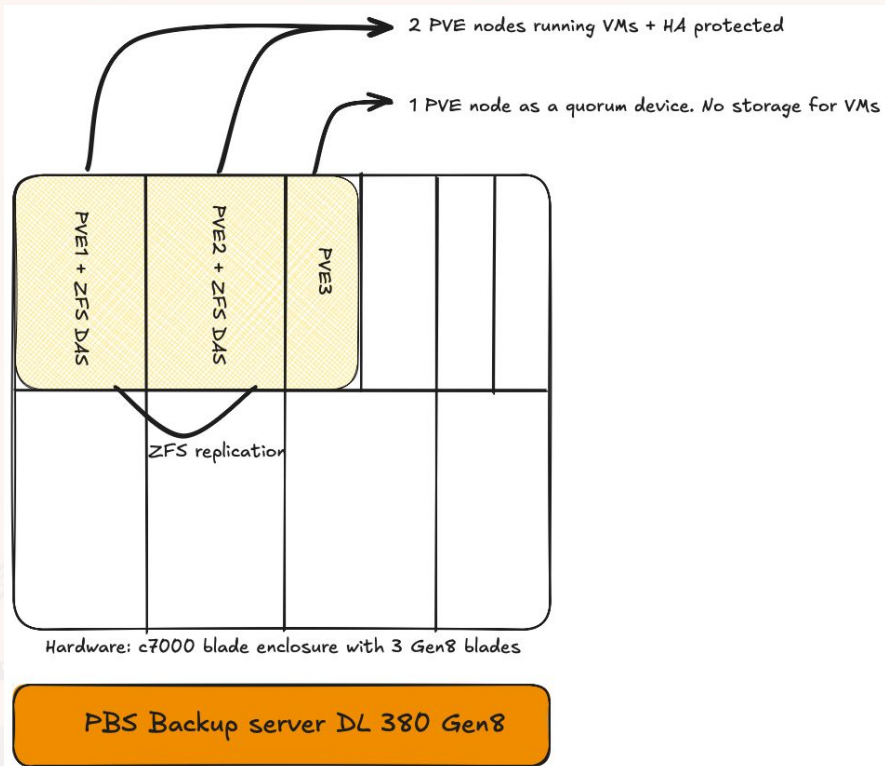
A faint, light gray silhouette of the London skyline is visible in the background. It includes the Tower Bridge, the Shard, the Gherkin, Big Ben, and St. Paul's Cathedral. A large, semi-transparent red Ceph logo is overlaid on the skyline, positioned behind the main title.

Enter Proxmox

Proxmox was high on our list from the start

- It fits our needs: run VMs, HA, backups, snapshots.
- Free You can test all you want.
- Runs on decommissioned simulation servers.
- Supports Ceph as an enterprise grade storage back-end 👍

2023: POC Proxmox cluster



Basic setup

- Decommissioned HPE c7000 blade enclosure
- BL460c gen8 blade (Xeon E5 Gen2)
- 2 nodes with replicated DAS ZFS pools
- 1 node as a “quorum device”, no storage
- HPE DL380 gen8 as a Proxmox Backup server (PBS)
- Ran some non-critical VMs.
- It just ... worked. Nothing to write home about during 2023-2024

Let's talk Ceph!



Why Ceph is a good fit (or a better fit than ZFS)

- Ceph is Free and Open Source, no vendor lock-in, full control
- It runs on our decommissioned simulation servers as well, no HCL 👍
- Ceph= “true shared storage”, ZFS can offer “pseudo shared storage”.
- Ceph’s failure handling is elegant
- Ceph grows with our company.
- You can shape it to your needs.
- CephFS might also replace OpenAFS and NFS in our infrastructure

- I dove into the docs, so much going on, couldn't wrap my head around "Ceph".
- We invested in a 3-day Ceph training (highly recommended).
- The training really helped to make sense of the documentation.
- One insight gathered is to run an external Ceph cluster.
 - Avoid potential pitfalls of HCI (compute resource contention)
 - We want better knowledge of Ceph in case things go wrong.
- Time to bootstrap a POC Ceph cluster! 🎉



ceph days LONDON 2026

Ceph PoC stage 1



CPU utilization

CPU	Users	Free	Mem%	Idle
1	0.0	5.0	0.0	95.0
2	0.0	0.0	0.0	98.0
3	1.0	2.0	0.0	97.0
4	1.0	5.7	0.0	92.3
5	1.0	1.0	0.0	98.0
6	20.4	2.0	0.0	77.6
7	2.0	1.0	0.0	97.0
8	0.0	1.0	19.8	79.2
9	1.0	3.0	0.0	96.0
10	0.0	1.0	0.0	99.0
11	0.0	1.0	51.5	47.5
12	1.0	2.0	0.0	97.0
13	0.0	5.0	0.0	100.0
14	0.0	1.0	12.1	86.9
15	1.0	3.0	0.0	98.0
16	1.0	1.0	36.6	59.4
17	0.0	0.0	0.0	100.0
Avg	1.9	1.8	7.6	86.7

Disk I/O /proc/diskstats mostly in KB/s Warning:contains duplicate

DiskName	Busy	Read	Write	IO
sdc	0%	0.0	0.0	0.0
sdc1	0%	0.0	0.0	0.0
sdc2	0%	0.0	0.0	0.0
sdc3	0%	0.0	0.0	0.0
sdd	54%	8.0	19.3	27.3
sdc	73%	0.0	25.0	25.0
sda	0%	0.0	37.4	37.4
dm-0	87%	0.0	28.1	28.1
dm-1	73%	0.0	32.0	32.0
dm-2	54%	0.0	76.0	76.0
dm-3	0%	0.0	0.0	0.0
dm-4	0%	0.0	0.0	0.0
dm-5	0%	0.0	0.0	0.0
dm-6	0%	0.0	0.0	0.0
dm-7	0%	0.0	0.0	0.0
Totals		Read-MB/s=0.0	Write-MB/s=163.1	Transfers/sec=2037.2



CPU utilization

CPU	Users	Free	Mem%	Idle
1	0.0	0.0	0.0	100.0
2	1.0	0.0	0.0	99.0
3	0.0	3.0	0.0	97.0
4	0.0	1.0	0.0	99.0
5	1.0	1.0	0.0	98.0
6	0.0	1.0	0.0	99.0
7	1.0	1.0	0.0	98.0
8	2.0	1.0	1.3	93.1
9	1.0	1.0	0.0	98.0
10	1.0	1.0	0.0	98.0
11	1.0	2.1	0.0	96.9
12	0.0	0.0	0.0	100.0
13	0.0	4.0	0.0	96.0
14	0.0	0.0	0.0	100.0
15	1.0	1.0	5.2	93.1
16	0.0	1.0	0.0	99.0
17	0.0	0.0	0.0	100.0
Avg	0.6	1.1	0.8	97.8

Disk I/O /proc/diskstats mostly in KB/s Warning:contains duplicate

DiskName	Busy	Read	Write	IO
sda	0%	0.0	0.0	0.0
sda1	0%	0.0	0.0	0.0
sda2	0%	0.0	0.0	0.0
sda3	0%	0.0	0.0	0.0
sdb	3%	0.0	12.0	12.0
sdc	41%	0.0	16.2	16.2
sdd	47%	0.0	15.0	15.0
dm-0	3%	0.0	12.0	12.0
dm-1	43%	0.0	24.0	24.0
dm-2	43%	0.0	16.0	16.0
dm-3	0%	0.0	0.0	0.0
dm-4	0%	0.0	0.0	0.0
dm-5	0%	0.0	0.0	0.0
dm-6	0%	0.0	0.0	0.0
dm-7	0%	0.0	0.0	0.0
Totals		Read-MB/s=0.0	Write-MB/s=96.3	Transfers/sec=1286.3

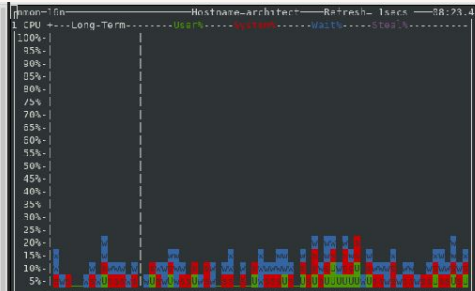


CPU utilization

CPU	Users	Free	Mem%	Idle
1	1.0	0.0	0.0	99.0
2	3.0	1.0	0.0	98.0
3	0.0	0.0	0.0	100.0
4	0.0	3.0	0.0	97.0
5	1.0	1.0	76.2	21.8
6	0.0	2.0	0.0	98.0
7	0.0	0.0	0.0	100.0
8	1.0	0.0	0.0	99.0
9	0.0	1.0	0.0	99.0
10	0.0	0.0	0.0	100.0
11	1.0	0.0	0.0	99.0
12	0.0	2.0	0.0	98.0
13	1.0	2.0	0.0	97.0
14	1.0	1.0	65.3	33.0
15	0.0	1.0	0.0	99.0
16	1.0	1.0	0.0	98.0
17	0.0	0.0	0.0	100.0
Avg	0.8	1.0	9.8	89.4

Disk I/O /proc/diskstats mostly in KB/s Warning:contains duplicate

DiskName	Busy	Read	Write	IO
sda	0%	0.0	0.0	0.0
sda1	0%	0.0	0.0	0.0
sda2	0%	0.0	0.0	0.0
sda3	0%	0.0	0.0	0.0
sdb	3%	0.0	16.0	16.0
sdd	100%	0.0	35.2	35.2
sdc	63%	0.0	20.5	20.5
dm-0	83%	0.0	36.0	36.0
dm-1	4%	0.0	16.0	16.0
dm-2	100%	0.0	32.0	32.0
dm-3	0%	0.0	0.0	0.0
dm-4	0%	0.0	0.0	0.0
dm-5	0%	0.0	0.0	0.0
dm-6	0%	0.0	0.0	0.0
dm-7	0%	0.0	0.0	0.0
Totals		Read-MB/s=0.0	Write-MB/s=166.5	Transfers/sec=2071.4



CPU utilization

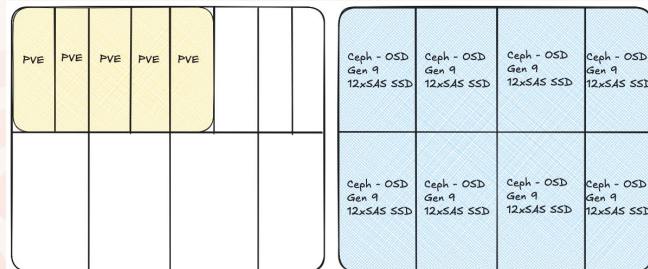
CPU	Users	Free	Mem%	Idle
1	1.0	11.2	0.0	87.8
2	4.0	8.0	0.0	88.0
3	2.0	6.1	0.0	91.8
4	2.1	4.1	0.0	93.0
5	4.0	4.0	0.0	92.1
6	3.0	0.0	2.0	24.9
7	0.0	0.0	0.0	100.0
8	2.0	1.0	0.0	97.0
9	2.0	5.1	0.0	92.9
10	3.0	0.0	0.0	100.0
11	0.0	2.0	0.0	98.0
12	2.0	0.1	2.0	96.9
13	4.1	2.1	69.0	25.8
14	3.0	5.9	2.0	89.1
15	1.0	1.0	0.0	98.0
16	1.0	1.0	1.0	97.0
17	0.0	0.0	0.0	100.0
Avg	2.0	3.7	4.5	89.8

Disk I/O /proc/diskstats mostly in KB/s Warning:contains duplicate

DiskName	Busy	Read	Write	IO
sda	2%	0.0	0.0	0.0
sda1	3%	0.0	0.0	0.0
sda2	3%	0.0	0.0	0.0
sda3	1%	0.0	0.0	0.0
sda	74%	0.0	78.0	78.0
sdc	83%	0.0	32.0	32.0
sdb	44%	0.0	16.3	16.3
dm-0	88%	0.0	48.0	48.0
dm-1	74%	0.0	28.0	28.0
dm-2	44%	0.0	70.0	70.0
dm-3	3%	0.0	0.0	0.0
dm-4	2%	0.0	0.0	0.0
dm-5	3%	0.0	0.0	0.0
dm-6	3%	0.0	0.0	0.0
dm-7	0%	0.0	0.0	0.0
Totals		Read-MB/s=0.0	Write-MB/s=172.4	Transfers/sec=2206.9

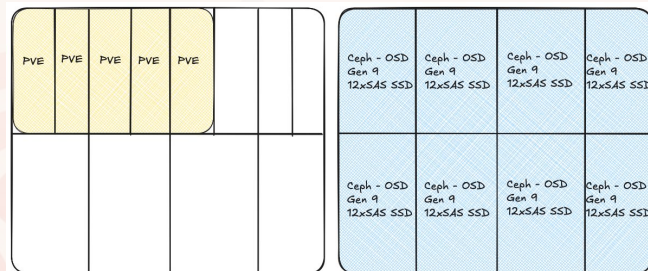
I scale up and out.

- Carefully started migrating more low impact production VMs to the new setup.
- PVE from 3 to 5 nodes
- Removed the ZFS data stores in PVE
- 8 Ceph nodes, each 12 x 3.84TiB SAS SSDs
- 256GiB RAM, 16 cores@3.2GHz per node.
- 96 OSDs, 368TiB raw capacity



Ceph has failure in mind. Our hardware/network stack had to match that:

- Upgrade from 10GbE to 20GbE
- Added a second NIC to each node, total of 4 “physical interfaces”
- 8 NICs, 2 Linux bonds for ceph-public and ceph-cluster
- CPU tuning got ping latencies from ~0.2ms to 0.04ms avg



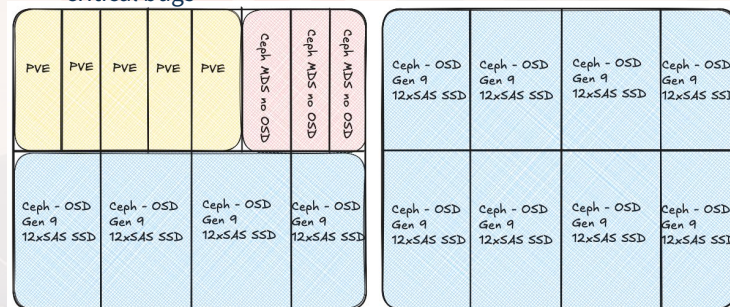
Main idea: I need many components to fail before production is affected.

File sharing OpenAFS:

- OpenAFS file share is a single large “directory”. My initial idea was to migrate to a similar large and monolithic CephFS filesystem.
- As my understanding of CephFS evolved, I shifted towards smaller, separate filesystems. (to contain blast radius)
- consequence: 2 pools/CephFS file system
- More pools > more PGs > more OSDs > need more block devices.
- “Intended” to add 48 extra SSDs
- 550TiB raw capacity
- Added 3 dedicated MDS nodes

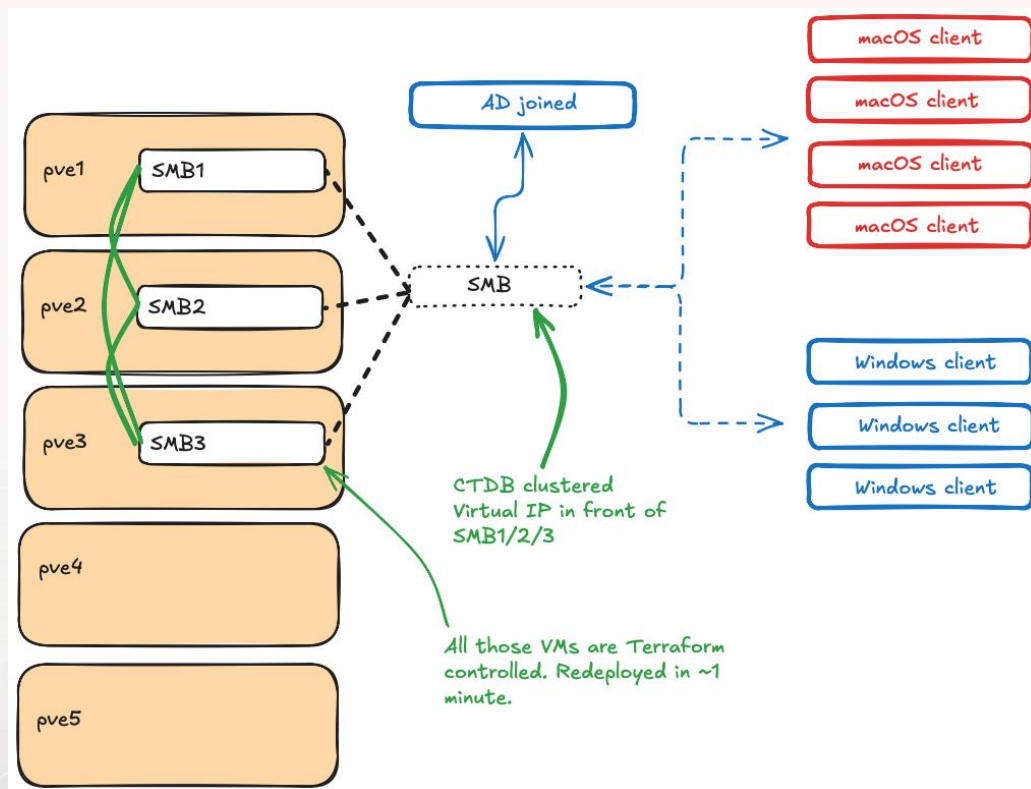
Added a support contract

- During cluster intake we were made aware of the `bluestore_elastic_shared_blobs` bug in Ceph Squid
- Mitigated by redeploying all OSDs
- Nevertheless, postponed cluster expansion to X-mas maintenance.
- Main key take-away: If you start out with Ceph: Look up “known critical bugs”



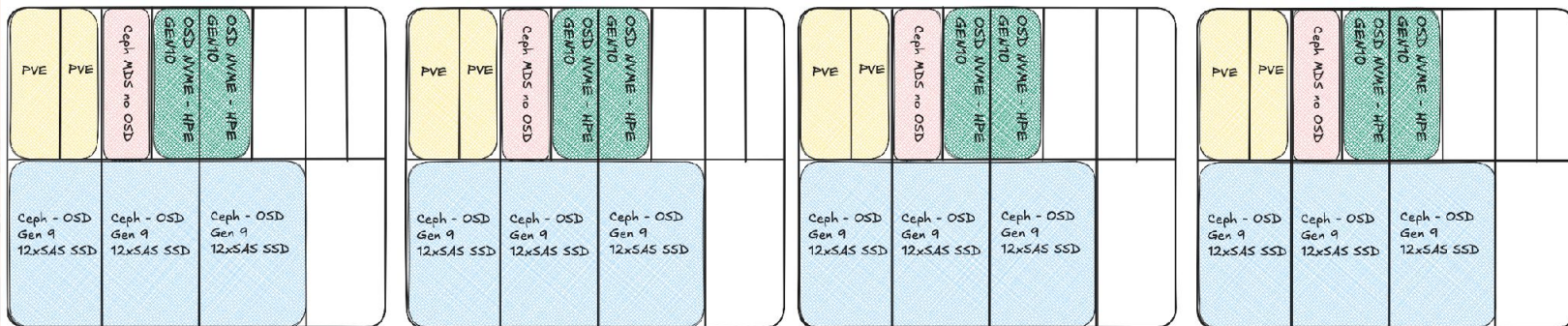
- macOS + Windows clients (native support)
- Samba `vfs_fruit` module is more than desirable for performance on macOS.
- I did not use the SMB manager module in Tentacle: it's brand new. We want more feedback from a larger install base.
- One EC-pool for archival data
- First user tests are promising!

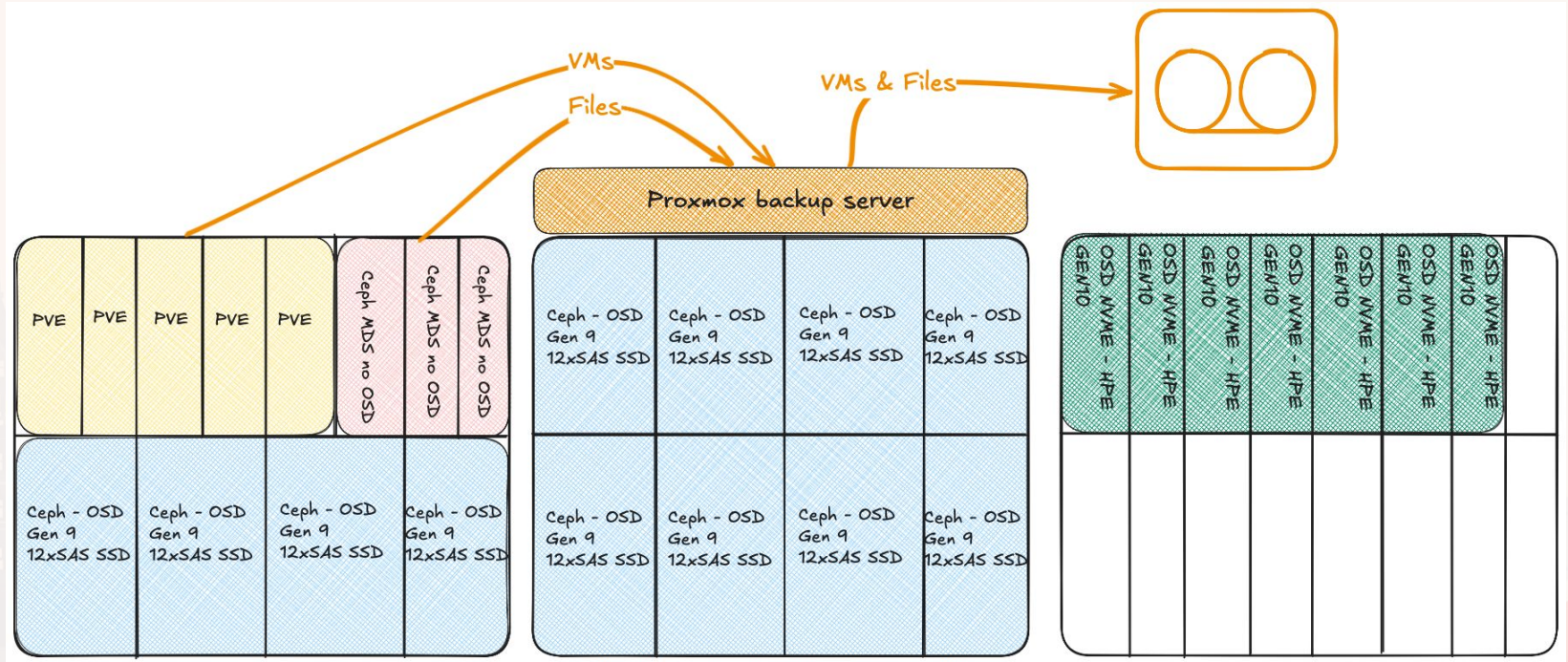
Q1 2026 – SMB setup



Scaling challenges: Failure domain & Layout

- Any enclosure crash is a recipe for disaster.
- To improve resilience, we want to add a 4th chassis, reshuffle hosts and bump the failure domain to “chassis”
- Challenges:
 - Injecting a new crushmap introduces a lot of IO





Any Mishaps?



Some noteworthy mishaps

- **ceph orch host maintenance enter** does not check if there's enough available RAM on other nodes to restart an MDS elsewhere. Result: had to remount affected share on all clients.
- Uncontrolled cluster shutdown
 - Root cause: Power outage
 - Result: next slide

```
root@persephone:~# ceph -s
cluster:
  id:          e9020818-2100-12a0-8aa2-9cdc71770100
  health: HEALTH_ERR
             unable to send alert email
             Invalid grafana certificate on host persephone: Invalid certificate key: [('PEM routines', '', 'no start line')]
             1 hosts fail cephadm check
             clock skew detected on mon.architect, mon.dujour, mon.seraph, mon.persephone
             2 osds(s) are not reachable

services:
  mon: 5 daemons, quorum apoc,architect,dujour,seraph,persephone (age 4h)
  mgr: architect.nrvjah(active, since 4h), standbys: niobe.yhkgpe, seraph.guwsy
  mds: 17/17 daemons up, 3 standby
  osd: 156 osds: 156 up (since 4h), 156 in (since 11d)
       flags noautoscale

data:
  volumes: 11/11 healthy
  pools:  24 pools, 8225 pgs
  objects: 327.05M objects, 103 TiB
  usage:   280 TiB used, 240 TiB / 520 TiB avail
  pgs:    8224 active+clean
          1   active+clean+scrubbing+deep

io:
  client:  0 B/s rd, 1.4 KiB/s wr, 0 op/s rd, 0 op/s wr
```

```
root@persephone:~# ceph health detail
...
[ERR] OSD_UNREACHABLE: 2 osds(s) are not reachable
      osd.53's public address is not in '192.168.11.0/24' subnet
      osd.86's public address is not in '192.168.11.0/24' subnet
```

Some final thoughts

What might make you think twice

- Ceph is not a set-and-forget
- The learning curve is real
- More in-house responsibility, rather than relying on your MSP.

Things we like!

- Disk space galore! 520TiB
- Future expansions are cost effective
- Ceph is flexible, we can shape/tune it to our needs
- We're free from vendor lock-in!
- Investment shifts from hardware people, knowledge & skill.
- Active and helpful community: Slack, mailing lists, Reddit, Online AMA,