



# SURFin' the Ceph wave:

How the Dutch NREN leverages Ceph.

Jean-Marie de Boer, SURF



# About me

- Jean-Marie de Boer, from Amsterdam, the Netherlands
- In the internet industry since 1995, always in a technical role
- Mainly focused on web-technologies
- Started as PHP developer/linux sysadmin
- Shifted to operations/HA scale-out solutions
- Joined SURF 9 years ago

# About SURF

- SURF is the Dutch NREN
- Not-for-profit co-operative working for members:
  - Universities
  - Higher education
  - Research institutes
- Oldest entities date back to 1971
- 800 employees, wide portfolio of services
- Many on-prem solutions in storage and compute

# Ceph usecases

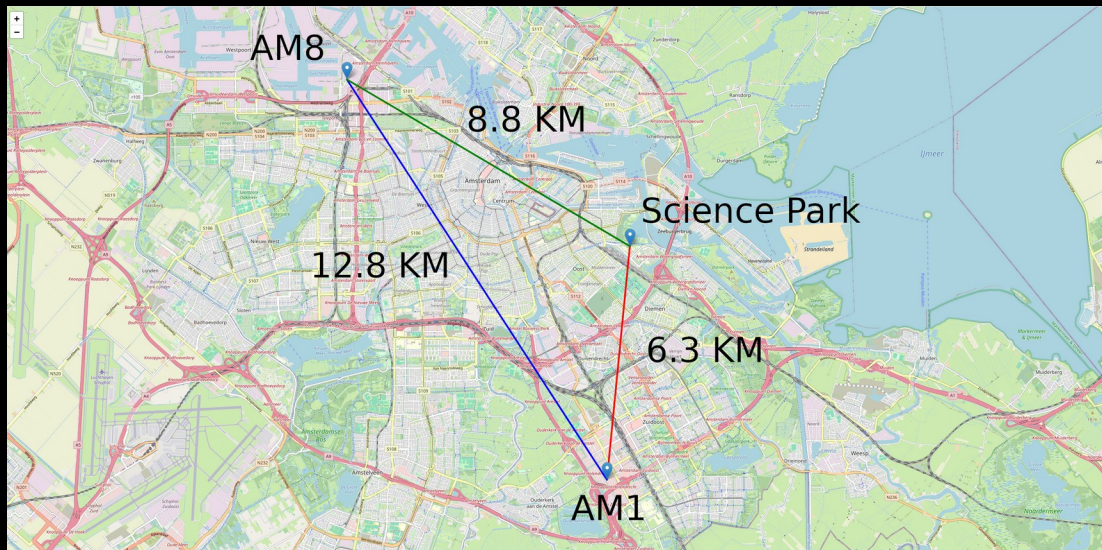
- First Ceph solutions implemented some 10 years ago
- Used for OpenStack/OpenNebula compute clusters
- RBD, replacing GlusterFS
- Initially < 1PB net storage, positive results, growth
- Since 2022 we converged on Ceph even more
- CephFS now POSIX storage for (Nextcloud) sync & share
- RGW replacing OpenStack Swift

# Current scale

- 5 production clusters:
  - RBD, 65 OSD nodes, 966 OSD's, 9PB raw
  - RBD, 64 OSD nodes, 857 OSD's, 10PB raw
  - CephFS/RGW, 66 OSD nodes, 1038 OSD's, 13PB raw
  - CephFS, 36 OSD nodes, 648 OSD's, 8.3PB raw
  - CephFS (ec), 24 OSD nodes, 288 OSD's, 3.5PB raw
- That's 3797 OSD's, 44PB raw
- In addition, a few test setups in virtual machines
- Versions range from Pacific to Squid

# Pragmatic geo-distribution

- EFSS and Object Storage are located in 3 datacenters, all active
- High bandwidth, low latency connections
- Data should be available when losing a datacenter



# But, it works!

```
[ceph: root@ceph-mon-b-01 /]# ceph -s
cluster:
  id:      1038152c-35d3-11ee-bad7-123456789012
  health: HEALTH_WARN
          12 hosts are in maintenance mode
          1/5 mons down, quorum ceph-mon-b-01,ceph-mon-b-02,ceph-mon-c-01,ceph-mon-c-02
          noout,nobackfill,norebalance,norecover flag(s) set
          1 datacenter (148 osds) down
          148 osds down
          10 OSDs or CRUSH {nodes, device-classes} have {NOUP,NODOWN,NOIN,NOOUT} flags set
          11 hosts (148 osds) down
          Degraded data redundancy: 606844769/1820534307 objects degraded (33.333%), 11501 pgs degraded, 13568 pgs undersized
```

- Of course, we are taking a risk with availability
- Long-term outages need to be considered

# Lessons learned

- Ceph is good! (but maybe don't mix workloads)
- It likes to scale
- Requires a stable network (network partitions do self-heal)
- Monitor your hardware, especially 'sick' disks
- Be patient!
- Look far ahead (esp. monitor fill rate & plan your capacity)
- If it's mission-critical, consider a consultancy contract (we have one)



# Future plans

- Scaling out the current CephFS/RBD solutions
- Object Storage seen as a growing demand
  - IAM policies
  - Replacing keystone auth
  - Implement UI/dashboard functionality
- NVMe-only cluster
- AI workloads
- National supercomputer
- European federated Object Storage

# Thank you!

(Any questions?)

