



redhat.

# Distributed File Storage in Multi-Tenant Clouds using CephFS

Openstack Vancouver 2018 May 23

**Patrick Donnelly**  
CephFS Engineer  
Red Hat, Inc.

Tom Barron  
Manila Engineer  
Red Hat, Inc.

Ramana Raja  
CephFS Engineer  
Red Hat, Inc.



# How do we solve this?

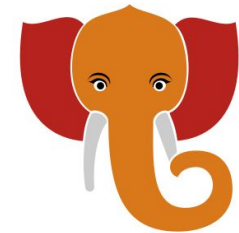


**MANILA**

*an OpenStack Community Project*



**ceph**



**GANESHA**

# Ceph Higher Level Components

OBJECT



**RGW**

S3 and Swift compatible object storage with object versioning, multi-site federation, and replication

BLOCK



**RBD**

A virtual block device with snapshots, copy-on-write clones, and multi-site replication

FILE



**CEPHFS**

A distributed POSIX file system with coherent caches and snapshots on any directory

**LIBRADOS**

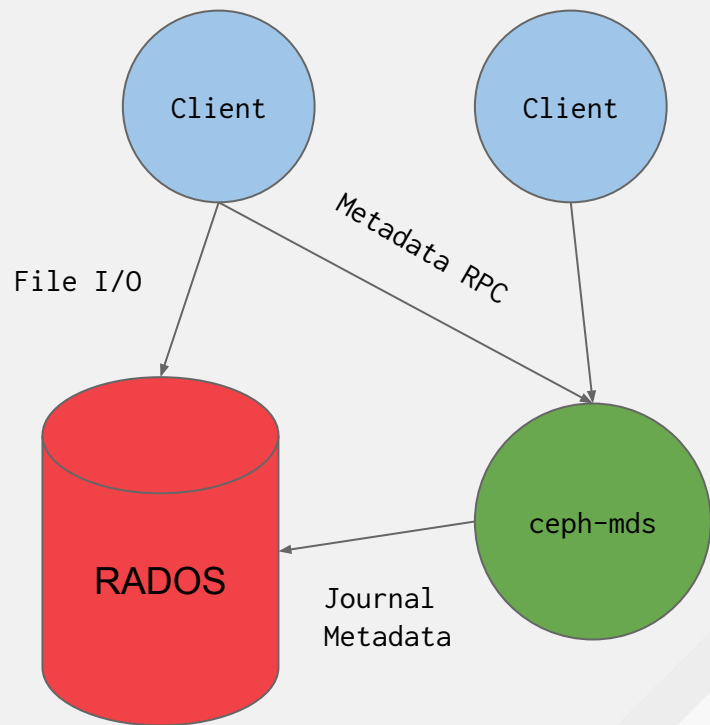
A library allowing apps to direct access RADOS (C, C++, Java, Python, Ruby, PHP)

**RADOS**

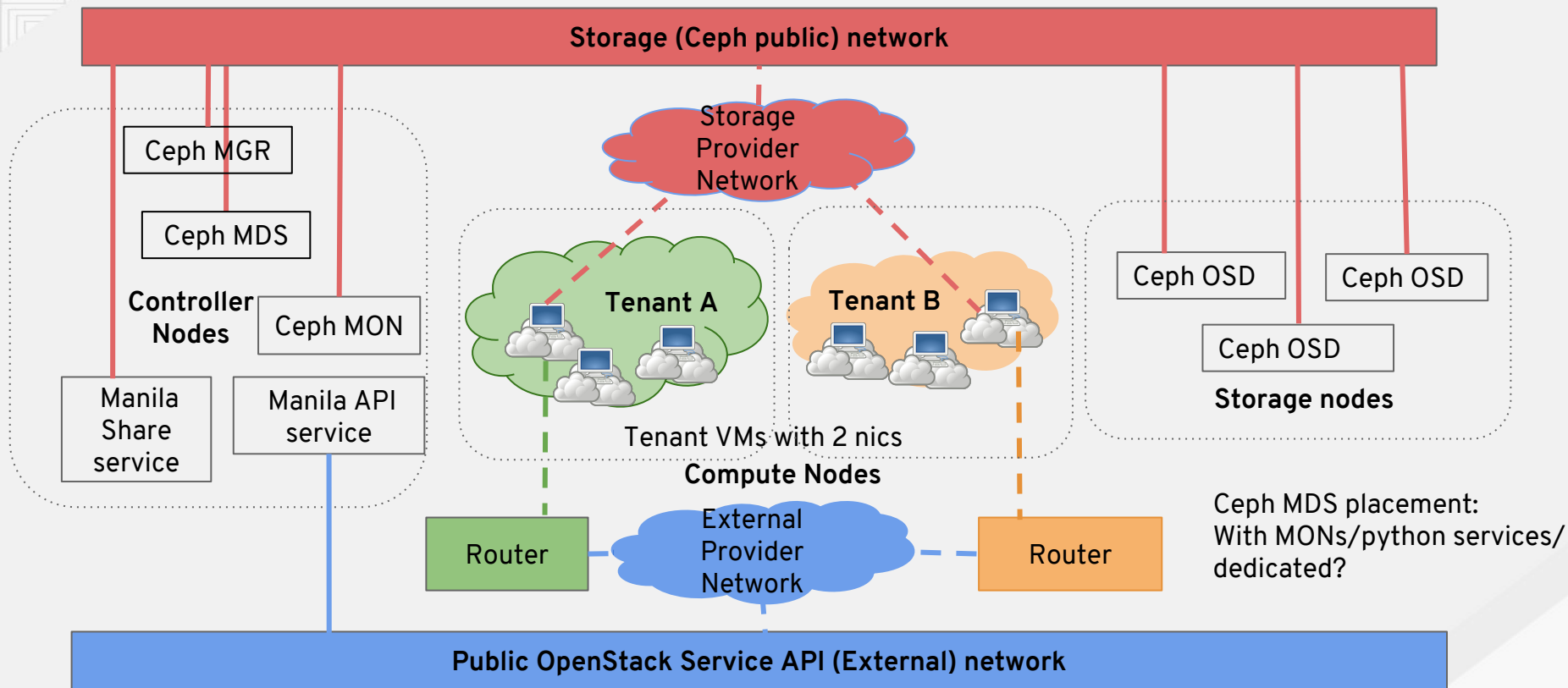
A software-based, reliable, autonomic, distributed object store comprised of self-healing, self-managing, intelligent storage nodes (OSDs) and lightweight monitors (Mons)

# What is CephFS?

- CephFS is a POSIX-compatible distributed file system!
- 3 moving parts: the MDS(s), the clients, and RADOS.
- Files, directories, and other metadata are stored in RADOS. The MDS has no local state.
- Mount using:
  - FUSE: `ceph-fuse ...`
  - Kernel: `mount -t ceph ...`
- Coherent caching across clients. MDS issues inode capabilities which enforce synchronous or buffered writes.
- Clients access data directly via RADOS.



# CephFS native driver deployment

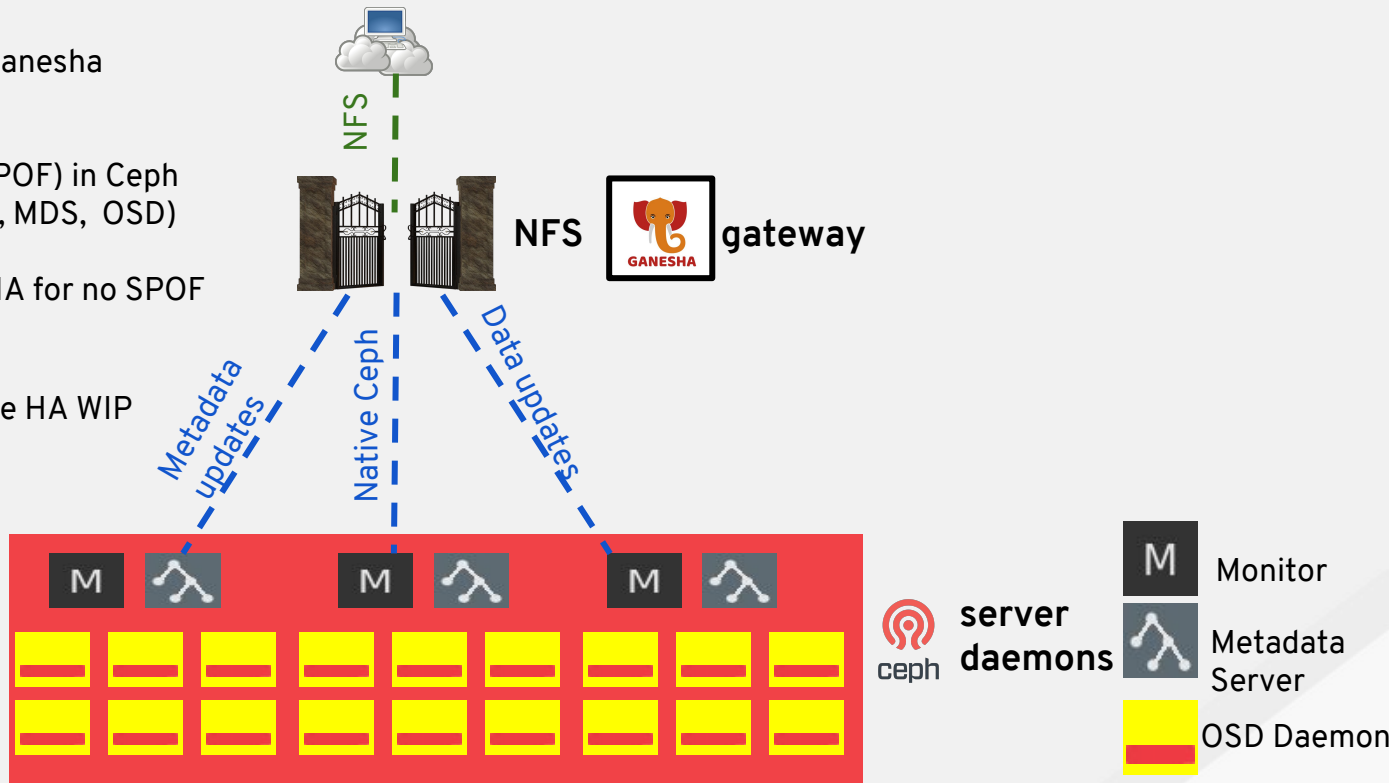


[https://docs.openstack.org/manila/ocata/devref/cephfs\\_native\\_driver.html](https://docs.openstack.org/manila/ocata/devref/cephfs_native_driver.html)

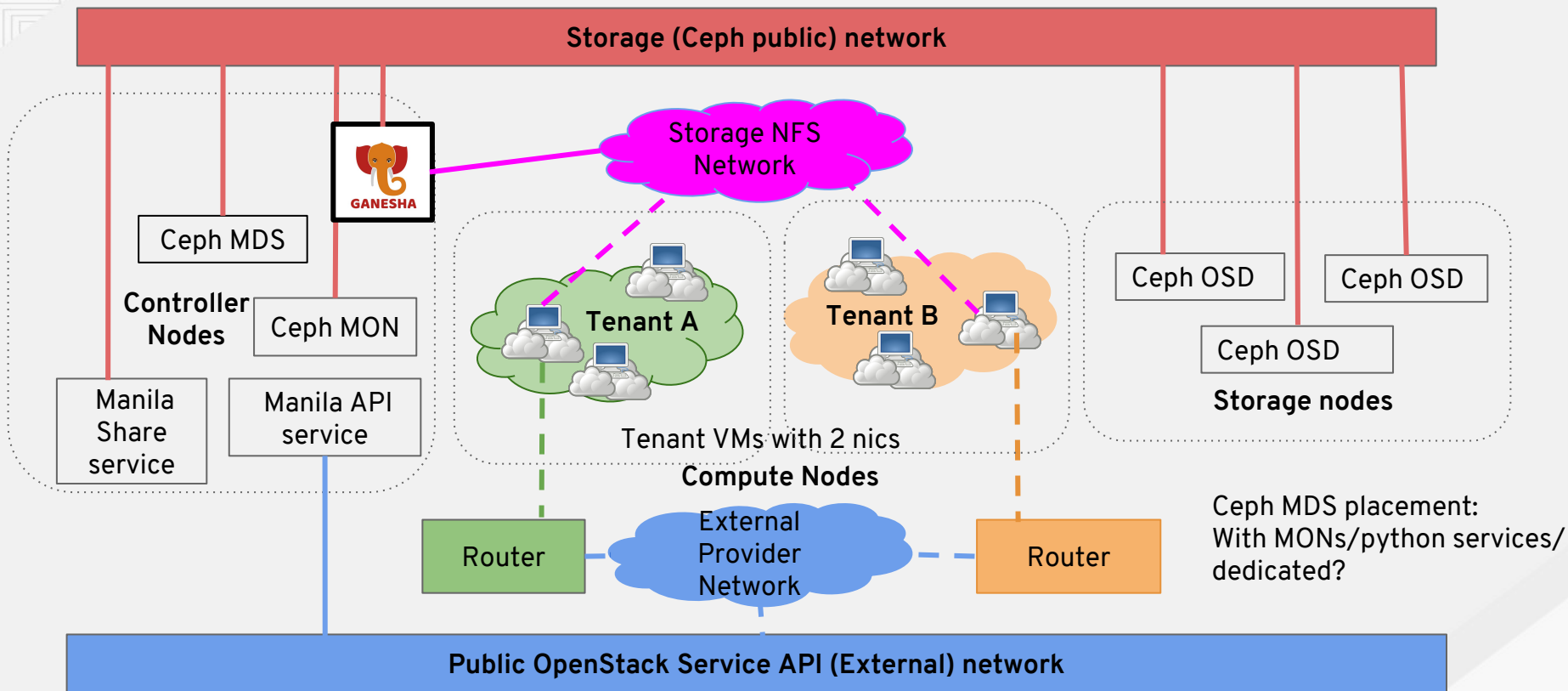
# CephFS NFS driver (in data plane)

OpenStack client/Nova VM

- Clients connected to NFS-Ganesha gateway. Better security.
- No single point of failure (SPOF) in Ceph storage cluster (HA of MON, MDS, OSD)
- NFS-Ganesha needs to be HA for no SPOF in data plane.
- NFS-Ganesha active/passive HA WIP (Pacemaker/Corosync)

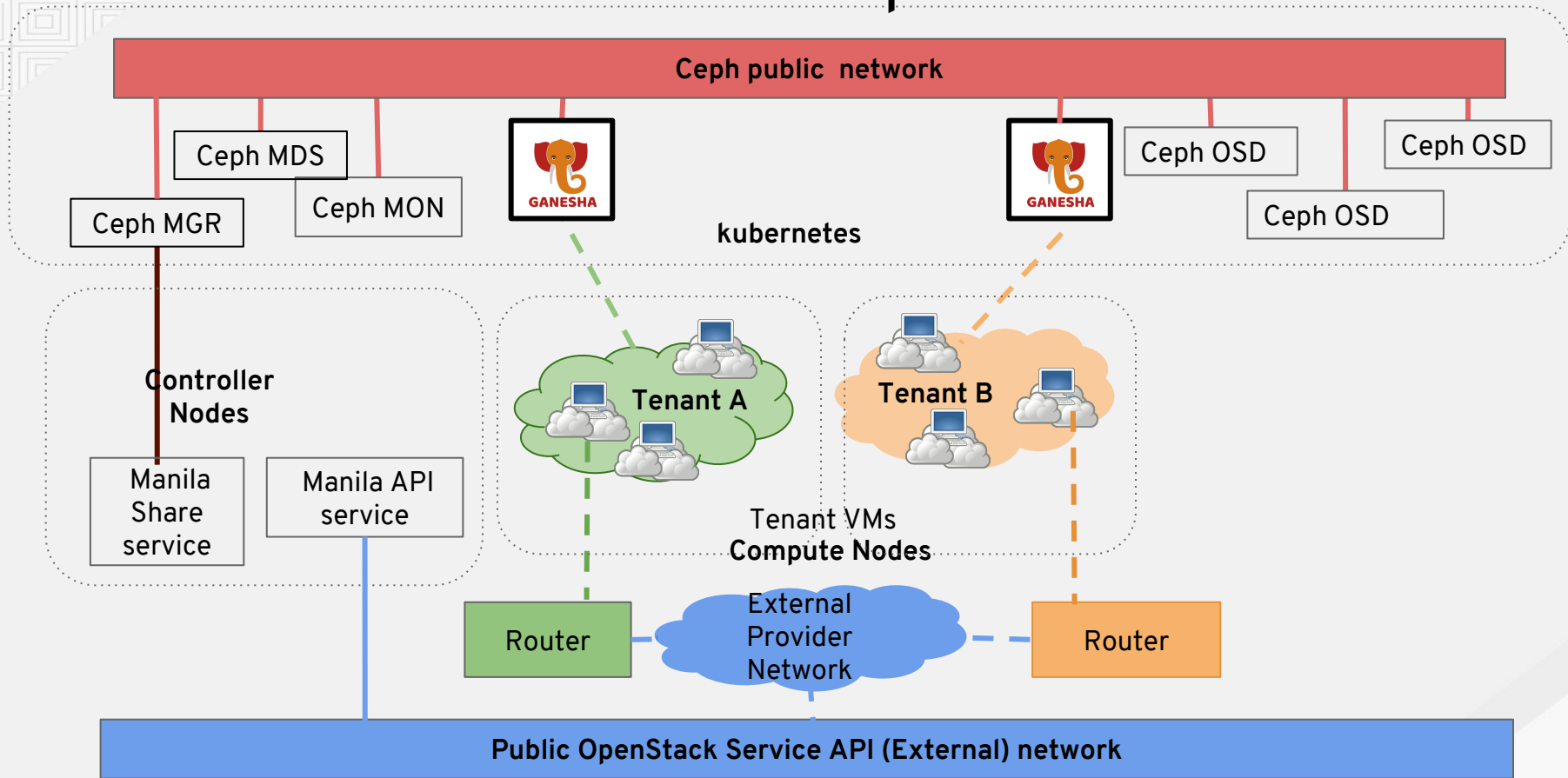


# CephFS NFS driver deployment

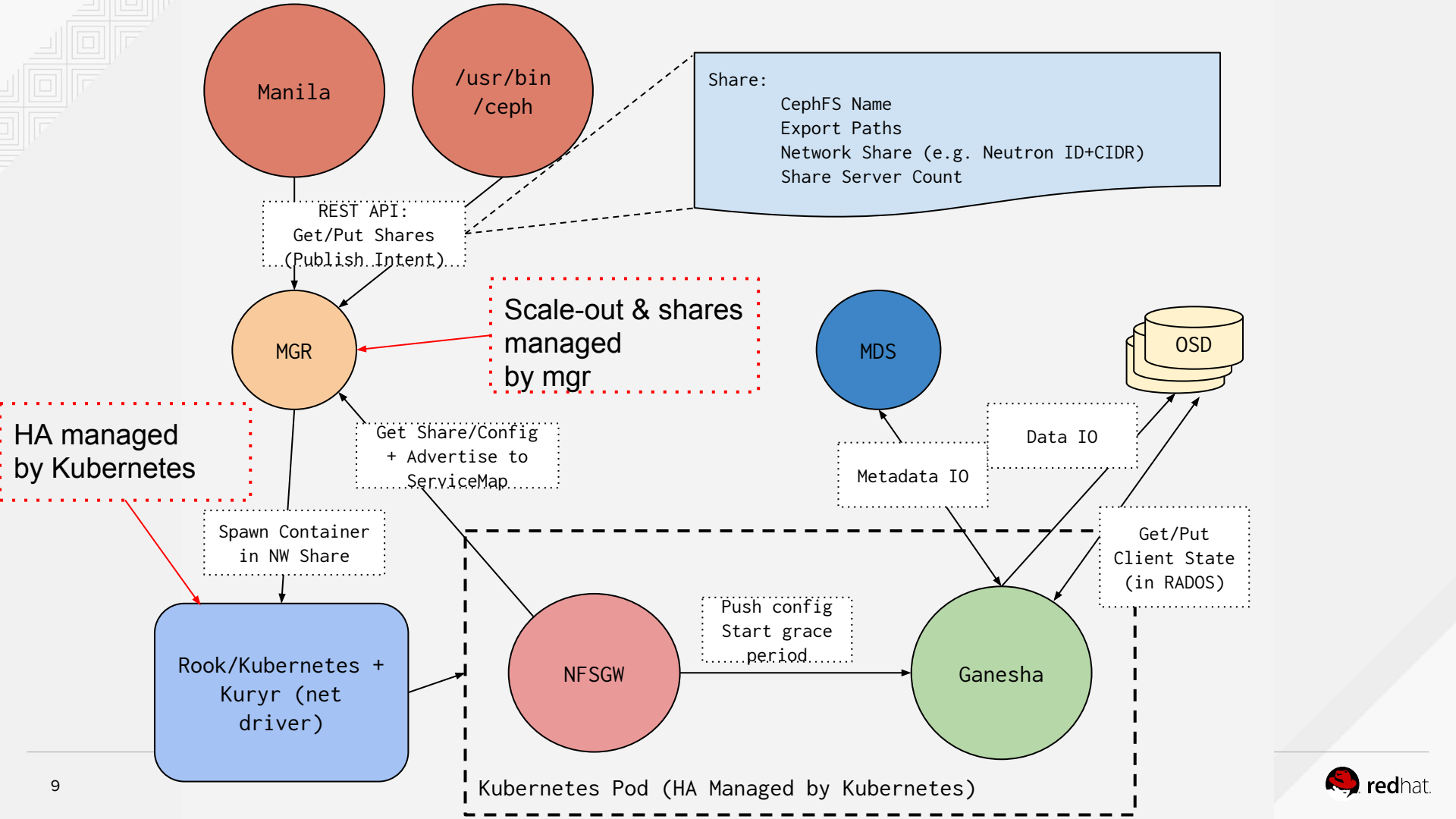


[https://docs.openstack.org/manila/queens/admin/cephfs\\_driver.html](https://docs.openstack.org/manila/queens/admin/cephfs_driver.html)

# Future: Ganesha per Tenant





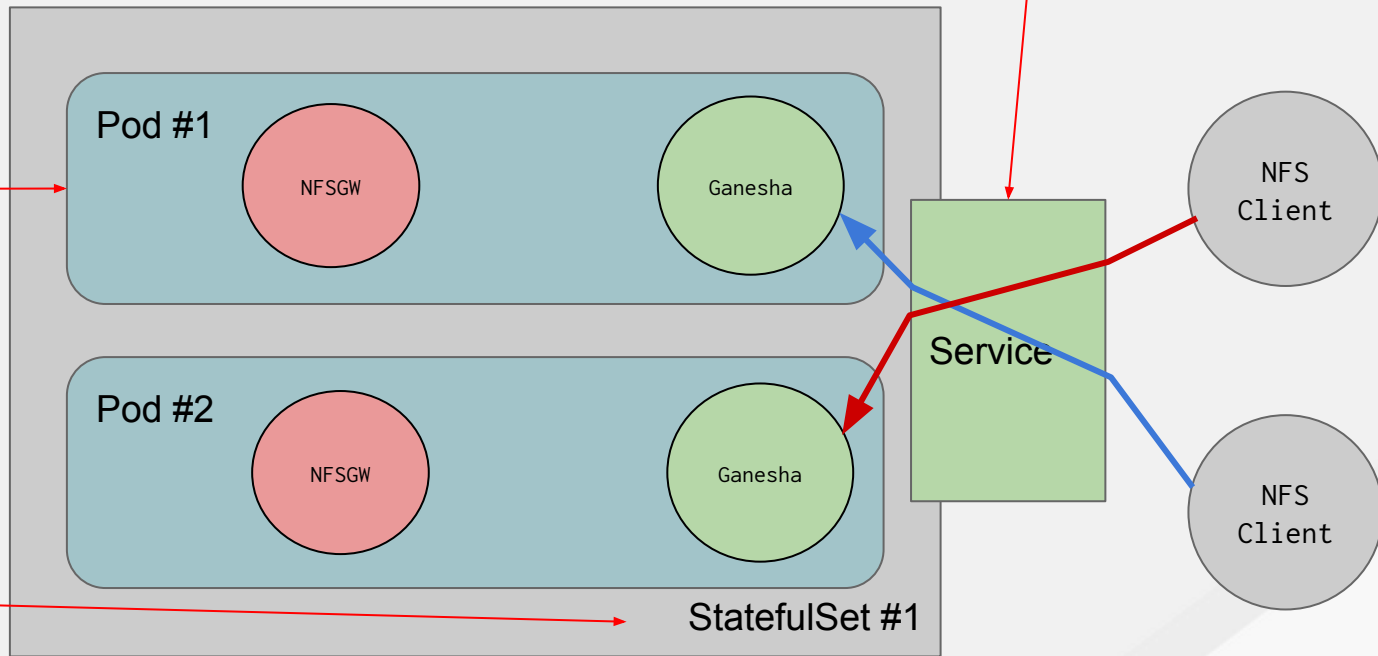


# Managing Scale-Out: “CephFSShareNamespace”

One IP to access NFS  
Cluster for tenant; only via  
tenant network

# of pods equal to  
scale-out for  
“CephFSShareNames  
pace”; dynamically  
grow/shrink??

Stable network  
identifiers for pods!



# Thanks!

Patrick Donnelly  
**`pdonnell@redhat.com`**

Thanks to the CephFS team: John Spray, Greg Farnum, Zheng Yan, Ramana Raja, Doug Fuller, Jeff Layton, and Brett Niver.

Homepage: <http://ceph.com/>

Mailing lists/IRC: <http://ceph.com/IRC/>