

Introduction to CEPH

CEPH- What?

Distributed storage platform which implements object storage without a single point of failure, scalable to exabyte level and freely available

CEPH - Why?

Traditional storage systems are not equipped for the amount of data generated nowadays . Software defined storage is the best alternative.

CEPH - What makes Ceph unique?

Design Principles :

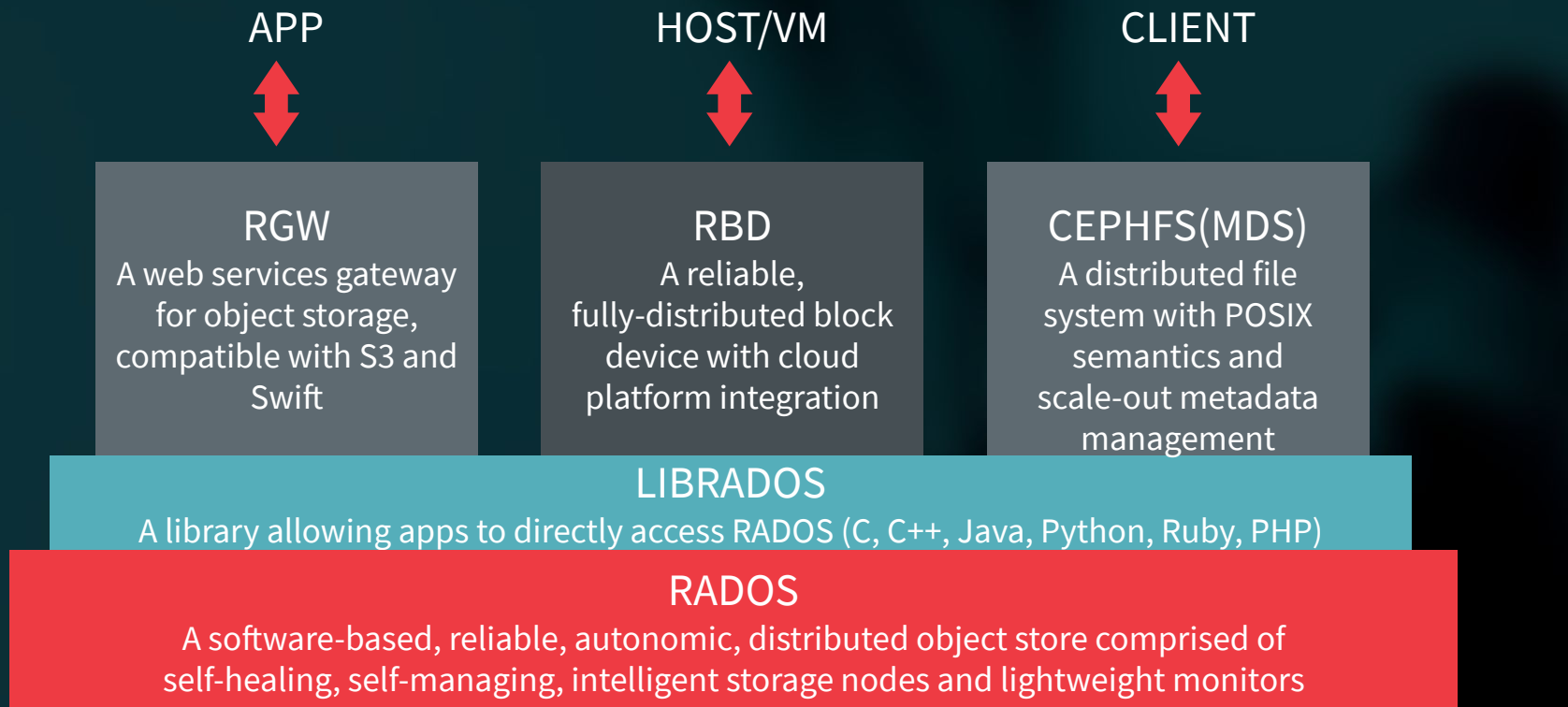
- Scalable
- No Single point of failure
- Software based
- Self Managing

Philosophy :

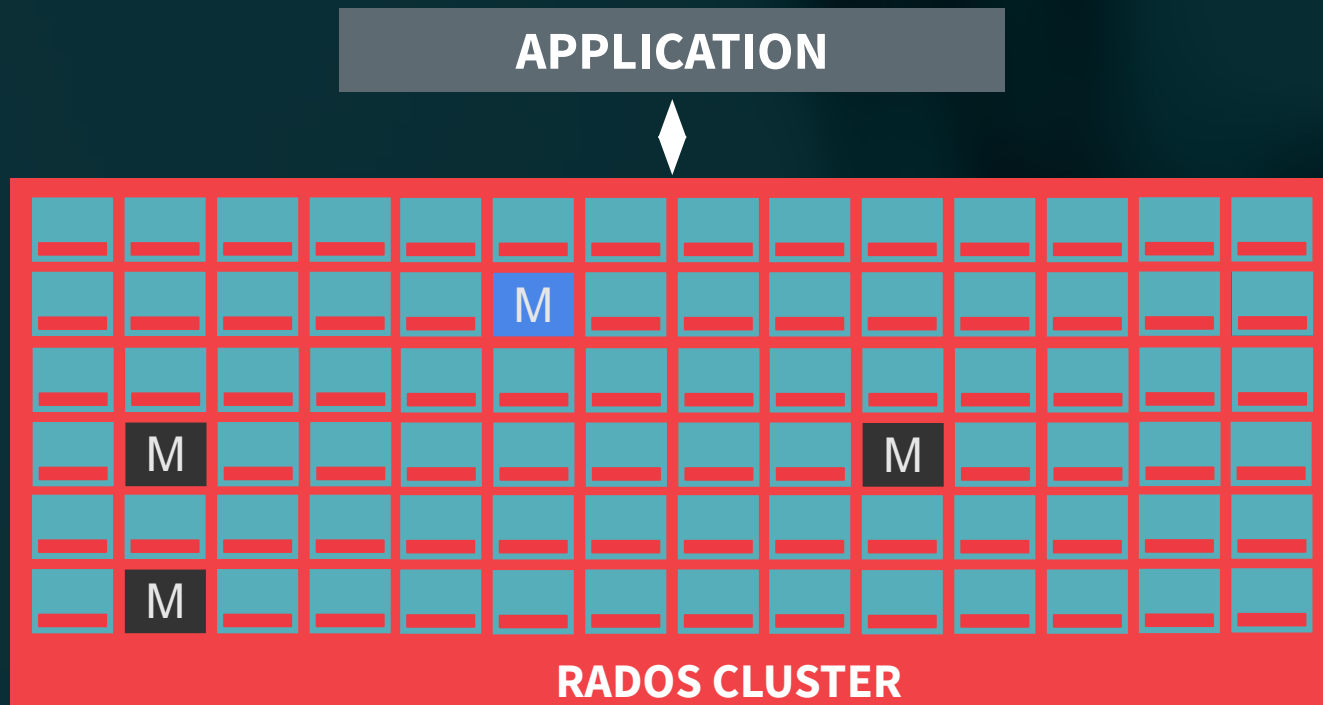
- Open source
- Community based

CEPH ARCHITECTURE

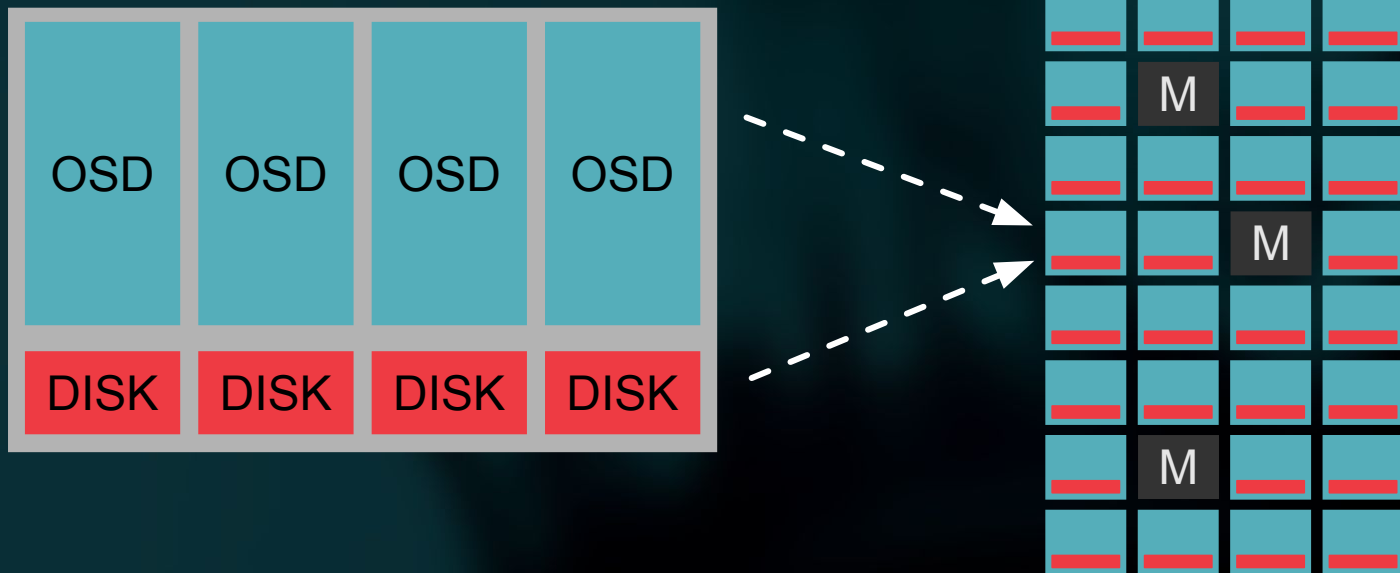
ARCHITECTURAL COMPONENTS



RADOS CLUSTER



OBJECT STORAGE DAEMONS



RADOS COMPONENTS



OSDs:

- 10s to 10000s in a cluster
- One per disk (or one per SSD, RAID group...)
- Serve stored objects to clients
- Intelligently peer for replication & recovery



Monitors:

- Maintain cluster membership and state
- Keeps the state of all MAPS.
- Small, odd number
- Do not serve data

Ceph Manager



Managers:

- Tightly coupled with Monitor
- Manage cluster “logistic” PG and maps
- Provide additional monitoring and interfaces to external monitoring and management systems
- Pluggable python interface to develop modules
- Some modules:
 - Balancer
 - Dashboard
 - RESTful
 - Prometheus

Responsibilities of RADOS Cluster :

- Write and read data
- Ensure durability by replicating or erasure coding data
- Monitor and report on cluster health—also called 'heartbeating'
- Redistribute data dynamically—also called 'backfilling'
- Ensure data integrity
- Recover from failures

Some Concepts of RADOS Cluster :

- Pools
 - Replicated
 - EC
- PG's
- Authentication
- CRUSH

OBJECT PLACEMENT WITH CRUSH

Client Communication with Ceph Storage Cluster

- Ceph Configuration File
- Pool Name
- User name and the path to the secret key

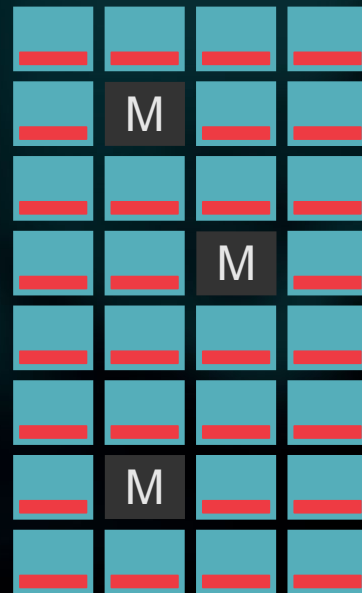
WHERE DO OBJECTS LIVE?

APPLICATION

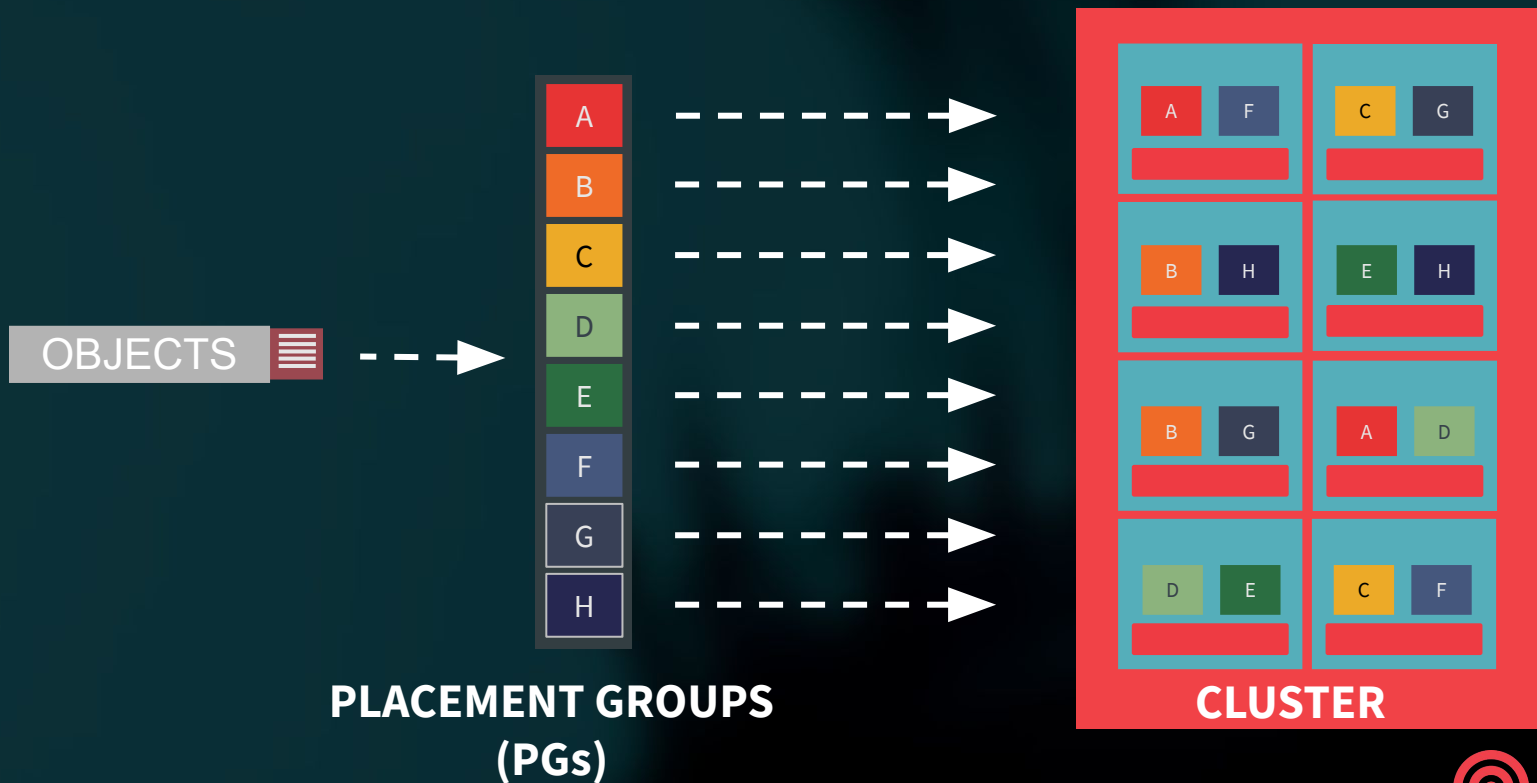
OBJECT



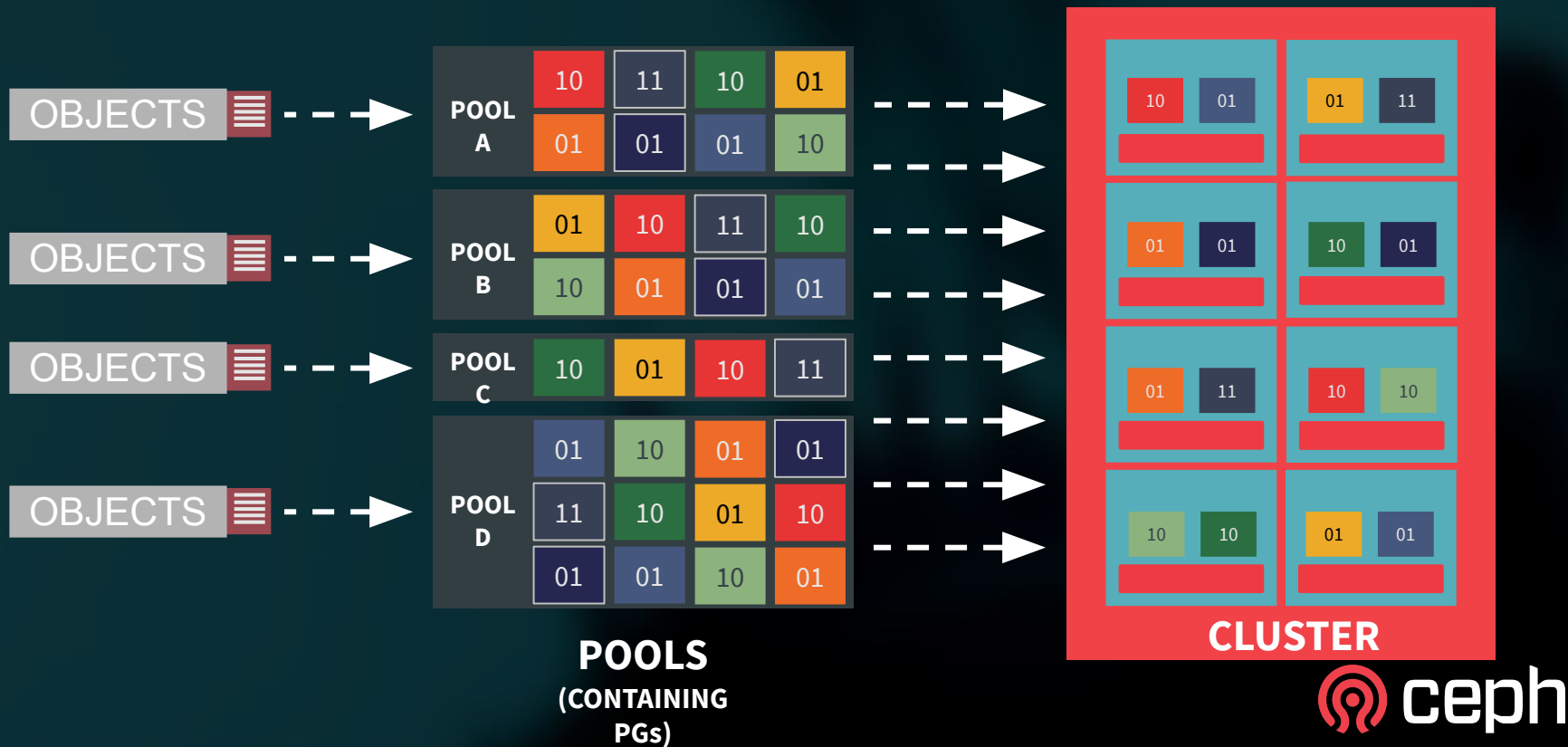
?
?



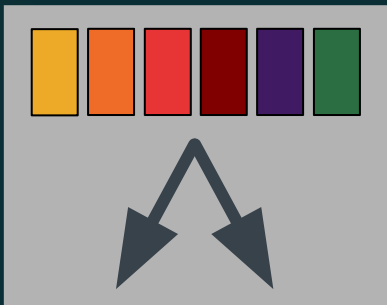
CRUSH: placement algorithm



DATA IS ORGANIZED INTO POOLS



CRUSH: DYNAMIC DATA PLACEMENT



CRUSH (Controlled Replication Under Scalable Hashing)

- Pseudo-random placement algorithm
 - Fast calculation, no lookup
 - Repeatable, deterministic
- Statistically uniform distribution
- Stable mapping
 - Limited data migration on change
- Rule-based configuration
 - Infrastructure topology aware
 - Adjustable replication
 - Weighting

 uestions ???