# Flexible + Scalable Storage

Red Hat's software-defined storage portfolio:
Ceph & Gluster

Patrick Ladd

Technical Account Manager, FSI

pladd@redhat.com

https://people.redhat.com/pladd/

# Red Hat Storage Overview

- ## Software Defined Storage

  - ### What and Why?

- ## Red Hat's Portfolio

  - ### Red Hat Ceph Storage

  - ### Red Hat Gluster Storage
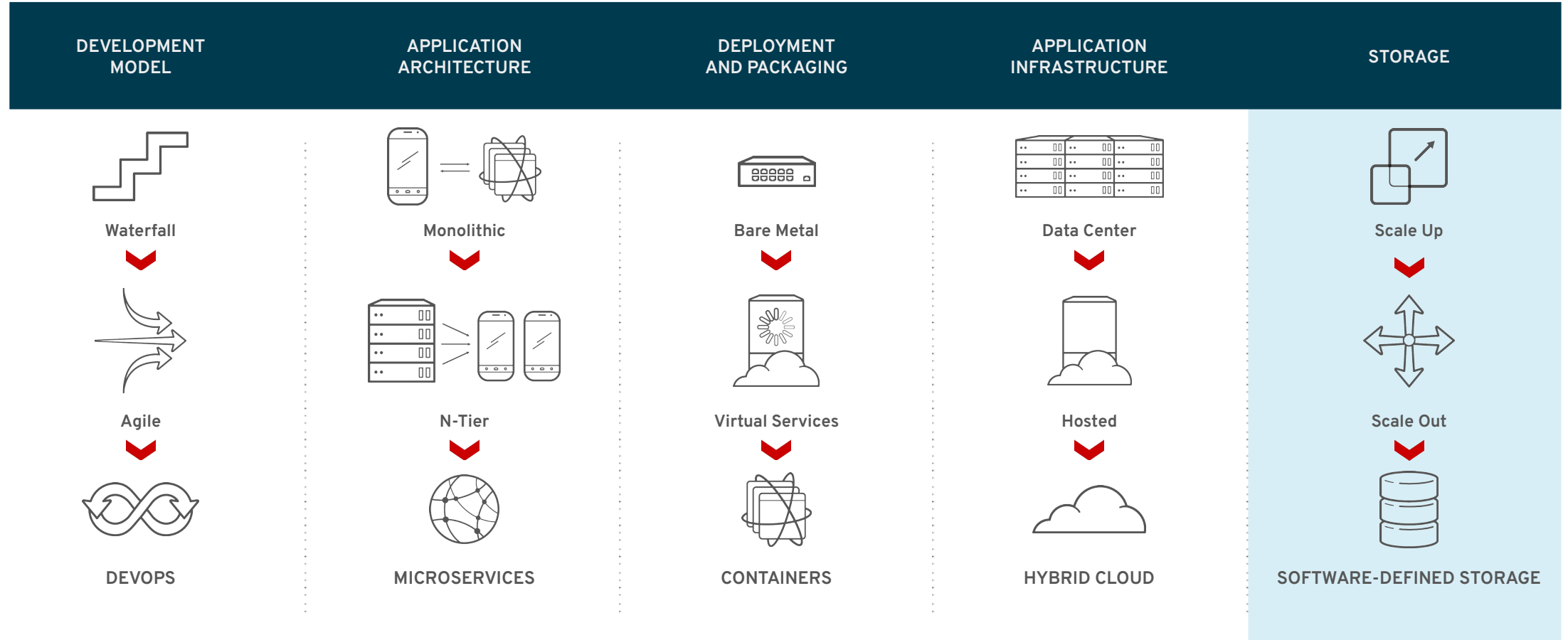
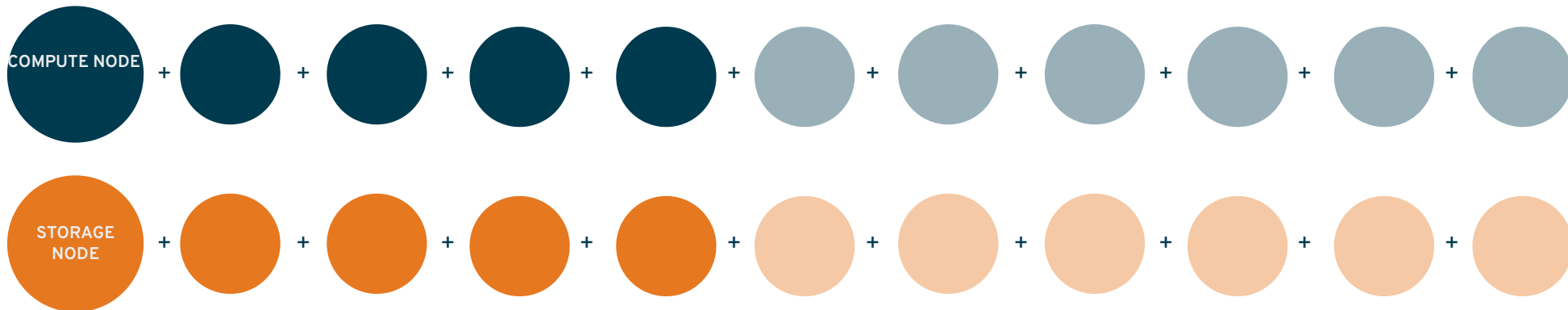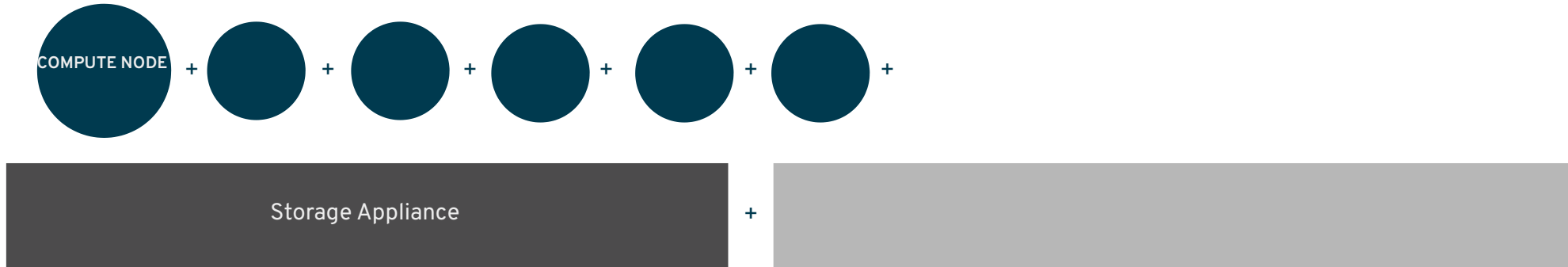# WHAT IS SOFTWARE-DEFINED STORAGE?

**SERVER-BASED**

**CENTRALIZED CONTROL**

**OPEN ECOSYSTEM**

Red Hat

# THE ROAD TO SOFTWARE-DEFINED STORAGE

| DEVELOPMENT MODEL | APPLICATION ARCHITECTURE | DEPLOYMENT AND PACKAGING | APPLICATION INFRASTRUCTURE | STORAGE |
|---|---|---|---|---|
| Waterfall | Monolithic | Bare Metal | Data Center | Scale Up |
| Agile | N-Tier | Virtual Services | Hosted | Scale Out |
| DEVOPS | MICROSERVICES | CONTAINERS | HYBRID CLOUD | SOFTWARE-DEFINED STORAGE |

Red Hat

# DISRUPTION IN THE STORAGE INDUSTRY

| PUBLIC CLOUD STORAGE | ← | TRADITIONAL APPLIANCES | → | SOFTWARE-DEFINED STORAGE |
|:---:|:---:|:---:|:---:|:---:|
| better | | COST EFFICIENCY | | better |
| faster | | PROVISIONING | | faster |
| more | | VENDOR LOCK-IN | | less |
| less | | SKILL REQUIRED | | more |
| weaker | | GOVERNANCE | | stronger |
| limited | | DEPLOYMENT OPTIONS | | broad |

Red Hat

# Virtualized Storage Scales Better

# Comparing Throughput and Costs at Scale
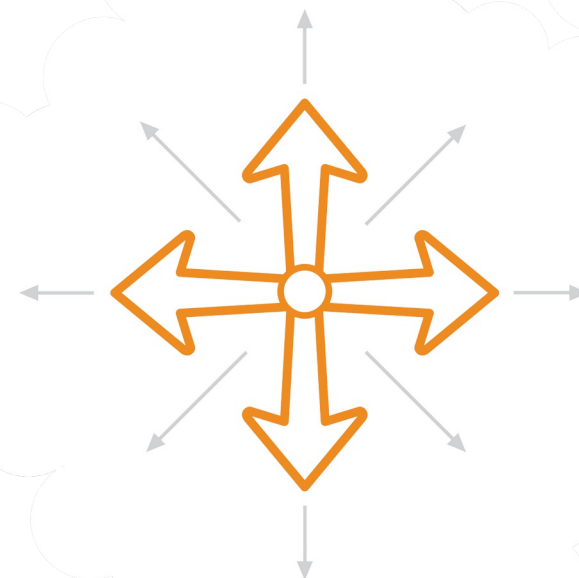
## STORAGE PERFORMANCE SCALABILITY

READS/WRITES THROUGHPUT (MBPS)

Software Defined Scale-out
Storage (GlusterFS)

Traditional Enterprise
NAS Storage

NUMBER OF STORAGE NODES

## STORAGE COSTS SCALABILITY

TOTAL STORAGE COSTS ($)

Traditional Enterprise
NAS Storage

Software Defined Scale-out
Storage (GlusterFS)

NUMBER OF STORAGE NODES

Red Hat

# The Robustness of Software

Software can do things hardware can't

Storage services based on software are  more flexible than hardware-based implementations

- Can be deployed on bare metal, inside containers, inside VMs, or in the public cloud

- Can deploy on a single server, or thousands, and can be upgraded and reconfigured on the fly

- Grows and shrinks programmatically to meet changing demands

Red Hat

# How Storage Fits

**RED HAT® STORAGE**

| PHYSICAL | VIRTUAL | PRIVATE CLOUD | CONTAINERS | PUBLIC CLOUD |
|---|---|---|---|---|
| **RED HAT®** CEPH STORAGE | **RED HAT®** CEPH STORAGE | **RED HAT®** CEPH STORAGE | **RED HAT®** CEPH STORAGE | **RED HAT®** CEPH STORAGE |
| **RED HAT®** GLUSTER STORAGE | **RED HAT®** GLUSTER STORAGE | **RED HAT®** GLUSTER STORAGE | **RED HAT®** GLUSTER STORAGE | **RED HAT®** GLUSTER STORAGE |
| **RED HAT®** ENTERPRISE LINUX® | **RED HAT®** ENTERPRISE LINUX®<br><br>**RED HAT®** ENTERPRISE VIRTUALIZATION | **RED HAT®** OPENSTACK® PLATFORM | **OPENSHIFT ENTERPRISE** by Red Hat | **RED HAT®** ENTERPRISE LINUX® |

Red Hat

# DIFFERENT KINDS OF STORAGE



## BLOCK STORAGE

Physical storage media appears to computers as a series of sequential blocks of a uniform size.
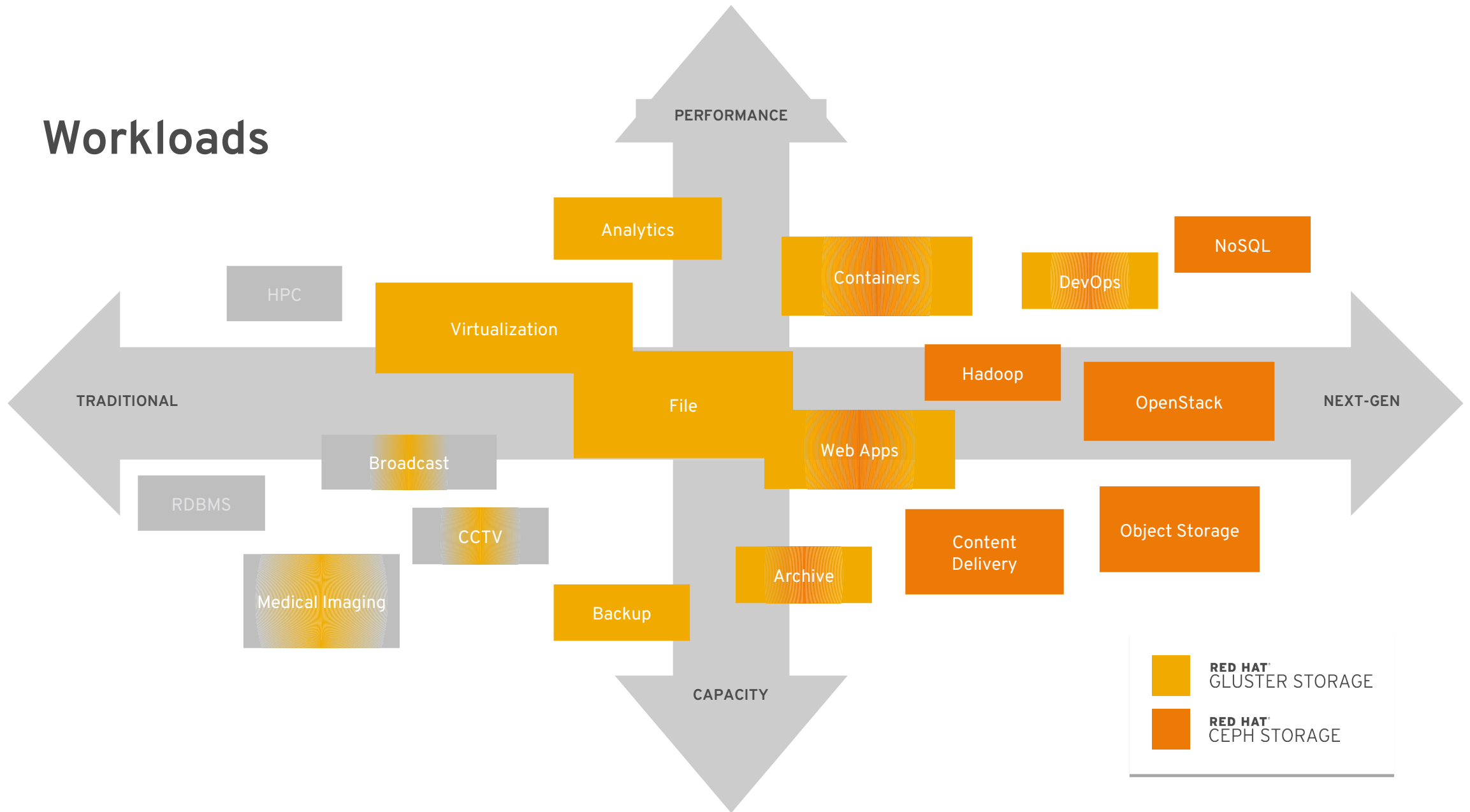
## FILE STORAGE

File systems allow users to organize data stored in blocks using hierarchical folders and files.

## OBJECT STORAGE

Object stores distribute data algorithmically throughout a cluster of media, without a rigid structure.
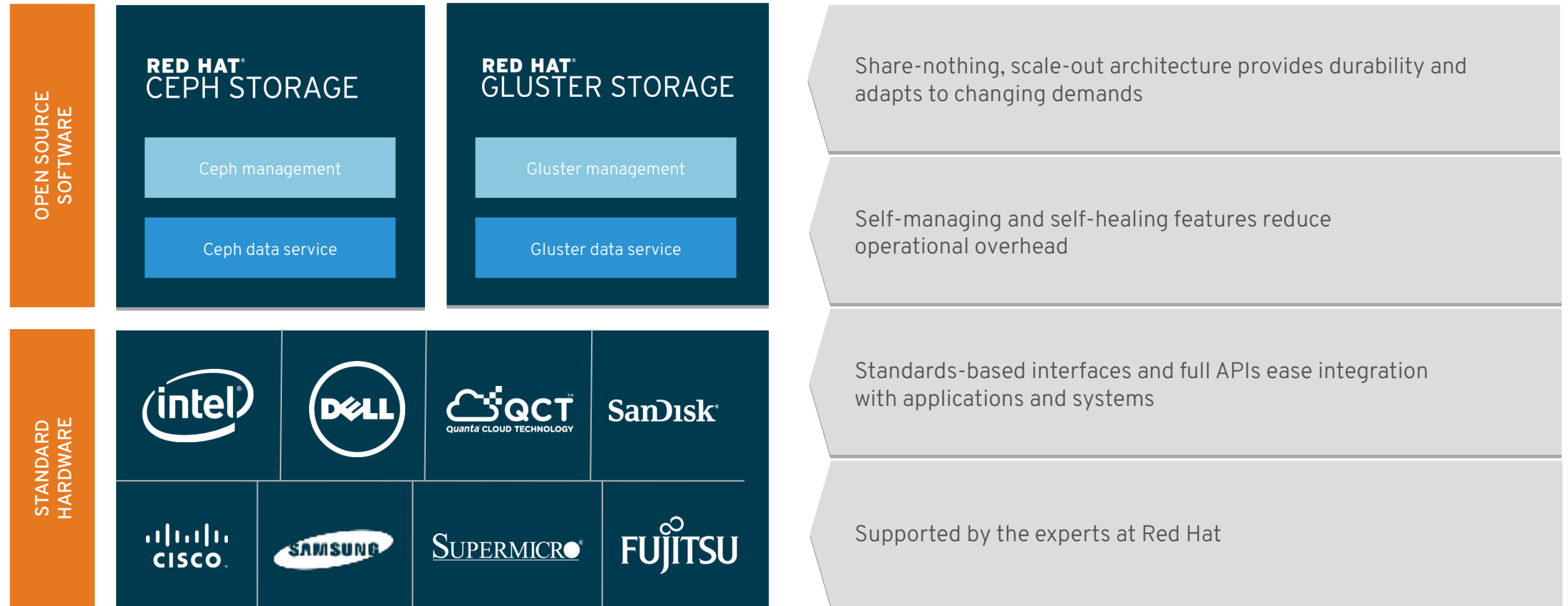
# Red Hat Storage Overview

# THE RED HAT STORAGE PORTFOLIO

**RED HAT®**
CEPH STORAGE

Ceph management

Ceph data service

**RED HAT®**
GLUSTER STORAGE

Gluster management

Gluster data service

intel

DELL

QCT
Quanta CLOUD TECHNOLOGY

SanDisk®

CISCO

SAMSUNG

SUPERMICRO®

FUJITSU

Share-nothing, scale-out architecture provides durability and adapts to changing demands

Self-managing and self-healing features reduce operational overhead

Standards-based interfaces and full APIs ease integration with applications and systems

Supported by the experts at Red Hat

Red Hat

# OVERVIEW: RED HAT CEPH STORAGE

## RED HAT® CEPH STORAGE

**TARGET USE CASES**

**Cloud Infrastructure**
- VM storage with OpenStack® Cinder, Glance Keystone, Manila, and Nova
- Object storage for tenant apps

**Rich Media and Archival**
S3-compatible object storage

## Powerful distributed storage for the cloud and beyond

- Built from the ground up as a next-generation storage system, based on years of research and suitable for powering infrastructure platforms

- Highly tunable, extensible, and configurable, with policy-based control and no single point of failure

- Offers mature interfaces for block and object storage for the enterprise

---

**CUSTOMER HIGHLIGHT: CISCO**

**CISCO**™

Cisco uses Red Hat Ceph Storage to deliver storage for next-generation cloud services

**Red Hat**

# OVERVIEW: RED HAT GLUSTER STORAGE

**RED HAT® GLUSTER**
**STORAGE**

**TARGET USE CASES**

**Enterprise File Sharing**
- Media streaming
- Active Archives

**Enterprise Virtualization**

**Rich Media and Archival**

## Agile file storage for petabyte-scale workloads

- Purpose-built as a scale-out file store with a straightforward architecture suitable for public, private, and hybrid cloud

- Simple to install and configure, with a minimal hardware footprint

- Offers mature NFS, SMB and HDFS interfaces for enterprise use

**CUSTOMER HIGHLIGHT: INTUIT**

**intuit.**

Intuit uses Red Hat Gluster Storage to provide flexible, cost-effective storage for its industry-leading financial offerings

Red Hat

# Red Hat
# Gluster Storage

# GLUSTER FUNDAMENTALS

- Clustered Scale-out General Purpose Storage Platform

- Fundamentally File-Based & POSIX End-to-End
  - Familiar Filesystems Underneath (EXT4, XFS)
  - Familiar Client Access (NFS, Samba, FUSE)
- No Metadata Server

- Standards-Based – Clients, Applications, Networks

- Modular Architecture for Scale and Functionality

# GLUSTER TERMINOLOGY

**Cluster:** Collection of peer systems

**Node:** System Participating in Cluster

**Brick:** Any Linux Block Device

**Volume:** Bricks taken from one or more hosts presented as a single unit

Cluster

Node

Node

Brick

Brick

Brick

Brick

Volume

Clients

# GLUSTER ELASTIC HASH ALGORITHM

| 0 | 1 | 1 | 0 | 1 | 0 | 1 | 0 |

**#**

BRICK        BRICK

## No Central Metadata Server

- Suitable for unstructured data storage
- No single point of failure

## Elastic Hashing

- Files assigned to virtual volumes
- Virtual volumes assigned to multiple bricks
- Volumes easily reassigned on-the-fly

## Location Hashed on Filename

- No performance bottleneck
- Eliminates risk scenarios

Red Hat

# TRANSLATION LAYERS

Translation layers handle:

- Data resilience scheme is maintained (replication, erasure coding)
- Metadata is stored and tracked with the object
- Dynamic mapping from virtual volumes to data volumes
- Heal, Rebalance, Bitrot Detection, Geo-Replication, …
- Data translation hierarchy (protocols, encryption, performance, …)
- Health monitoring, alerting, and response

Translator 1

Translator 2

...

Translator *n*

Red Hat

# SERVER- AND CLIENT-SIDE TRANSLATORS

Translations layers may be distributed!

- Some layers in the translator stack may be implemented on the client
- Higher performance and efficiency

# GLUSTER DEFAULT DATA PLACEMENT

# Distributed Volume



**MOUNT POINT**

**Uniformly distribute files across all bricks in the namespace**

**DISTRIBUTED VOLUME**

server1:/exp1          server2:/exp2

BRICK                  BRICK

FILE 1    FILE 2       FILE 3

# GLUSTER FAULT-TOLERANT DATA PLACEMENT

## Distributed-Replicated Volume



**MOUNT POINT**

DISTRIBUTED-REPLICATED VOLUME

Replicated Vol 0

Replicated Vol 1

**Creates an fault-tolerant distributed volume by mirroring same file across 2 bricks**

server1

server2

server3

server4

BRICK
(exp 1)

BRICK
(exp 2)

BRICK
(exp 3)

BRICK
(exp 4)

FILE 1

FILE 1

FILE 2

FILE 2

# GLUSTER ERASURE CODING

## Storing more data with less hardware



**MOUNT POINT**

**Disperses fragments of data and data parity across volume bricks**

DISPERSED VOLUME

FILE 1

| server1 | server2 | server3 | server4 | server5 | server6 |
| --- | --- | --- | --- | --- | --- |
| BRICK | BRICK | BRICK | BRICK | BRICK | BRICK |
| 1 | 2 | 3 | 4 | X | Y |

Red Hat

# GLUSTER CLIENT ACCESS

Multi-protocol distributed file system access with optional Swift object translator

| FUSE | NFS | SMB | API | Swift |
|------|-----|-----|-----|-------|

**FILE**

CLUSTERS OF SERVER HOSTS

HOSTS = ON-PREMISE OR PUBLIC CLOUD HOST INSTANCES

Red Hat

# GLUSTER GEO-REPLICATION

# Multi-site content distribution

Asynchronous across LAN, WAN, or Internet

Master-slave model, cascading possible

Continuous and incremental

Multiple configurations

- One to one
- One to many
- Cascading



One to One replication

Site A

Site B

Cascading replication

Site B

Site A

Site C

# SELF HEALING

- Automatic Repair of Files
  - As they are accessed
  - Periodic via Daemon

## Scenarios:
- Node offline
  - Bricks on node need to be caught up to current
- Node or brick loss
  - New brick needs to be completely rebuilt

# BIT ROT DETECTION

- Scans data periodically for bit rot

- Check sums are computed when files are accessed and compared against previously stored values

- On mismatch, an error is logged for the storage admin

ADMIN

| 0 | 0 | 0 | 0 | 0 | X | 0 | 0 |

# SNAPSHOTS

- Volume level, ability to create, list, restore, and delete

- LVM2 based, operates only on thin-provisioned volumes

- User serviceable snapshots

- Crash consistent image

**BEFORE SNAPSHOT**

**AFTER SNAPSHOT**

CURRENT FILE SYSTEM

SNAPSHOT

CURRENT FILE SYSTEM

A B C D

A B C D

# TIERING

- Automated promotion and demotion of data between "hot" and "cold" sub volumes

- Based on frequency of access

- Cost-effective flash acceleration

# QUOTAS

- Control disk utilization at both directory and volume level

**Quota Limits**

- Two levels of quota limits: Soft (default) and hard

- Warning messages issued on reaching soft quota limit

- Write failures with EDQUAT message after hard limit is reached

**Global vs. Local Limits**

- Quota is global (per volume)

- Files are psuedo-randomly distributed across bricks



VOLUME

BRICK 0     BRICK 1     BRICK n

Red Hat

# Red Hat Gluster Storage Demo

# INSTALL

# INSTALL

# Red Hat
# Ceph Storage

# CEPH FUNDAMENTALS

- Single, efficient, unified storage platform (object, block, file)
- User-driven storage lifecycle management with 100% API coverage

- Integrated, easy-to-use management console

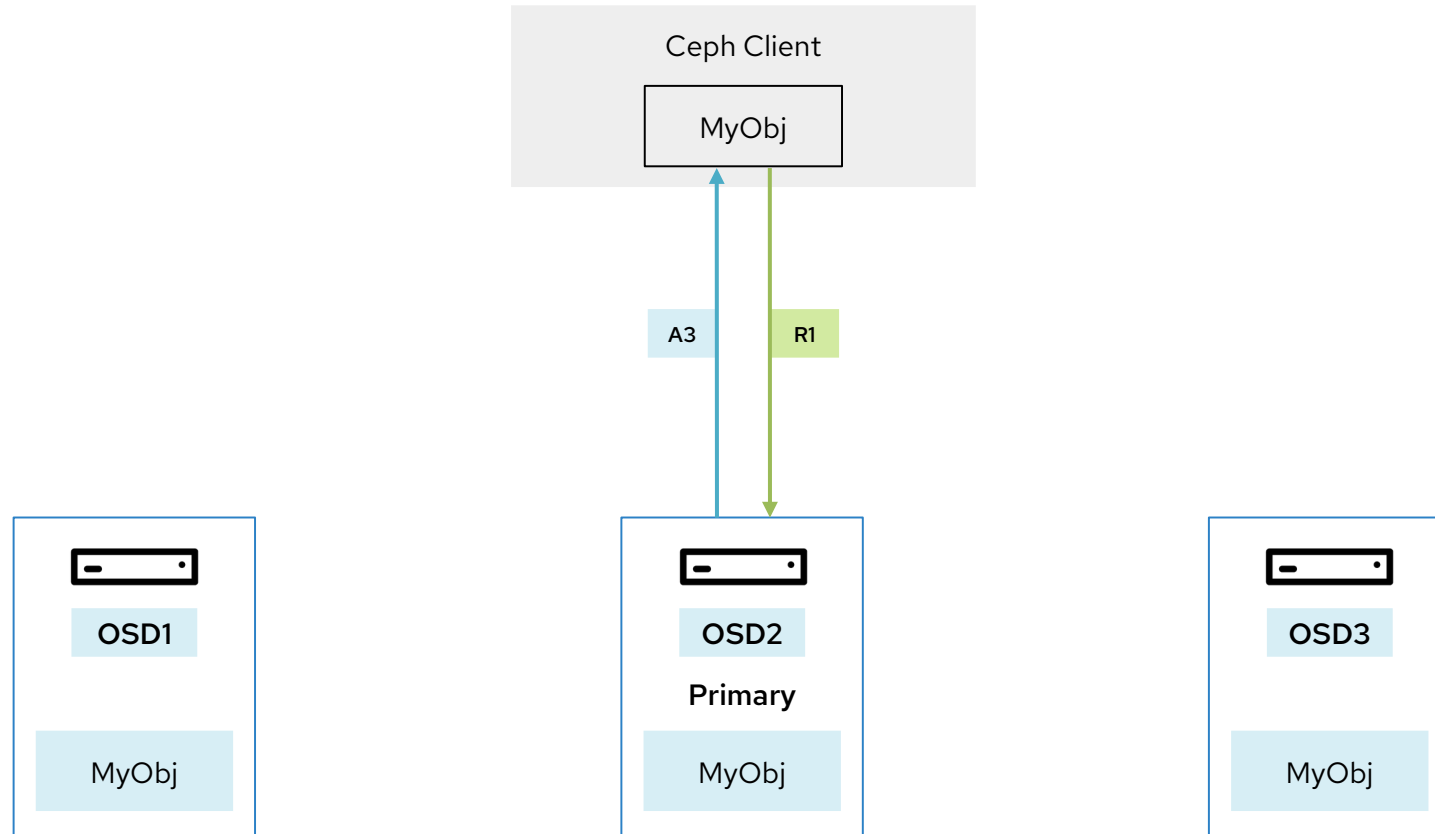- Designed for cloud infrastructure and emerging workloads
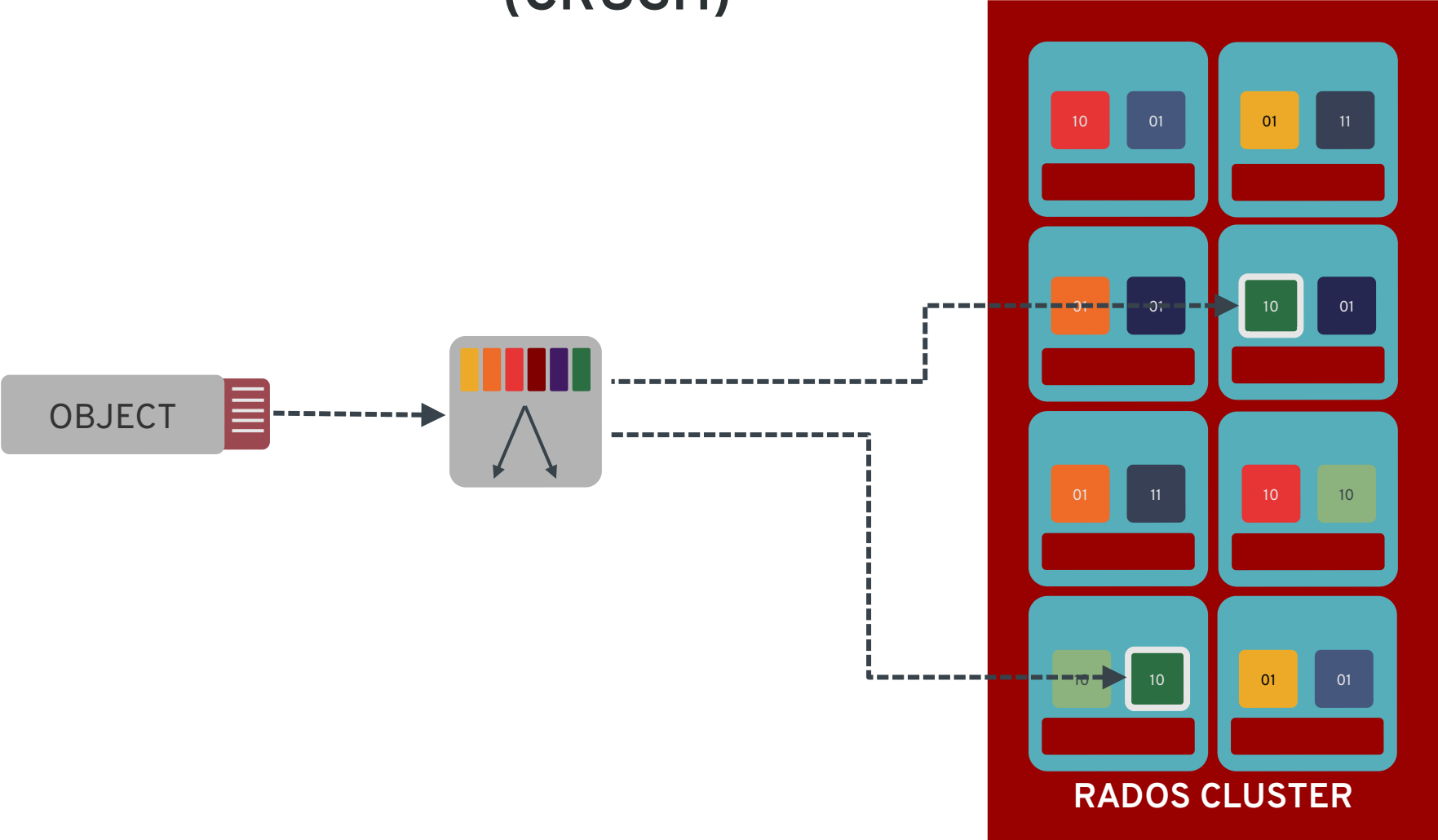
# CEPH ARCHITECTURE

APP                     HOST/VM                 CLIENT

### RGW
A web services gateway for object storage, compatible with S3 and Swift

### RBD
A reliable, fully-distributed block device with cloud platform integration

### CEPHFS
A distributed file system with POSIX semantics and scale-out metadata management

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

### RADOS
A software-based, reliable, autonomous, distributed object store comprised of self-healing, self-managing, intelligent storage nodes and lightweight monitors

Red Hat

# DETAILED ARCHITECTURE

**Client interface layer**

LIBRADOS | RADOSGW | RBD

RADOS

**Objects in pools**

Obj Obj Obj Obj Obj Obj Obj Obj

Obj Obj Obj Obj Obj Obj Obj Obj

**CRUSH ruleset**

Pool ID. (Hash(Object) %Num of PGs)

Pool ID. (Hash(Object) %Num of PGs)

**Placement groups**

PG PG PG PG PG PG PG PG

**CRUSH Map**

**Ceph nodes:**
**OSD hosts**
**monitors (mons)**

OSD1 | OSD2 | OSD3

OSD4 | OSD5 | OSD6

Mon1    Mon2    Mon3

Red Hat

# WRITES

# READS

# CEPH DATA PLACEMENT (CRUSH)



OBJECT

RADOS CLUSTER

# STORAGE CONSOLE

- An easy to use interface for managing cluster lifecycles

- Ansible-based deployment tools for driving granular configuration options from CLI or GUI

- Monitoring and graphs for troubleshooting with statistical information about components

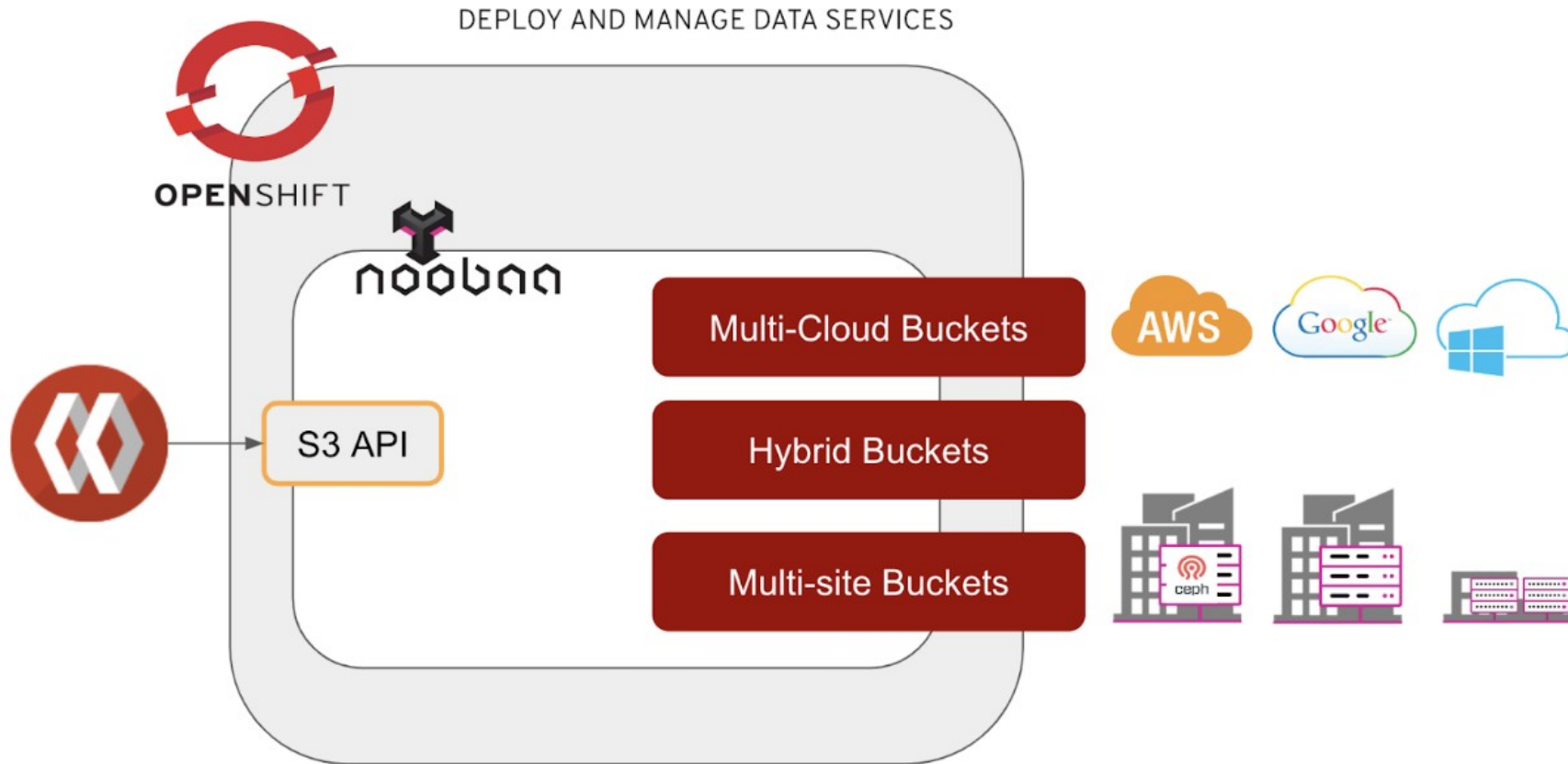# Red Hat OpenShift Container Storage

# HYPERCONVERGED STORAGE

# OPENSHIFT INTEGRATION

# MULTI-CLOUD WITH NOOBAA

# Thank you

Red Hat is the world's leading provider of enterprise

open source software solutions. Award-winning

support, training, and consulting services make

Red Hat a trusted adviser to the Fortune 500.

linkedin.com/company/red-hat

youtube.com/user/RedHatVideos

facebook.com/redhatinc

twitter.com/RedHat

**Red Hat**