# Let's decompose storage (again)

*Why? How? Huh?*

## MSST May 17 2017

**Evan Powell**

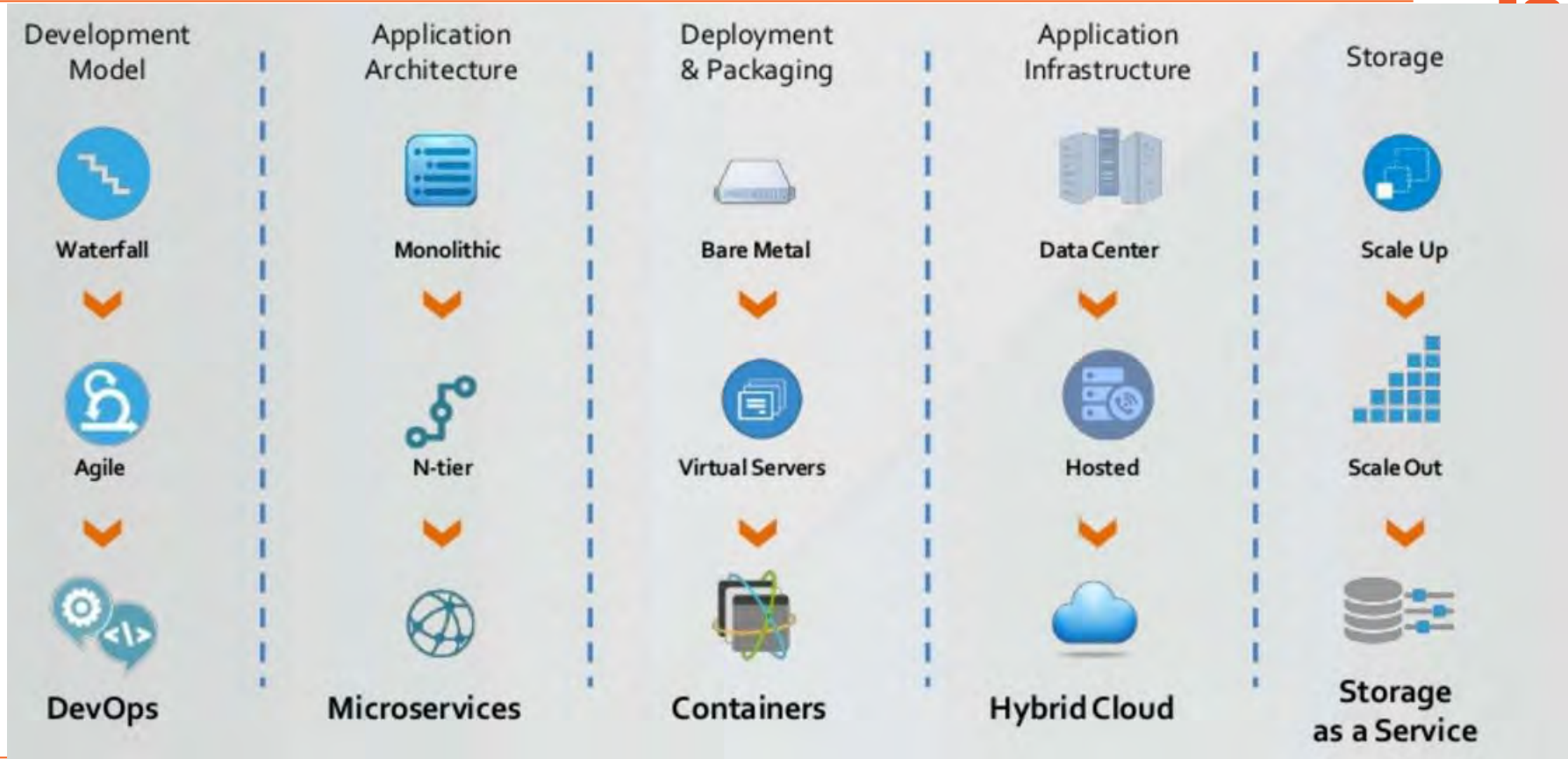blog.openebs.io

https://github.com/openebs

Join the community #slack slack.openebs.io

@openebs

# What's new?

| Development Model | Application Architecture | Deployment & Packaging | Application Infrastructure | Storage |
|---|---|---|---|---|
| Waterfall | Monolithic | Bare Metal | Data Center | Scale Up |
| Agile | N-tier | Virtual Servers | Hosted | Scale Out |
| **DevOps** | **Microservices** | **Containers** | **Hybrid Cloud** | **Storage as a Service** |

# What's new?



| Development Model | Application Architecture | Deployment & Packaging | Application Infrastructure | Storage |
|---|---|---|---|---|
| Waterfall | Monolithic | Bare Metal | Data Center | Scale Up |
| Agile | N-tier | Virtual Servers | Hosted | Scale Out |
| **DevOps** | **Microservices** | **Containers** | **Hybrid Cloud** | **Storage as a Service** |

# Layering

(3-TIER APPs)

CERTIFIED
SYSADMIN(S)

Supervised Provisioning
Manage Storage Upgrades!
Manage 100s of Volumes
Managed Upgrades

**NFS,iSCSI**

Enterprise Storage
(SAN, NAS, ScaleOut,
DP, DR, Backup,
Compression, Dedup)

(LOCAL APPs)

GEEK

Format Disks and Use.
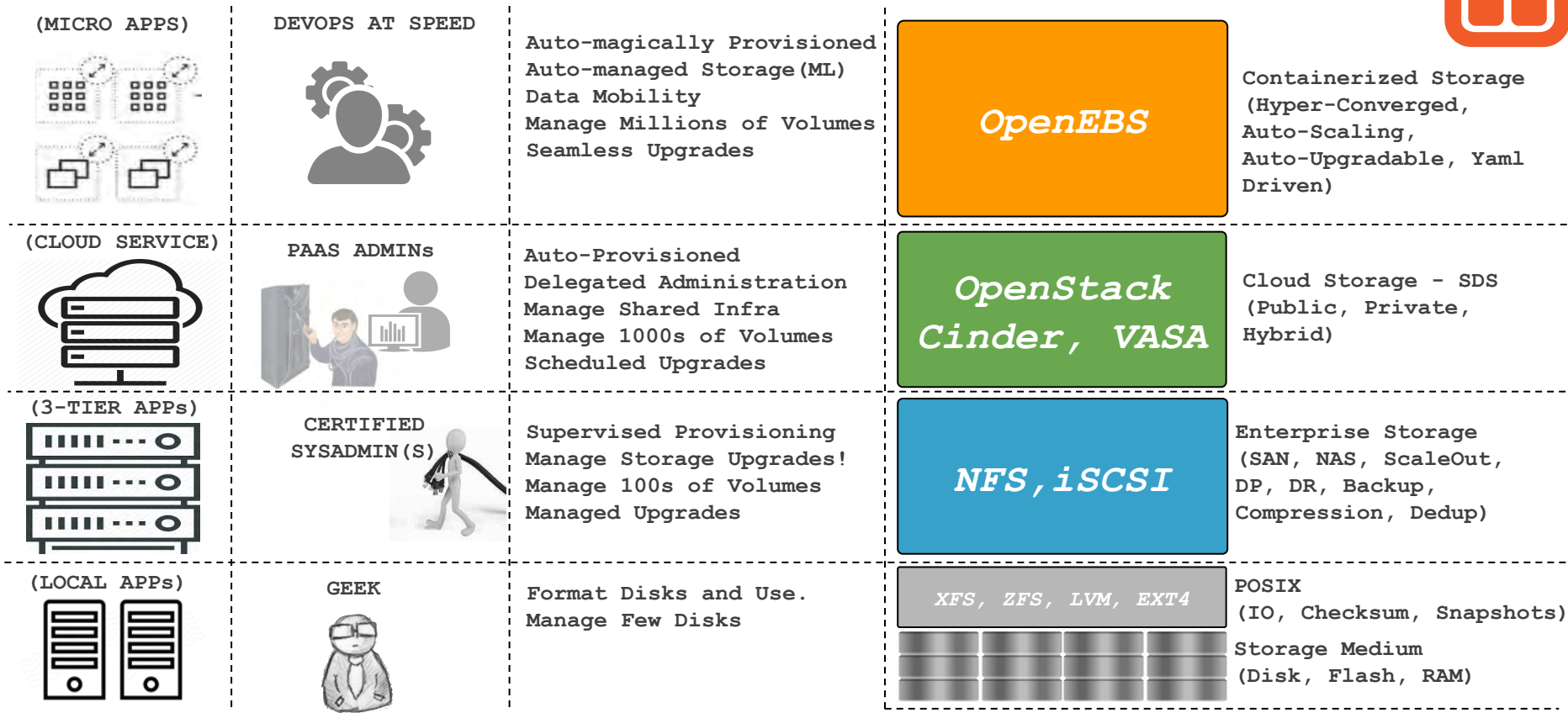Manage Few Disks
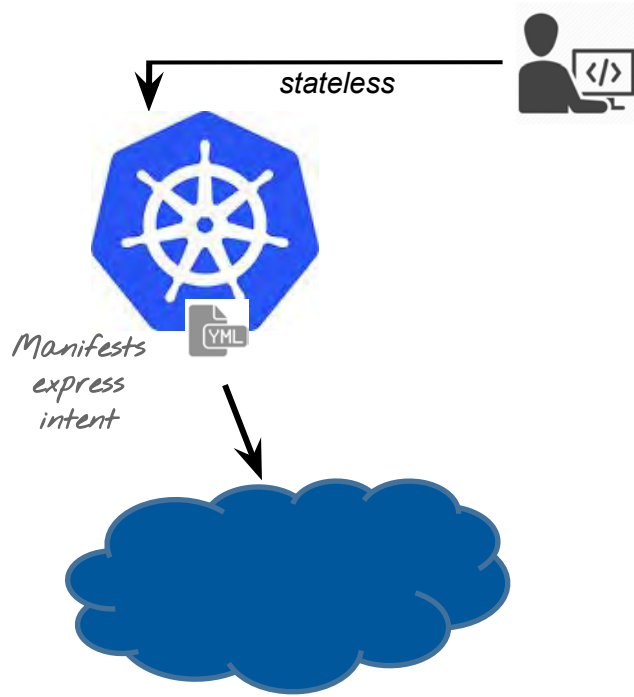
*XFS, ZFS, LVM, EXT4*

POSIX
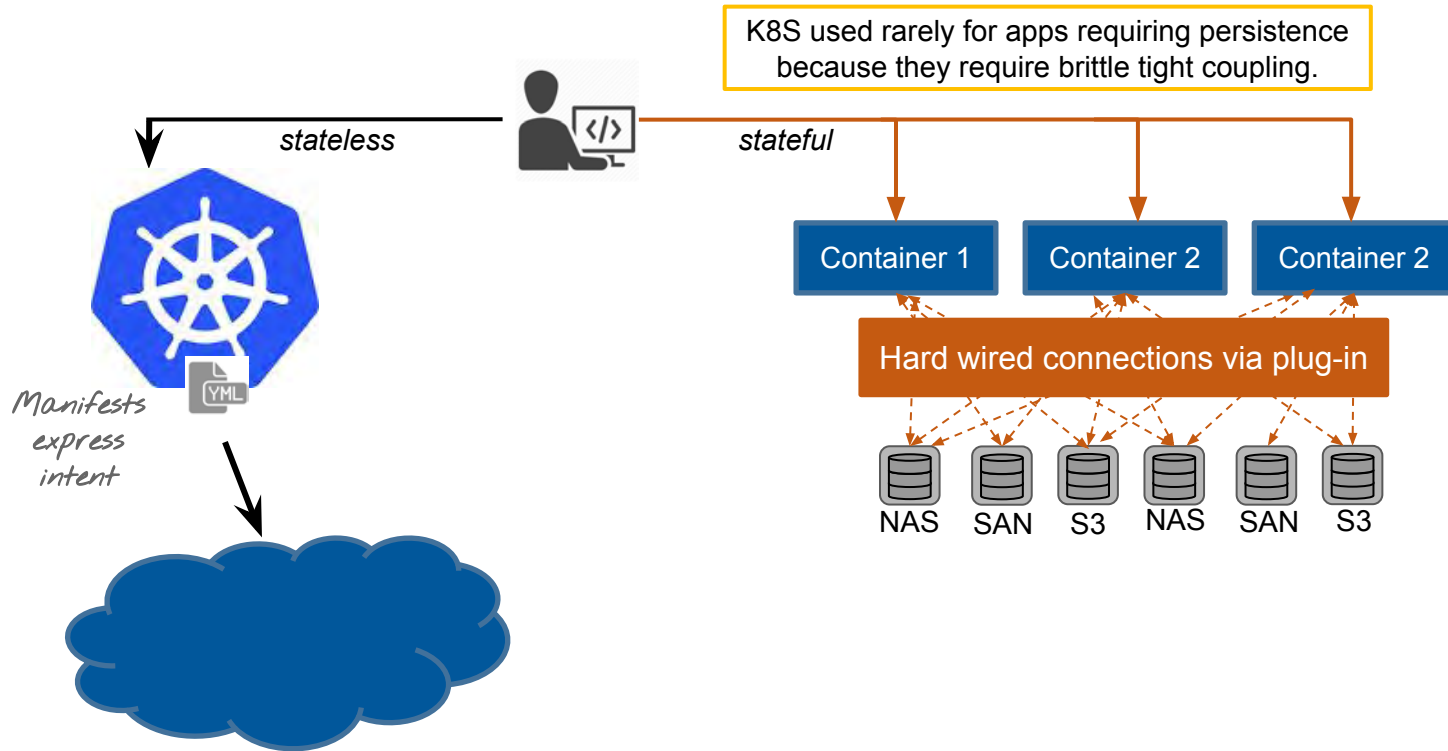(IO, Checksum, Snapshots)

Storage Medium
(Disk, Flash, RAM)

# Layering

| (MICRO APPS) | DEVOPS AT SPEED | Auto-magically Provisioned<br>Auto-managed Storage(ML)<br>Data Mobility<br>Manage Millions of Volumes<br>Seamless Upgrades | **OpenEBS** | Containerized Storage<br>(Hyper-Converged,<br>Auto-Scaling,<br>Auto-Upgradable, Yaml<br>Driven) |
|---|---|---|---|---|
| (CLOUD SERVICE) | PAAS ADMINs | Auto-Provisioned<br>Delegated Administration<br>Manage Shared Infra<br>Manage 1000s of Volumes<br>Scheduled Upgrades | **OpenStack<br>Cinder, VASA** | Cloud Storage - SDS<br>(Public, Private,<br>Hybrid) |
| (3-TIER APPs) | CERTIFIED<br>SYSADMIN(S) | Supervised Provisioning<br>Manage Storage Upgrades!<br>Manage 100s of Volumes<br>Managed Upgrades | **NFS,iSCSI** | Enterprise Storage<br>(SAN, NAS, ScaleOut,<br>DP, DR, Backup,<br>Compression, Dedup) |
| (LOCAL APPs) | GEEK | Format Disks and Use.<br>Manage Few Disks | *XFS, ZFS, LVM, EXT4* | POSIX<br>(IO, Checksum, Snapshots)<br>Storage Medium<br>(Disk, Flash, RAM) |

*OpenSource Technology Stack*

# Happy days!

stateless

Manifests express intent

# Painful persistence



K8S used rarely for apps requiring persistence because they require brittle tight coupling.

*stateless*

*stateful*

Manifests express intent

Container 1

Container 2

Container 2
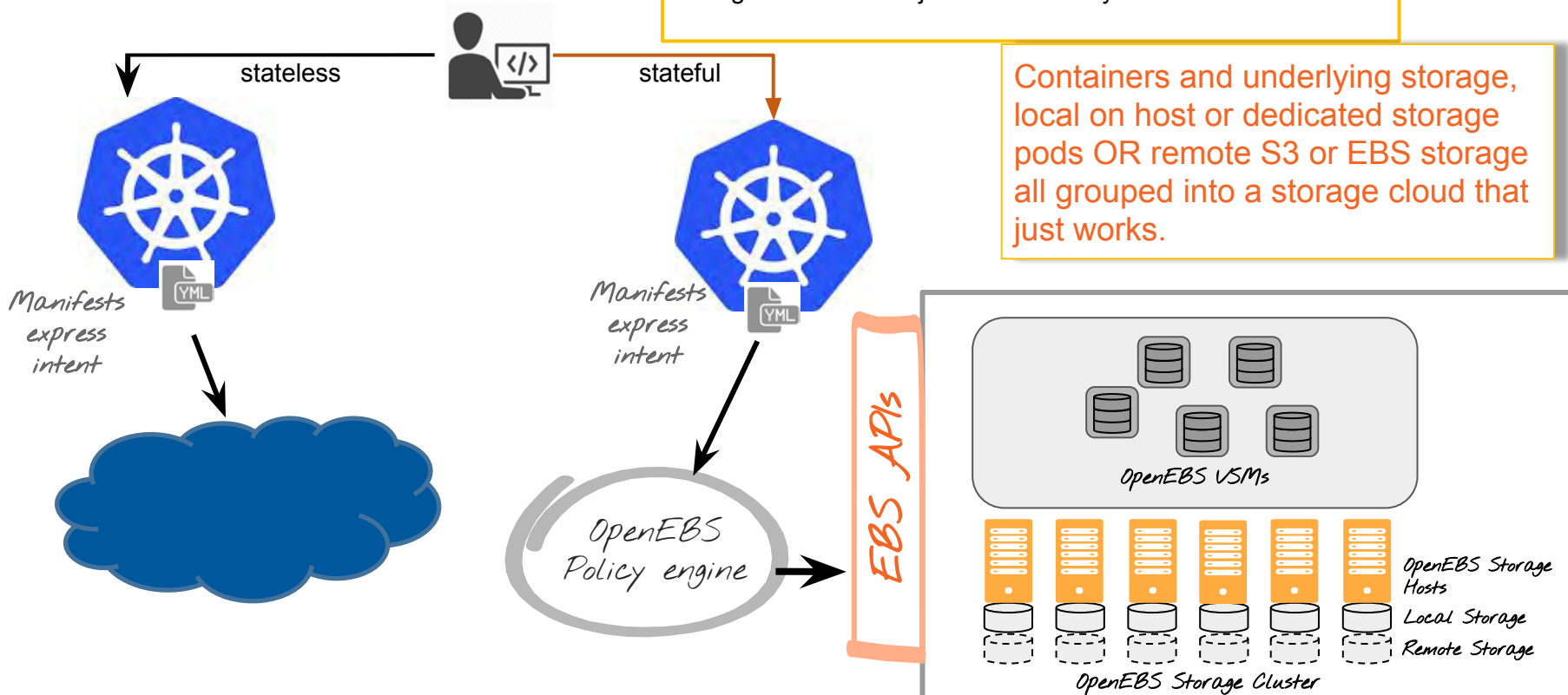
Hard wired connections via plug-in

NAS    SAN    S3    NAS    SAN    S3

# Desired state of state

No changes to DevOps workflow even for containers requiring persistence. Users manifest their intent and the storage and storage controllers adjust automatically as needed.

stateless

stateful

Containers and underlying storage, local on host or dedicated storage pods OR remote S3 or EBS storage all grouped into a storage cloud that just works.

Manifests express intent

Manifests express intent

OpenEBS Policy engine

EBS APIs

OpenEBS VSMs

OpenEBS Storage Hosts

Local Storage

Remote Storage

OpenEBS Storage Cluster

# Architecture and Design

- Powered by Linux, Go and OpenSource
- Built and Delivered as Containers / Micro-services
- Longhorn, Gotgt, Kubernetes, Consul

# Design Goals and Constraints

Fault tolerant and secure by default

Low entry barrier, easy to setup

Storage optimized for Containerized Applications

Horizontally scalable to millions of Containers

Seamless integration into existing private and public cloud environments

Non-disruptive upgrades

Developer and Operators Friendly

Completely OpenSource (Apache license)

Microservices based

DevOps architecture

# Overview & Terminology

**K8s master**

**K8s minions**

**Pod**

Storage
Driver

HTTPS
(manage)

Network (Flannel)

Network (Flannel*)

Data

**OpenEBS VSMs / Storage Pod**

**OpenEBS
Maya
master**

Local Storage

Remote Storage

**OpenEBS Storage Hosts**

# Deployment - Hyper-Converged



**K8s master**

OpenEBS
Maya-K8s
Adaptors

Network (Flannel)

**K8s minions**

OpenEBS Maya Storage
Orchestrator

**Pod**

TCMU

TCP

**Storage Pods(3)**

# VSM - Storage in Containers



**VSM / Storage Pod**

- Data (iSCSI/TCMU)
- Frontend Containers
- Inline Replication
- Backend Containers to Persist Data (Cached, Protected)

**OpenEBS Storage Hosts**

- Container (Docker)
- Maya Storage Orchestrator
- NVMe Flash

**Storage**

- Multiple Storage Backends

Local Disks   NAS or SAN   Cloud Storage

# Jiva - Containerized Storage Image

# Maya - Container Storage Orchestration

**Integrations**

mTerraform  mDriver  mAWS-MP

**OpenEBS Maya Master**

mAPI  mConnect  mGUI

mDB  mAIEngine  mAnalytics

maya  mSCH  mCluster

**OpenEBS Storage Host**

mAgent  mCluster

mTelmetry  mJRunner

maya  mStorageInterface

# Storage Internals

- Capacity Management
- QoS
- Access - iSCSI, TCMU
- Snapshot / Restore (S3)
- Backup / Migration
- Caching/Tiering
- Replication / Rebuild

# OpenEBS - Core differentiations

- The block storage software is made into a micro service
- The 'micro service' has its own block protocol stack, tiering engine, QoS engine and ML prediction capability
- The block storage knowledge is maintained on a per-volume basis. The data of each volume is divided into cold-data and hot-data. Cold-data resides on NVMe-Flash or on 3DX-Memory. Hot-data resides in slower disks / SAN/ Cloud-Storage/S3
- The metadata knowledge also is maintained at a volume level (not the entire storage). This saves us from the issue of huge-metadata-sifting at scale. The traversal through meta-data depends on the "size of the volume" and not on the "number of volumes".
- Within the volume the meta data is not managed at "block-level" but at "chunk-level". Typical block-size is 4KB and Typical chunk-size is 4 MB. This results in the huge reduction of metadata size of the block-volume that needs to be maintained.
- Checksum - One of the important metadata is checksum. OpenEBS guarantees bit-rot protection through the use of checksums. The checksums are managed at a chunk level only on Cold-Data. The checksums are not managed on hot-data, the blocks go in and out of chunks on the hot-data without the need of checksum calculation on the fly.
- Deduplication-while-tiering: Deduplication has capacity benefits but kills performance (either inline or offline). But in OpenEBS, we do this while moving the data from hot-to-cold tiers. In effect, the benefits of deduplication without the performance penalty.
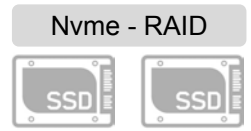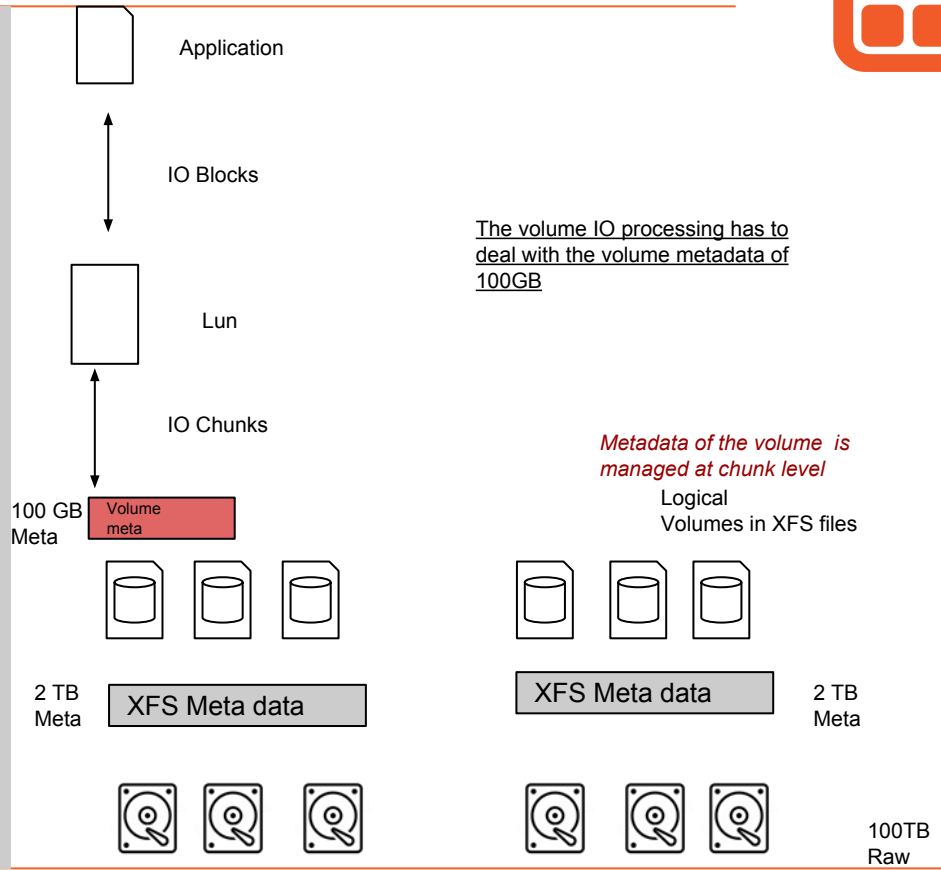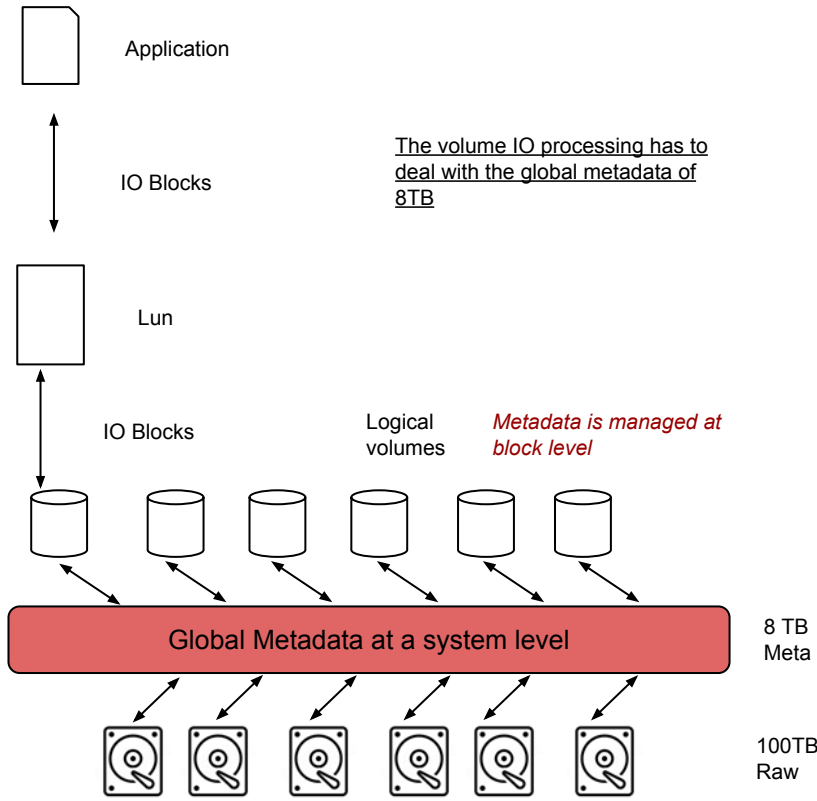
# Cchchchcunking

Data blocks in NVRAM (NVDIMM)

Blocks to chunks (coalescing)

Data block to chunkbook mapper

lun1.1.metabook
lun1.2.metabook

lun chunkbooks

lun1 data chunks

lun2 data chunks

lun3 data chunks

Read - uncompress

Write Dedup and compress

User land

Data in nvme Flash

Kernel

Nvme - RAID

SSD
SSD

XFS

LVM RAID

# OpenEBS - Metadata at scale is not an issue

Application

IO Blocks

The volume IO processing has to deal with the global metadata of 8TB

Application

IO Blocks

Lun

The volume IO processing has to deal with the volume metadata of 100GB

Lun

IO Blocks

Logical volumes

*Metadata is managed at block level*

IO Chunks

*Metadata of the volume is managed at chunk level*

Logical Volumes in XFS files

100 GB Meta

Volume meta

| Global Metadata at a system level | 8 TB Meta |

2 TB Meta

XFS Meta data

XFS Meta data

2 TB Meta

100TB Raw

100TB Raw

# Storage Data format

Meta table of the lun
(Fixed size)

Two dimensional array indexed chunk number.

Chunk number (computed with block byte range)

- Location (fast memory/chunk)
- A
- B
- C
- D
- ?

Chunk table1 (Fixed size)

Chunk table2 (Fixed size)

Block layout of the lun inside the xfs file

/xfs/outside/lun1 (sparse file)

# Storage Interface

- HardDisks
- SAS/SATA Flash
- NVMe Flash
- PCIe Flash
- S3
- Cloud Block Storage

# VSM Network Interface

- Host Networking
- VLANs / IPSpaces

# Ease of Configuration

- VSM Configuration Spec
- Infra Spec
- Integrate into K8s / EBS Compatible

# Integration to Orchestration

Options to consume the storage by containers:

- iSCSI Driver ( Pre-provisioned)
- Maya Volume Driver
- Integrated Orchestration

# Storage Connectivity - iSCSI Claims



**K8s master**

iSCSI Driver

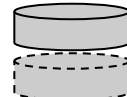HTTPS (manage)

**OpenEBS Maya master**

**K8s minions**

**Pod**

Network (Flannel)
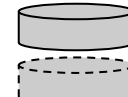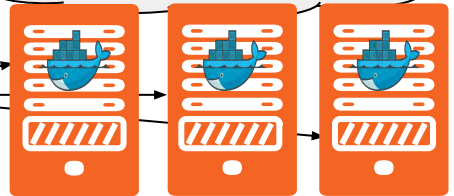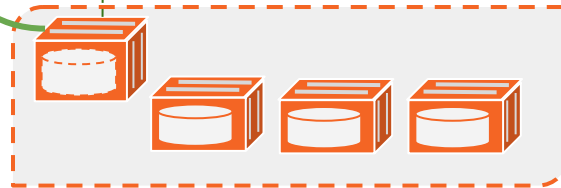
Network (Flannel*)

iSCSI

**OpenEBS VSMs / Storage Pod**

Local Storage

Remote Storage

**OpenEBS Storage Hosts**

# Storage Connectivity - Maya Driver

# Storage Connectivity - Shared Orchestration

# Resiliency and Fault Tolerance

- Scaleout
- Blue-Green Upgrades - Infra
- Rolling Upgrades - VSMs
- High-Availability

# Security

- Data Security
- Encryption
- Secure Delete

# Telemetry

- Monitoring and Troubleshooting
- Analytics

# Performance

- IO Latency
- Provisioning
- Analytics

# Scale

- Capacity
- Number of Volumes

# Deployment Flexibility

OpenEBS Deployment Options for:

- Dedicated Storage (External)
- Hyper-converged
- Hybrid-Cloud (AWS)

# OpenEBS Roadmap

- K8s Provisioning via EBS-like driver
- S3 building blocks
- Complete longhorn integration
- First usable release for community consumption

0.2 - **Basic k8s integration**

- Tiering and QoS demonstrated
- Building blocks of ML for storage analytics
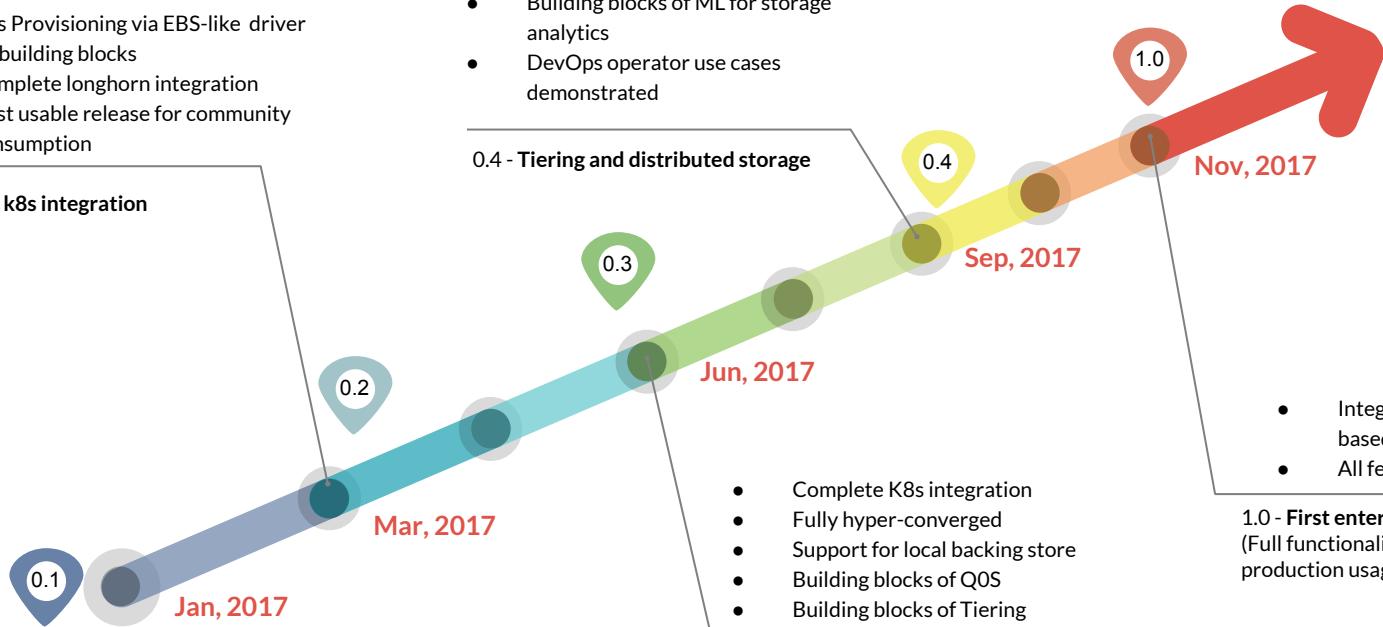- DevOps operator use cases demonstrated

0.4 - **Tiering and distributed storage**

**0.1**

**0.2**

**0.3**

**0.4**

**1.0**

**Jan, 2017**

**Mar, 2017**

**Jun, 2017**

**Sep, 2017**

**Nov, 2017**

- **Soft launch / Basic version**
- Containerized controller
- Longhorn integration basics

- Complete K8s integration
- Fully hyper-converged
- Support for local backing store
- Building blocks of Q0S
- Building blocks of Tiering

0.3 - **Full k8s integration**
(Hyper-Converged)

- Integration into Enterprise LDAP w/ role based access control
- All features supportable at scale

1.0 - **First enterprise edition**
(Full functionality for basic production usage)

# Stateful containers?!



"For which workloads or application use cases have you used/do you anticipate to use containers?"

Data Apps 77%
Cloud Apps 71%
Systems of Engagement 62%
Systems of Record 62%
Web and Commerce Software 57%
Mobile Apps 52%
Social Apps 46%

# Stateful containers?!



"For which workloads or application use cases have you used/do you anticipate to use containers?"

Data Apps 77%
Cloud Apps 71%
Systems of Engagement 62%
Systems of Record 62%
Web and Commerce Software 57%
Mobile Apps 52%
Social Apps 46%