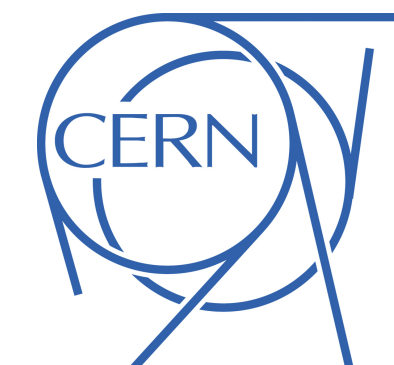# The alternative to tape based archive storage at KISTI

Sang Un Ahn[1], Latchezar Betev[2], Eric Bonfillou[2], Heejune Han[1], Jeongheon Kim[1], Seung Hee Lee[1], Bernd Panzer-Steindel[2], Andreas Joachim Peters[2], Heejun Yoon[1]

[1]KISTI, Daejeon, South Korea
[2]CERN, Geneva, Switzerland

*WLCG GDB*
*11 March 2020*

# Motivation

- Why we do?

  - Shrinking market: Tape technology mono(or bi-)poly

    ▸ One enterprise tape drive manufacturer; Two tape cartridge manufacturers

  - High cost of operating HSM for tape storage

    ▸ Commercial licenses for Spectrum Protect (TSM) and Spectrum Scale (GPFS)

    ▸ Expensive to update or upgrade - .5 Million USD @ KISTI

  - Tape operation requires own experts, not easy to find and train

- Goal

  - Replace the existing 3+ PB tape archive system with equally data-secure alternative

  - Use cheap off-the-shelf equipments and open-source storage solution
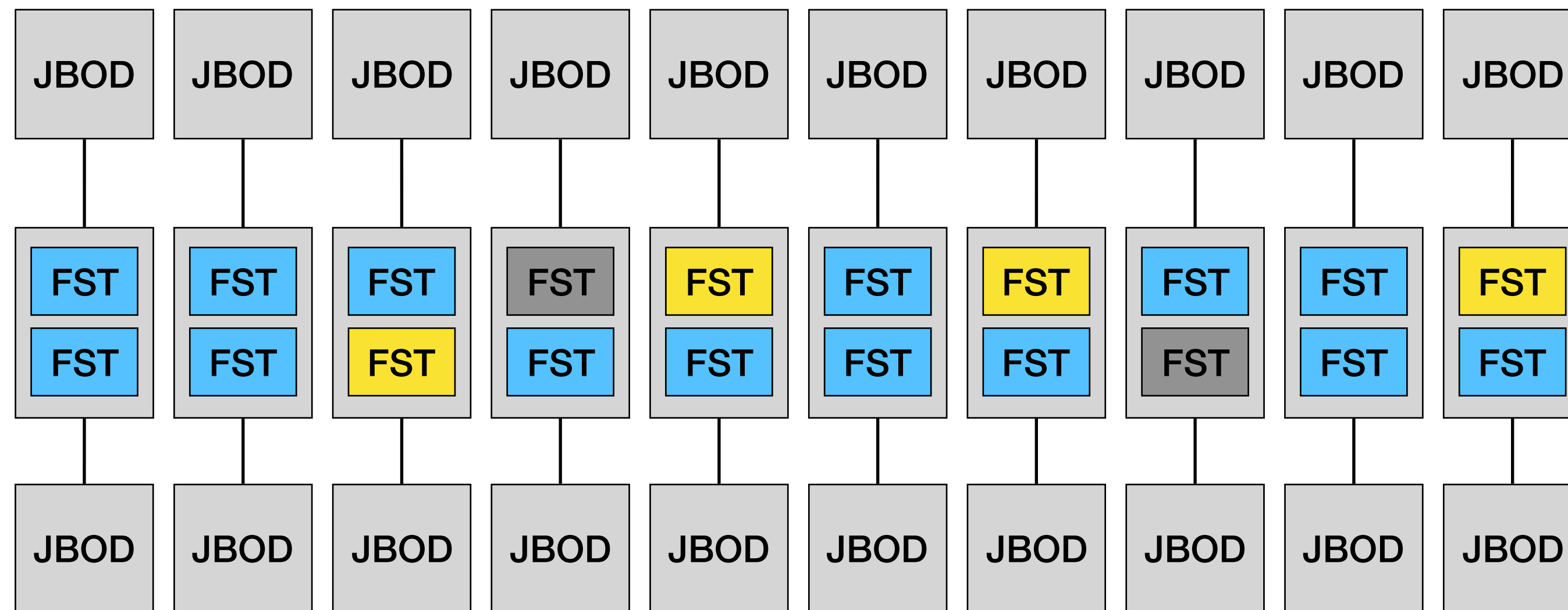
# ATAS Project

- A proposal on seeking an alternative to tape archive system approved by WLCG Overview Board (30 Nov 2018) and endorsed by ALICE

- Expert meetings in mid of February @ KISTI and in end of March @ CERN

  - Focus on design of disk-based custodial storage system

    ▸ Latest model JBODs with high density (up to 102 HDDs), 12Gb/s SAS HBAs

    ▸ Storage management through EOS

    ▸ Data protection through erasure coding RAIN

    ▸ Project budget ~1M USD

# Initial System Design

- 10 EOS front-end node, each hosts 2 EOS FSTs, each EOS FST serves 1 JBOD box

  - EOS EC (M, K) = (14, 4) to balance between usable space (77.7% of physical capacity) and data security

  - data loss probability ~ 0.000000005% (acceptable for ALICE)

- Each front-end node equipped with 2 SAS HBA cards (2 ports for each)

  - 1 HBA = 1 JBOD, SAS multi-path configuration to be tested for HA
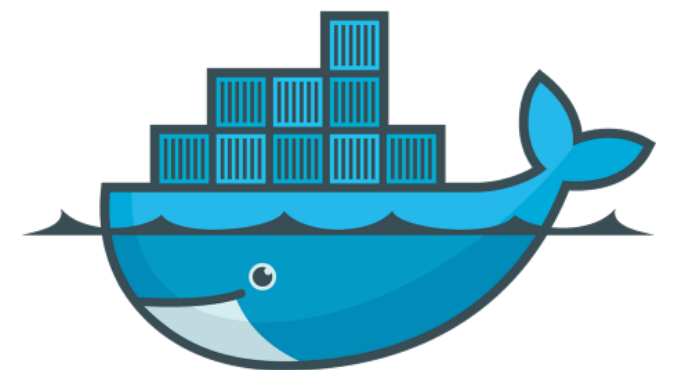
M = data node
K = parity node

## EOS RAIN6 (14,4)

| JBOD | JBOD | JBOD | JBOD | JBOD | JBOD | JBOD | JBOD | JBOD | JBOD |
|------|------|------|------|------|------|------|------|------|------|
| FST  | FST  | FST  | FST  | FST  | FST  | FST  | FST  | FST  | FST  |
| FST  | FST  | FST  | FST  | FST  | FST  | FST  | FST  | FST  | FST  |
| JBOD | JBOD | JBOD | JBOD | JBOD | JBOD | JBOD | JBOD | JBOD | JBOD |

- (x2) EOS FSTs based on Docker container
- EOS decides where to store data fragments across FST nodes randomly (no fixed scheme)
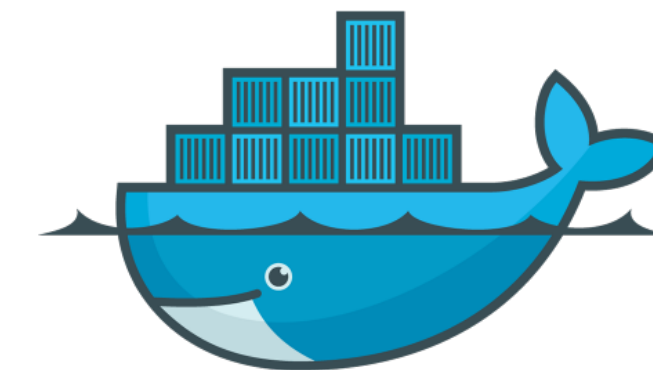
4

# EOS-Docker (1/2)

- Based on CERN EOS Docker Project + AARNet EOS Docker

  - Dumb-init + bash-script-with-exec to forward properly SIGTERM signal when docker stop command issued

- Ansible configuration in "2-pass"

  - Pass 1: Pre run to auto-generate inventory, group_vars, host_vars from config.yaml

  - Pass 2: Full run using the auto-generated configs

- Not ready for orchestration (e.g. Kubernetes); optimal network equipments and configurations are required

- Host container network to share the host IP among containers

- Available on GitHub

  - https://github.com/jeongheon81/gsdc-eos-docker

- Consul : health check and DNS service, event trigger (experimental)

- Simple log collection : Grafana + Loki + Promtail (all-in-one setup)

- Simple monitoring : Cockpit (cluster mode)

- Others : chrony, firewalld, network (nmcli)

Consul



Grafana + Loki



Promtail



Cockpit

# DEMO Equipment & Setup

- JBOD: DELL PowerVault ME484

  - Disk: 70EA (HGST 12TB 7.2k NL-SAS), 840 TB

- Front-end Server: DELL PowerEdge R640

  - CPU: Intel Xeon Scalable 6150 2.7GHz 18 core * 2EA

  - Memory: DDR4 16GB 2666MHz * 24EA

  - HBA: DELL PowerEdge 12Gbps SAS HBA (FW version: 16.17.00.03)

  - NIC: QLogic 4x10GE QL41164HMCU CNA

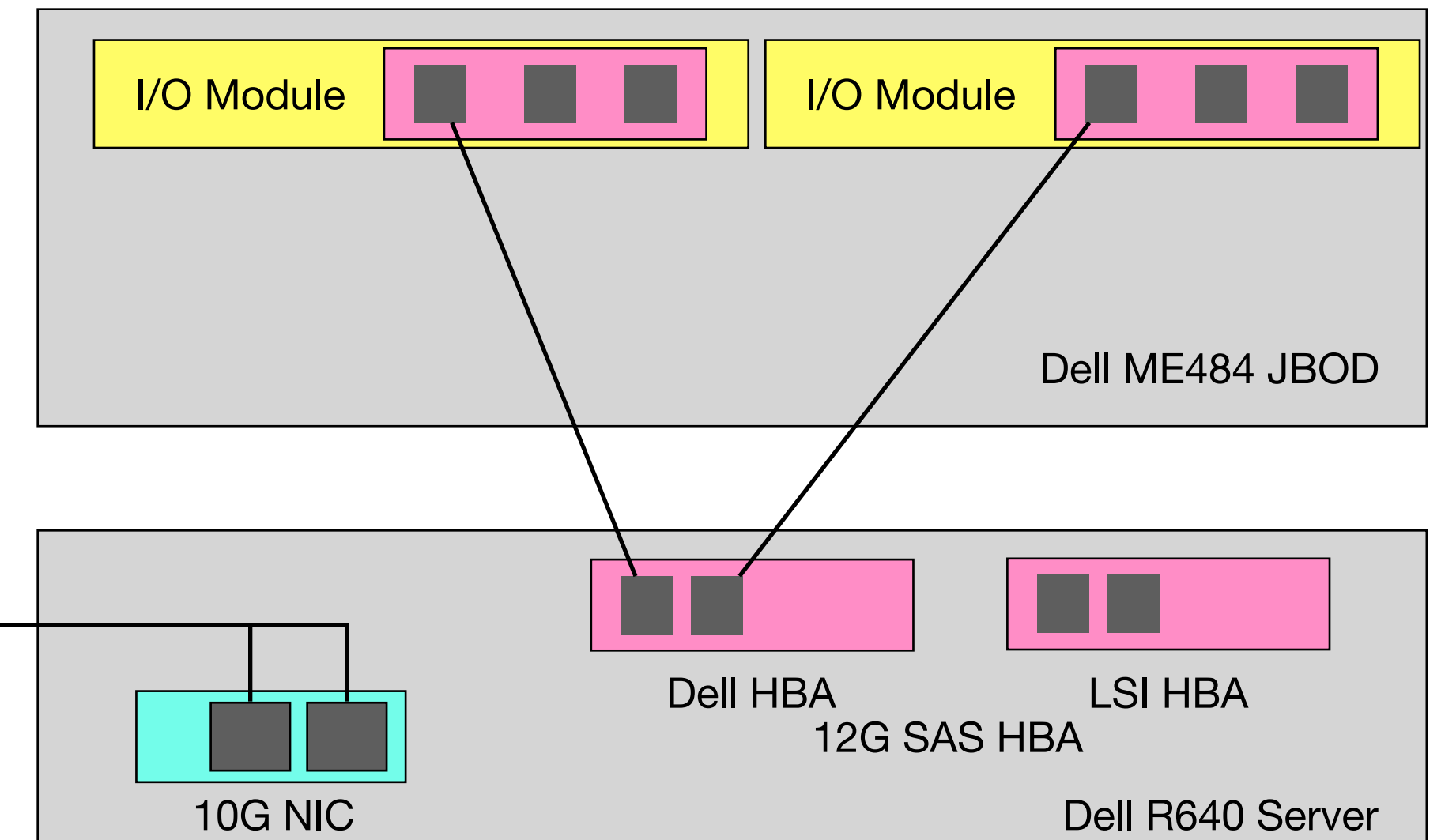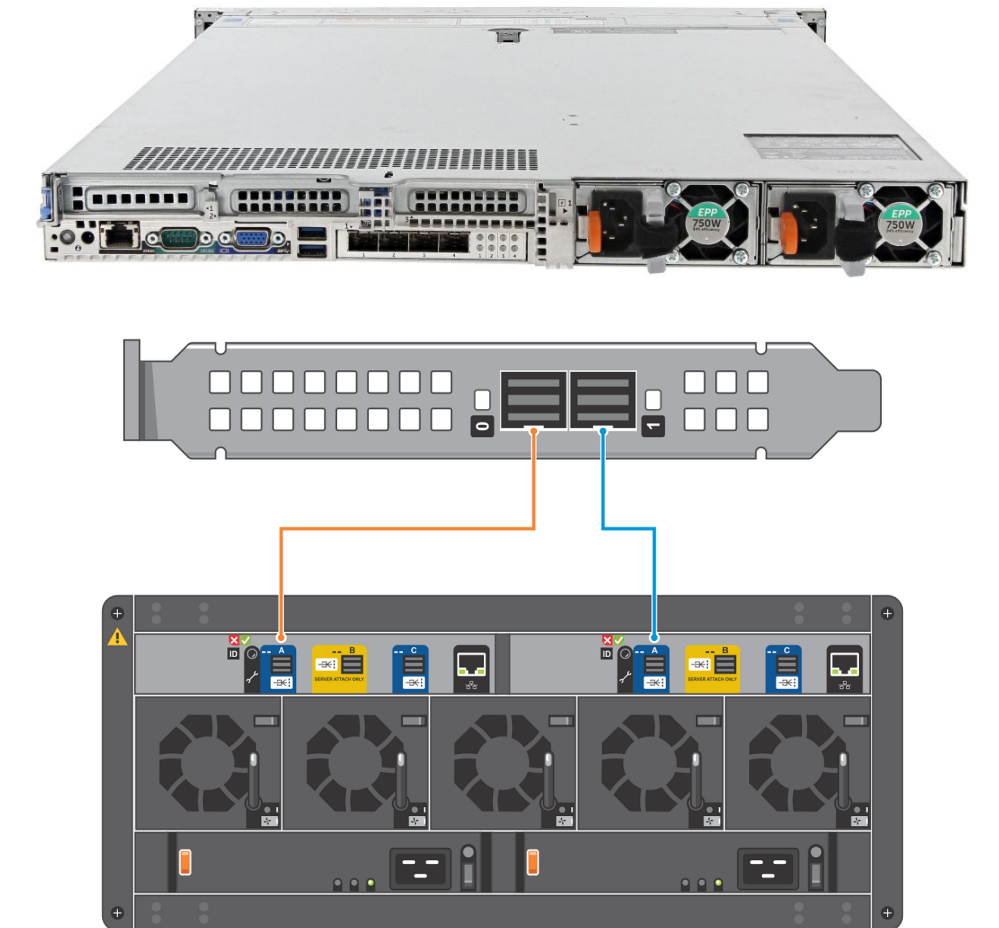**System Information**

**Operating System**
CentOS Linux 7 (Core)

**Operating System Kernel Version**
7 (Core) Kernel 3.10.0.-957.el7.x86_64

**BIOS Version**
1.5.6

**Filesystem: XFS (Default EL7 Distribution)**



I/O Module

I/O Module

Dell ME484 JBOD

UpLink

10G NIC

Dell HBA

LSI HBA

12G SAS HBA

Dell R640 Server

# Performance Test Results

## I/O Test: Multipath Mode

- Multipath mode: failover (active-standby) vs. multibus (active-active)

  - **multibus** mode showed the maximum I/O speed up to 6GB/s for read/write

    ▸ Bottleneck on PCIe 3.0 (6400MB/s)

  - **failover** could not fulfill the available bandwidth, limited under 1 SAS port (48Gb) pipe

I/O Test: VDBench

**XFS, failover, VDBench**

4GB/s ~ 1 SAS port (4800MB/s)

Read/Write 4GB/s

5000 / 3750 / 2500 / 1250 / 0

4K 8K 16K 32K 64K 128K 256K 512K 1024K 2048K 4096K

◇ VDBench Read   ◇ VDBench Write

**XFS, multibus, VDBench**

Write 6.2GB/s
Read 5.8GB/s

7000 / 5250 / 3500 / 1750 / 0

4K 8K 16K 32K 64K 128K 256K 512K 1024K 2048K 4096K

◇ VDBench Read   ◇ VDBench Write

18

## I/O Test: Read/Write

- XFS read/write performance (simultaneous read and/or write from 70 disks)

  - **VDBench** shows full read/write transfer performance @ transfer size >= 2048k (6GB/s)

  - **IOZone** shows full read/write transfer performance @ transfer size ~ 2048k (6GB/s)

Disk: 70EA
Filesize: 2GB

**IOZone & VDBench Read/Write Tests**

Write 6.2GB/s
Read 5.8GB/s

7000 / 5250 / 3500 / 1750 / 0

4K 8K 16K 32K 64K 128K 256K 512K 1024K 2048K 4096K

◇ IOZone Read   ◇ IOZone Write   ◇ VDBench Read   ◇ VDBench Write

**Read:Write = 95:5 Scenario**

Read 5.2GB/s
Write 0.3GB/s

7000 / 5250 / 3500 / 1750 / 0

4K 8K 16K 32K 64K 128K 256K 512K 1024K 2048K 4096K

◇ IOZone Read   ◇ IOZone Write   ◇ IOZone Random-Read
◇ IOZone Random-Write   ◇ VDBench Read   ◇ VDBench Write

* IOZone tests with different Read/Write ratio Scenario did not much affect on the performance

19

## Power Consumption

- JBOD Test Equipment (70 Disks)

  - JBOD (DELL ME484): idle = 830W; load = 860W (Max 960) **(1.12W/TB)**

  - Server: idle = 200W; load = 270W

  - Switch: idle = 246W; load = 246W

  - **1.75W/TB** including JBOD, Server and Switch

- Disk Storages (Full Load)

  - DellEMC SC7020, 2.5PB  - 12,120W **(4.8W/TB)**

  - EMC Isilon, 16 Nodes, 2.95 PB- 13,730W **(4.6W/TB)**

  - EMC VNX, 12 Nodes, 2.36 PB - 5,100W **(2.2W/TB)**

  - HITACHI VSP, 2 PB - 18,300W **(9.15W/TB)**

  - EMC Isilon, 15 Nodes, 1.43 PB - 12,880W **(9W/TB)**

  - EMC CX4-960, 1.5PB - 14,900W **(9.9W/TB)**

- Tape Library (Full Load)

  - **IBM TS3500 5-Frame (3.2PB) - 1,600W (0.5W/TB)**

L3(P3)

Rack Switch

[1] PDU Unit.Power.Active.Value 246 W
Timestamp        05.10.

L2(P2)

x86 Server

[1] PDU Unit.Power.Active.Value 272 W
Timestamp        05.10.

L1(P1)

JBOD

[1] PDU Unit.Power.Active.Value 874 W
Timestamp        05.10.

20

- Confirmed the upper cap of read/write performance ~ 6GB/s (intrinsic limit by PCIe 3.0)

- Power consumption shown ~ 1.75W/TB, not uncomfortably higher than Tape (0.5W/TB)

  - High-end Enterprise Disk Storage 5 ~ 9W/TB

# Procurement Schedule



| | | | | | |
|---|---|---|---|---|---|
| 1 week | 2 weeks | 1.5 weeks | 10 weeks | | 2 weeks |

**Week**

Pre-annoucnement
19' Sep 9

Main announcement
19' Sep 24

Bid
19' Oct 9

Delivery & Installation
19' Oct 18

Finished
Install & BMT
19' Dec 26

Setup
EOS & Docker
20' Jan 10

Finished
EOS & Docker
Setup
20' Jan 23

# Delivery & Installation

- Dec 17th ~ 27th

# Delivered Systems

**5U Storage Enclosure**

SAS Ports          SAS Ports

I/O Module      I/O Module

Uplink ◄

40G NIC

**5U Storage Enclosure**

I/O Module      I/O Module

SAS Ports          SAS Ports

## Specifications

- x9 Servers: Dell PowerEdge R730
    - Intel(R) Xeon(R) CPU E5-2637 v4 @ 3.50GHz * 2EA
    - DDR4 16GB 2,400MHz * 12EA (192GB)
    - Dell PowerEdge 12Gbps SAS HBA * 4EA
    - MLNX 40Gb 2P ConnectX3Pro Adpt * 4EA
- x18 JBOD: Dell PowerVault ME484
    - 84EA HGST or Seagate 12TB 7.2K NL-SAS)
- x2 40G network switches

ATAS Storage
18PB (RAW)
14PB (Usable)

- Confirmed the upper cap of read/write performance ~ 6GB/s (intrinsic limit by PCIe 3.0)
- Power consumption will be measured throughout the whole testing & commissioning periods

# JBOD Cabling

- Recommended

  - One Server / Four HBA per server / Single path

- Target

  - One Server / Four HBA per server / Dual path



- Disk Access, Recognition Test via Multipath
- Data Consistency, Corruption Test
- Read/Write Performance Test

# Network Topology

# Schedule

| Tasks | 2019 | | | | | | | | | | | | 2020 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 1 | 2 | 3 |
| Technology Search | ■ | | | | | | | | | | | | | | |
| Product Survey | | ■ | | | | | | | | | | | | | |
| Design & Specification | | | ■ | ■ | | | | | | | | | | | |
| Testing | | | | | ■ | ■ | ■ | | | | | | | | |
| Procurement | | | | | | ■ | ■ | ■ | ■ | ■ | ■ | | | | |
| Installation & Validation | | | | | | | | | | | ■ | ■ | | | |
| EOS Deployment Test | | | | | | | | | | | | | ■ | ■ | ■ |

KISTI-CERN Expert Meeting @ KISTI

KISTI-CERN Expert Meeting @ CERN

EOS Workshop @ CERN

Call for tender (delayed) → Delivery

Today

RAIN Layout

- Change of procurement planning had been approved in May by National Facility & Equipment Committee
- Call for tender delayed due to change of procurement procedure (technical pre-estimation included)
- Commissioning and production test in 2020 targeting the production service before RUN3

# Conclusions

- We are investigating a disk-based storage, using standard JBODs and EOS with erasure coding, as an alternative to tape-based custodial storage

- Obvious benefits: avoid single-vendor dependency, common expertise for all storage systems across the computing centre

- A final system unit I/O tests show ~6GB/s read/write performance, as expected from the limits of the PCIe 3.0 and SAS 12Gb/s HBA

- Power consumption is shown to be 1.75W/TB, not uncomfortably higher than a tape library

- Procurement and system deliveries finished in November 2019, installation and validation on delivered systems finished in January 2020

- EOS deployment with RAIN layout is being applied and tested after the recent EOS workshop in February, special thanks to EOS developers!!

- During the whole year of 2020, this disk-based custodial storage will be tested and verified repeatedly with ALICE targeting the production service before the start of RUN3

# Questions?

# Backup

# Concerns about Tape Market

- One enterprise tape drive manufacturer, two tape cartridge manufacturers

- Oracle enterprise tape drive

  - https://www.theregister.co.uk/2017/02/17/oracle_streamline_tape_library_future/

- Concerning steady tape cartridge supply, tape suppliers shrunk over the past three years from six to two - Sony, Fujifilm

  - https://www.bloomberg.com/news/articles/2018-10-17/the-future-of-the-cloud-depends-on-magnetic-tape

- Patent dispute between Sony and Fujifilm => No LTO-8 supply available globally

  - https://www.theregister.co.uk/2019/05/31/lto_patent_case_hits_lto8_supply/

  - https://www.theregister.co.uk/2019/08/06/sony_fujifilm_storage_patent_lawsuit_settled/

- Sony, Fujifilm stopped patent dispute (however not officially announced from both sides) at the end of July, starting production of LTO-8 media

- Disk = $25/TB, Tape = $10/TB, SSD $100/TB (QLC), SpectraLogic 2019 Report

# Data Loss Probability

Data loss probability

$$p = e^{-\lambda}\frac{\lambda^k}{k!}$$

where

$$\lambda = \frac{AFR \times (Number\ of\ Disks)}{365 \times 24 \div MTTR}$$

MTTR = Mean Time To Repair
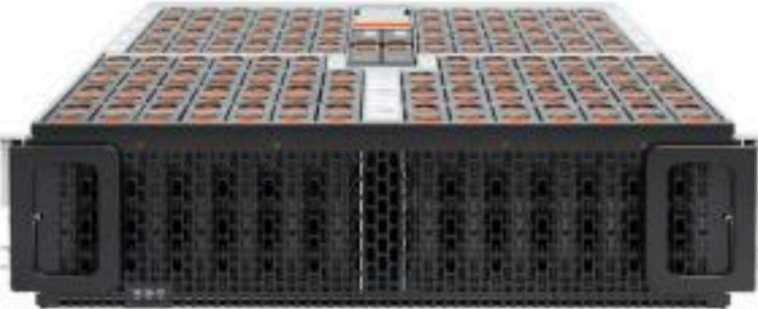AFR = Annualized Failure Rate

Assuming 1680 disks, 2% of AFR and 24h of MTTR, one can have λ = 0.092 so with 4 parity disks the data loss probability $p$ gives,

$$p = e^{-0.092}\frac{0.092^5}{5!} = 0.000000050242575 = 5.02 \times 10^{-9}$$

# High Density JBOD Products

| Image |  |  |  |  |  |  |
|---|---|---|---|---|---|---|
| Model | Dell EMC PowerVault ME484 | HPE D6020 | QCT JB4602 JB9T | WD Ultrastar Data102 H4102-J SE4U102-102 | WD Ultrastar Data60 H4060-J SE4U60-60 | Promise VTrak J5800S |
| Unit | 5U | 5U | 4U | 4U | 4U | 4U |
| Disk | 12TB | 12TB | 12TB | 12TB | 12TB | 12TB |
| # Disks | **84** | **70** | **60** | **102** | **60** | 24 |

- Note that each JBOD enclosure has different dimensions depending on its unit and the number of disk drives to mount
- Proprietary SAS HBA cards shipped with x86 server may not provide enough compatibility to other JBOD products
- JBOD enclosures with RAID controller to provide hardware-level data protection are available in the market

# State-of-the-art SAS HBA

<u>3rd Generation</u>

● Broadcom (Avago, LSI) SAS 9300 16(8)-port 12Gb/s SAS HBA

- IO Controller: Two I/O controller

- PCI Data Burst Transfer Rates: Half Duplex, 19200MB/s

- Device support: <u>1024 non-RAID devices</u>

In case of 4 ports

Allowing the transmission of signals in both directions but not simultaneously

# Design Limitation Study

- In case of direct attached storage, PCIe 3.0 is the bottleneck

  - Third generation 12Gb/s SAS

  - Typical HDD transfer rate : 230MB/s for 15k, 100MB/s ~ 170MB/s for slower

  - Theoretical burst of PCIe 3.0 is about 8000MB/s while typical number is 6400MB/s (80% efficiency)

SAS Two Ports
4 Lane each port
1 Lane = 12Gb/s
∴ 48Gb/s or 4800MB/s (per port)
Total bandwidth = 9600MB/s

| Configuration | Bottleneck (MB/s) | # of HDDs | # of SSDs |
|---|---|---|---|
| 6Gb/s SAS x4 / PCIe 2.x | SAS (2200) | 9 | 4 |
| 6Gb/s SAS x8 / PCIe 2.x | PCIe (3200) | 14 | 6 |
| 12Gb/s SAS x4 / PCIe 2.x | PCIe (3200) | 14 | 6 |
| 12Gb/s SAS x4 / PCIe 3.0 | SAS (4400) | 19 | 8 |
| 12Gb/s SAS x8 / PCIe 3.0 | PCIe (6400) | 28 | 12 |

For 15k HDD (~230MB/s)

56 slower disks can fulfill
the bandwidth provided by
Two port 12Gb/s SAS HBA card
connected to a PCIe 3.0 slot

**Table 4 – Sample storage configurations showing each one's bottleneck and the number of drives supported at their peak throughput**

*https://docs.broadcom.com/docs/12353459*

# I/O Test: Multipath Mode

- Multipath mode: failover (active-standby) vs. multibus (active-active)

  - **multibus** mode showed the maximum I/O speed up to 6GB/s for read/write

    ‣ Bottleneck on PCIe 3.0 (6400MB/s)

  - **failover** could not fulfill the available bandwidth, limited under 1 SAS port (48Gb) pipe

I/O Test: VDBench



XFS, failover, VDBench

4GB/s ~ 1 SAS port (4800MB/s)

Read/Write
4GB/s

VDBench Read    VDBench Write

XFS, multibus, VDBench

Write
6.2GB/s

Read
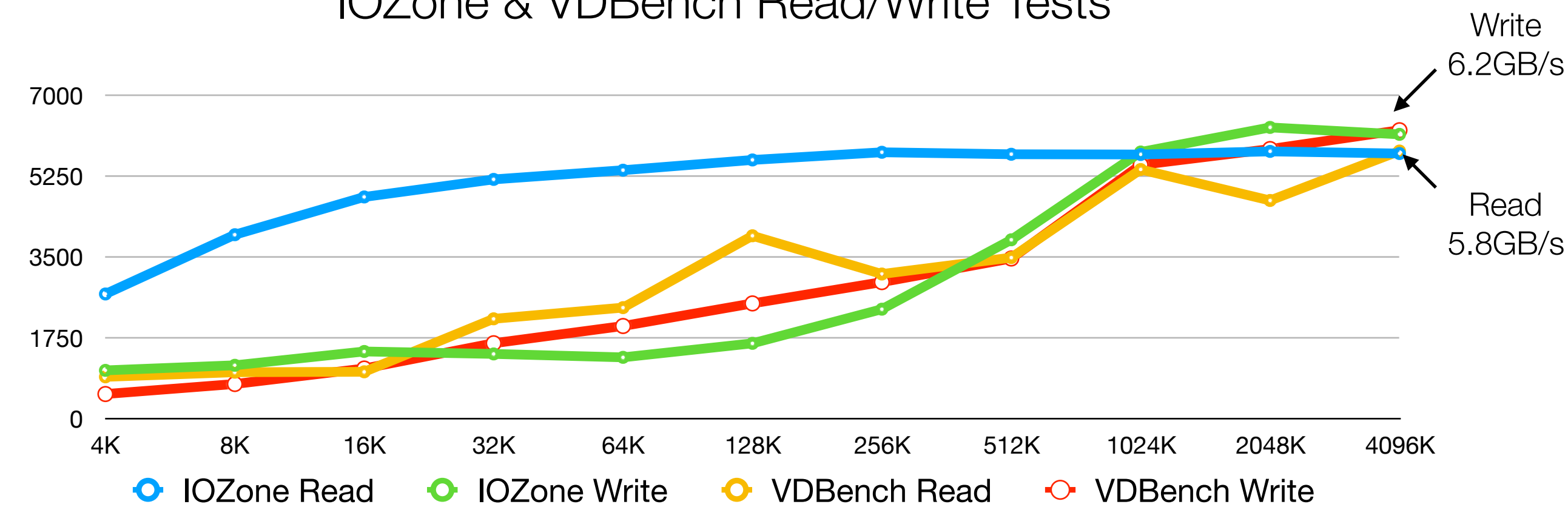5.8GB/s

VDBench Read    VDBench Write
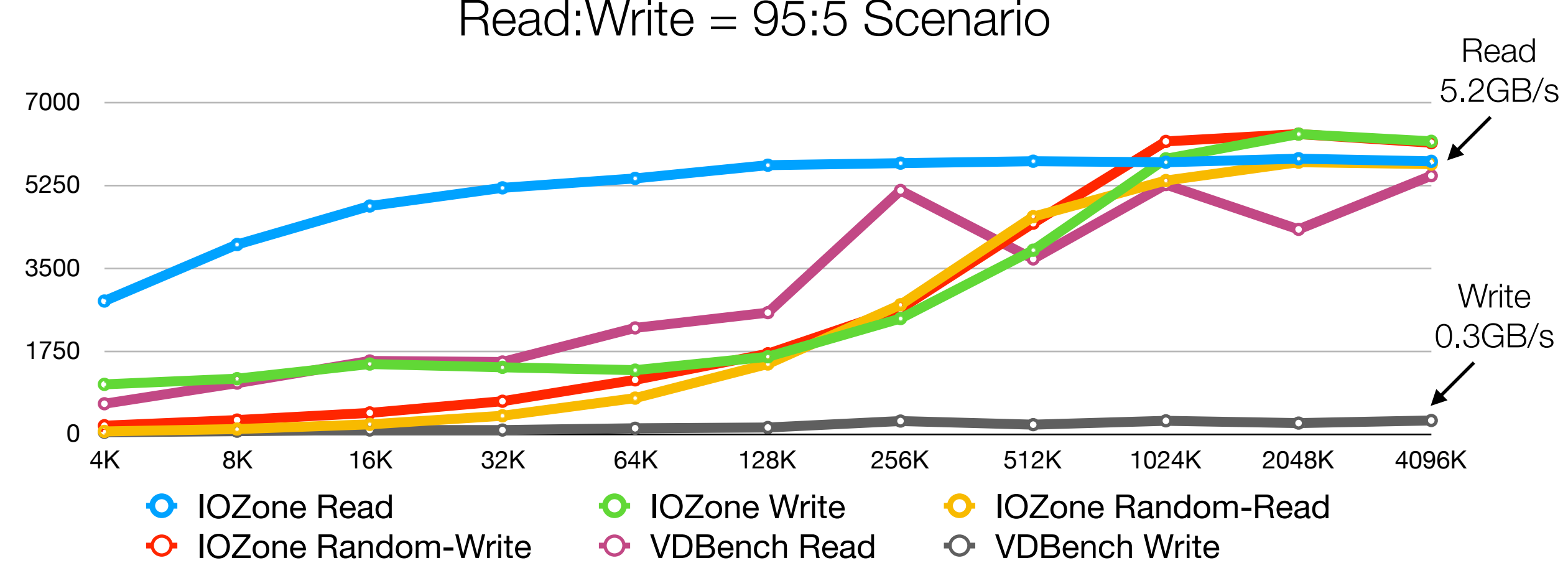
# I/O Test: Read/Write

- XFS read/write performance (simultaneous read and/or write from 70 disks)

  - **VDBench** shows full read/write transfer performance @ transfer size >= 2048k (6GB/s)

  - **IOZone** shows full read/write transfer performance @ transfer size ~ 2048k (6GB/s)

Disk: 70EA
Filesize: 2GB



IOZone & VDBench Read/Write Tests

Write 6.2GB/s
Read 5.8GB/s

- IOZone Read
- IOZone Write
- VDBench Read
- VDBench Write

Read:Write = 95:5 Scenario

Read 5.2GB/s
Write 0.3GB/s

- IOZone Read
- IOZone Write
- IOZone Random-Read
- IOZone Random-Write
- VDBench Read
- VDBench Write

* IOZone tests with different Read/Write ratio Scenario did not much affect on the performance

# Power Consumption

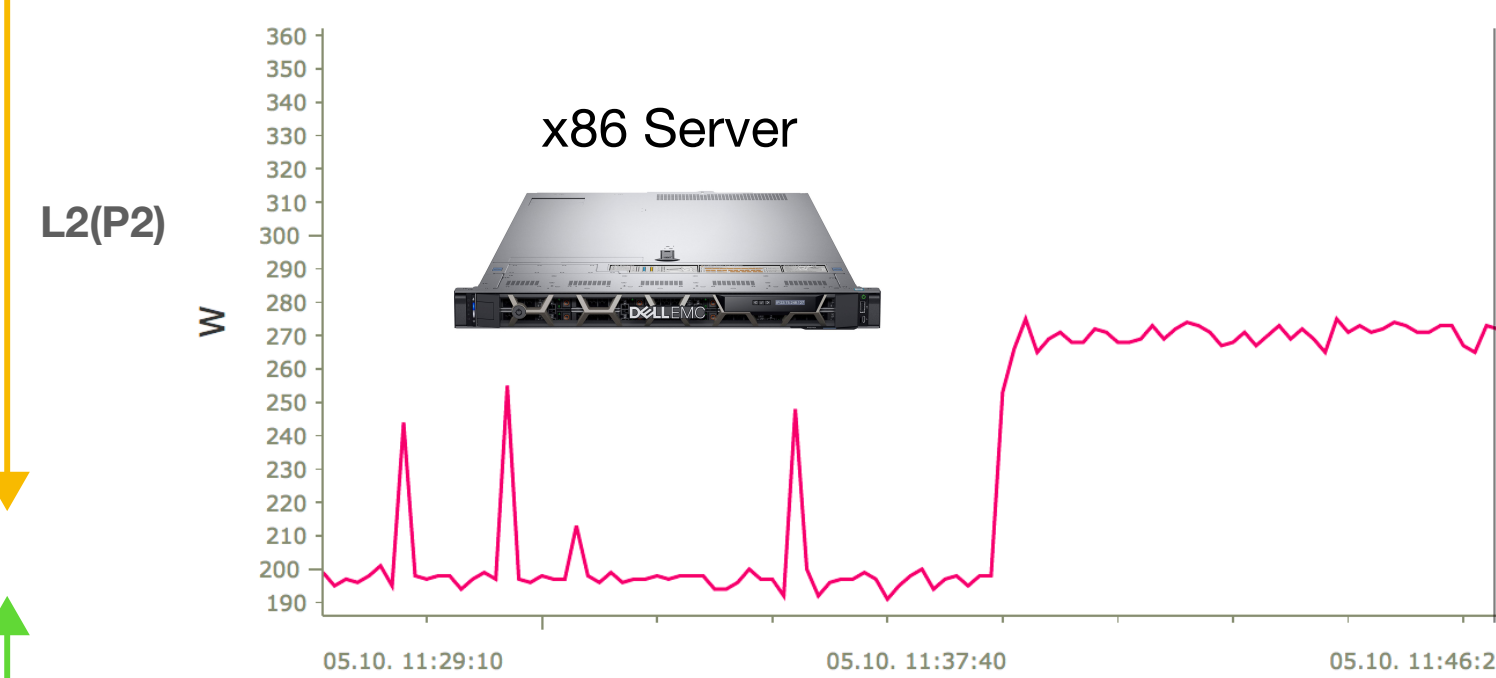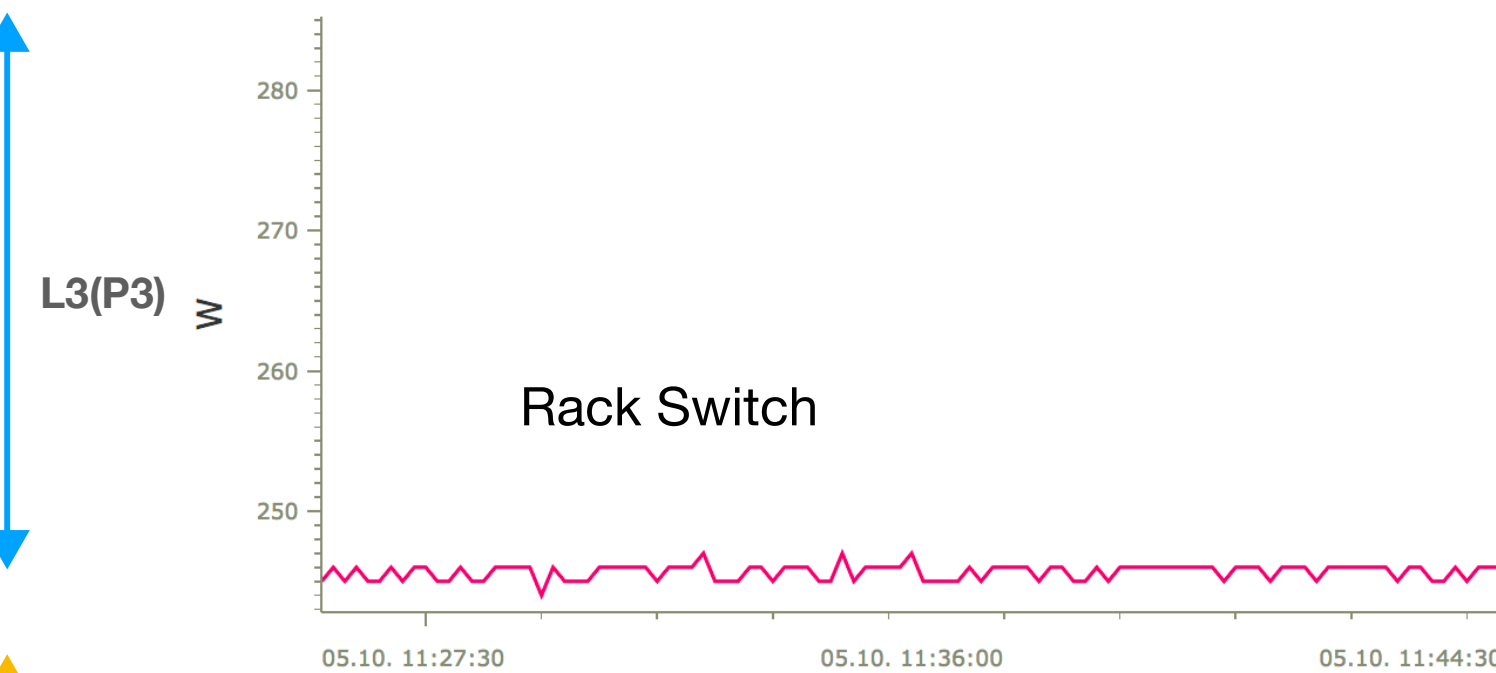- JBOD Test Equipment (70 Disks)

  - JBOD (DELL ME484): idle = 830W; load = 860W (Max 960) **(1.12W/TB)**

  - Server: idle = 200W; load = 270W

  - Switch: idle = 246W; load = 246W

  - **1.75W/TB** including JBOD, Server and Switch

- Disk Storages (Full Load)

  - DellEMC SC7020, 2.5PB - 12,120W **(4.8W/TB)**

  - EMC Isilon, 16 Nodes, 2.95 PB- 13,730W **(4.6W/TB)**

  - EMC VNX, 12 Nodes, 2.36 PB - 5,100W **(2.2W/TB)**

  - HITACHI VSP, 2 PB - 18,300W **(9.15W/TB)**

  - EMC Isilon, 15 Nodes, 1.43 PB - 12,880W **(9W/TB)**

  - EMC CX4-960, 1.5PB - 14,900W **(9.9W/TB)**

- Tape Library (Full Load)

  - **IBM TS3500 5-Frame (3.2PB) - 1,600W (0.5W/TB)**



L3(P3)

Rack Switch

[1] PDU Unit.Power.Active.Value 246 W

Timestamp                  05.10.

L2(P2)

x86 Server

[1] PDU Unit.Power.Active.Value 272 W

Timestamp                  05.10.

L1(P1)

JBOD

[1] PDU Unit.Power.Active.Value 874 W

Timestamp                  05.10.

# Present Network Diagram