# Optimizing FUSE for Cloud Storage

**|| Parallels**™

*Profit* from the Cloud

## Maxim Patlasov
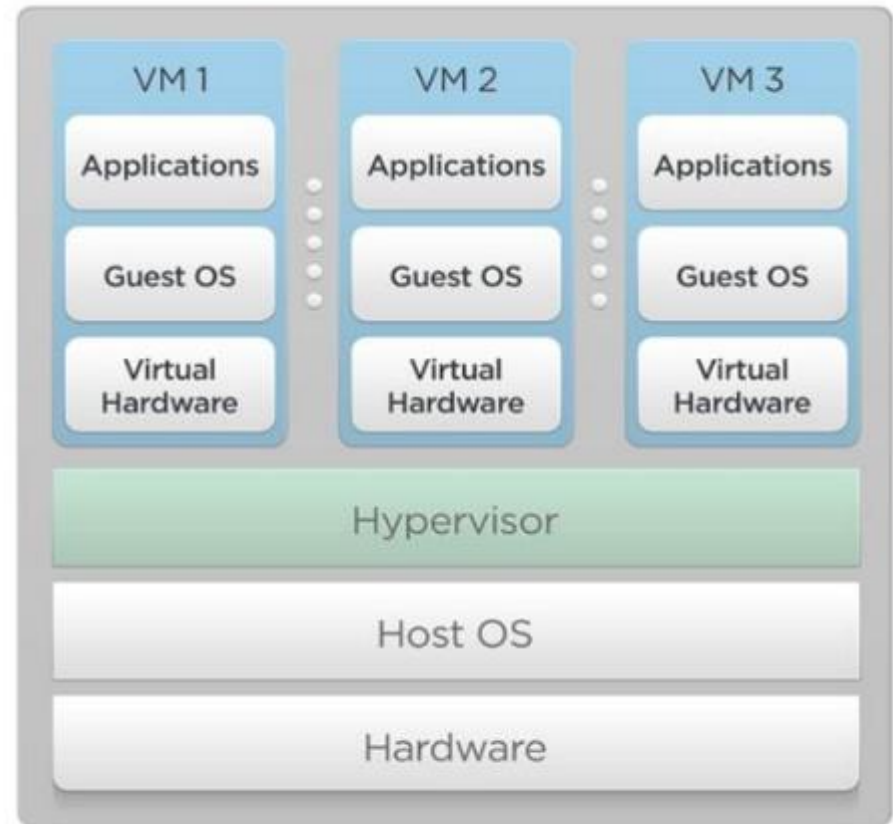Linux Kernel Developer, Parallels Inc.

Vault   2015

# Agenda

1. Parallels Cloud Storage
2. FUSE concept
3. FUSE optimizations
4. Performance achieved
5. Future improvements

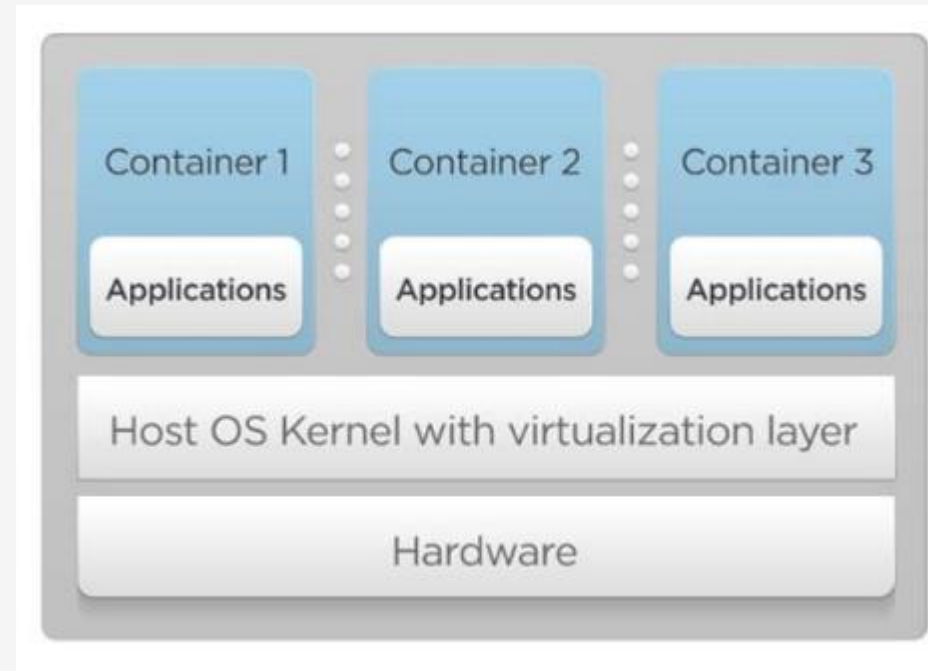# Parallels Cloud Storage

# **Parallels Hypervisor Virtualization**

- OS flexibility
- HW emulation
- Bare metal installation

# Parallels Containers

- More efficient memory management
- More efficient caching reduces I/O
- CT resource management
- Easy migration
- Easy backups&snapshots

| Container 1 | Container 2 | Container 3 |
| --- | --- | --- |
| Applications | Applications | Applications |

Host OS Kernel with virtualization layer

Hardware

# Storage requirements
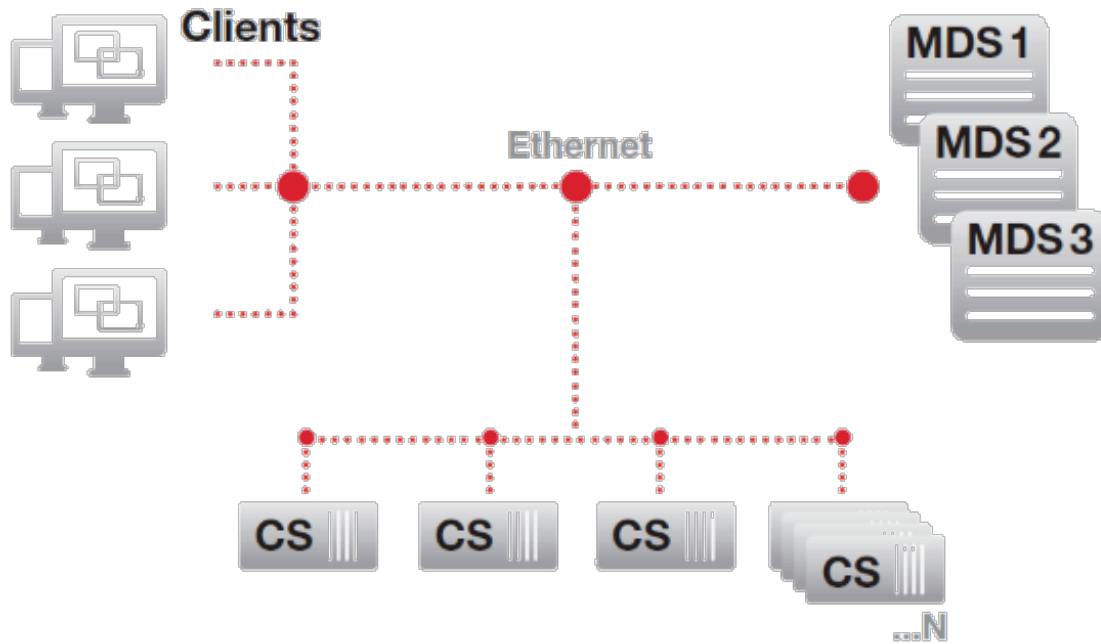
Key requirements for VM and Container's needs:

- Strong consistency
- High performance
- Fault tolerance
- Fast recovery
- Address all space from any node
- Commodity hardware
- In-flight reconfiguration and update

# Parallels Cloud Storage solution

Key decisions made:

- Optimize for big files
- Union of all local storages
- Replication for fault tolerance
- Keep data and metadata separately
- Multiple metadata servers

# PStorage architecture



**Clients**

**Ethernet**

**MDS 1**
**MDS 2**
**MDS 3**

**CS**  **CS**  **CS**  **CS**
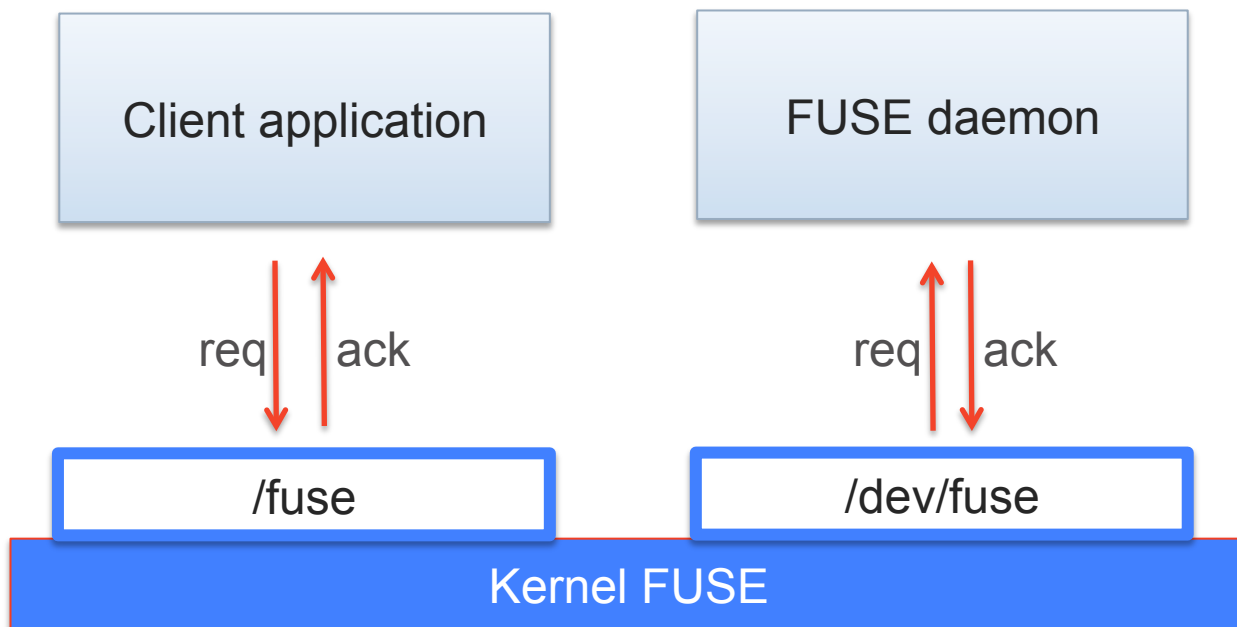...N

**Meta Data Server (MDS)**
- Stores metadata in memory
- Tracks data chunks and their versions
- Is highly available
- Can run on the same server as the chunk server and client

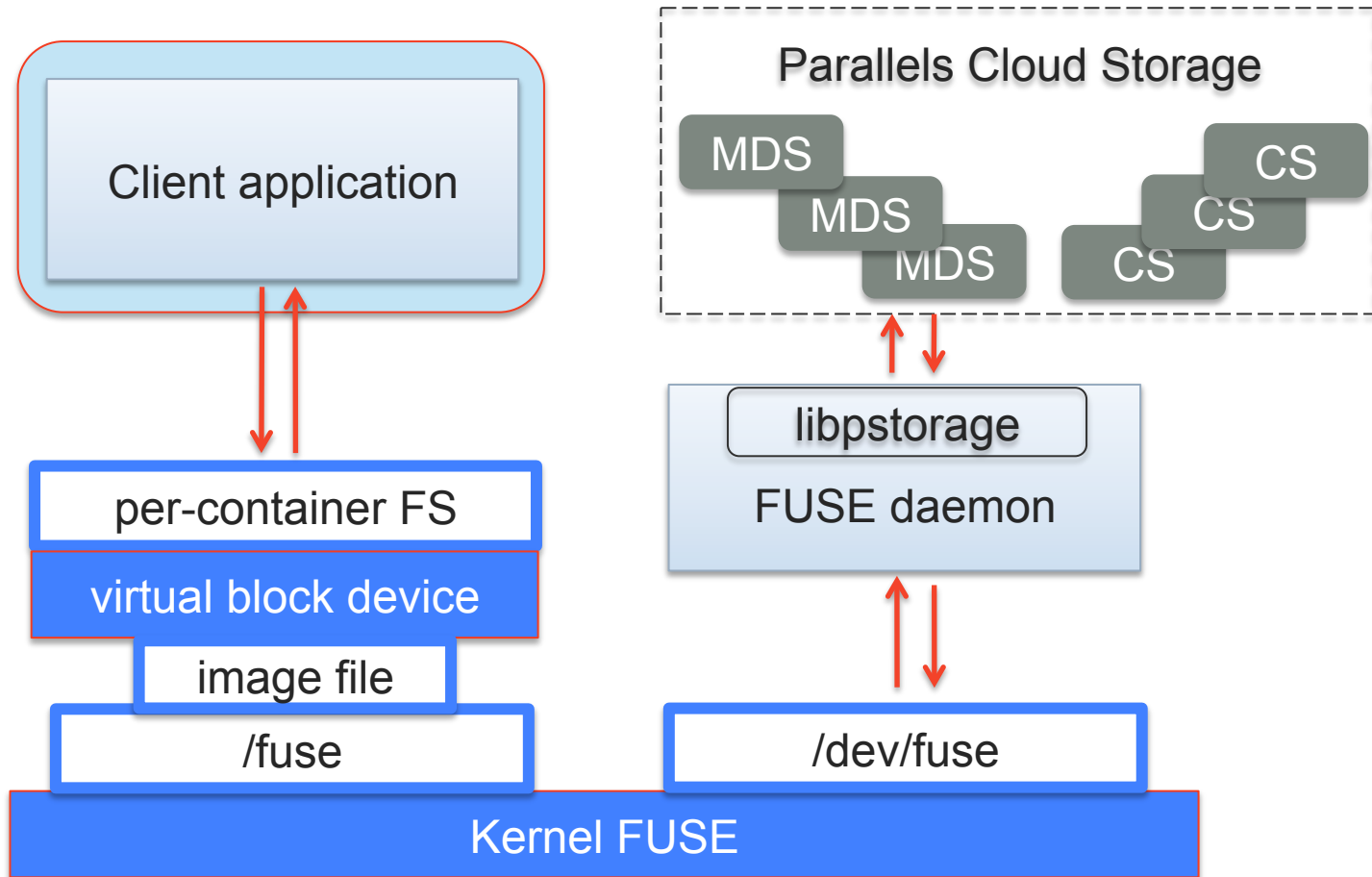**Chunk Server (CS)**
- Stores data chunks
- Manages data chunks
- Performs read/write operations on data chunks
- Can run on the same server as the client

|| Parallels™

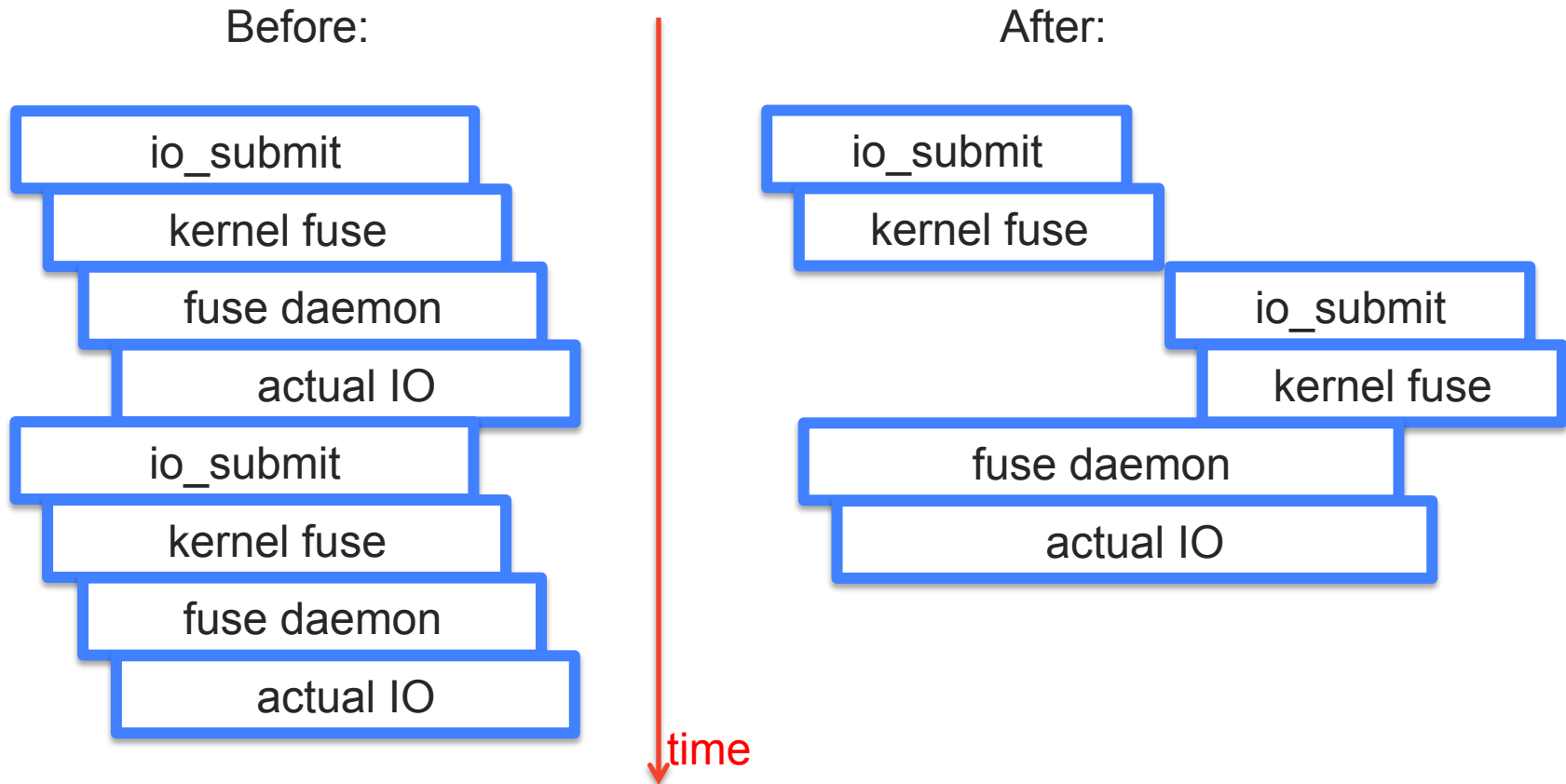*Profit* from the Cloud

# FUSE

# FUSE framework

# FUSE: Containers on PStorage

# FUSE optimizations

# Asynchronous direct IO

Application: io_submit(&iocb1); io_submit(&iocb2);

Before:

After:



Before (left column):
- io_submit
- kernel fuse
- fuse daemon
- actual IO
- io_submit
- kernel fuse
- fuse daemon
- actual IO

After (right column):
- io_submit
- kernel fuse
- io_submit
- kernel fuse
- fuse daemon
- actual IO

time

Profit from the Cloud

# Synchronous direct IO

Application: fd = open(O_DIRECT); write(fd, buf, 1<<20);

Before:

8x:

| write |
| kernel fuse: 128K |
| fuse daemon |
| actual IO |

. . .

| kernel fuse: 128K |
| fuse daemon |
| actual IO |

After:

| write |
| kernel fuse: 8 x 128K |
| fuse daemon |
| actual IO |

time

# Writeback cache

Application: buffered write(fd, buf, 1<<20);

Before:

After:

8x:

| write |
| kernel fuse: 128K |
| fuse daemon |
| actual IO |

. . .

| kernel fuse: 128K |
| fuse daemon |
| actual IO |

| write |
| kernel fuse: populate page cache |
| kernel writeback |
| kernel fuse: 8 x 128K |
| fuse daemon |
| actual IO |

time

Key benefits:
- Lower latency of write(2)
- Parallel processing writeback

# Performance Comparison :: HW



## iSCSI SAN Storage
### DELL EqualLogic PS6510E

x1 HW SAN EQL PS6510E

**48 SATA Disks**: 1TB 7200rpm
(Seagate ST31000524NS)

Network: **10Gbit**
(Dell Force10 S4810)

vs.

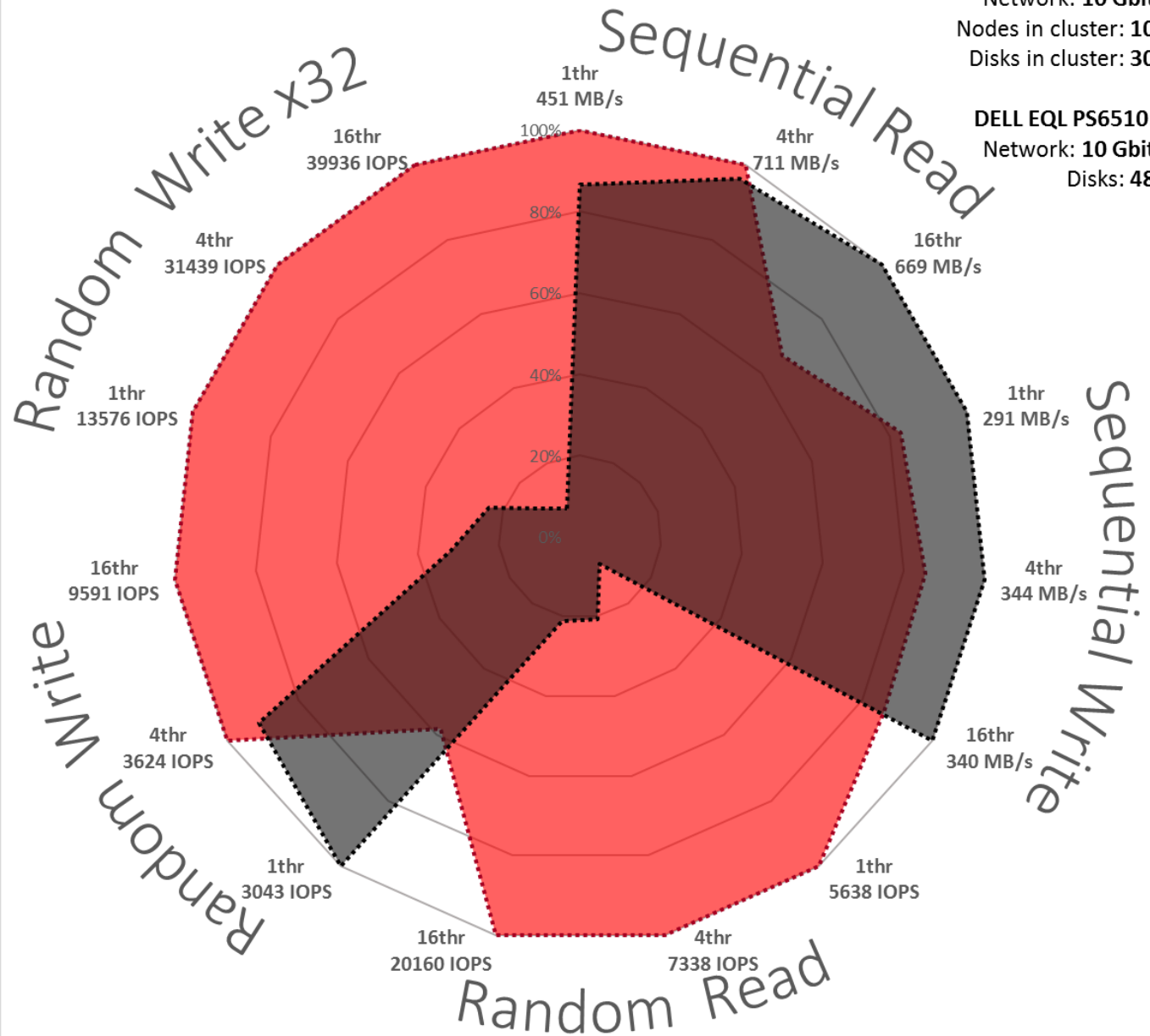## Parallels Cloud Storage
### (FUSE based)

x10 compute nodes

**30 SATA Disks**: 2TB 7200rpm
(Seagate ST2000DM001)
+ 10 SSD for caching
(Intel SSD 520)

Network: **10Gbit**
(Brocade FastIron SuperX SX-F42XG)

Profit from the Cloud

PERFROMANCE MAP

PCS 6.0:
Network: **10 Gbit**
Nodes in cluster: **10**
Disks in cluster: **30**
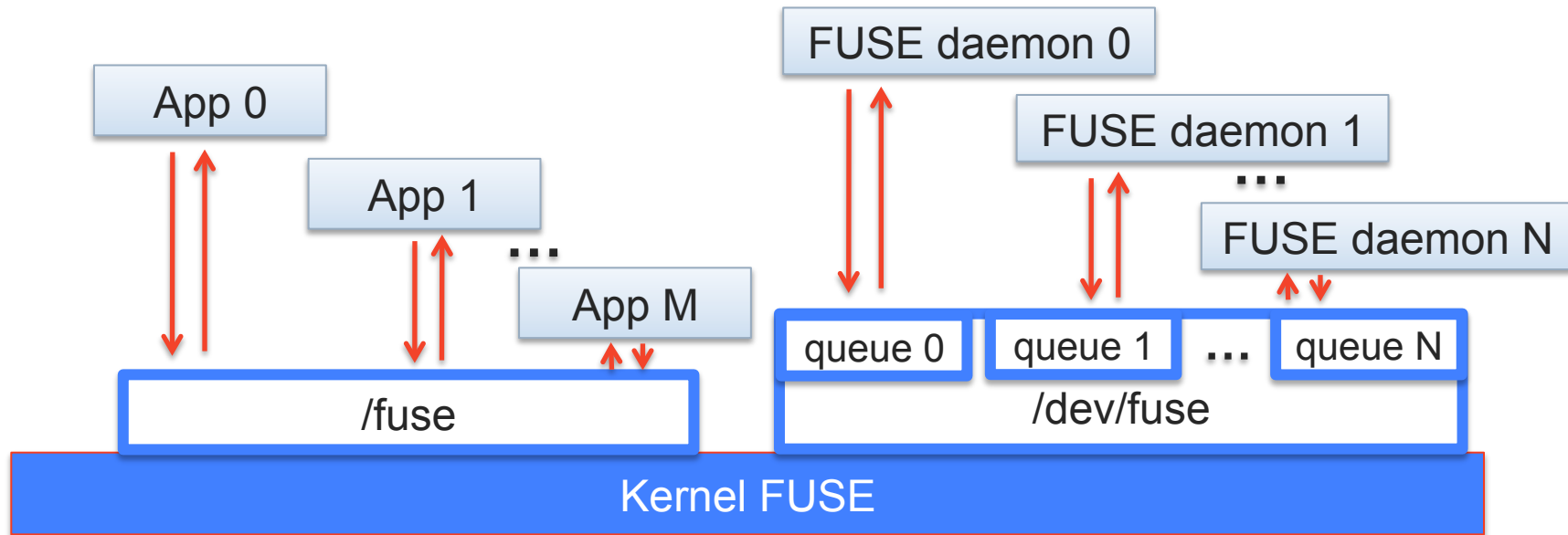
DELL EQL PS6510:
Network: **10 Gbit**
Disks: **48**

**PCS FASTER**
**than**
**HW SAN**

Just 10 nodes PCS cluster faster than DELL EQL SAN ($97000) in most workloads

Sequential Read
1thr 451 MB/s
4thr 711 MB/s
16thr 669 MB/s

Sequential Write
1thr 291 MB/s
4thr 344 MB/s
16thr 340 MB/s

Random Write x32
16thr 39936 IOPS
4thr 31439 IOPS
1thr 13576 IOPS

Random Write
16thr 9591 IOPS
4thr 3624 IOPS
1thr 3043 IOPS

Random Read
16thr 20160 IOPS
4thr 7338 IOPS
1thr 5638 IOPS

100%
80%
60%
40%
20%
0%

▦ **Parallels Cloud Storage 6.0**     ▦ **HW SAN (DELL EQL PS6510)**

‖ Parallels™

# FUSE: what's next?

# FUSE: future improvements

- Variable message size (currently 128K)
- Eliminate global lock
- Multi-queue
- CPU and NUMA affinity

> Q&A