# NVRAM & Software Defined Storage

## Fabian Trumper

Senior Product Marketing Manager, Performance Solutions Group
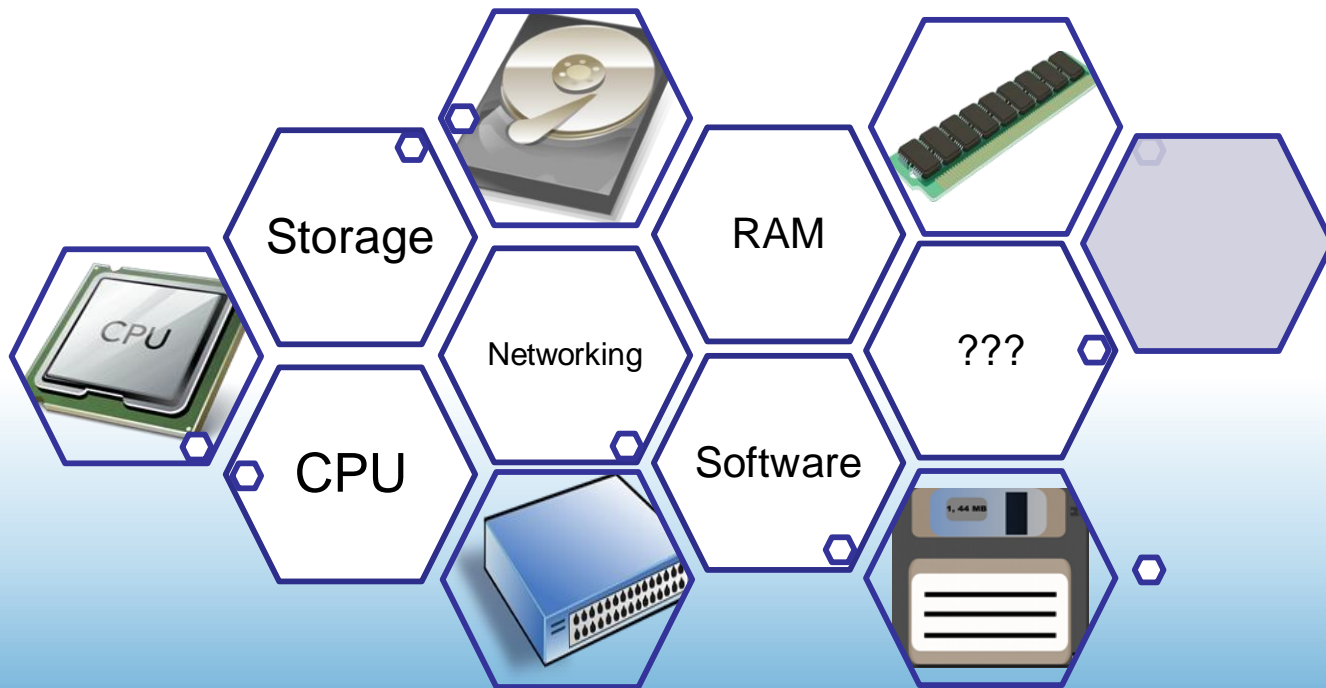
# What is Software-Defined-Storage?

- Logical Storage Layer

- Abstracts HW differences

- Pools together resources

- Uses commercially available HW

- Can be Block, File or Object
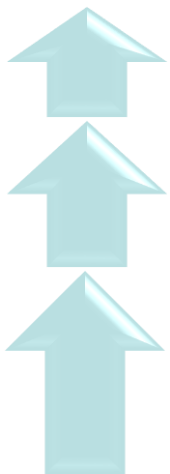
- Can be Scale-Out or Scale-Up

# Building Blocks for SDS

Storage

RAM

CPU

Networking

Software

???

Are you making the most of it?

# Scale Out Approach

Availability? More Replicas!

Throughput? More Nodes!

Capacity? More Nodes!

Simple, but it works…

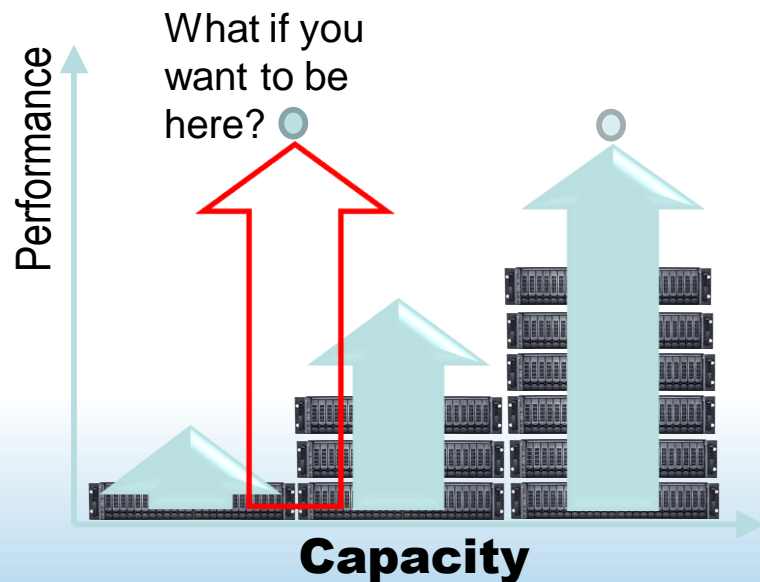(if you have very smart software to make it simple)

# Can Scale-Out solve everything?

## Over Provisioning Costs
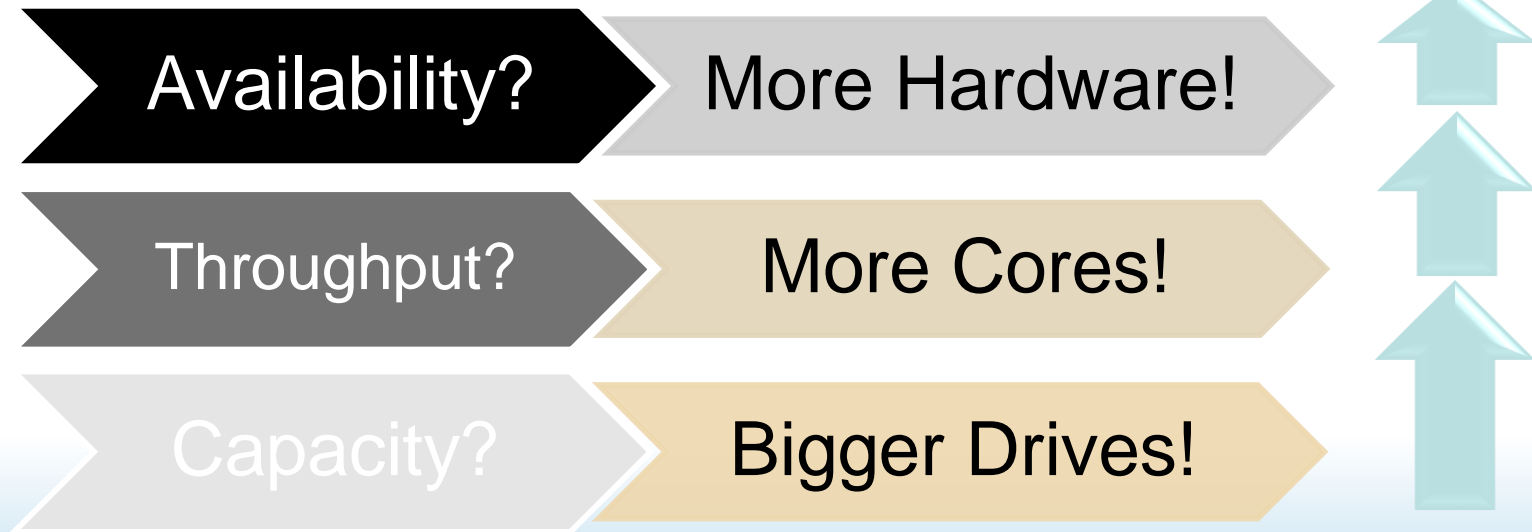
- Racks
- Network
- Power
- Cooling
- Rent, etc

## Diminishing returns on investment

- Difficult to get it just right; efficiency decreases
- Some tasks aren't easily parallelized!

Performance

What if you want to be here?

Capacity

Amdahl's Law: $\frac{1}{\alpha + \frac{1-\alpha}{P}}$

# Scale up Approach

Availability? → More Hardware!

Throughput? → More Cores!

Capacity? → Bigger Drives!

Mostly Hardware Options

Not very popular among the Software Define Storage folks

But can we take a page of that book to the SDS world?

# What Defines Performance?

*Throughput = Bandwidth / Latency*

*What if…*

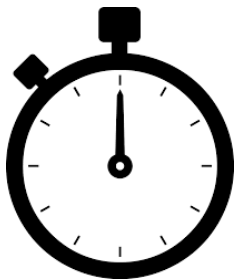Double your bandwidth → 2X Throughput!

Cut your latency by half → 2X Throughput!

# Can We reduce Latency?

## Caching!

Back End Processes

Network Delays in between

Front End Processes

TARGET

NETWORK

INITIATOR

Read Cache?

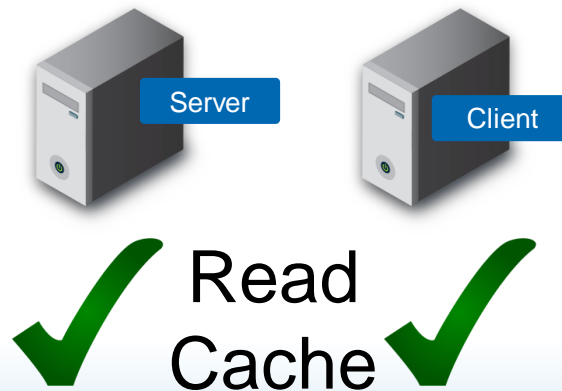Write Cache?

Read Cache?

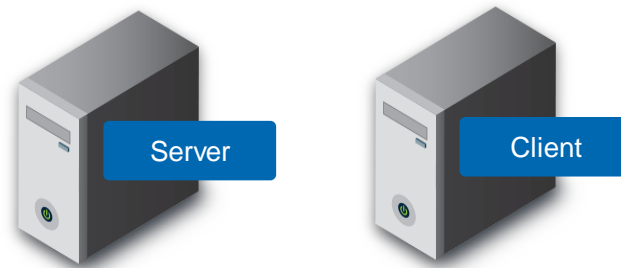Write Cache?

# Read Caching works quite well

SSD's are perfect for this!

- Large Capacity → Better Hit Ratio

- Low Latency → Better Throughput

- Works for both Client and Server

- DRAM is used extensively as well

Server

Client

✔ Read ✔
Cache

# Write Caching is More challenging
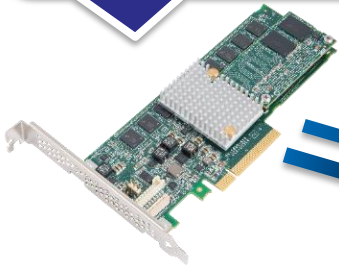
## What about a Write Cache?

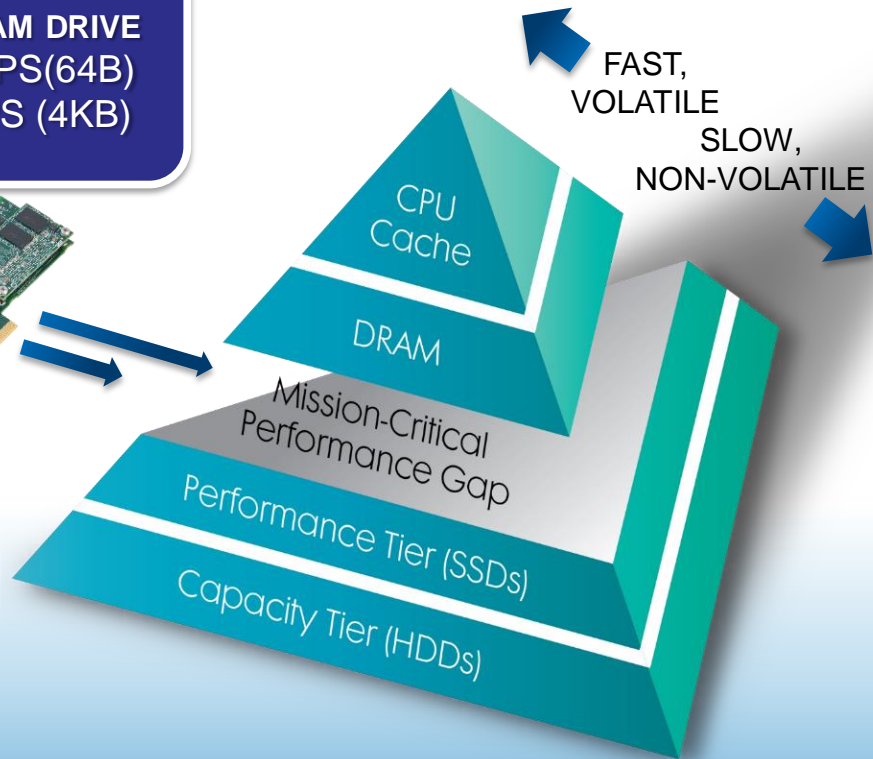- Adoption of SSD's as primary storage calls for a new layer to handle Write Caching!

? Write Cache ✗

Server

Client

# FLASHTEC™ NVRAM DRIVES
# ESTABLISHING A NEW STORAGE TIER

**Flashtec NVRAM DRIVE**
10 Million IOPS(64B)
1 Million IOPS (4KB)

FAST, VOLATILE

SLOW, NON-VOLATILE

CPU Cache

DRAM

Mission-Critical Performance Gap

Performance Tier (SSDs)

Capacity Tier (HDDs)

# Adding NVRAM to the mix

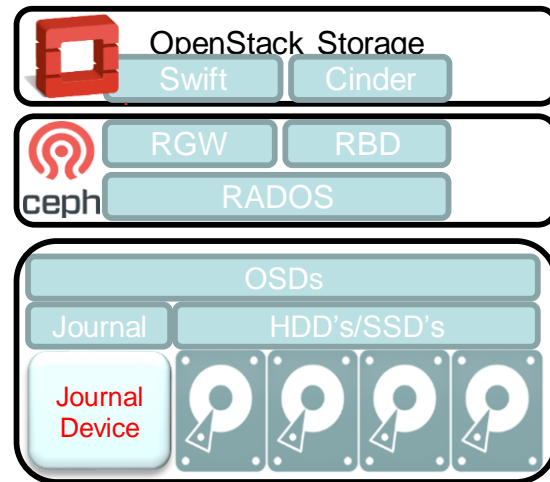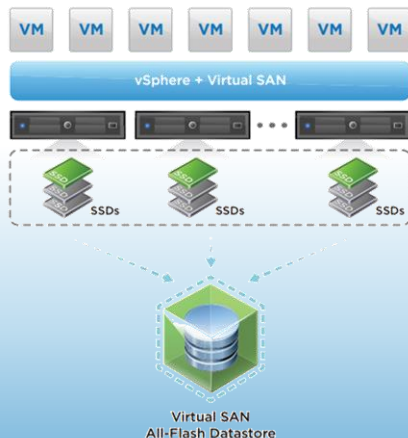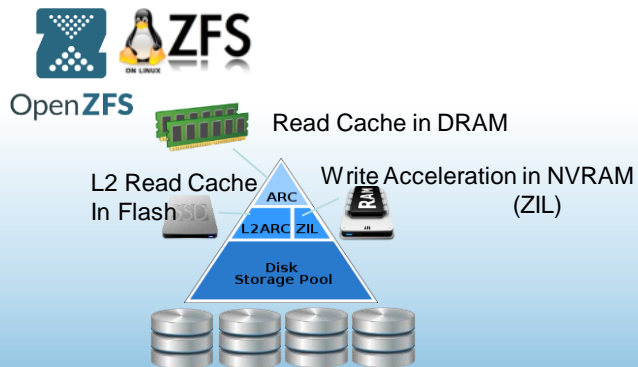Storage

RAM

Networking

NVRAM

CPU

Software

What are the benefits?

# Write Caching Implementations

- ■ Examples of Write Cache Usage
  - ZFS uses ZIL (ZFS Intent Log)
  - Ceph OSD's Journaling Device
  - VMWare VSAN Write Cache

OpenStack Storage

| Swift | Cinder |

| RGW | RBD |
| RADOS | |

ceph

| OSDs | |
| Journal | HDD's/SSD's |
| Journal Device | |

OpenZFS
ZFS ON LINUX

Read Cache in DRAM

L2 Read Cache
In Flash

ARC
L2ARC  ZIL

Write Acceleration in NVRAM
(ZIL)

Disk
Storage Pool

VM VM VM VM VM VM VM

vSphere + Virtual SAN

SSDs    SSDs    SSDs

Virtual SAN
All-Flash Datastore

PMC

# How much does NVRAM benefit Ceph?

# Journaling @ Ceph

Every CEPH OSD Writes Twice!



DATA

Journal

DATA

Replicas

DATA*    DATA*

+ Replica Copies

2X Writes!

½ Throughput!

# NVRAM to the rescue!

Min Latency
No Wear

PMC | flashtec™
NVRAM

DATA

Journal

DATA

Replicas

DATA*

DATA*

LAST ACTION HERO

BAM!

2X
Throughput!

PMC

# Can we do even better?

| | Samsung PM1633 |
|---|---|
| Form Factor | 2.5" |
| Capacity (GB) | 480/960/ 1920/3840 |
| Host Interface | SAS 3 (12 Gb/s) |
| MTBF | 2.0 Million Hours |
| Uncorrectable Bit Error Rate (UBER) | 1 in $10^{17}$ |
| Power Consumption (Active) | 11.0W |
| Power Consumption (Idle) | 4.0W |
| Random Reads (up to) | 180,000 IOPS |
| Random Writes (up to) | 15,000 IOPS |
| Sequential Reads (up to) | 1,100 MB/s |
| Sequential Writes (up to) | 1,000 MB/s |
| Endurance (up to) | 1 DWPD |

- 3.84TB Capacity – GREAT! ✔

- 180K IOPS Random Read – GREAT! ✔

- 15K Random Write (60 MB/s) – FAIL! ✘

!?

# Write Caching + Re-Ordering

| | Samsung PM1633 |
|---|---|
| Form Factor | 2.5" |
| Capacity (GB) | 480/960/ 1920/3840 |
| Host Interface | SAS 3 (12 Gb/s) |
| MTBF | 2.0 Million Hours |
| Uncorrectable Bit Error Rate (UBER) | 1 in $10^{17}$ |
| Power Consumption (Active) | 11.0W |
| Power Consumption (Idle) | 4.0W |
| Random Reads (up to) | 180,000 IOPS |
| Random Writes (up to) | 15,000 IOPS |
| Sequential Reads (up to) | 1,100 MB/s |
| Sequential Writes (up to) | 1,000 MB/s |
| Endurance (up to) | 1 DWPD |

**BAM!**

1 GB/s Sequential Write! – GREAT!

Can we turn Random Workloads to Sequential?

HDD's have been doing it for years!

With an NVRAM memory layer, SDS Storage Stacks could use these age old techniques to improve throughput and reduce costs at the same time