# GLUSTER CAN DO THAT!
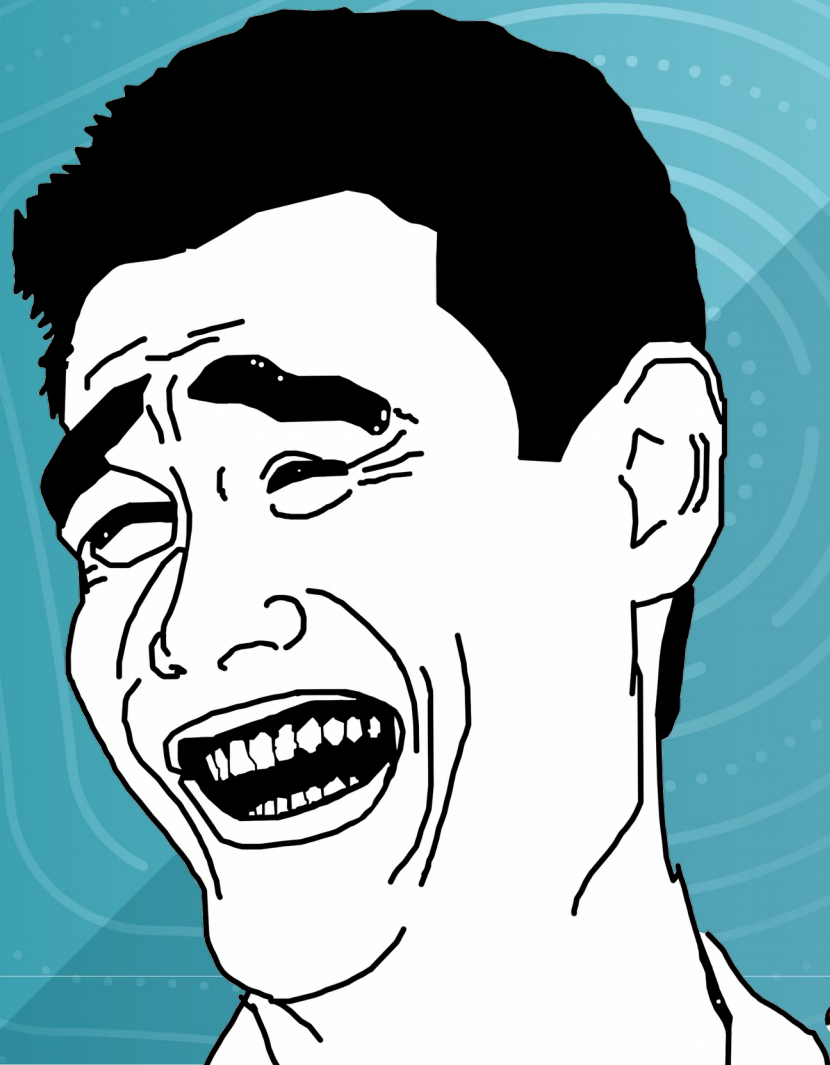
Architecting and Performance Tuning
Efficient Gluster Storage Pools

**Dustin Black**
**Senior Architect, Software-Defined Storage**
**@dustin1black**

**Ben Turner**
**Principal Quality Engineer**
**@bennyturns**

2017-05-02

GLUSTER 101
IN 5 SECONDS

# THE DATA EXPLOSION

**WEB, MOBILE, SOCIAL MEDIA, CLOUD**
Our digital assets have grown exponentially due to web scale services like Facebook, Flickr, Snapchat, YouTube, and Netflix.

**VIDEO ON-DEMAND SERVICES**
Rapid growth of video on-demand has culminated in 50% of households using this service.

**MEDIA AND ENTERTAINMENT INDUSTRIES**
A staggering amount of content is created during today's optimized production processes.

**MEDICAL INDUSTRY**
Medical imaging needs are vast, and regulatory requirements can be demanding.

redhat.

# NOT SURE IF YOU GOT IT?
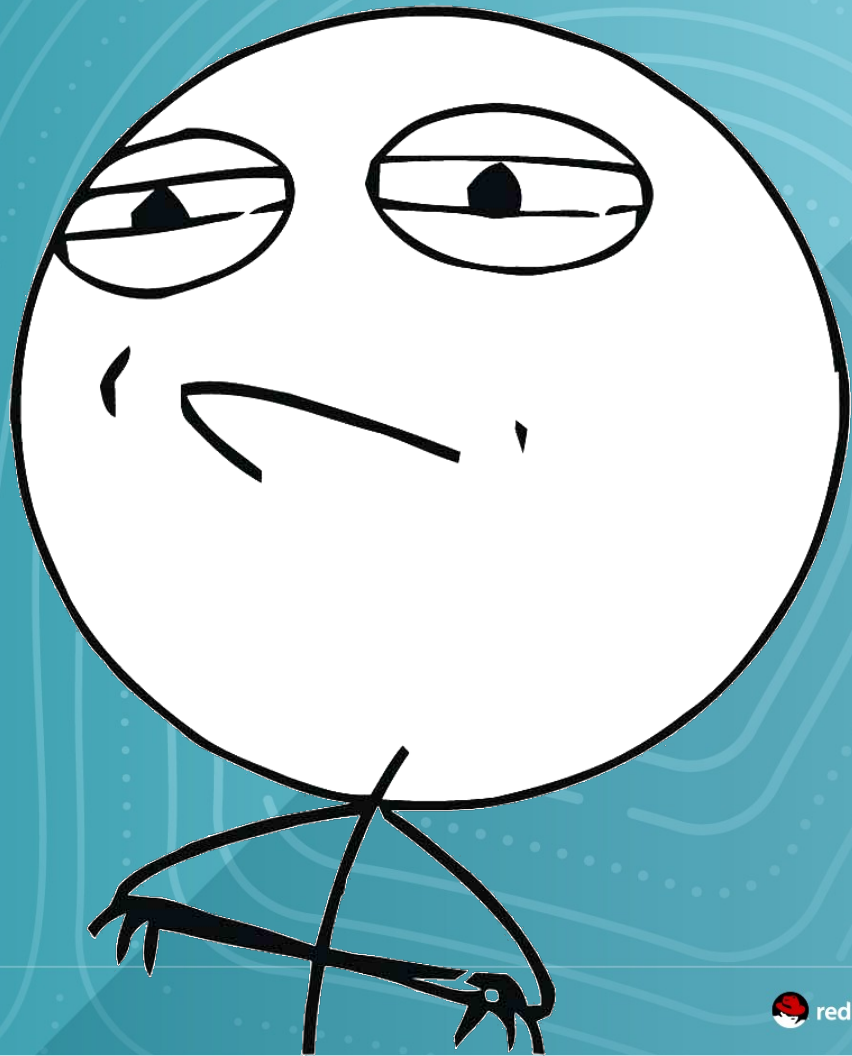
https://people.redhat.com/dblack/summit2017

# GLUSTER
# CAN DO THAT!*

*If you build it right

# A SIX-NODE POOL CAN PROCESS…

JPEG Web
Image Files
(32KB)

72x 7.2K HDD

**1,700** JPEGs
per second

or

Optimized
72x 7.2K HDD

**12,000** JPEGs
per second

or

72x SSD

**23,000** JPEGs
per second

redhat.

# OR...

DVD
Movie Files
(4GB)

72x 7.2K HDD

**Optimized**
**72x 7.2K HDD**

72x SSD

**1** DVD
per second

or

**2** DVDs
per second

or

**4** DVDs
per second

redhat.

# OR...



High-Def
CCTV Camera
Recording Streams

72x 7.2K HDD

**200** CCTV streams
within latency threshold

or

Optimized
72x 7.2K HDD

**500** CCTV streams
within latency threshold

or

72x SSD

**?** CCTV streams
within latency threshold

# Keep It Simple, Stupid

redhat.

# SWTWD

# START WITH THE WORKLOAD, DUMMY

redhat.

**WHY DO YOU ASK THE WRONG QUESTIONS?**

Delivered-To: dblack@redhat.com
From: ████████████████████████████
Date: Fri, 31 Mar 2017 11:59:29 +0200
Subject: Pererformance testing with fio


…
One of the things ██████████ wants is see that gluster performs similarly to the ████████ NFS
system it is intended to replace.


Now I noticed the following:


- Doing a **simple test with dd** yields a write throughput of around 500MB/s, which for a rep 2
volume on a 10Gb connection is quite good.
- Doing a **read with dd** strangely yields slower throughput....
…

Delivered-To: dblack@redhat.com
From: ▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓
Date: Sun, 5 Feb 2017 20:16:40 +0900
Subject: RHGS scale-out options

…
▓▓▓▓▓▓▓▓▓▓▓▓ plans to **add physical nodes to increase "performance"**
(currently ▓▓▓▓▓ is experiencing performance problem)
…
Current Env : 80 X 2-way distributed replicated vols on 6 nodes
To-Be : add 6 more nodes... becomes 80 X 2-way distributed replicated vols on 12 nodes


I'm not sure which one is the best way to increase performance.


1. extend current cluster from 6 to 12 nodes and add bricks from new 6 nodes into existing 80 vols
2. extend current cluster from 6 to 12 nodes and migrate some vols to new new 6 nodes.
3. create another RHGS gluster cluster with new 6 nodes and migrate some vols to new RHGS cluster
4. ??
…

redhat.

Delivered-To: dblack@redhat.com
From: ████████████████████████
Date: Mon, 6 Mar 2017 10:54:17 -0800
Subject: Fwd: ██████████████████ server quote ███████████████████████

…
What are your calculations for the ██████████ NAS storage RFP?

██████████████ is asking for the **IOPS per drive / Raid Volume** for the design?

They would like to make sure they are getting **28,000 IOPs per site**.
…


---------- Forwarded message ----------
From: ██████████████████████████
Date: Mon, Mar 6, 2017 at 10:45 AM


…
Thank you.  The next question that I have is how many IOPS per drive (or per RAID volume, or per server), for 3.5" 7200RPM SATA drives, are you assuming.  The requirement is for 28,000 IOPS at each site. Thanks.
…

THE WORKLOAD
IS COMING

# SMALL FILE JPEG WORKLOAD
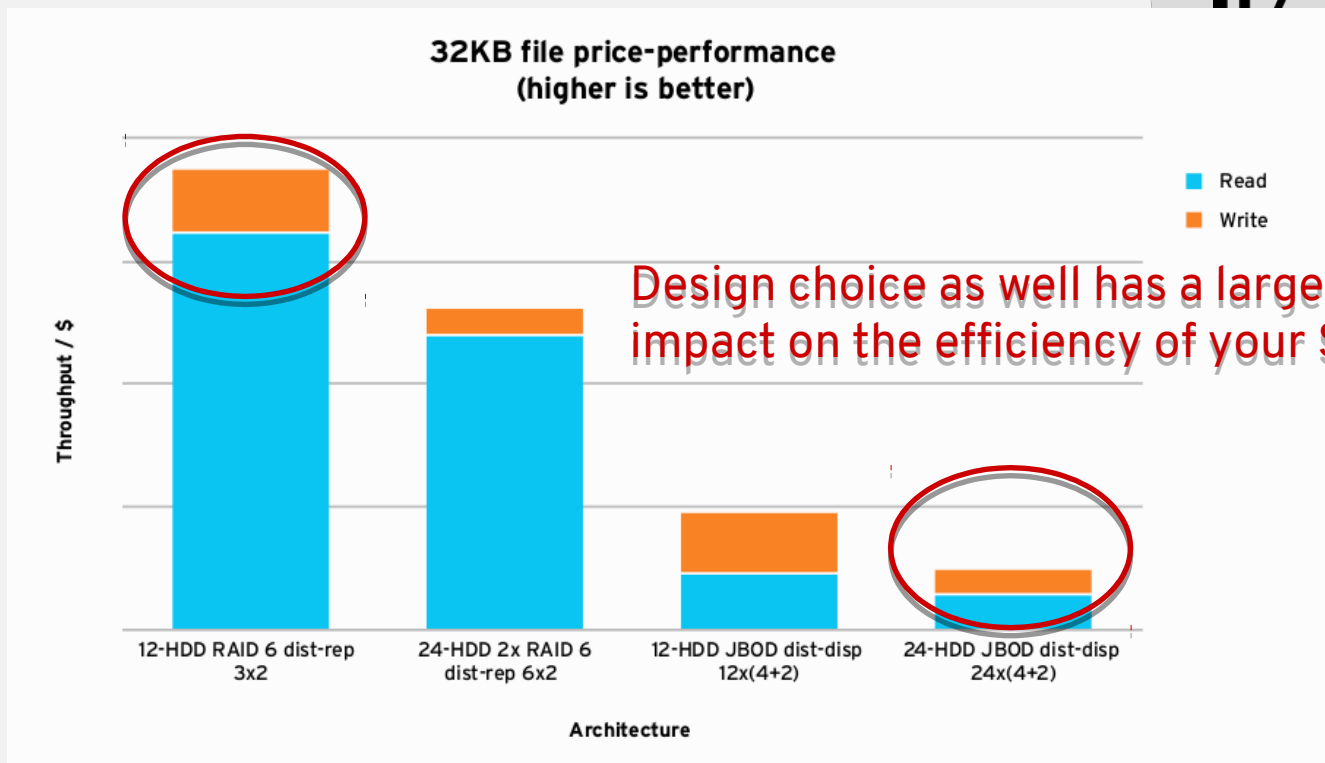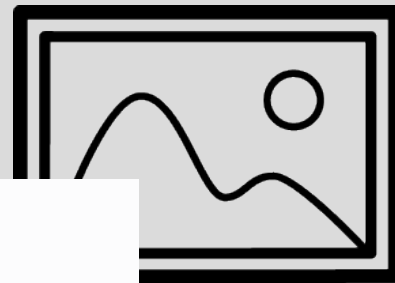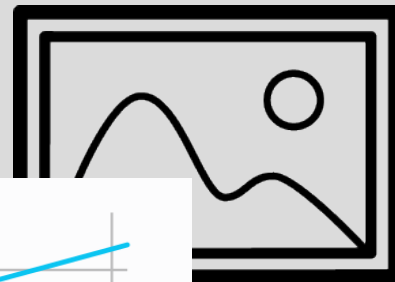
**Standard servers - 32KB file throughput by architecture**

Client / architecture

- Native / 12-HDD RAID 6 dist-rep 3x2 — Read: 170.85, Write: 26.96
- Native / 12-HDD RAID 6 disp 1x(4+2) — Read: 70.54, Write: 36.19
- Native / 12-HDD JBOD dist-disp 12x(4+2) — Read: 23.76, Write: 26.69
- NFS / 12-HDD RAID 6 dist-rep 3x2 — Read: 184.86, Write: 12.15
- NFS / 12-HDD RAID 6 disp 1x(4+2) — Read: 35.66, Write: 4.53
- NFS / 12-HDD JBOD dist-disp 12x(4+2) — Read: 47.48, Write: 10.61

Legend: ■ Read  ■ Write

X-axis: Files/s/drive (0, 50, 100, 150, 200)

**Same Hardware**

**Design choice dramatically affects system capabilities**

redhat.

# SMALL FILE JPEG WORKLOAD



**32KB file price-performance (higher is better)**

Read
Write

Throughput / $

Architecture

12-HDD RAID 6 dist-rep 3x2
24-HDD 2x RAID 6 dist-rep 6x2
12-HDD JBOD dist-disp 12x(4+2)
24-HDD JBOD dist-disp 24x(4+2)

Design choice as well has a large impact on the efficiency of your $$

# SMALL FILE JPEG WORKLOAD



files/s vs Workers (Threads)

Client concurrency is important for maximizing system throughput

Write — Read

# SMALL FILE JPEG WORKLOAD

**Server Aggregate Network Utilization**

10% of Theoretical Maximum

| | |
|---|---|
| 1.3 GBps | |
| 1.0 GBps | |
| 750 MBps | |
| 500 MBps | |
| 250 MBps | |
| 0 Bps | |

11:40   11:45   11:50   11:55   12:00   12:05   12:10   12:15   12:20   12:25   12:30   12:35

— Outgoing  — Incoming

**Server Aggregate CPU Utilization**

25% Consumption

| | |
|---|---|
| 50% | |
| 40% | |
| 30% | |
| 20% | |
| 10% | |
| 0% | |

11:40   11:45   11:50   11:55   12:00   12:05   12:10   12:15   12:20   12:25   12:30   12:35

— Wait  — User  — System

redhat.

# SMALL FILE JPEG WORKLOAD



**Server Aggregate Memory Utilization**

768 GiB Maximum

- Cached
- Used

**Server HDD Busy**

Maximum Utilization

We are reaching a disk bottleneck on reads

IF A FILE IS VERY VERY SMALL

IS IT STILL A FILE?

IT'S LIKE A SERIES OF TUBES
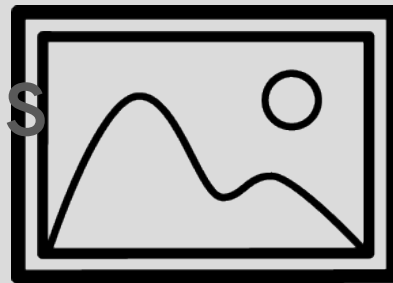
# GLUSTER TRANSLATORS

```c
struct xlator_fops fops = {
        .open        = ra_open,
        .create      = ra_create,
        .readv       = ra_readv,
        .writev      = ra_writev,
        .flush       = ra_flush,
        .fsync       = ra_fsync,
        .truncate    = ra_truncate,
        .ftruncate   = ra_ftruncate,
        .fstat       = ra_fstat,
        .discard     = ra_discard,
        .zerofill    = ra_zerofill,
};

struct volume_options options[] = {
        { .key  = {"force-atime-update"},
          .type = GF_OPTION_TYPE_BOOL,
          .default_value = "false"
        },
        { .key  = {"page-count"},
          .type = GF_OPTION_TYPE_INT,
          .min  = 1,
          .max  = 16,
```

# SMALL FILE AND METADATA WORKLOADS

**What the Gluster community is doing:**

Improve efficiency of individual calls

Store metadata in client cache

Prefetch metadata

Compound file operations

Coming Soon! Negative lookups and parallel readdirp

redhat.

# TUNING FOR SMALL FILE & METADATA

Since small file workloads are metadata intensive, I use the same tuning for both.
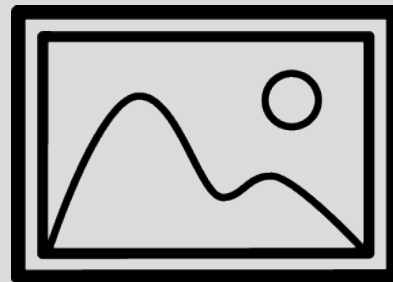
RAID 10 or RAID 6 are recommended for bricks

Tuned profile: rhgs-throughput-performance

Event Threads = 4

lookup-optimize = on

Features.cache-invalidation = on

Performance.stat-prefetch = on

# SMALLFILE CREATES & READS

Create & read of 32 KB files
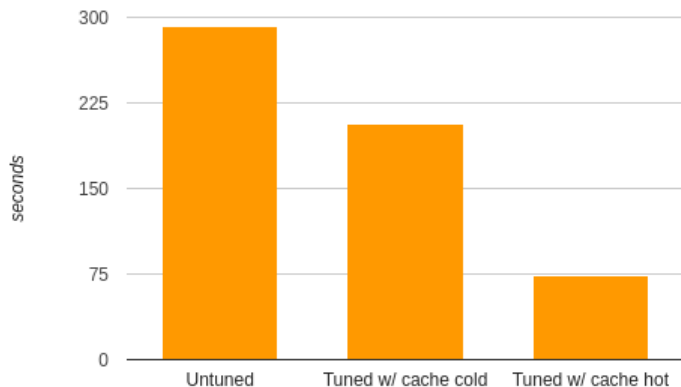untuned vs tuned w/ cold cache vs tuned w/ hot cache
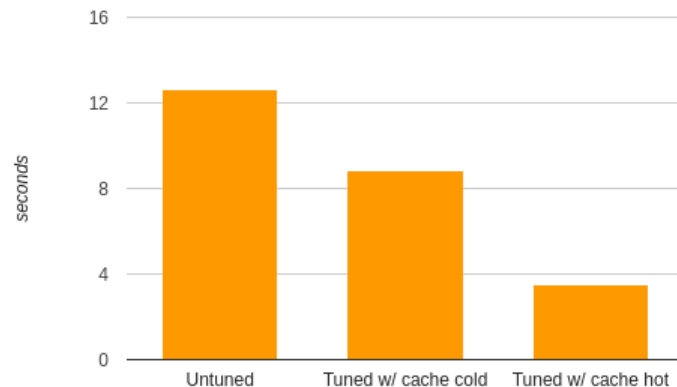
# SMALLFILE METADATA WORKLOAD

Single and multi-threaded ls -l workloads
untuned vs tuned w/ cold cache vs tuned w/ hot cache
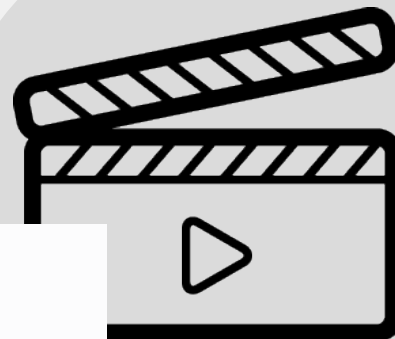


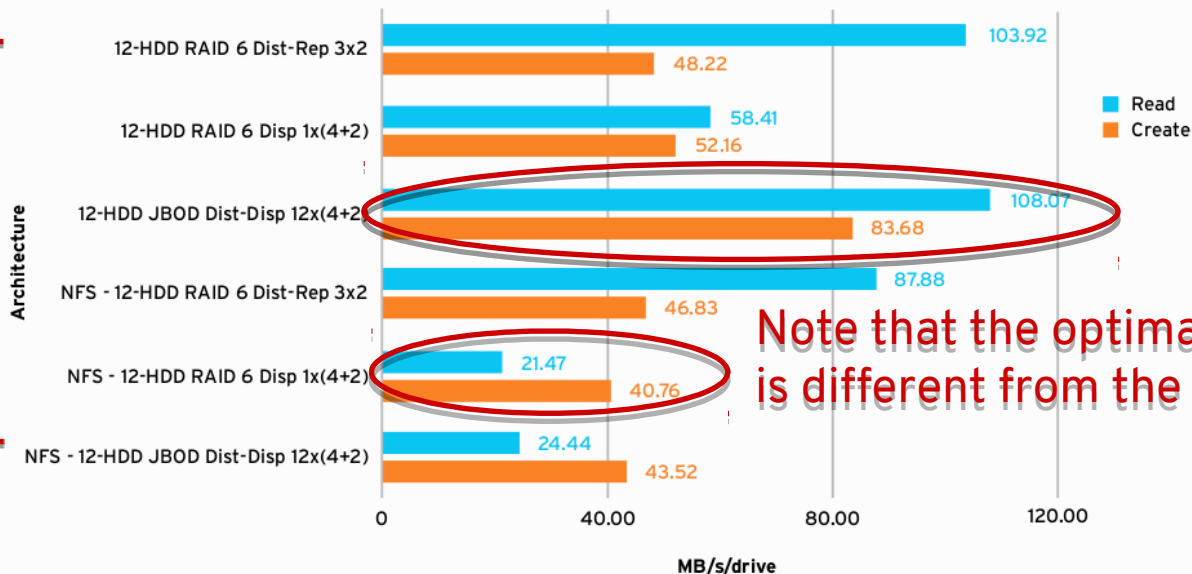320k files ls -laR (single client thread) ----
smaller = better



320k files ls -laR (4 clients, 8 threads/client) --
smaller = better
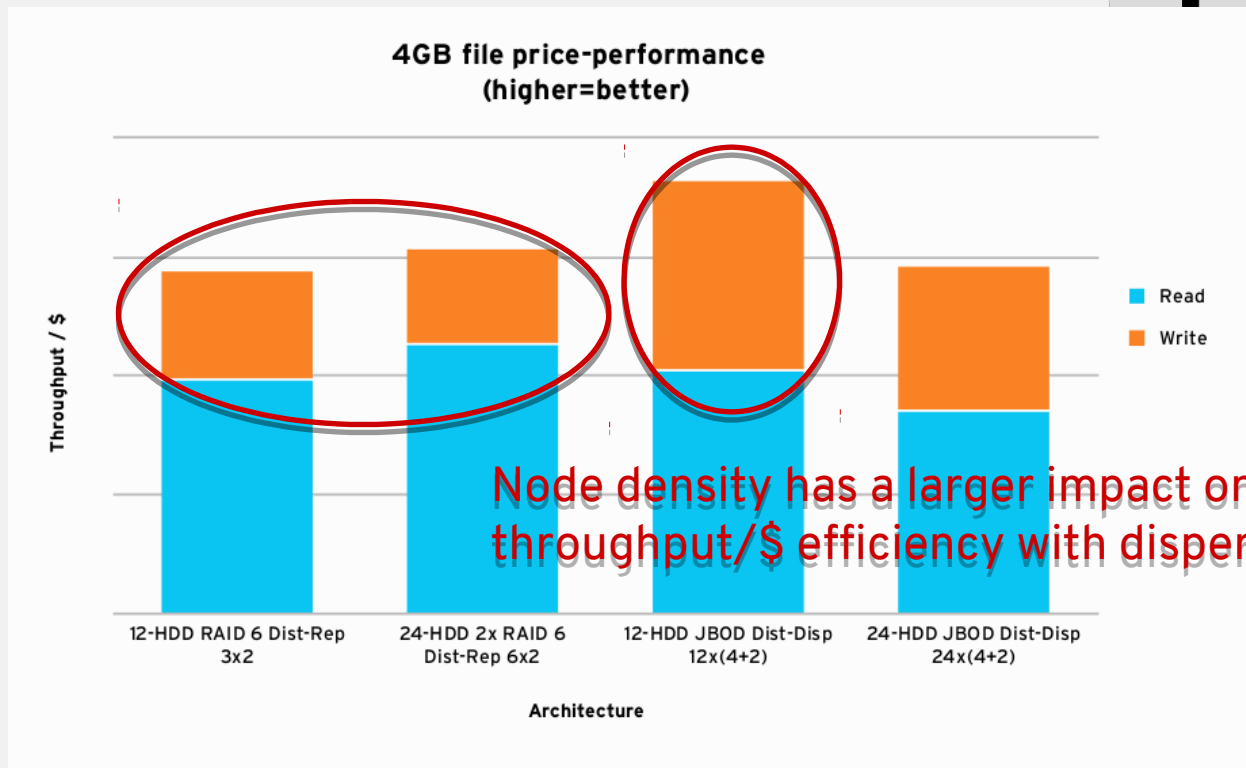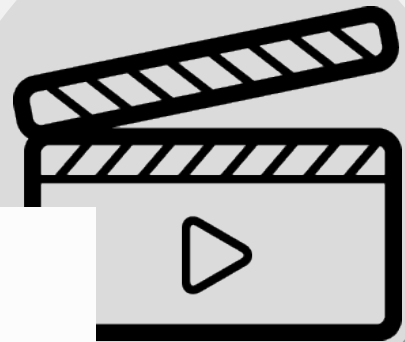
# LARGE FILE DVD WORKLOAD



**Standard servers - 4GB file throughput by architecture**

Same Hardware

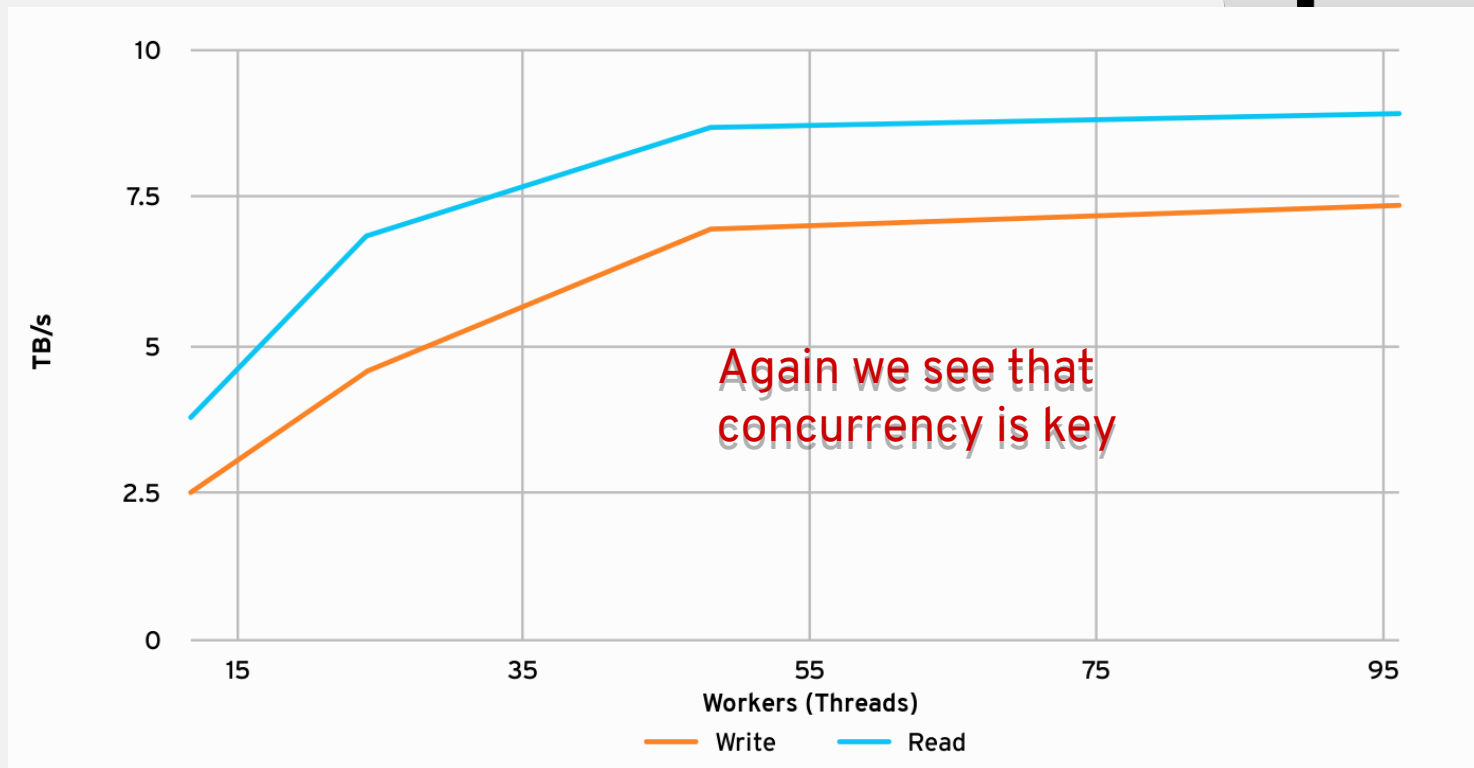Note that the optimal configuration is different from the small file results

# LARGE FILE DVD WORKLOAD



**4GB file price-performance (higher=better)**

Throughput / $

Read
Write

12-HDD RAID 6 Dist-Rep 3x2

24-HDD 2x RAID 6 Dist-Rep 6x2

12-HDD JBOD Dist-Disp 12x(4+2)
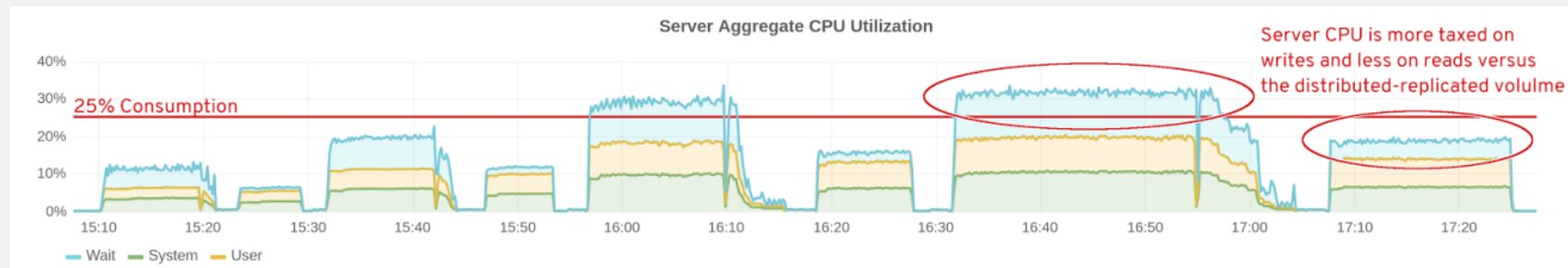
24-HDD JBOD Dist-Disp 24x(4+2)

Architecture
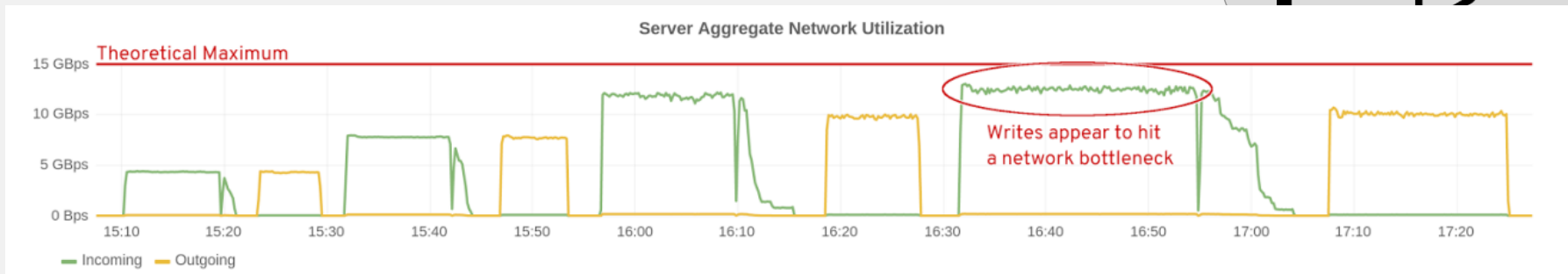
Node density has a larger impact on your throughput/$ efficiency with disperse volumes.

# LARGE FILE DVD WORKLOAD



Again we see that
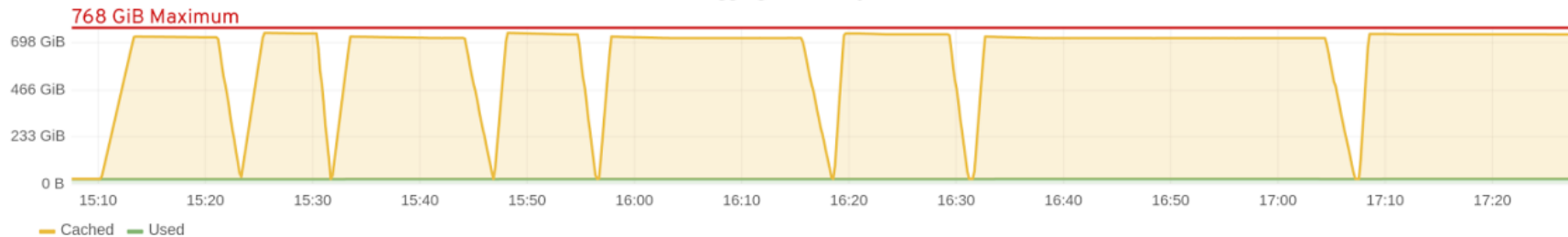concurrency is key

# LARGE FILE DVD WORKLOAD

## Server Aggregate Network Utilization

Theoretical Maximum

15 GBps

10 GBps

5 GBps

0 Bps

15:10  15:20  15:30  15:40  15:50  16:00  16:10  16:20  16:30  16:40  16:50  17:00  17:10  17:20

Writes appear to hit a network bottleneck

— Incoming  — Outgoing

## Server Aggregate CPU Utilization

Server CPU is more taxed on writes and less on reads versus the distributed-replicated volulme

40%

30%

25% Consumption

20%

10%

0%

15:10  15:20  15:30  15:40  15:50  16:00  16:10  16:20  16:30  16:40  16:50  17:00  17:10  17:20

— Wait  — System  — User

# LARGE FILE DVD WORKLOAD



Server Aggregate Memory Utilization

768 GiB Maximum

— Cached  — Used

Server HDD Aggregate Bandwidth

We appear to be reaching an aggregate HDD throughput limit for reads

— write  — read

# TUNING FOR LARGE FILE SEQUENTIAL

How Dustin got his performance gains from tuning!

RAID 6 or EC are recommended for bricks

Tuned profile: rhs-high-throughput

Read-ahead on bricks

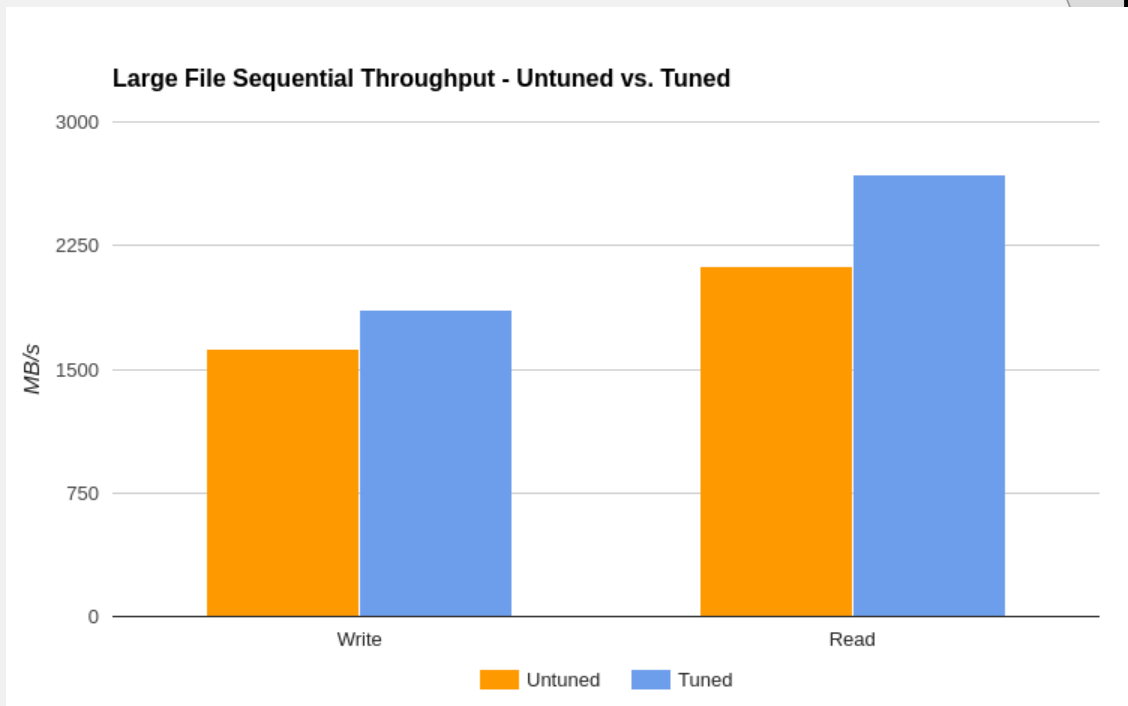Deadline scheduler
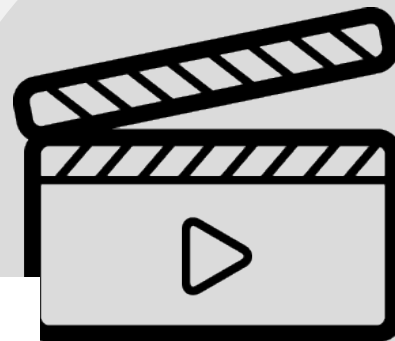
vm.dirty-ratio

Jumbo Frames

Event Threads = 4

Smallfile tuning may have some effect, especially with metadata operations.

# LARGE FILE SEQUENTIAL
## 4 Servers, 4 Clients, 4 Workers/Client, 16GB File/Worker



Large File Sequential Throughput - Untuned vs. Tuned

# SCOPING FOR LARGE FILE WORKLOADS

Now that you understand the workload, how can you size your cluster?

Formula for *guesstimating* large file performance:

Writes = (Slowest of NIC / DISK) / # replicas * .7(overhead)

1200 MB / 2 * .7 = 420 MB / sec

Reads = (Slowest of NIC / Disk ) * .6(overhead)

1200 * .6 = 720 MB / sec

*This is just a rule of thumb, actual results are highly dependant on hardware.*

# TAKEAWAYS FOR LARGE FILE WORKLOADS

EC on JBOD outperforms replica 2 on RAID 6 high worker concurrency workloads

Replica 2 on RAID 6 outperforms EC on JBOD when there are less files / clients / threads and on single threaded workloads

Read ahead on block devices as well as jumbo frames provide the most performance benefit of the tunables

Again, start with the workload when designing your storage cluster.  The proper brick architecture from the start will yield far better performance than any of the tunables mentioned.  Design in a way that avoids problems, don't try tune your way out of them.

redhat.

YOUR WORKLOAD CAN'T BE SLOW IF YOU NEVER RUN IT
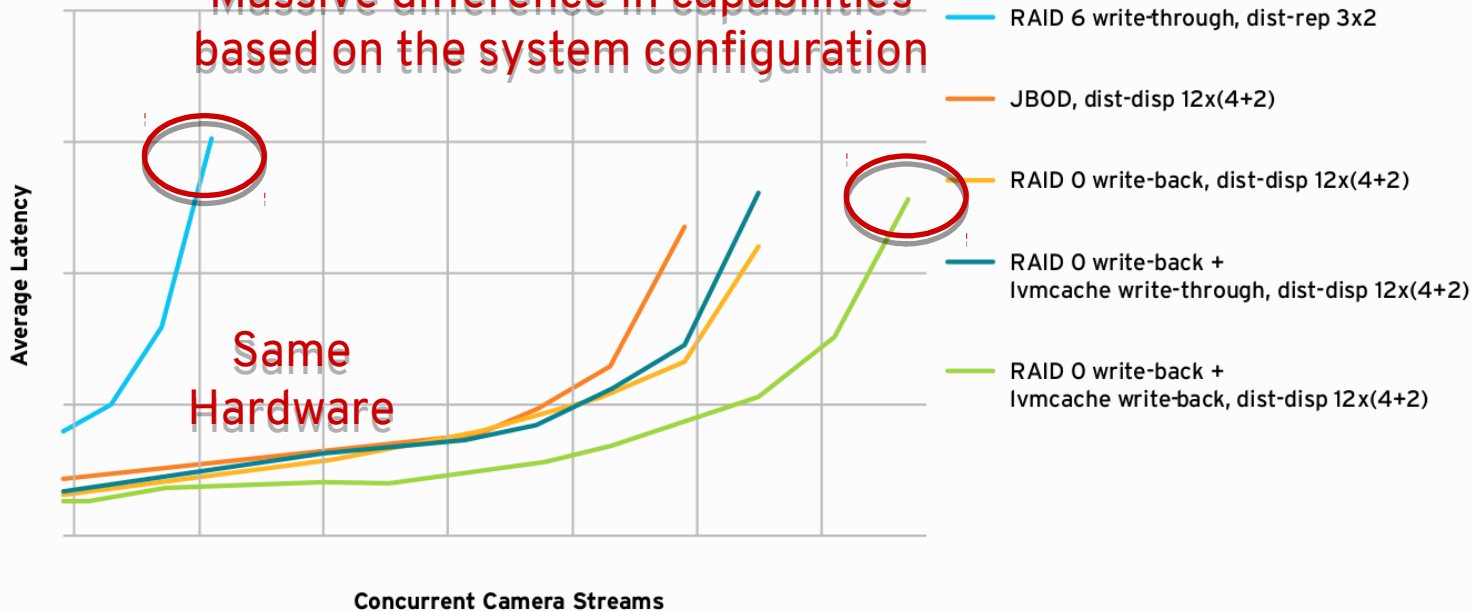
#redhat #rhsummit

# CCTV STREAMING WORKLOAD



**STREAMING VIDEO CAPTURE LIMIT PER GLUSTER CONFIGURATION**

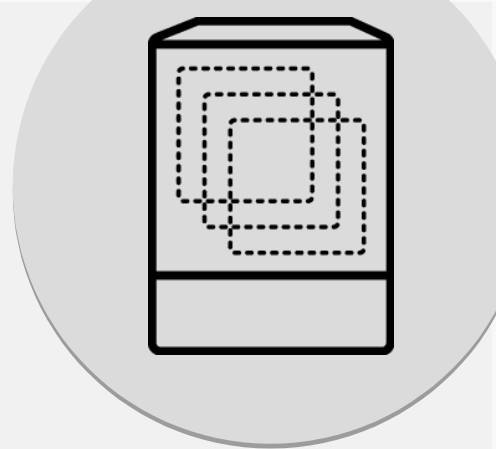Massive difference in capabilities based on the system configuration

Same Hardware

Average Latency

Concurrent Camera Streams

RAID 6 write-through, dist-rep 3x2

JBOD, dist-disp 12x(4+2)

RAID 0 write-back, dist-disp 12x(4+2)

RAID 0 write-back + lvmcache write-through, dist-disp 12x(4+2)

RAID 0 write-back + lvmcache write-back, dist-disp 12x(4+2)

# HYPERCONVERGED RHV / RHGS

Setup Details

Storage and compute on the same systems

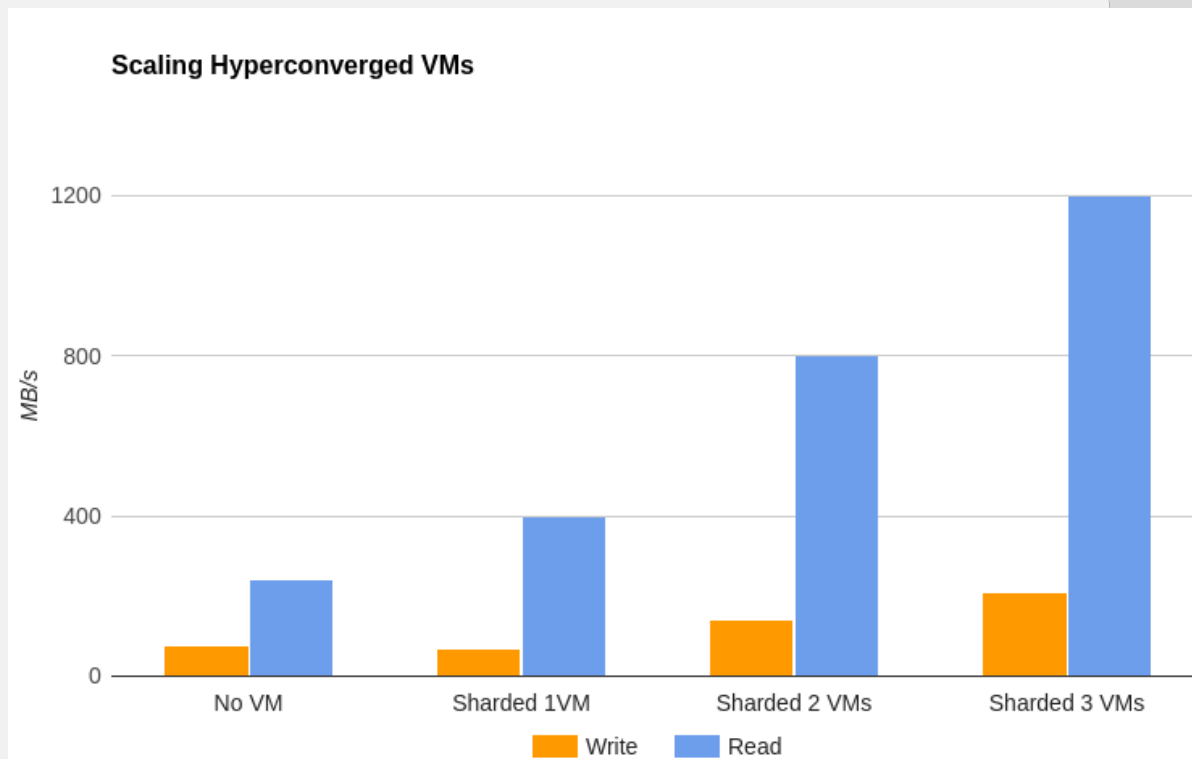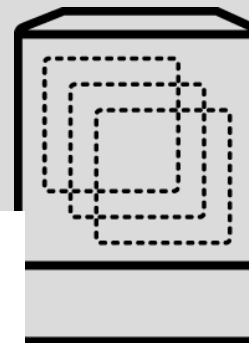Cost advantage

Management using the same linux based tools

```
# gdeploy -c robo.conf

# hosted-engine --deploy --config-append=<path to hosted engine answer
file>
```
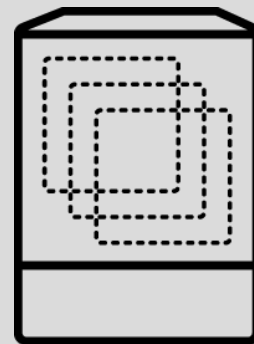
# Hyperconverged Infrastructure Example Arch



Hosted Engine

RHV

KVM

KVM

KVM

Gluster Volume

Stores sharded VM metadata and disk files

# VM PERFORMANCE



Scaling Hyperconverged VMs

# PERFORMANCE TEST TOOL - GBENCH

**Gbench was used to gather the performance data**

https://github.com/gluster/gbench

Wraps IOZone, smallfile, FIO

Run multiple iterations and averages it

Multi host capable

# Install The Summit Mobile App

https://install.events/reraa17

**PLEASE**

**Help us out by reviewing our session**

redhat.