# Best Practices for RHS Performance

**RED HAT SUMMIT**

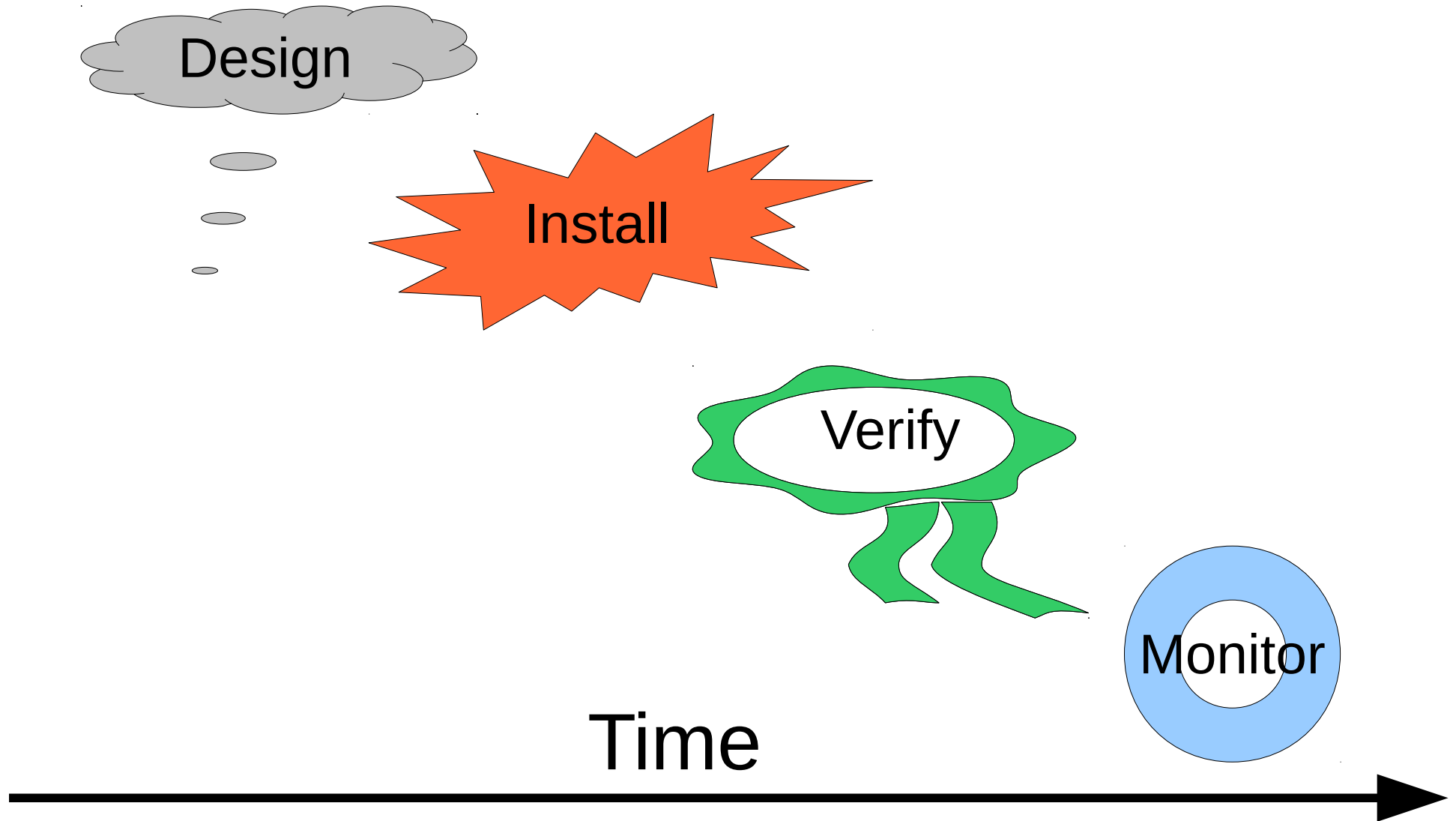**CONNECT TO THE INFORMATION REVOLUTION**

Ben England

Neependra Khare

June 2013

redhat

# Scope of this session

- Quick review of Gluster from perf. perspective

- Best practices to provision faster RHS storage

  - characterize use case (workload, config)

  - provision network + storage appropriately

- GUI and RHS performance

# Workflow for provisioning RHS

Design

Install

Verify

Monitor

Time

RED HAT®
STORAGE

# balancing RHS server configuration
# don't underfund network

## relative capacity – CPUs idle, network busy

Network ~1 GB/s

CPUs ~10 GB/s

Storage ~1-2 GB/s,

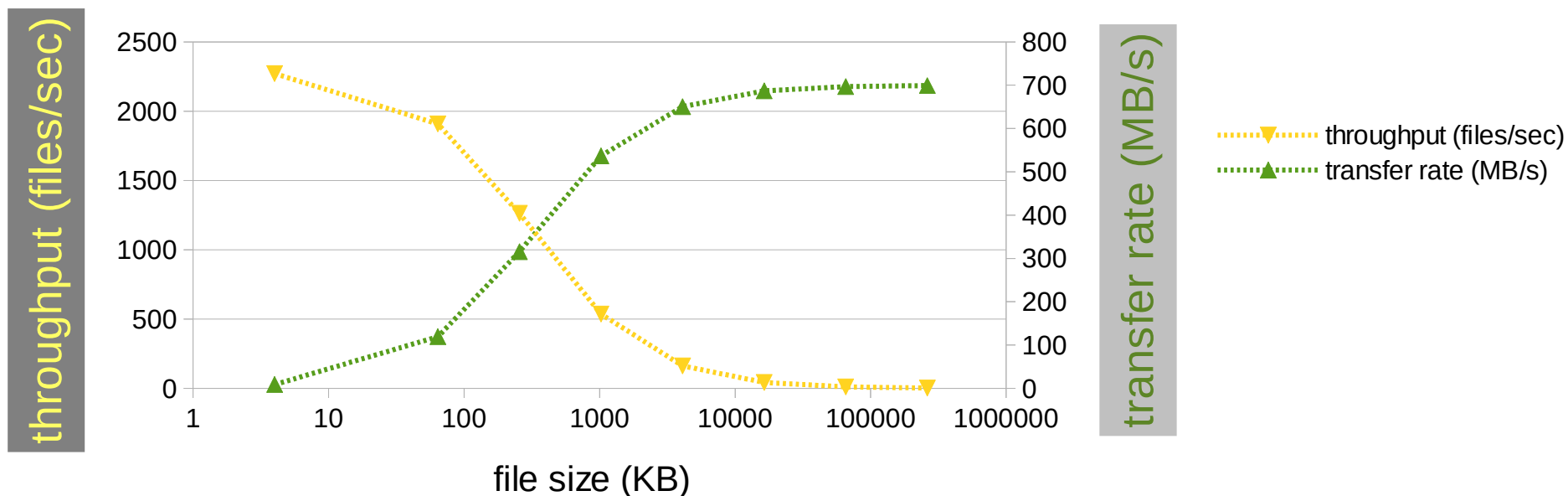## relative cost – storage is cost driver

Network 10%

CPU 30%

Storage 60%

**price-performance ratio**

**RED HAT®**
**STORAGE**

# file size distribution and effect on sizing

one set of performance measurements with different file sizes
looking at measurements with 2 y-axes

file create, 64 GB total data, 4 threads/client, 4 servers = 4 clients,
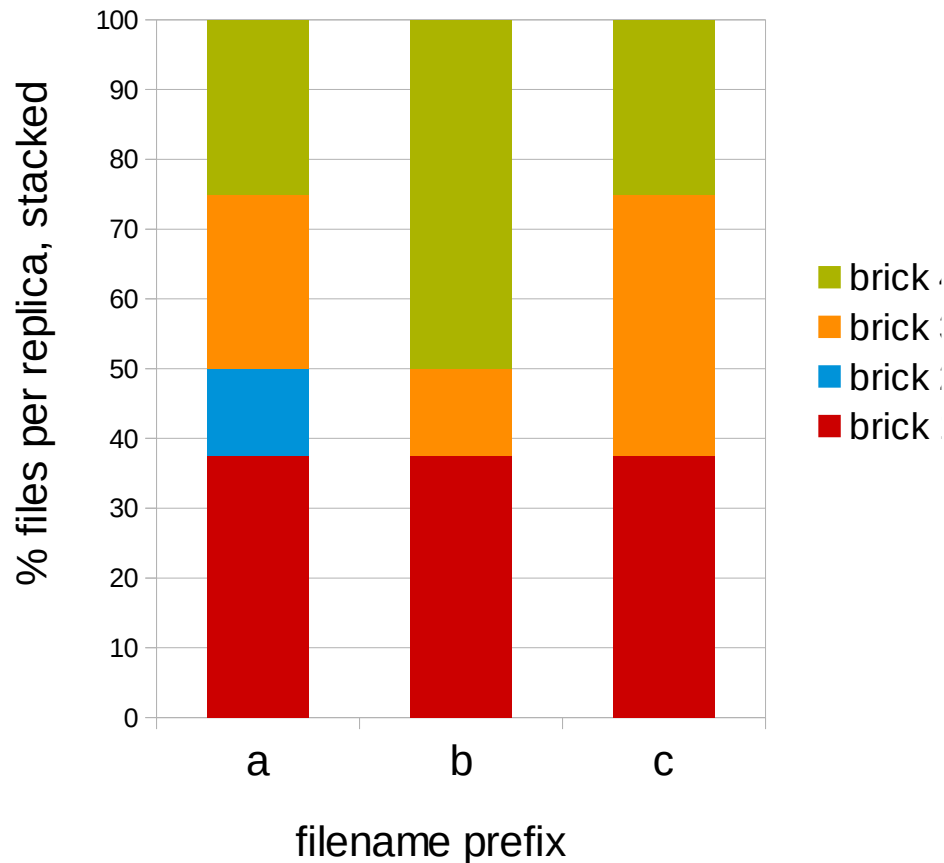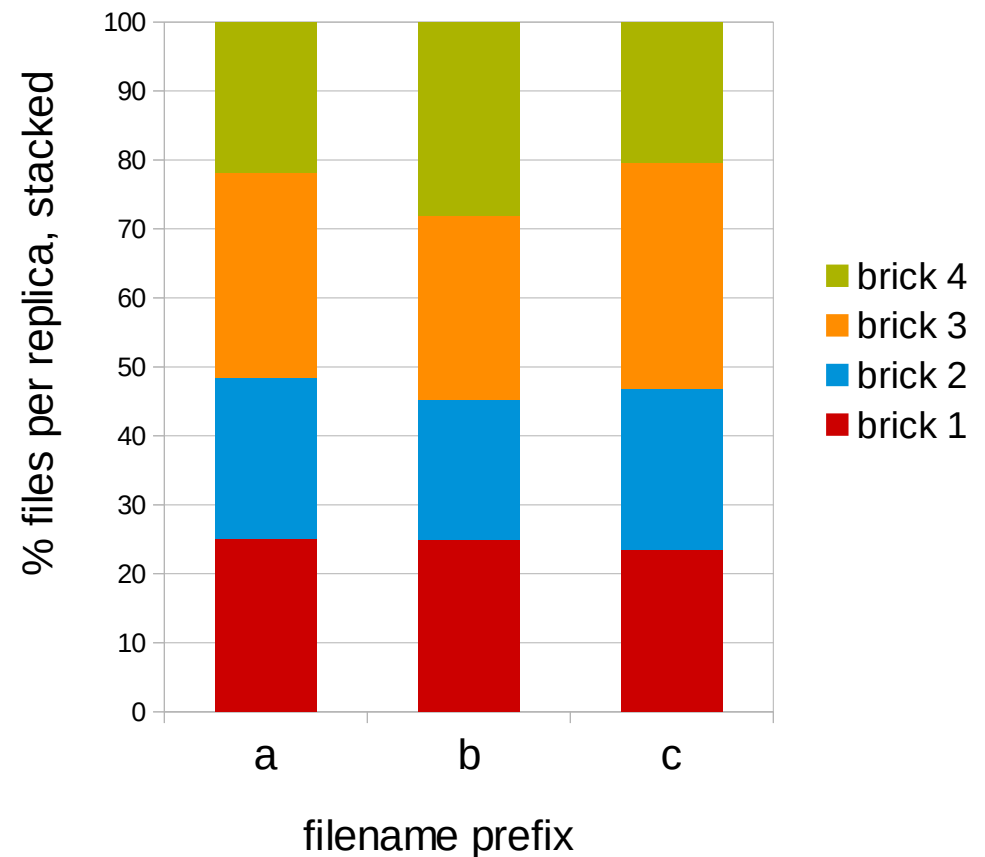


CPU,
Storage

*bottleneck*

network

RED HAT®
STORAGE

# streams:server ratio must be >> 1 for even load example: distribution of files across 4 bricks
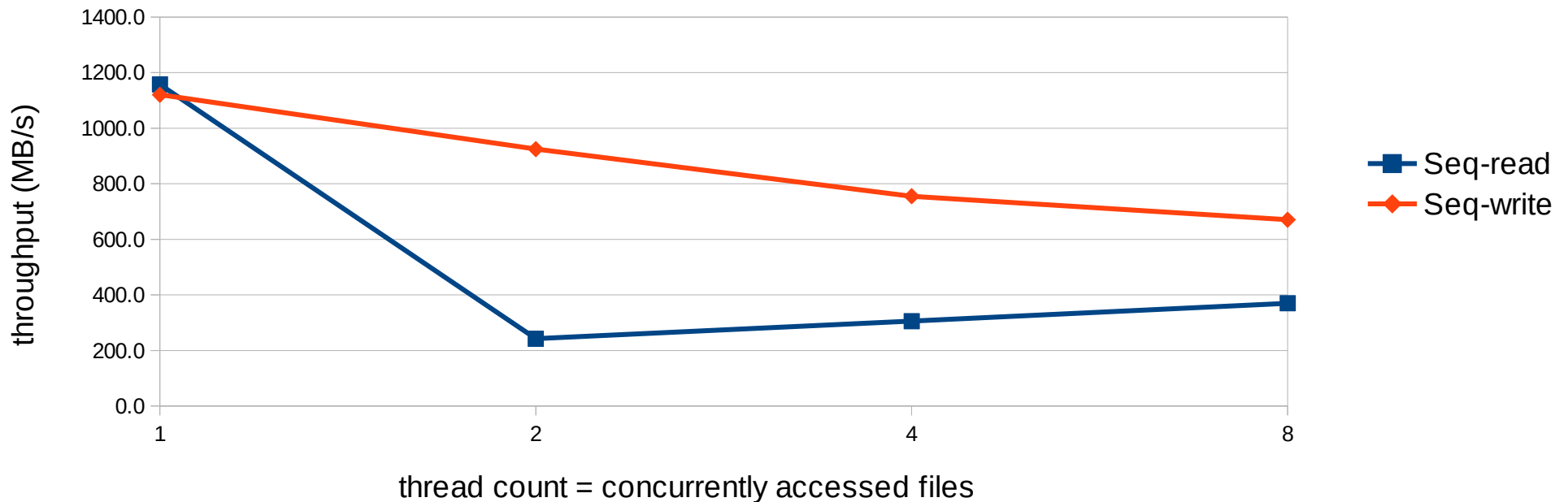
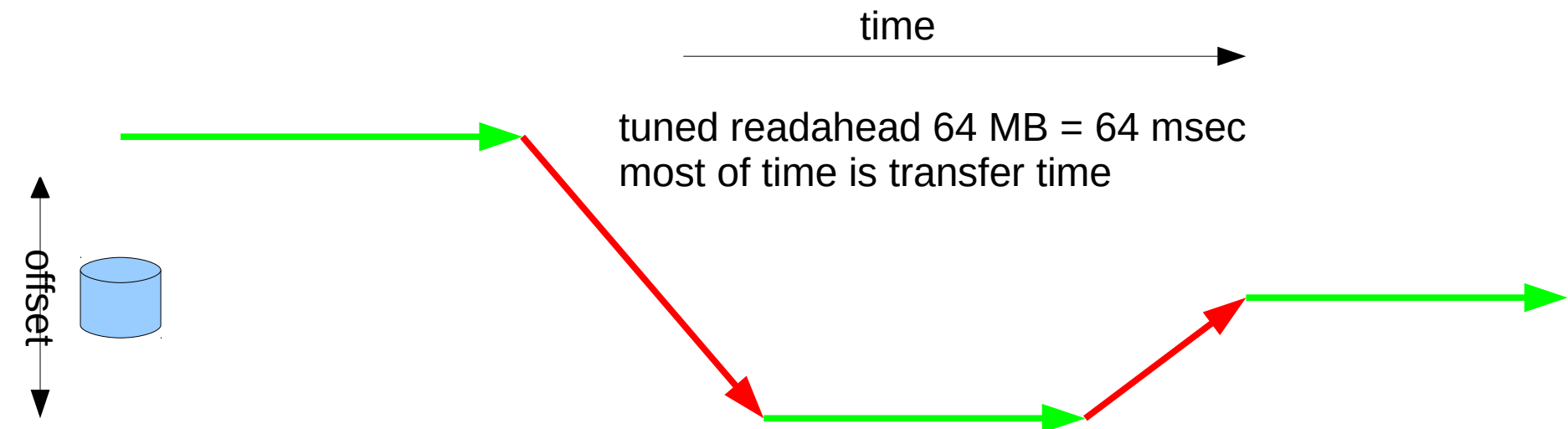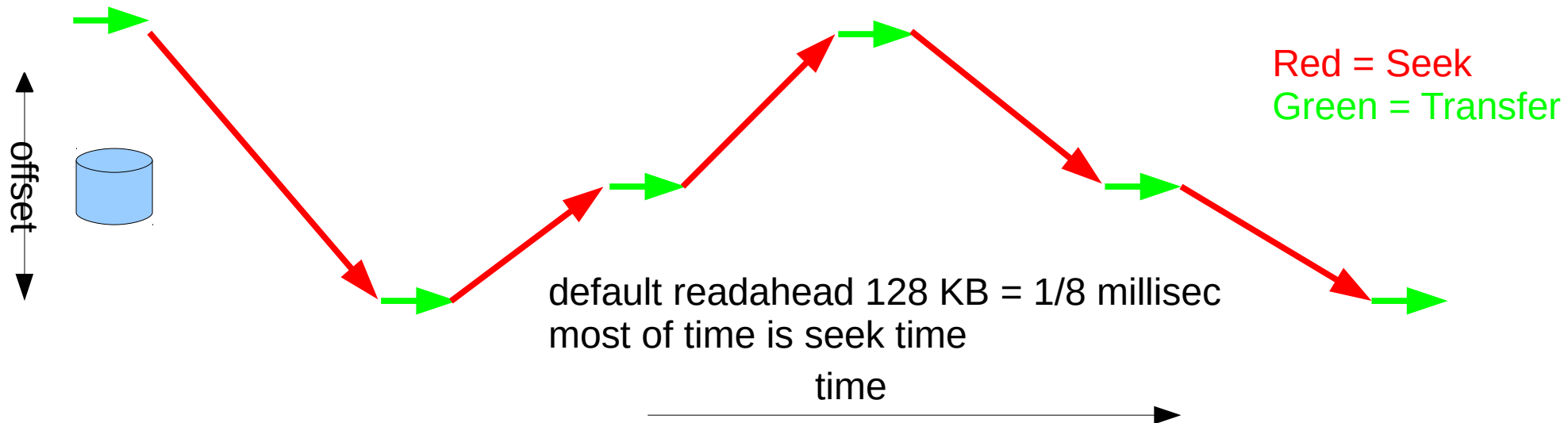## 2 files/brick



## 16 files/brick

# with >> 1 stream:brick, large-file sequential reads see I/O contention



XFS sequential large-file throughput vs threads

read_ahead_kb=512, deadline scheduler

RED HAT
STORAGE

# Eliminating seek time in multi-stream read

Red = Seek
Green = Transfer

offset

default readahead 128 KB = 1/8 millisec
most of time is seek time

time

time

tuned readahead 64 MB = 64 msec
most of time is transfer time

offset

RED HAT
STORAGE

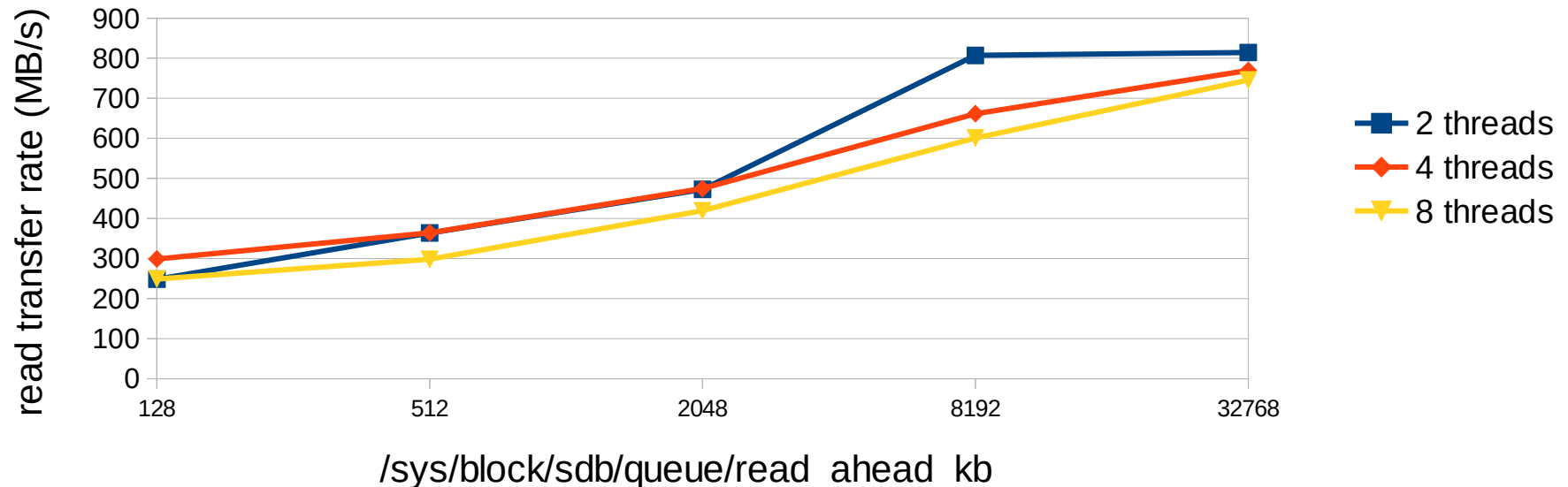# Solution: tuned profile "rhs-high-throughput"

uses 64-MB readahead in XFS bricks
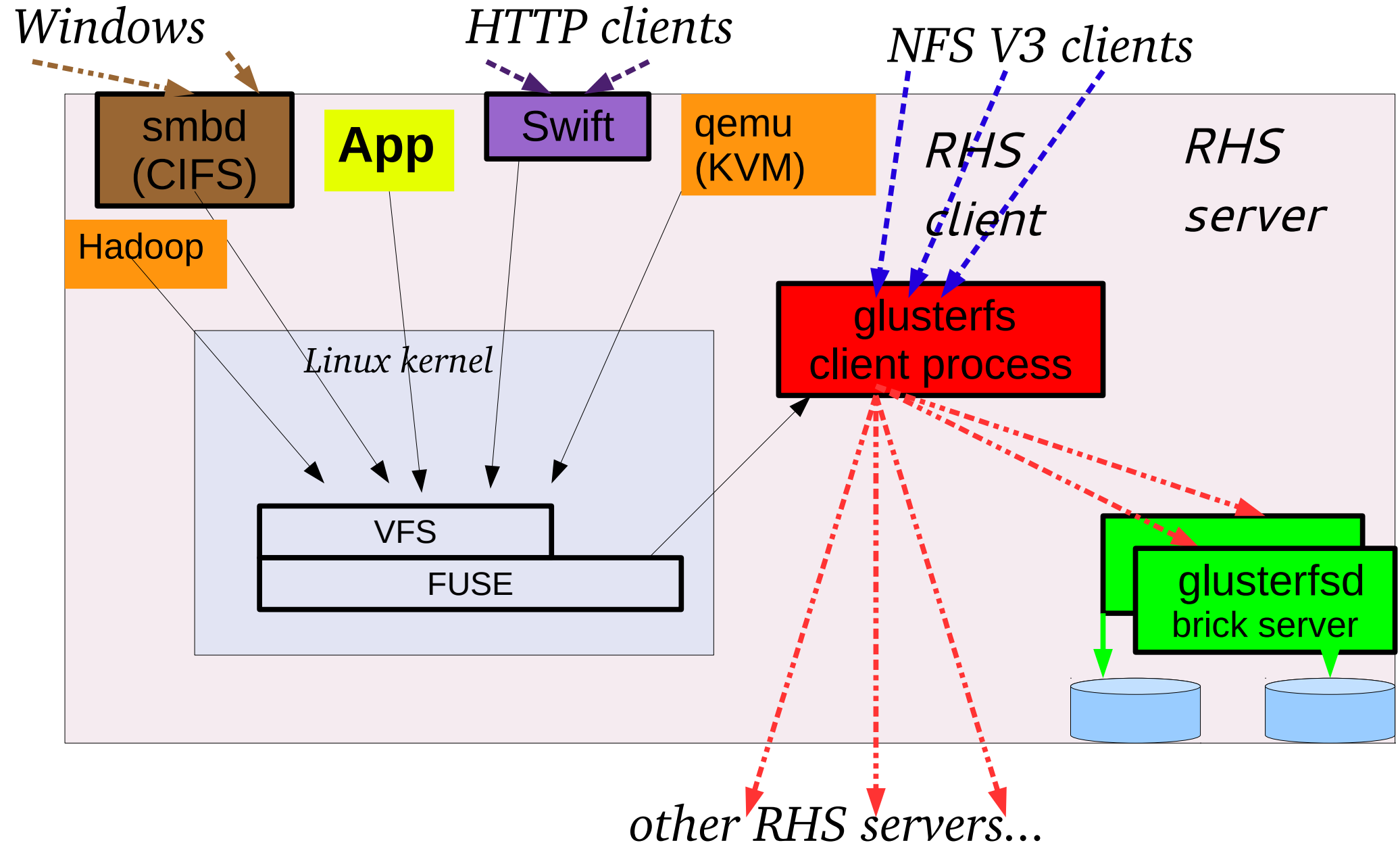
trades off latency for higher throughput

on servers: **tuned-adm profile rhs-high-throughput**

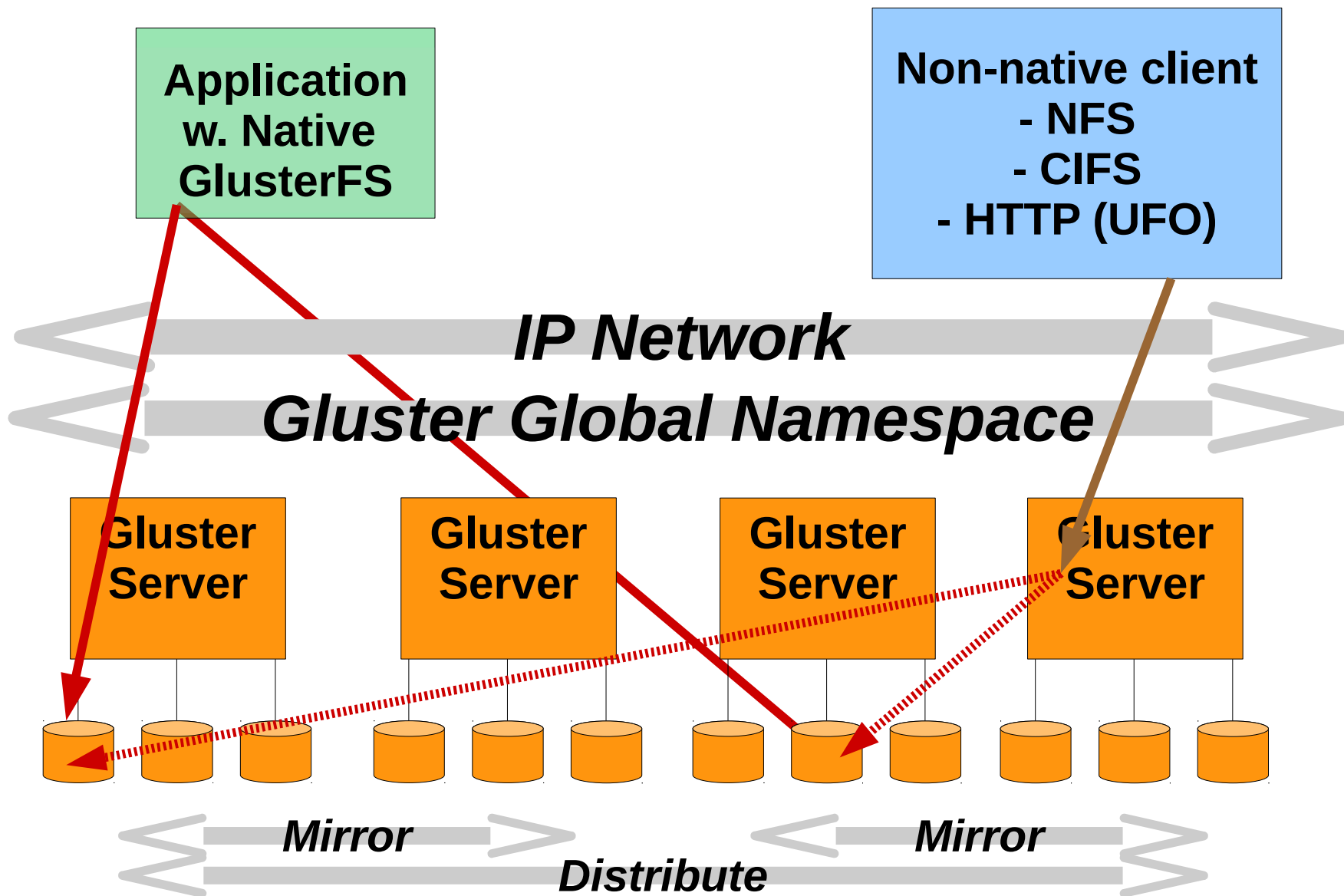effect of /dev/sdb readahead on Gluster read throughput

legend:  concurrently read files (threads)



- ■ 2 threads
- ◆ 4 threads
- ▼ 8 threads

y-axis: read transfer rate (MB/s), 0 to 900

x-axis: /sys/block/sdb/queue/read_ahead_kb — 128, 512, 2048, 8192, 32768

# Anatomy of Red Hat Storage (RHS)



*Windows*

*HTTP clients*

*NFS V3 clients*

smbd (CIFS)

**App**

Swift

qemu (KVM)

*RHS client*

*RHS server*

Hadoop

*Linux kernel*

VFS

FUSE

glusterfs client process

glusterfsd brick server

*other RHS servers...*

**RED HAT STORAGE**

# non-native protocol adds a network hop

**Application w. Native GlusterFS**

**Non-native client**
**- NFS**
**- CIFS**
**- HTTP (UFO)**

*IP Network*
*Gluster Global Namespace*

**Gluster Server**

**Gluster Server**

**Gluster Server**

**Gluster Server**

*Mirror*

*Mirror*

*Distribute*

RED HAT® STORAGE

# before deployment
# Is hardware sufficient?

capture network requirements

- – non-native reads –  2x application transfer rate
- – native writes – 2x app. transfer rate
- – non-native writes – 3x app. transfer rate
- – clients on same VLAN as Gluster volume?

capture storage requirements:

- – random reads – each disk = 100 IOPS
- – random writes – each disk = 15 IOPS
- – SSDs, flash can temporarily accelerate

# recommended storage brick configuration

12 drives/RAID6 LUN, 1 LUN / brick

    limited bricks/volume -> make bricks big

hardware RAID stripe size 256 KB (default 64 KB)

pvcreate –dataalignment 2560k

mkfs.xfs **-i size=512** -n size=8192 \

    -d su=256k,sw=10 /dev/vg_bricks/lv*N*

mount options: inode64

# Deploying network

if non-native protocol only, separate Gluster and non-Gluster traffic onto separate VLANs
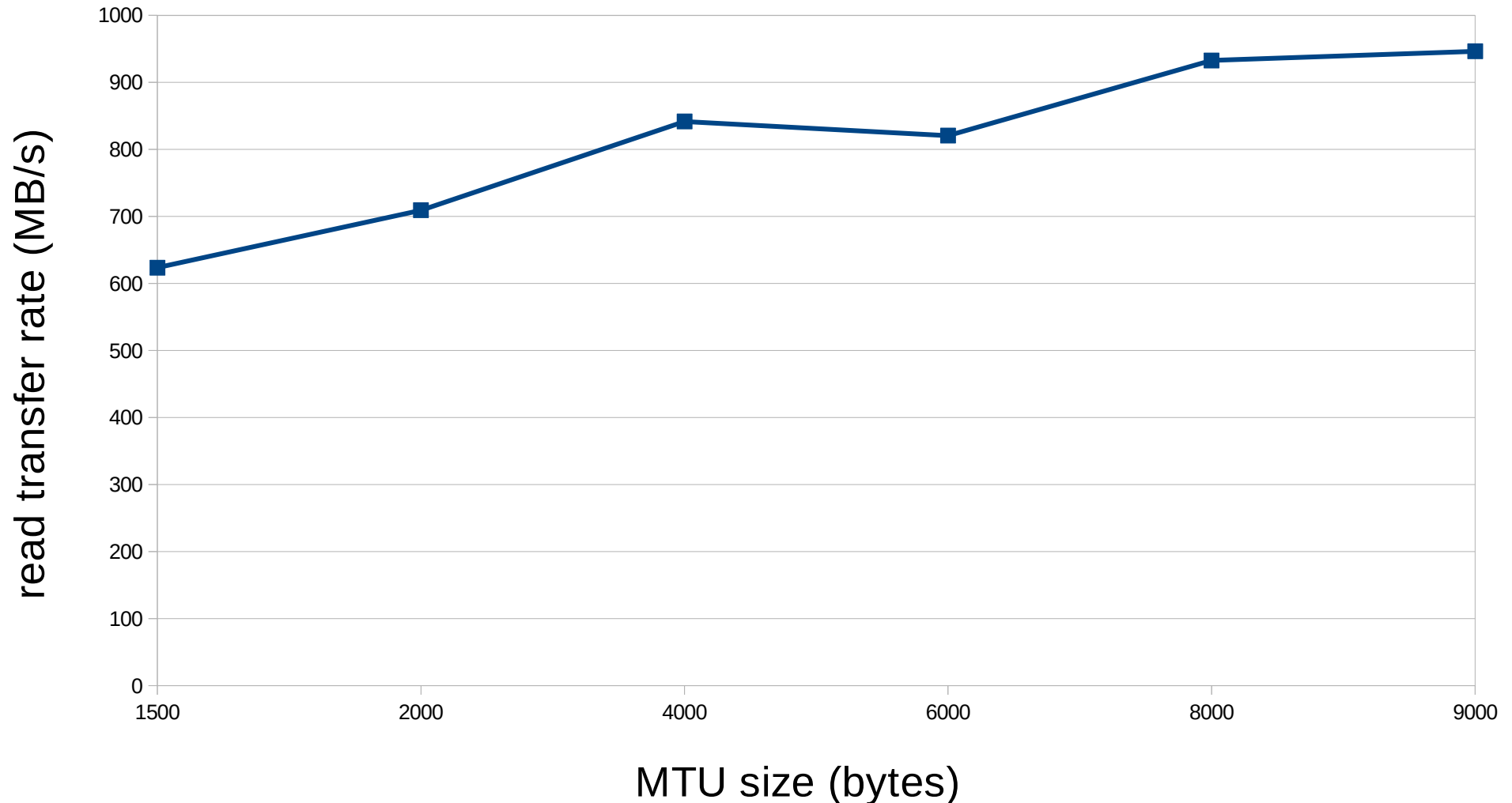
- separates replication traffic from user traffic

jumbo frames – improve throughput, but requires switch configuration

network design – Gluster doesn't know about switch boundaries

**RED HAT®
STORAGE**

# effect of jumbo frames

effect of MTU size on Gluster read throughput

1 client, 1 server

RED HAT®
STORAGE

# Network interconnection problem example top-of-rack switches

RED HAT
STORAGE

# tuning tips

**gluster volume set *your-vol* <u>eager-lock on</u>**

**per-server:**

- **tuned-adm profile rhs-high-throughput**

**native (glusterfs) client:**

- write transfer size > 32 KB

- avoid single-threaded apps
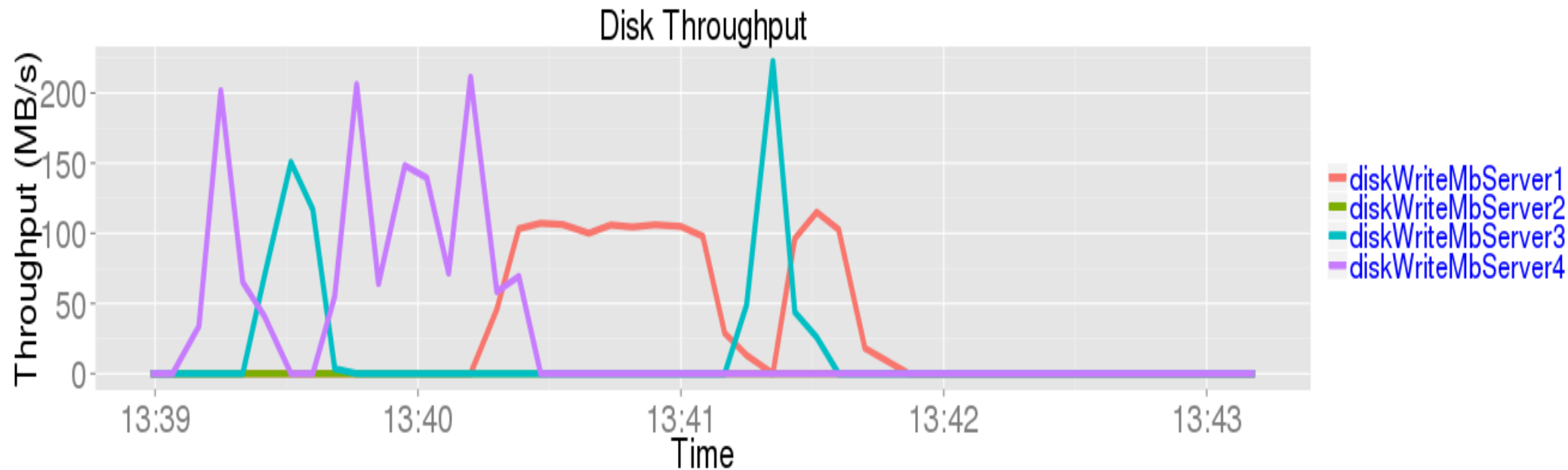
- more bricks, mountpoints help wi random I/O, small files

# Performance with and without brick recommendations

RED HAT®
STORAGE

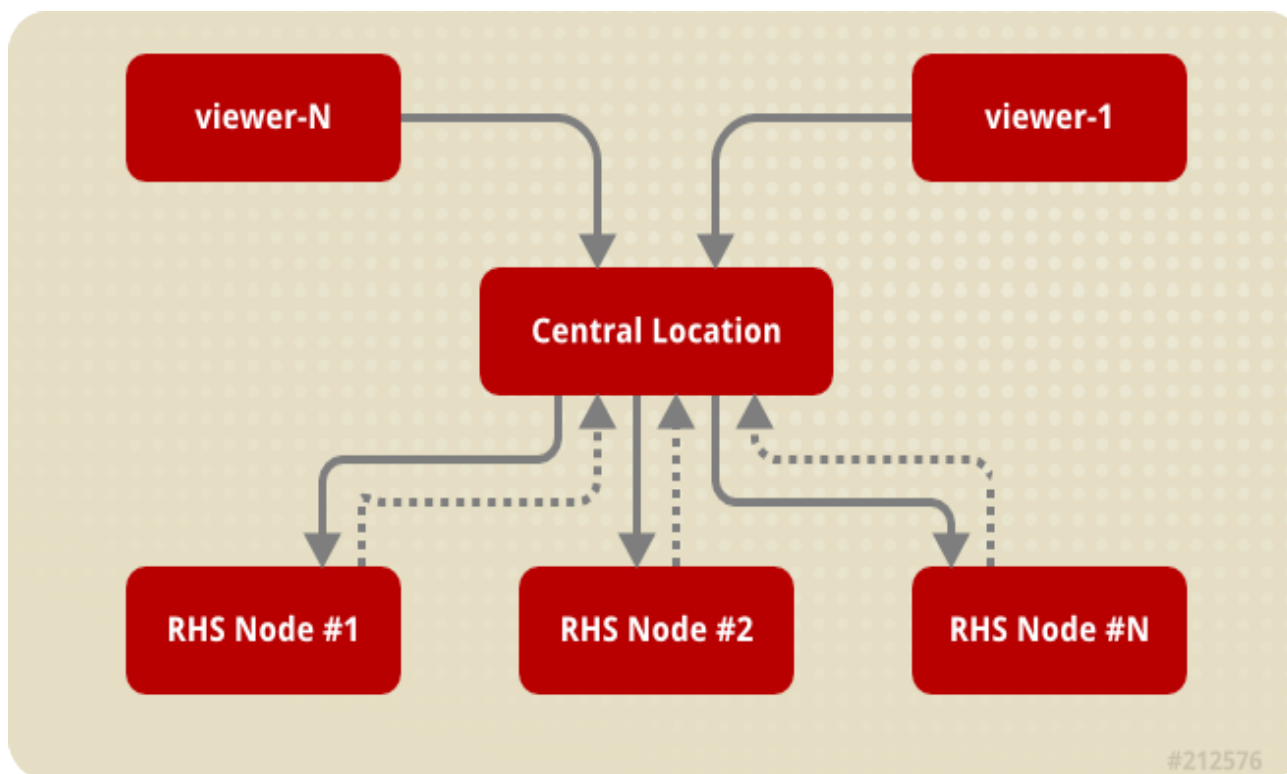# Performance with and without brick recommendations

**RED HAT® STORAGE**

# Hashing Behaviour for small number of files

**RED HAT® STORAGE**

# Hashing Behaviour for large number of files

**RED HAT® STORAGE**

# Real time data collection and analysis tool

RED HAT® STORAGE

# Real time data collection and analysis tool

# Real time data collection and analysis tool

- On RHS nodes
  - System tool to collect the stats
    - PCP  (Performance Co-Pilot)
      - http://oss.sgi.com/projects/pcp/
- On Central Location
  - Way to query specific performance stats from RHS nodes.
    - pmwebd
      - A web daemon of PCP to request the stats from PCP service running on RHS nodes.
  - Way to store the results from above query
    - Elasticsearch (http://www.elasticsearch.org/)
      - distributed RESTful search and analytics
  - Way to serve requested stats from the viewers
    - R
      - R is a free software environment for statistical computing and graphics.
    - Shiny
      - Turn analyses of R into interactive web applications that anyone can use.
- On Viewers
  - Web browser

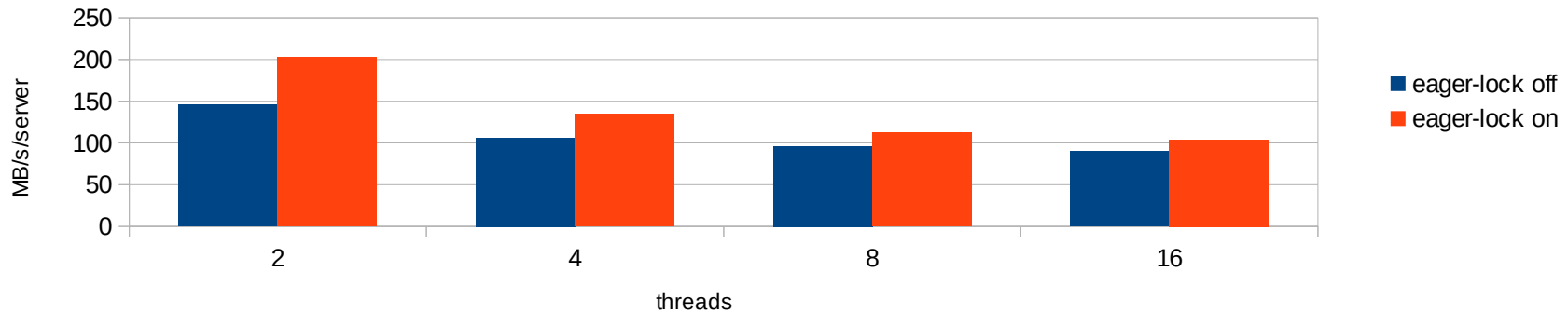# Real time data collection and analysis tool

https://forge.gluster.org/gluster-performance-stats-collection-and-analysis-tool

RED HAT®
STORAGE

# End of presentation

Additional material in slides that follow

**RED HAT®
STORAGE**

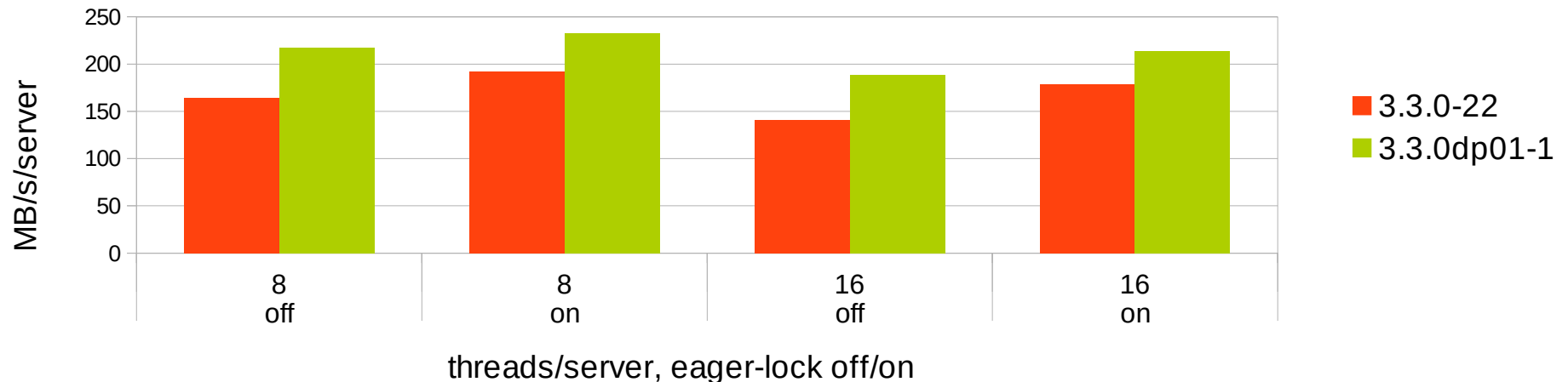# NFS large-file performance gains



effect of eager-lock on Gluster-NFS sequential write

RHS 2.0, 8 servers, 2 clients/server



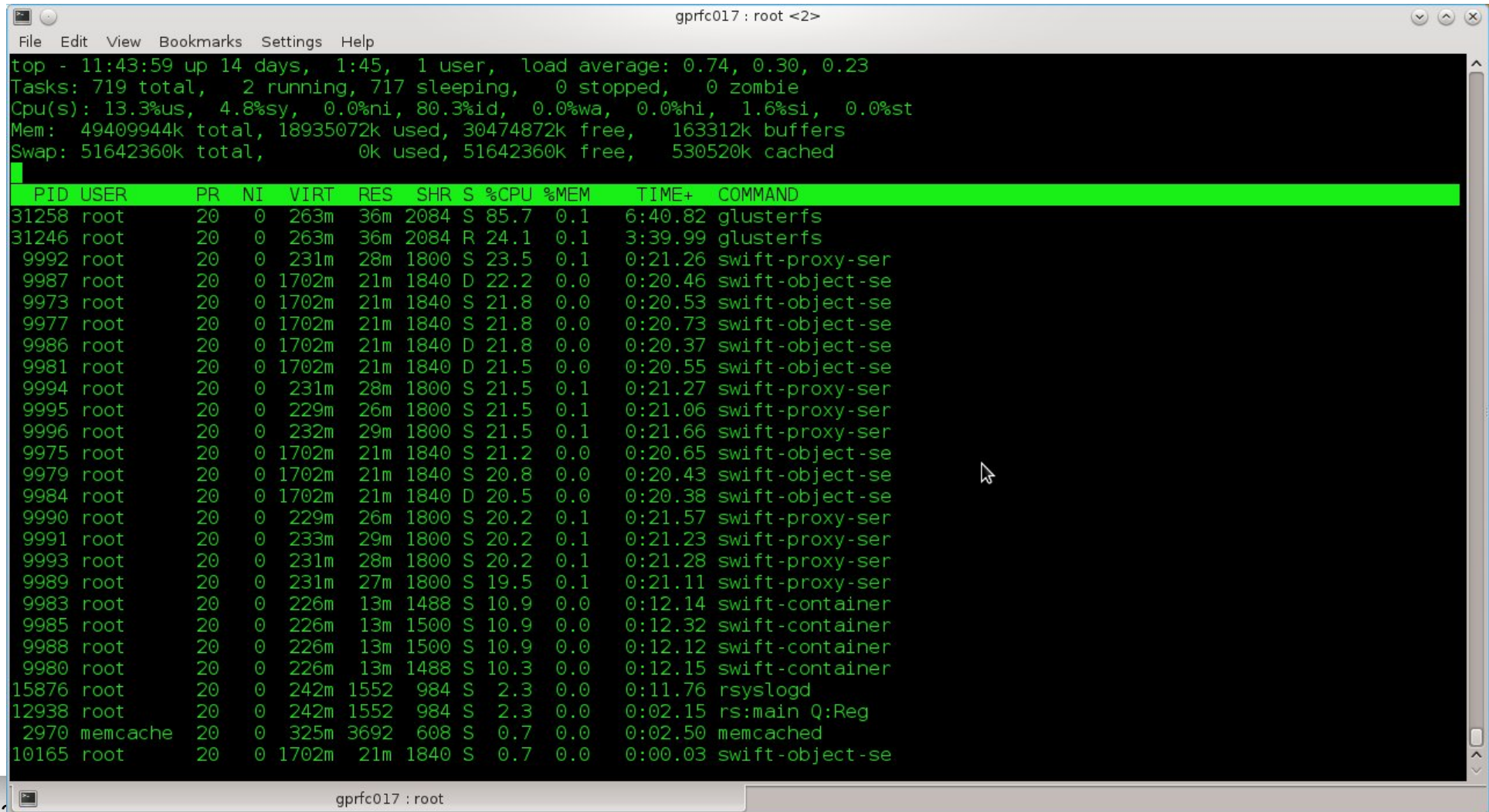Effect of deferred-unlock patch for higher thread/server counts

2 replicas, 4 servers, 4 clients/server

RED HAT
STORAGE

# IOPS – requests/sec bottleneck

In Object Server PUT workload

- – glusterfs thread is 85% of 1 core -- bottleneck
- – average system CPU core utilization 20%

STORAGE

# RPC read optimization in RHS 2.1

17:39:16.467980 epoll_wait(3, {{EPOLLIN, {u32=9, u64=4294967305}}}, 257, 4294967295) = 1 <0.000013>
17:39:16.468035 **readv**(9, [{"\200\0\20\214", 4}], 1) = 4 <0.000026>
17:39:16.468097 **readv**(9, [{"\0 \366$\0\0\0\1", 8}], 1) = 8 <0.000034>
17:39:16.468168 **readv**(9, [{"\0\0\0\0", 4}], 1) = 4 <0.000010>
17:39:16.468209 **readv**(9, [{"\0\0\0\0\0\0\0\0", 8}], 1) = 8 <0.000007>
17:39:16.468242 **readv**(9, [{"\0\0\0\0", 4}], 1) = 4 <0.000008>
17:39:16.468283 **readv**(9, [...], 1) = 116 <0.000015>
17:39:16.468359 **readv**(9, [..., 4096}], 1) = 4096 <0.000023>
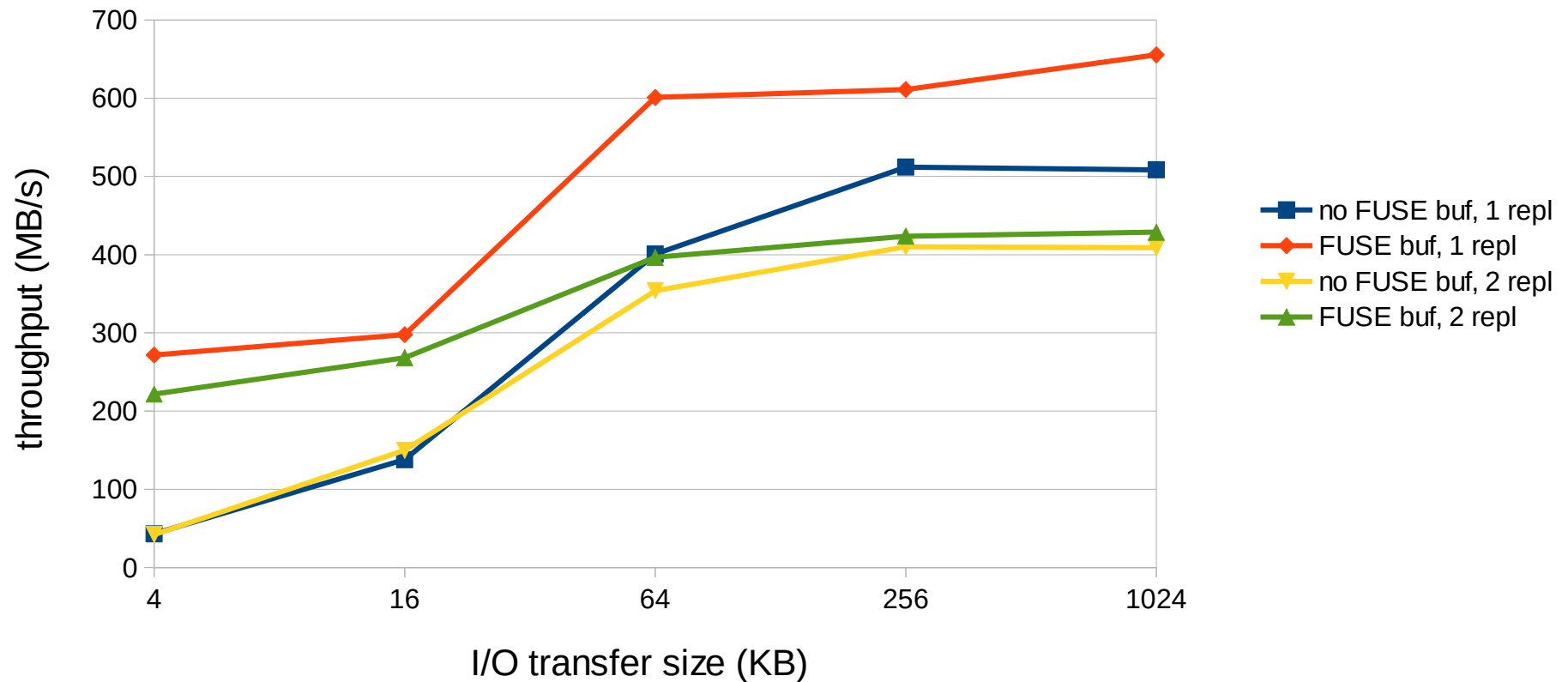
RED HAT®
STORAGE

# readdirplus - efficient metadata reads

- classic small-file app scans directory and performs operations on files in directory
  - example: **rsync -ravu  /home/ben /mnt/glusterfs/ben**
  - problem: each file's metadata requires *network round-trip*, so big directory listings take a long time
  - solution from NFS: **readdirplus** – return directory contents together with important metadata about each file in directory
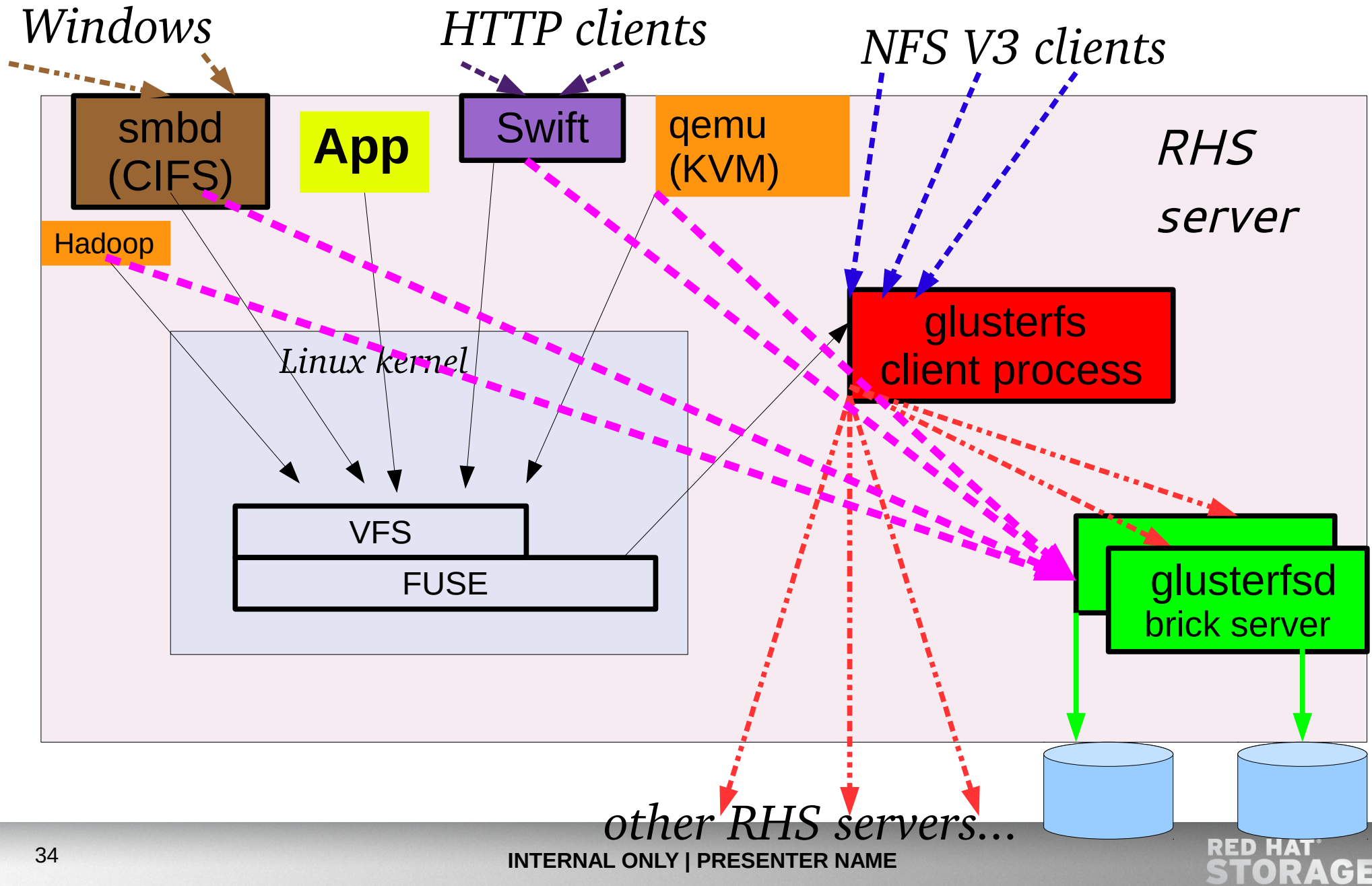  - upstream now, release in RHS 2.1

# Upstream Linux 3.7 FUSE enhancement for small-transfer writes

effect of FUSE write buffering on Gluster single-thread write performance

1 client, 1 thread,4-GB file size,  volume has eager-lock and client-io-threads enabled

RED HAT
STORAGE

# Anatomy of Red Hat Storage (RHS) – libgfapi

# The network IS the system

test between clients and servers if possible

test between servers 2nd best

network – scripts available

large-file – **iozone -+m (-+h)**

small-file – **smallfile_cli.py --host-set**

Swift – **ssbench** will exercise

**RED HAT®
STORAGE**

# potential  perf. problems in storage device

Check no failed disks  – RAID6 will continue but performance will degrade.

Check writeback caching is enabled.  If dependent on battery, check that battery is installed and working.

disconnected brick can cause perf. degradation (example: brick wasn't mounted)

**RED HAT®**
**STORAGE**

# capturing perf. problems onsite

**top** utility – press **H** to show per-thread CPU utilization, will detect "hot-thread" problems where thread is using up its core
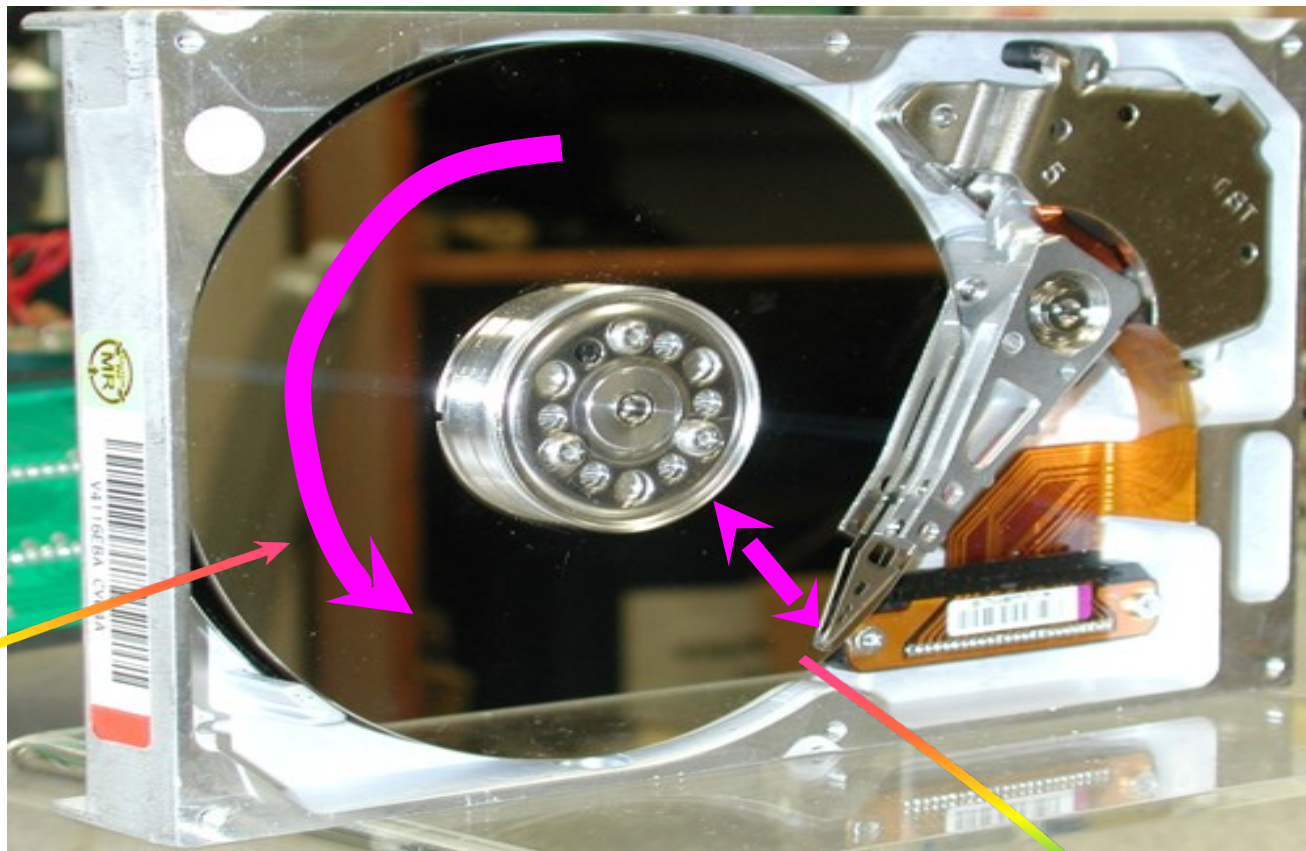
NFS: **nfsiostat** and **nfsstat** utilities

**gluster volume profile** –  shows latency, throughput for Gluster RPC operations

**gluster volume top** – shows which files, servers are hot

**Wireshark Gluster plug-in** – isolates problems

RED HAT® STORAGE

# effect of non-sequential I/O



7200 RPM → 4.15 ms average rotation time

Assume 5 ms average cylinder seek time

At 50 MB/s, 32 KB request size takes 0.625 ms per transfer

Total = about 10 ms / rq → 100 random IOPS/drive

divide by 6 for random writes to RAID6 LUN

•divide by 2 for random writes to RAID10 LUN