



CERN

European Organization for Nuclear Research  
Organisation Européenne pour la Recherche Nucléaire



# Large Storage Systems Present & Future

Dr. Andreas-Joachim Peters

CHEP 2012 - New York

[andreas.joachim.peters@cern.ch](mailto:andreas.joachim.peters@cern.ch)

CERN IT Department  
CH-1211 Genève 23  
Switzerland  
[www.cern.ch/it](http://www.cern.ch/it)



# Outline

- Why are large storage systems relevant?
- What are the challenges?
- How large are storage systems today?
- Which technologies are used?
- Which technologies are emerging?
- What can we expect from current and future technology?
- What can we improve in LHC storage?



# Why are large storage systems relevant?

According to [International Data Corporation](#), the total amount of global data is expected to grow to **2.7 zettabytes** during 2012.

This is **48%** up from 2011.

The storage industry has sold more than **350 exabyte** of storage in 2011.

**Multi-PB** storage systems are a **norm** and available by many vendors!



# Storage Market 2011

Technology	Situation	Units Sold [ Million ]	Volume Sold [ Exa Byte ]
<b>DRAM</b>	4 companies have 90% of market share: Samsung, Samsung, Hynix, Micron, Elpida (bankrupt)	800	2
<b>NAND</b>	4 companies have 99% market share: Samsung, Toshiba, Micron, Hynix	4000	19
<b>SSD</b>	> 50 companies	17	3*
<b>HDD</b>	3 companies only: Western Digital 50%, Seagate 39%, Toshiba 11%	630	330
<b>TAPE</b>	3 technologies: IBM, Oracle, LTO---consortium LTO has 90% market share	27	20

90% of capacity sold!

\* included in NAND numbers

From Bernd Panzer Steindel/CERN



# LSS Challenges and Implications

- **large** volume & meta data

- PB Stores with millions to billions of objects
  - can a hierarchical namespace be kept?

**scale-out** storage systems, aggregation & federation

new **redundancy** methods - RAIN, eventual consistency with quorum on object versions

- new scale of hardware **failures**

- e.g. 12k disks ~ 1 failing disk per day
- multi-PB RAID-6 probability of data loss within 5y at  $n\%$  level

**elastic** block storage w/o DHTs/dynamo principle  
virtualization

- flexible, sizable and administrable

- capacity growth and life-cycle management in production

**storage tiering**

- manifold **requirements** and tuning parameters

- bandwidth, IOPS, latencies, costs, volume/object/user scalability
- interfaces for objects, files, block devices, cloud infrastructure
- technology evolution decreases performance:capacity ratio

**virtualization**

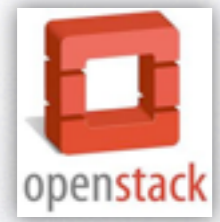
**volume** grows faster than **bandwidth** - MB/s per GB

Challenge

finally: costs matter more at large scale!



# Some Challengers ...

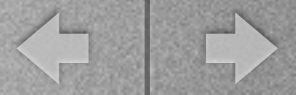


IBM Systems > Software > SONAS

## General Parallel File System

Efficient storage management for big data applications





# The world is divided

**POSIX**

**Cluster Filesystems**

**Commercial Products**

**Hardware Solutions**

“Storage in a Box”



**Cloud/MapReduce  
Storage**

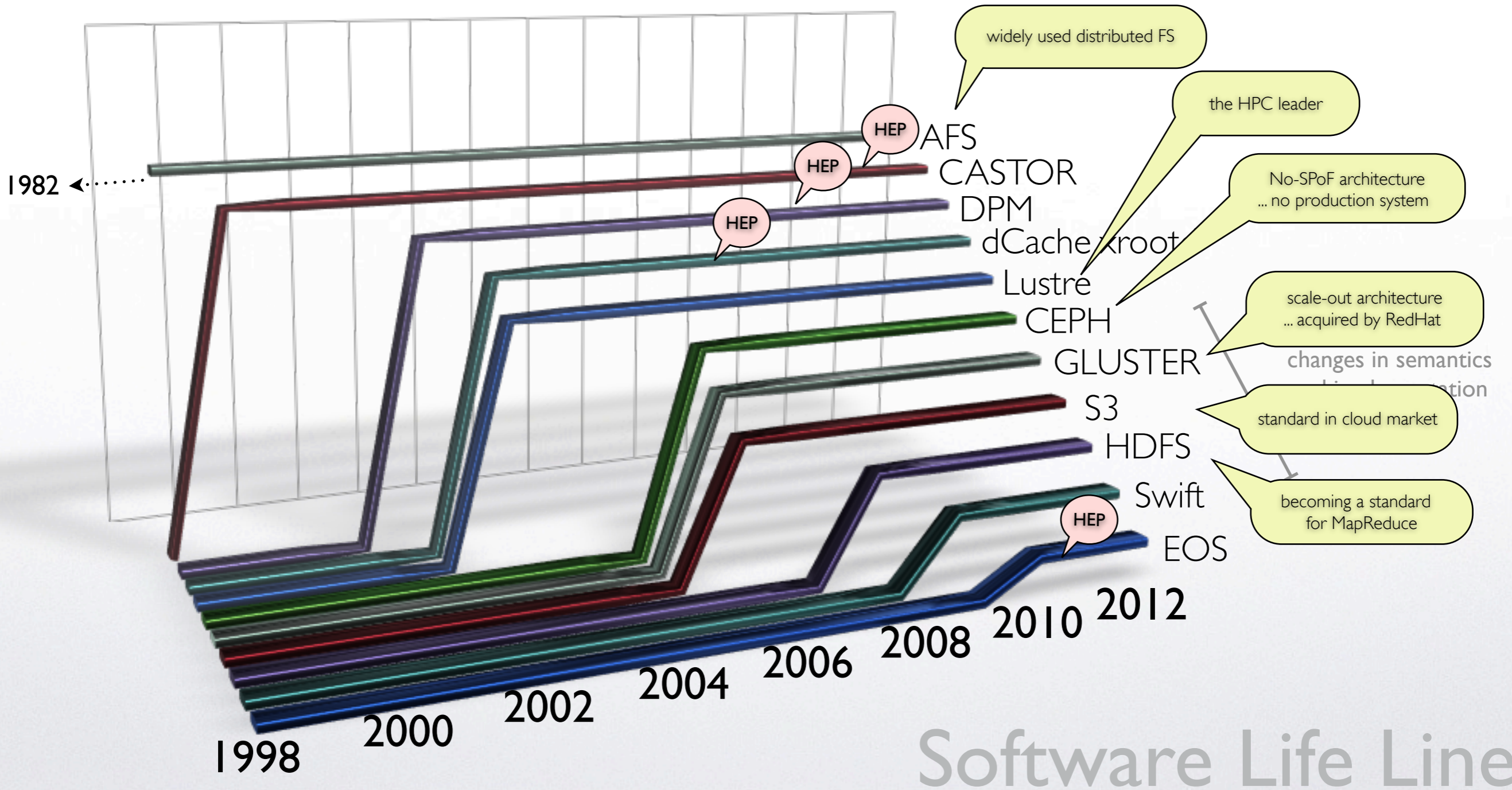
**Open Source**

**Software Solutions**  
“Buy & Build”

... but also coalescing  
into many hybrid storage systems ...



# (Open) Solutions have a long way in(-to) production ...





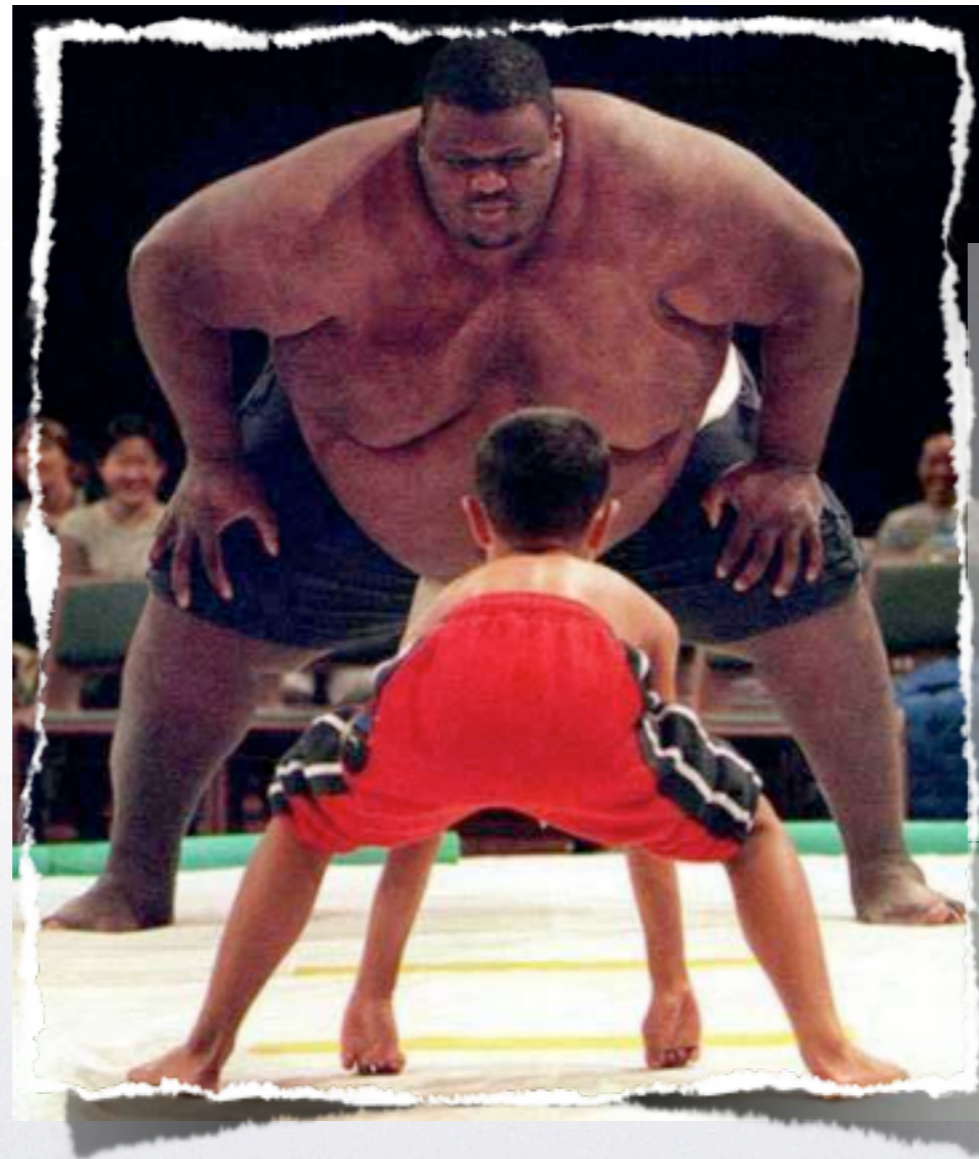


# What is a large storage system today?

What?



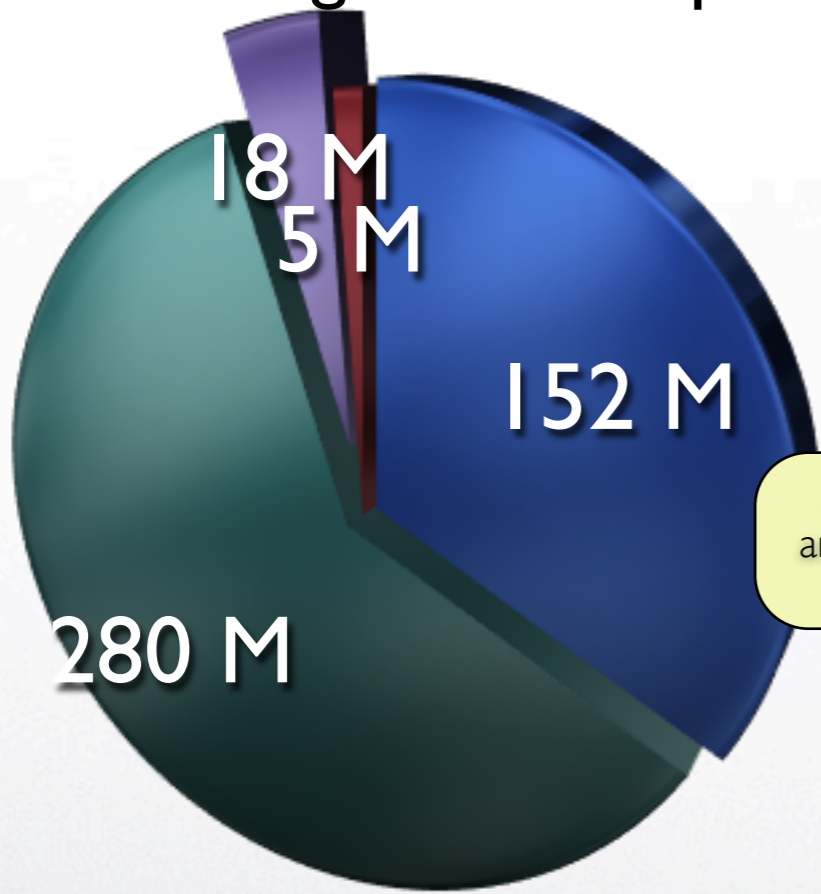
# Is LHC Storage large?





# LHC Storage

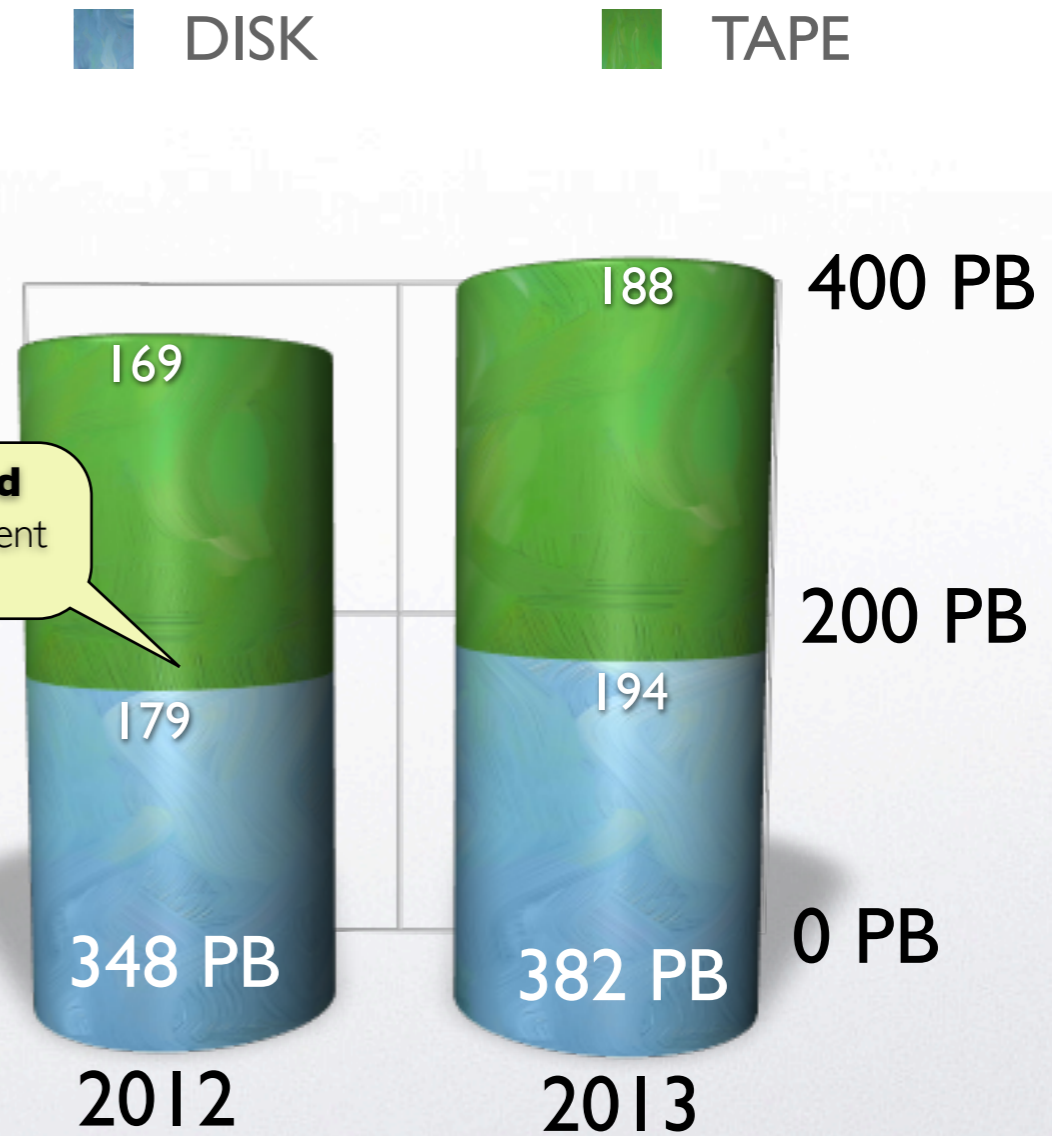
### Number of managed files - April 2012



This are only managed files (there are more user files)

- ALICE
- ATLAS
- CMS
- LHCb

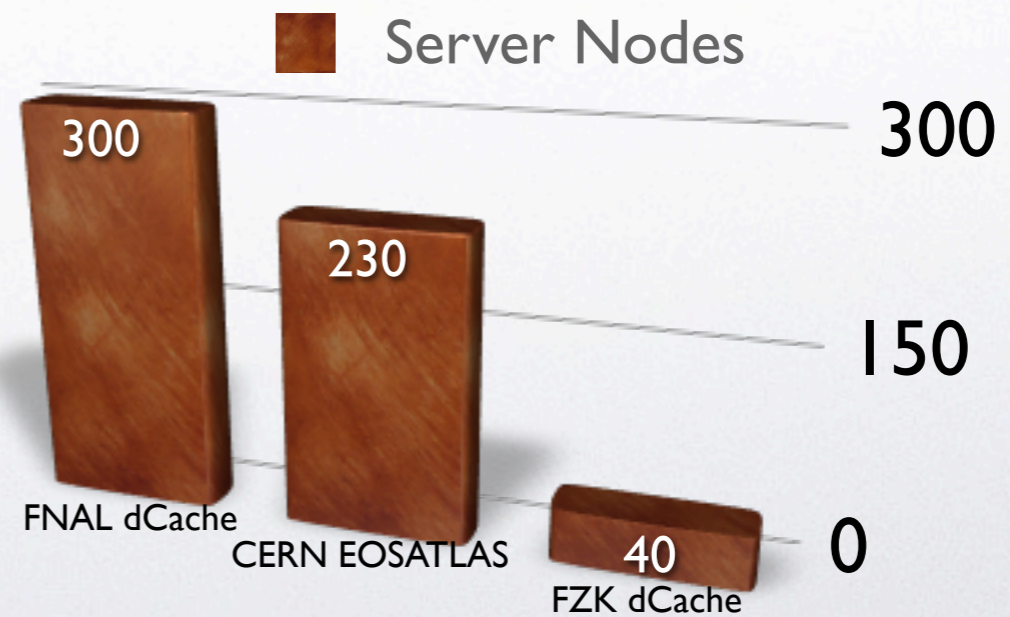
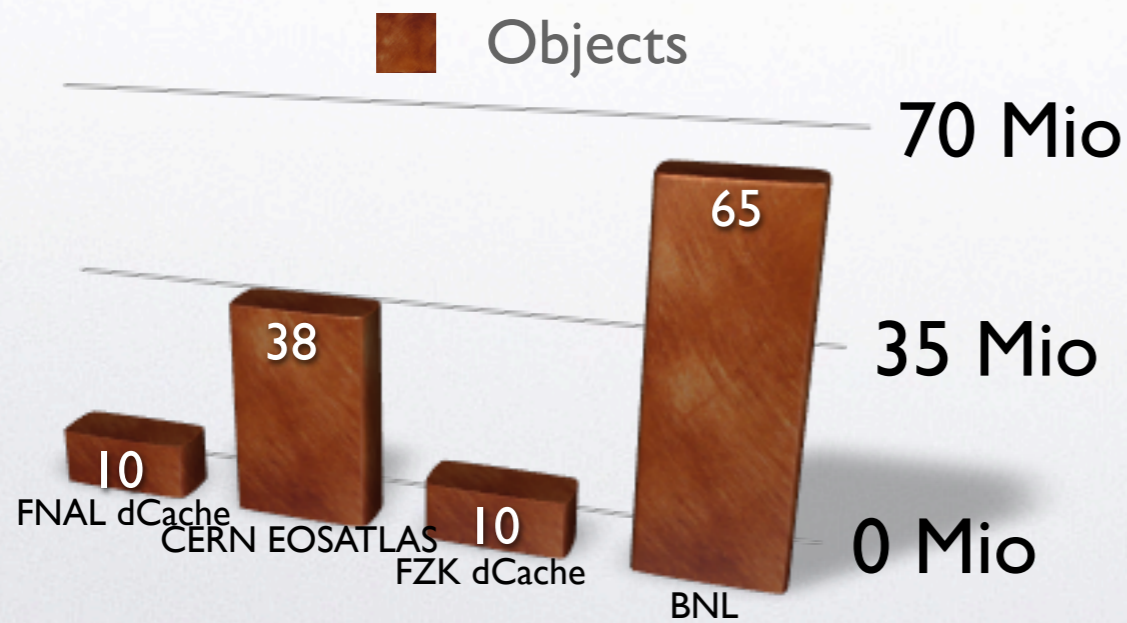
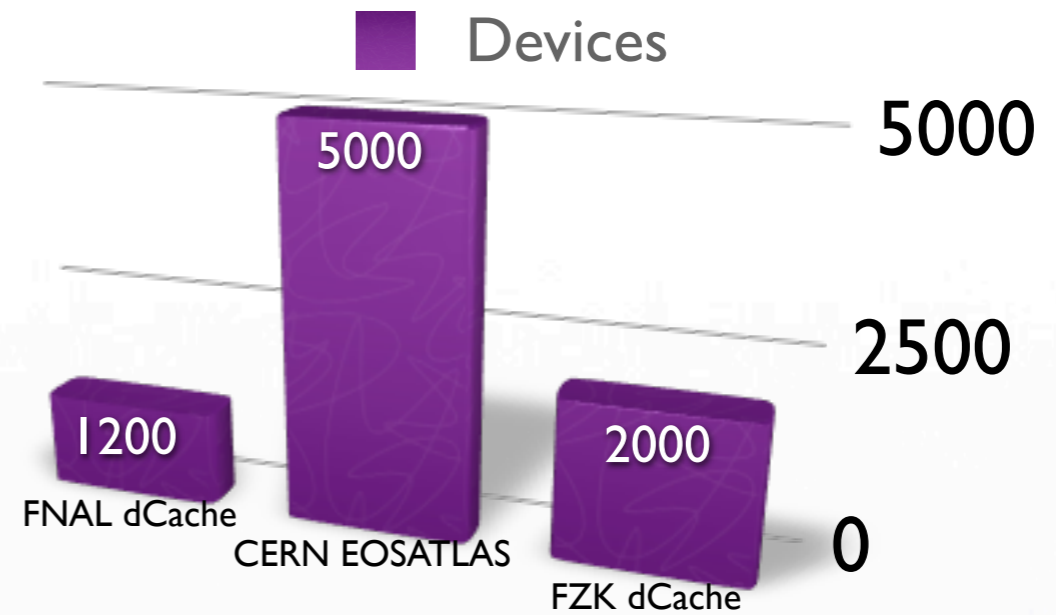
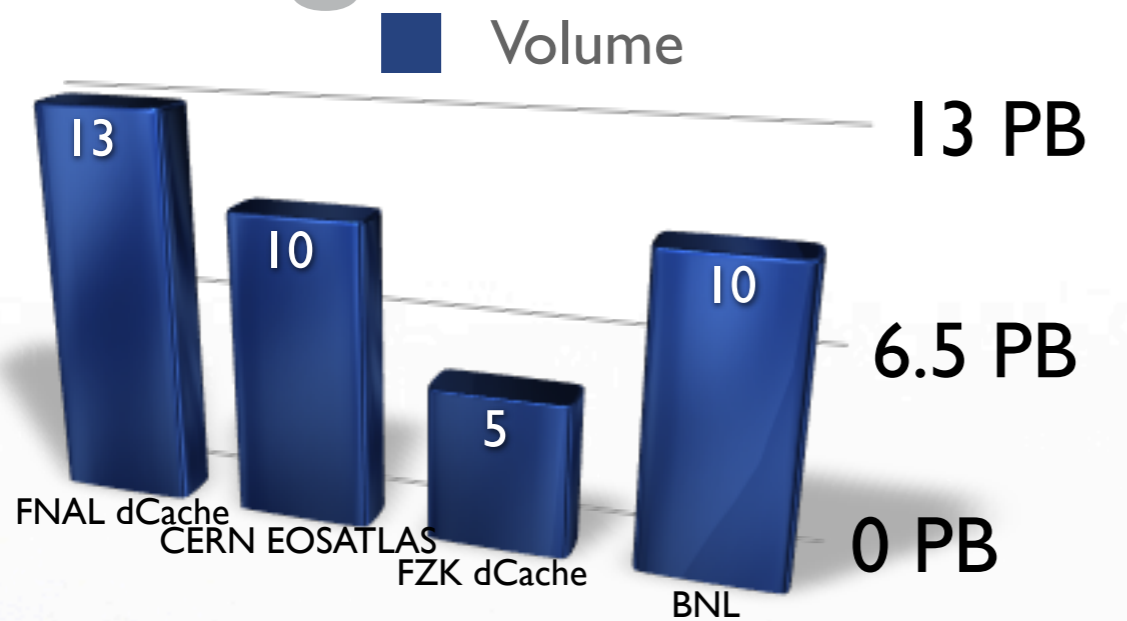
### CRSG Recommendations

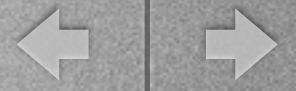


... the storage is **aggregated** and **virtualized** by experiment frameworks



# Large LHC Storage Instances ...

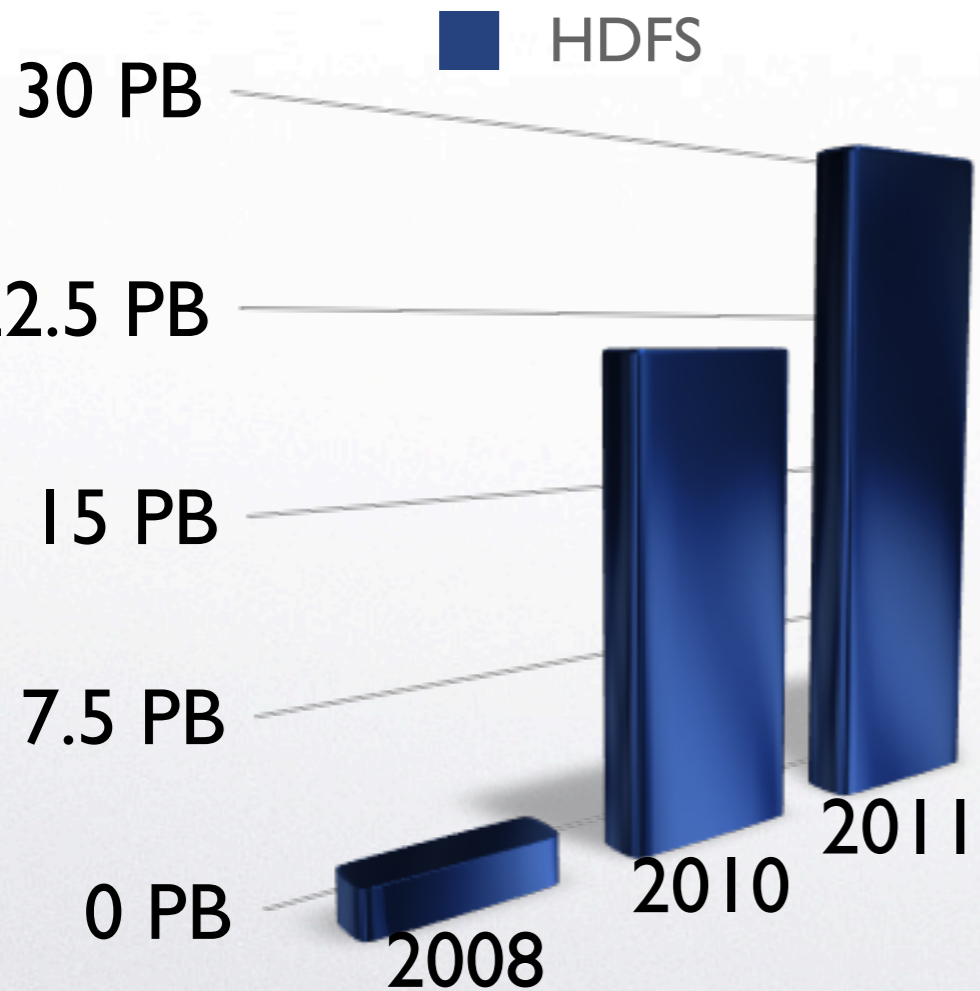




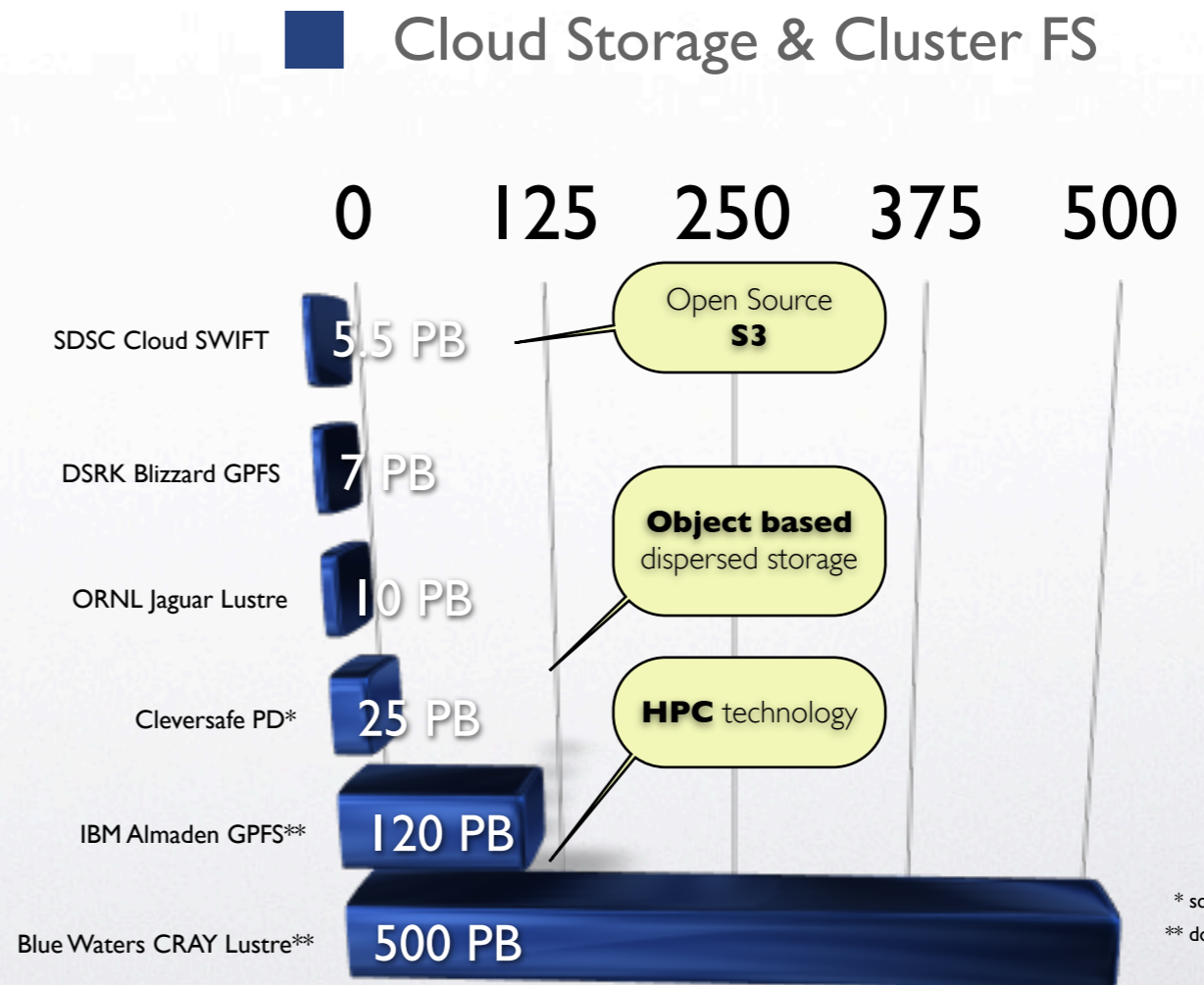
# Other Large Storage Systems

## Facebook

## Large Storage Installations



From facebook.com Paul Yang Facebook Engineering



\* scales to 10 EB  
\*\* does not exist yet!

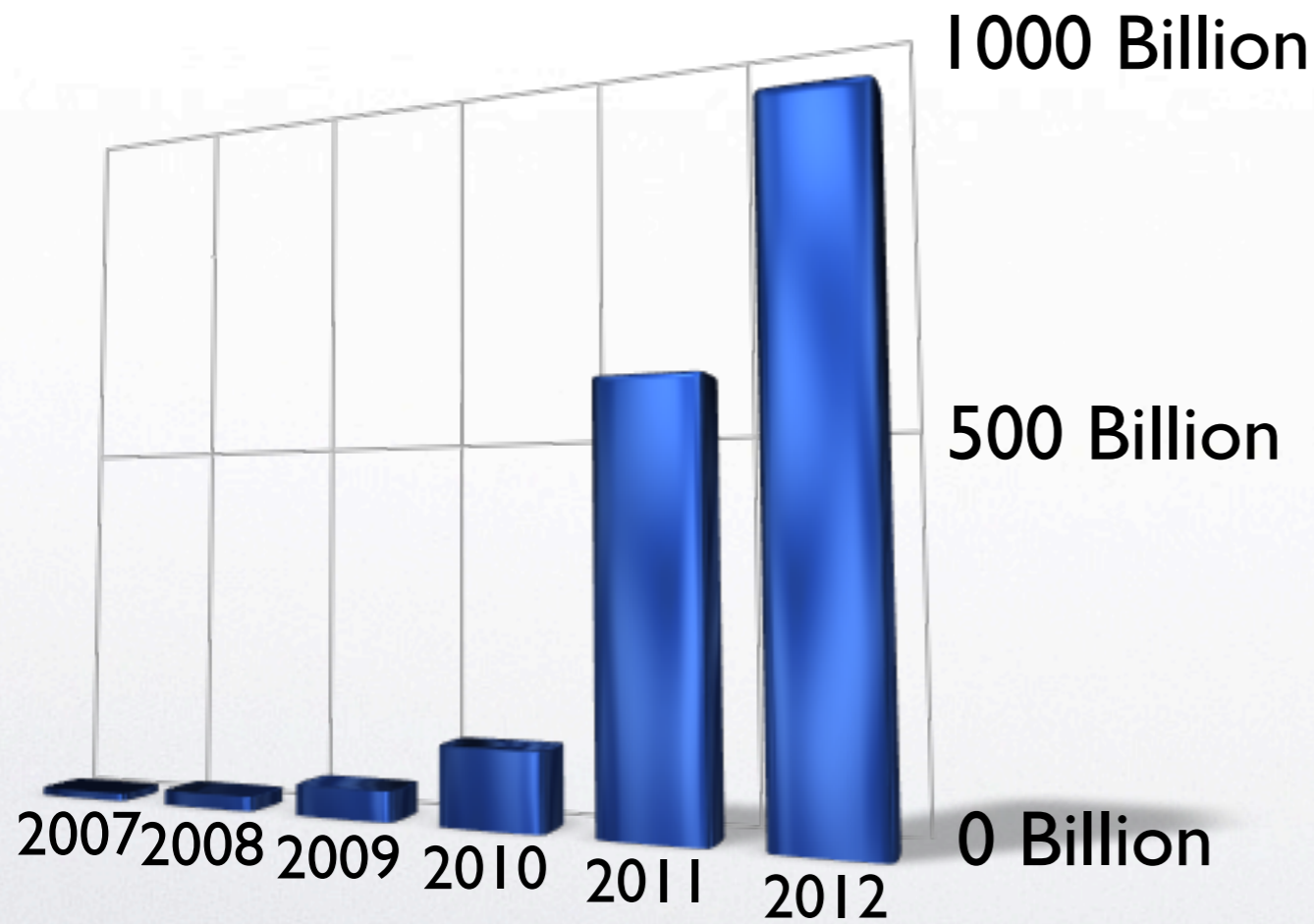
Today: **Cleversafe** 10 EB system would require 4.5M disks and cost several billion \$\$ !



# Amazon S3 & Yahoo!

■ Objects

■ Disk Storage



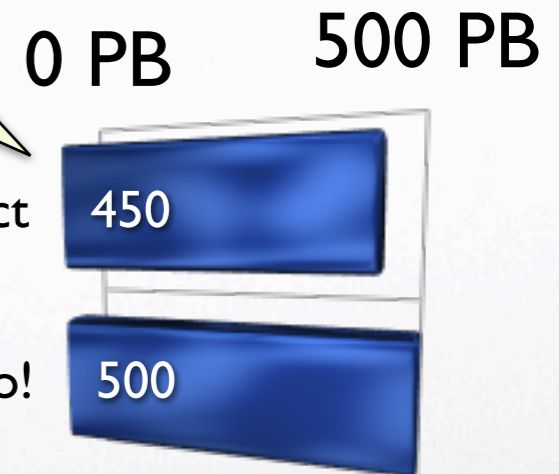
From Amazon Web Service Blog

650k Request/s

There is no information available about the average object size. This is just an exemplary assumption.

Amazon 512k per object

Yahoo!



From LTUG 2012

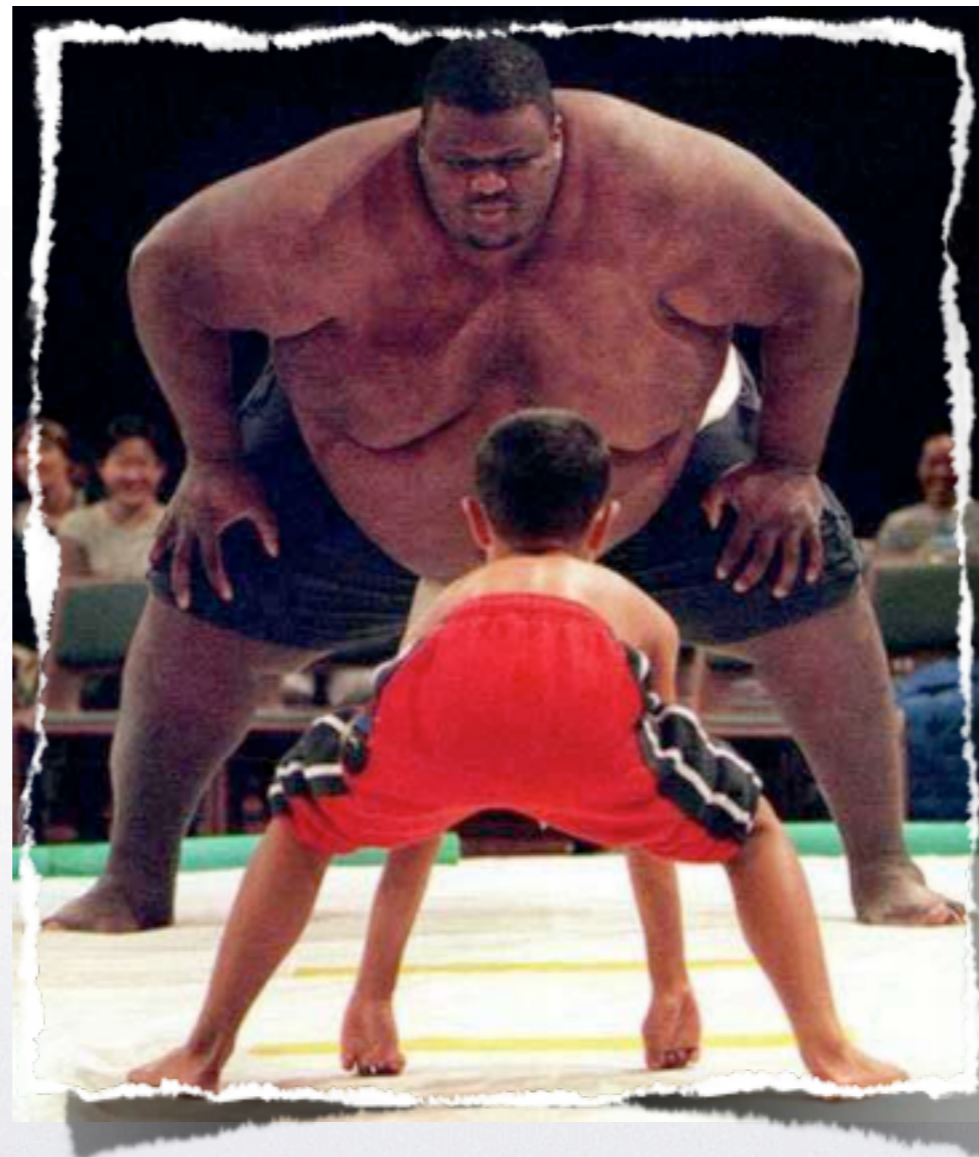
GOOGLE. does not publish these numbers ...

... there are many other Cloud Storage Provider: DropBox, iCloud, Google Drive ...

Other Storage



# Is LHC Storage large?



LHC Storage is large in volume - not in number of objects!



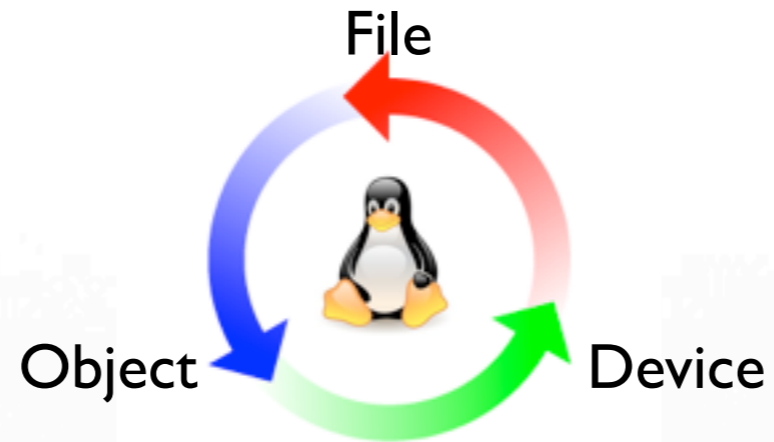
# Technology Trends







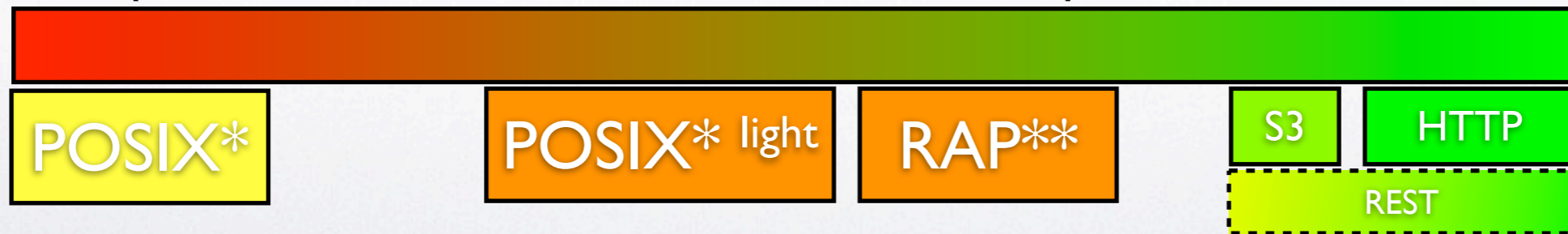
# Storage Interfaces



## Access Protocols

Complex/rich features set

Simple/reduced feature set



standard applications run out of the box with this interface

applications have to be enabled or use a download - process with POSIX - upload approach

\* POSIX via AFS, Lustre, GPFS, pNFS or FUSE client
\*\* RAP: remote access protocols like ftp, XRootD, DCap++, rfiio



# Storage Semantics

## 1. POSIX(-like) Storage



Based on **filesystems, RDBMS ...**

- GPFS, Lustre, AFS, pNFS, OrangeFS, GLUSTER **et.al.**
- CEPH, FUSE driver for `<xyz>` **et.al.**

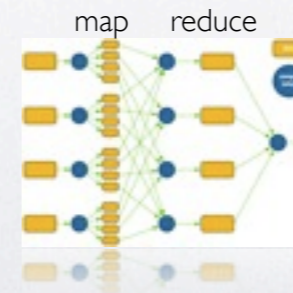
## 2. Cloud Storage



Based on **Object Stores, DHTs and key-value DB's**

- Amazon S3, Swift, Facebook Haystack **et al.**

## 3. Map-Reduce Storage



- GoogleFS, HDFS **et.al.**

*Some solutions mix semantics and technology*



# How many files can a Cluster filesystem have?

*“GPFS scans 10 billion files”*

Scalability is sufficient for designed use cases. ... still highlights also the problem of hierarchical namespaces.

Richard Freitas, Joseph Slember, Wayne Sawdon, and Lawrence Chiu. [GPFS Scans 10 Billion Files in 43 Minutes](#). 2011.

<http://www.almaden.ibm.com/storagesystems/resources/GPFS-Violin-white-paper.pdf>

- theoretical **exercise (policy scan)**
  - 0-size files and 6.5TB of meta data
- meta data on SSDs (violin memory)
  - > IMIOPS @ 4k
  - 4 GB/s
  - ‘would require hundreds of hard disks to reach SSD performance’



# From POSIX to Cloud API

Web Scale Problems ...

[access](#)  
[aio\\_cancel](#)  
[aio\\_error](#)  
[aio\\_read](#)  
  
[aio\\_return](#)  
[aio\\_suspend](#)  
[aio\\_write](#)  
[chdir](#)  
[chmod](#)  
[chown](#)  
[close](#)  
[closedir](#)  
[creat](#)  
[dup](#)  
[dup2](#)  
[fcntl](#)  
[fdatasync](#)  
[fdopen](#)  
[fstat](#)  
[fsync](#)  
[getcwd](#)

[link](#)  
[lio\\_listio](#)  
[lseek](#)  
[mkdir](#)  
[mkfifo](#)  
[msync](#)  
[open](#)  
[opendir](#)  
[read](#)  
[readdir](#)  
[rename](#)  
[rewinddir](#)  
[rmdir](#)  
[stat](#)  
[umask](#)  
[uname](#)  
[unlink](#)  
[utime](#)  
[write](#)

POSIX IO



**Service**  
[GET Service](#)  
**Bucket**  
[DELETE Bucket](#)  
[DELETE Bucket lifecycle](#)  
[DELETE Bucket policy](#)  
[DELETE Bucket website](#)  
[GET Bucket \(List Objects\)](#)  
[GET Bucket acl](#)  
[GET Bucket lifecycle](#)  
[GET Bucket policy](#)  
[GET Bucket location](#)  
[GET Bucket logging](#)  
[GET Bucket notification](#)  
[GET Bucket Object versions](#)  
[GET Bucket requestPayment](#)  
[GET Bucket versioning](#)  
[GET Bucket website](#)  
[HEAD Bucket](#)  
[List Multipart Uploads](#)  
[PUT Bucket](#)  
[PUT Bucket acl](#)  
[PUT Bucket lifecycle](#)  
[PUT Bucket policy](#)  
[PUT Bucket logging](#)  
[PUT Bucket notification](#)  
[PUT Bucket requestPayment](#)  
[PUT Bucket versioning](#)  
[PUT Bucket website](#)

**Object**  
[DELETE Object](#)  
[Delete Multiple Objects](#)  
[GET Object](#)  
[GET Object ACL](#)  
[GET Object torrent](#)  
[HEAD Object](#)  
[POST Object](#)  
[PUT Object](#)  
[PUT Object acl](#)  
[PUT Object - Copy](#)  
[Initiate Multipart Upload](#)  
[Upload Part](#)  
[Upload Part - Copy](#)  
[Complete Multipart Upload](#)  
[Abort Multipart Upload](#)  
[List Parts](#)

S3 API

- Simplifications mainly:
- Drop Namespace Hierarchy
  - WORM - write once read many
  - [GetObject \[Range\]](#) [PutObject](#), [Delete Object](#)
  - bucket handling

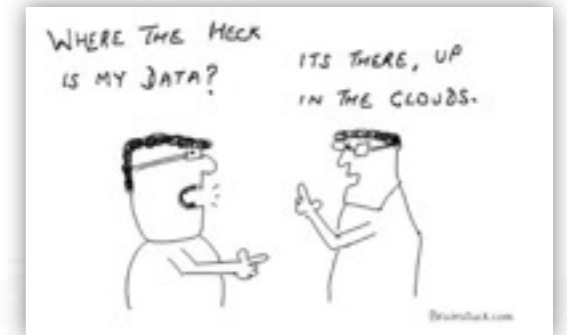
simplifications=> simple resilient **scale-out architecture**

interface changes => **compromises**



# The Cloud Compromise

- **increased** latency
- **eventual** consistency
- **reduced** but **simpler** storage interface
- goes along with **MapReduce** for efficient data access
  - move task where the data is, optimize for large IOs, **rewrite** your application “no POSIX”, WORM & append-only
- **proven** scalability **and** manageability e.g.
  - 900 Billion Objects in Amazon S3
  - 100 Billion Photos in Facebook 2011
- **can run in** no maintenance **mode**
  - no repair approach - a failed disk or node needs no intervention
  - exchange/migrate all after natural lifecycle of the whole system



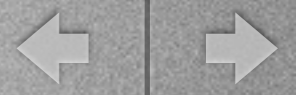
... solved the web problem

... is not optimal for non-sequential small reads

... allows large installations to be easy operable and cheap

**Cloud:**  
Overall performance and scalability

**HPC:**  
Single client performance



# Storage Tiering

Tiered storage is the assignment of different categories of data to different types of storage media in order to **reduce total storage cost**.

- **Examples**

- **HSM Systems**

- Lustre<sup>HSM</sup>, GPFS<sup>HSM</sup>, dCache, CASTOR

- **HHD/SSH = HD + SSD Cache**

- **VTL - Virtual Tape Libraries**

- **CAS - Content-addressed Storage**

LHC storage demonstrated that HSM is not a perfect solution  
**random user file recalls**

Large impact for the **OS, DB's** and **meta data** stores - not for volume

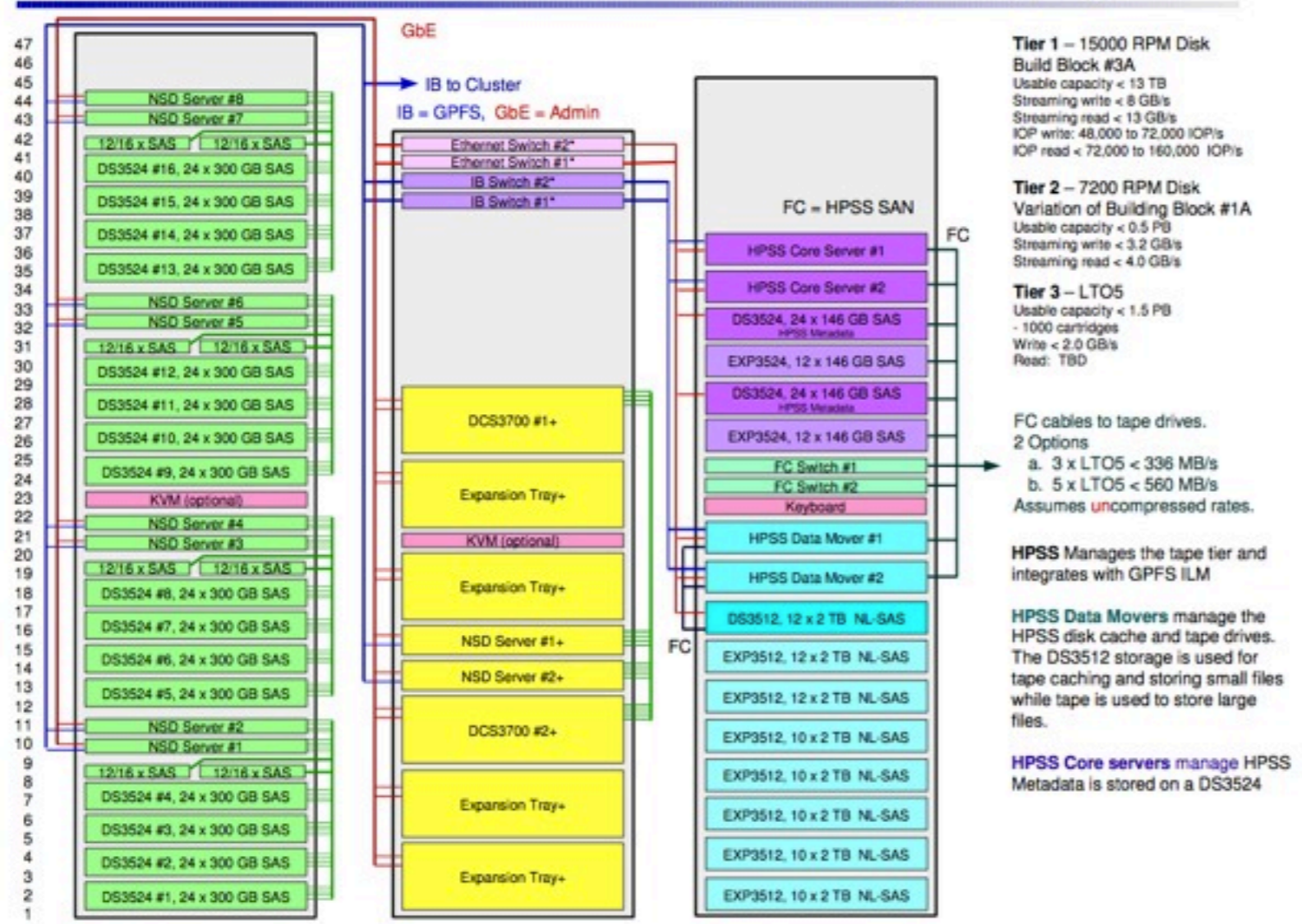
**hides** the **tape IO** characteristics.

**Cloud** storage with **POSIX** API.



# Multi Tier Storage Solutions

## Three-Tier Solution: Fast Disk, Capacity Disk, Tape

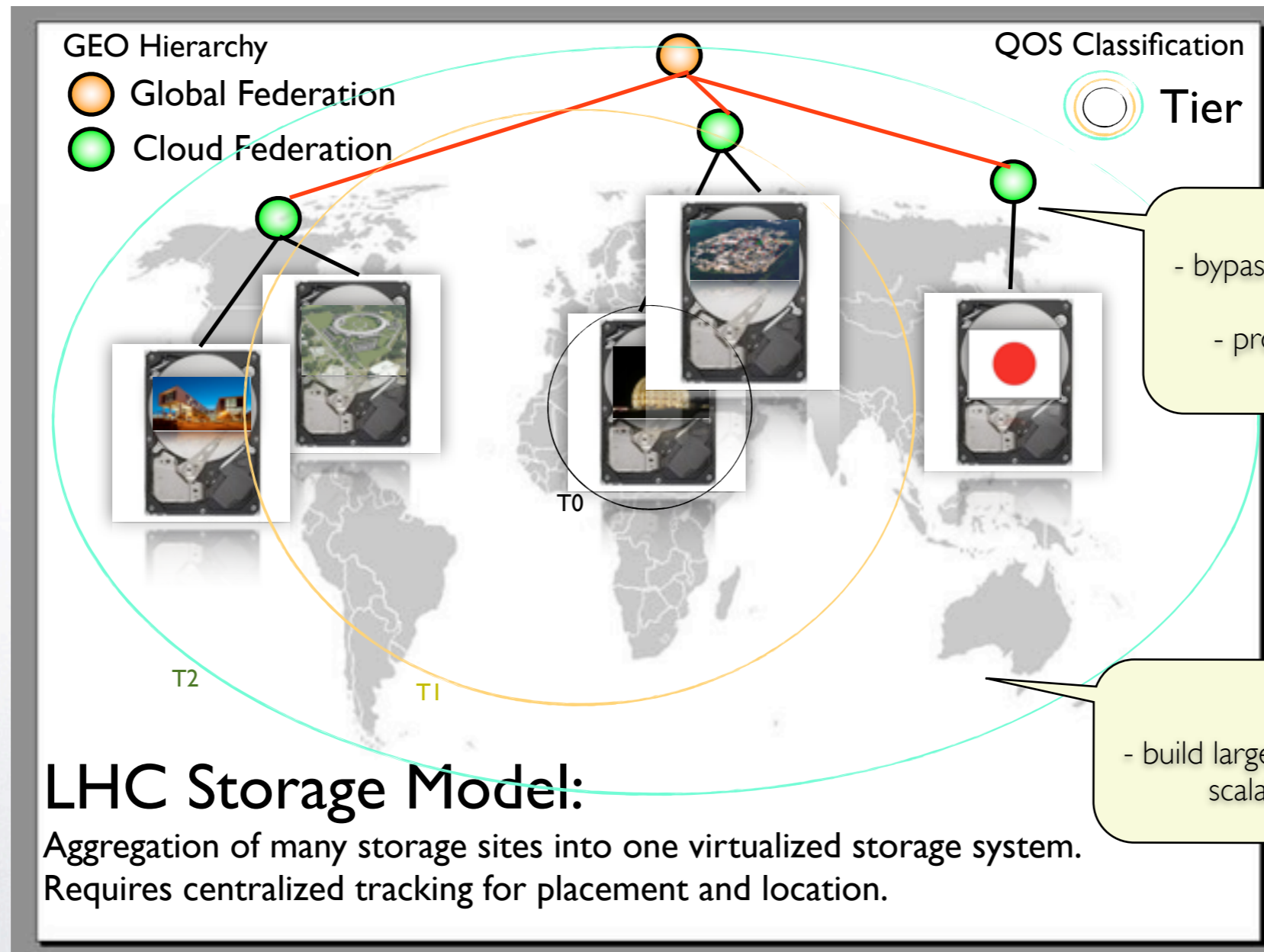


- allow to **shape/optimize performance**
  - e.g. convert stream perf. to IOPS
  - save money having the largest capacity on the cheapest media
- **does not work** for all work loads
  - the dimension and performance parameters of the tiers must meet usage pattern.
  - otherwise:
    - no guaranteed gain
    - no savings

Promising approach:  
**combination of cloud storage** as a capacity store **+**  
**front-end** with fine-grained and performant IO interface e.g. HSM enabled filesystems, dCache, XRootD FRM



# Scalability by Aggregation & Federation

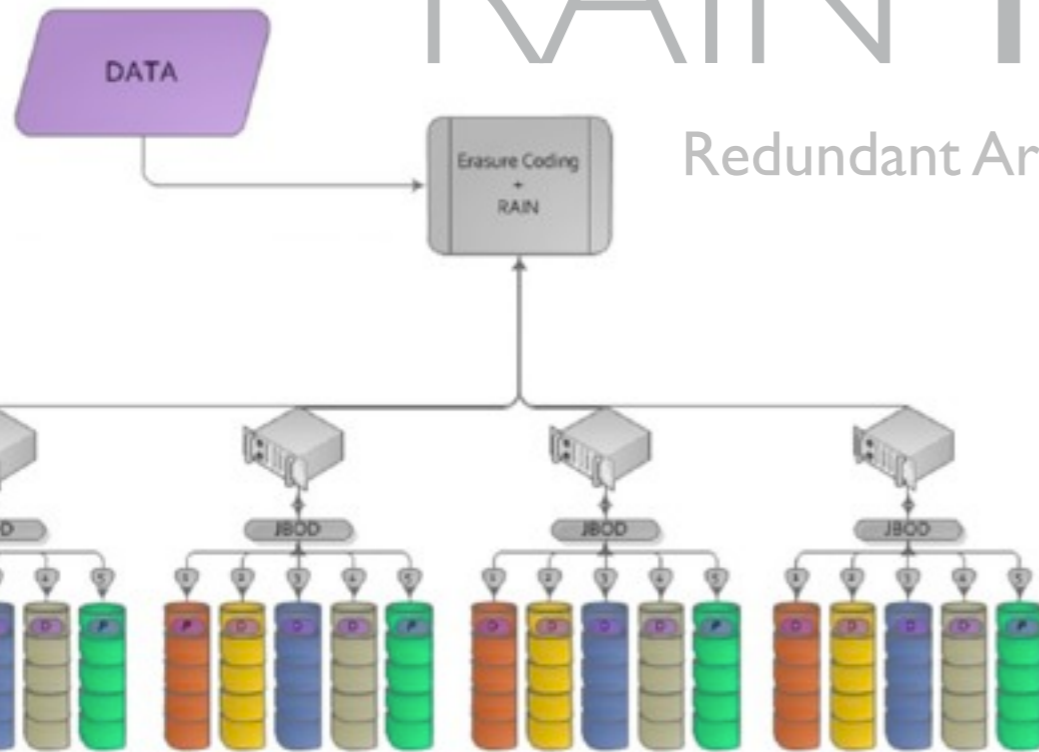






# RAIN Technology

Redundant Array of Inexpensive Nodes



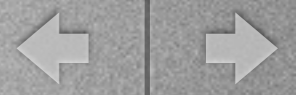
RAID: uses redundancy algorithm on device level

RAIN: uses redundancy algorithm on file/object level distributed over nodes



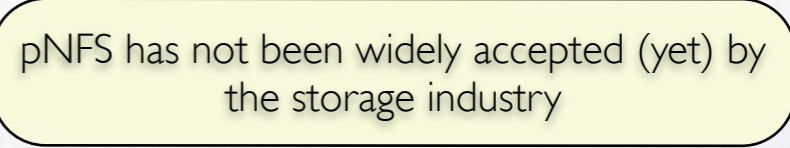
RAIN combined with block checksumming is a perfect match for very large storage systems

- **space** overhead **configurable** via algorithm e.g. Reed Solomon  $10+3 = 30\%$  overhead
- allows to **scale rebuild performance** independently of stripe widths
- allows to **repair** device failures and **silent corruptions**
- allows to **scale IO performance** per object
- **reduces** the **data loss** scenario compared to RAID
- **no** need for hardware RAID **controller** or **multi-path** storage

Available e.g. by NetApp, PanFS, Cleversafe ... GPFS implemented as native RAID (~RAIN) on AIX.



# Trends & Standards

- GLUSTER & CEPH are providing an **S3** and **object store** interface 
- LUSTRE & HDFS work on **distributed/ federated namespaces** 
- **pNFS** - client in RedHat 6.2
- **new NFS v4.2 standard** coming 
  - defines SSC (server side copy), application data blocks, space reservation, sparse files and IO advise, targeting virtualized data centers

Will pNFS ever displace native clients in Lustre & GPFS?



# Storage Virtualization

Storage as a Service





# Virtualization

## ● Virtualization of Storage

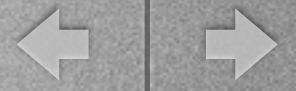
- **not new** concept
  - Cluster Filesystems use storage virtualization on block, disk, file and filesystem level and virtualization of tape via HSM etc.
  - important for fully virtualized data centers hosting VM and database images

## ● Storage in virtualized Environments

- run storage system as a **virtual appliance** in a virtual machine e.g. the **GLUSTER** Storage Appliance
- simplifies deployment and configuration '**Storage as a Service**'
- allows on-the-fly **deployment** of a distributed storage system in **cloud environments**
- allows performance **confinement** within a virtual machine
- **GLUSTER** reports only a 5% degradation in performance

Limitation: Storage is stateful and can not quickly be moved or exchanged once it is filled.

1 PB @ 1GB/s = 11,5 days



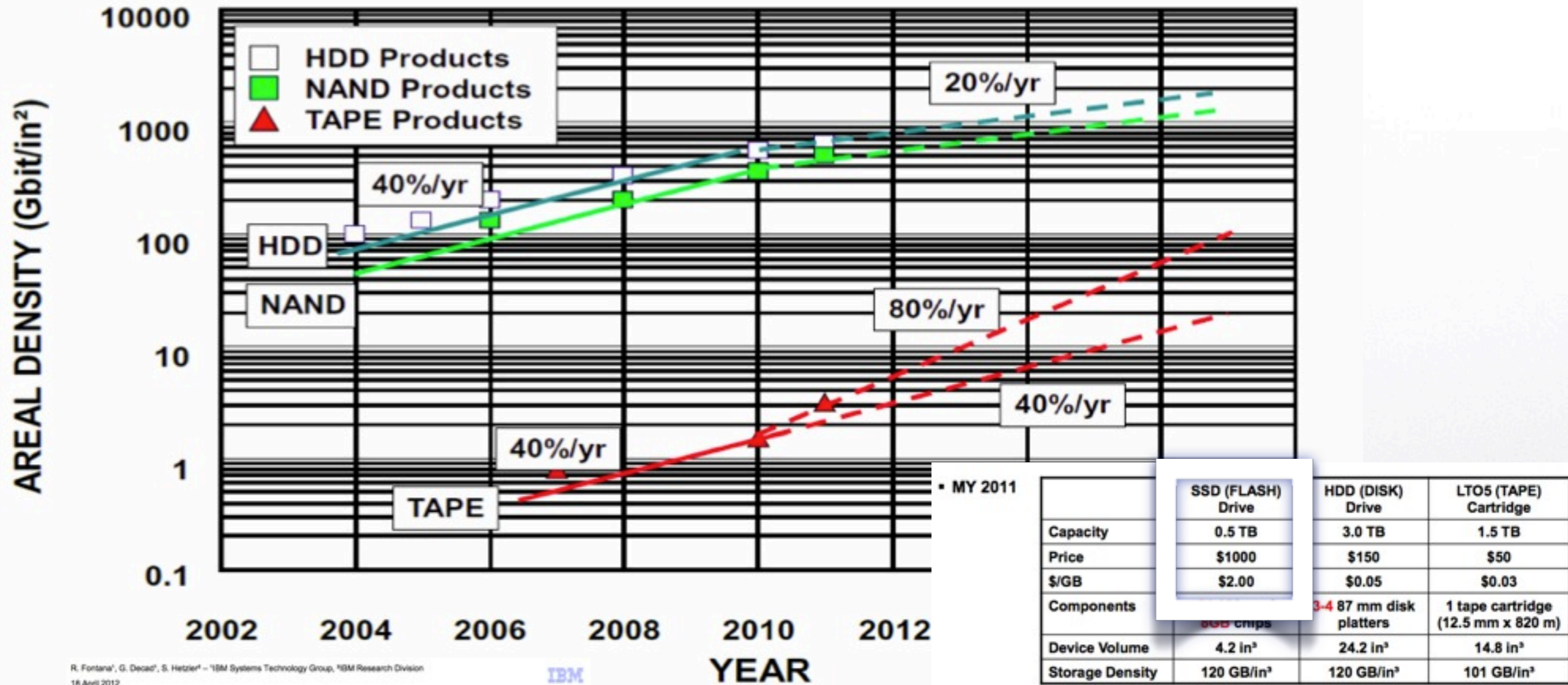
# Market Trends





# Storage Density & Volume Metrics

### Growth Rate Predictions



R. Fontana<sup>1</sup>, G. Decad<sup>2</sup>, S. Hetzler<sup>3</sup> - <sup>1</sup>IBM Systems Technology Group, <sup>2</sup>IBM Research Division  
18 April 2012





# 2014 Forecast

R. Fontana<sup>1</sup>, G. Decad<sup>2</sup>, S. Hetzler<sup>3</sup> – IBM Systems Technology Group, IBM Research Division  
18 April 2012



## Scenarios for 2014

	Historical	Conservative	Tape Aggressive
<b>Areal Density Growth (Specifics)</b>	<b>40%/yr--TAPE 40%/yr--HDD 40%/yr--NAND</b>	<b>40%/yr--TAPE 20%/yr--HDD 20%/yr--NAND</b>	<b>80%/yr--TAPE 20%/yr--HDD 20%/yr--NAND</b>
<b>TAPE</b>			
-- Areal Density	4.8 Gbit/in <sup>2</sup>	4.8 Gbit/in <sup>2</sup>	12.0 Gbit/in <sup>2</sup>
-- Minimum Feature			
-- Cartridge Capacity	6.0 TB	6.0 TB	15.0 TB
-- Volumetric Density			
<b>HDD</b>			
-- Areal Density	2500 Gbit/in <sup>2</sup>	1300 Gbit/in <sup>2</sup>	1300 Gbit/in <sup>2</sup>
-- Minimum Feature	0.010 um	0.018 um	0.018 um
-- HDD Capacity <sup>1</sup>	12.0 TB	6.0 TB	6.0 TB
-- Volumetric Density	480 GB/in <sup>3</sup>	240 GB/in <sup>3</sup> ✱	240 GB/in <sup>3</sup> ✱
<b>NAND Flash</b>			
-- Areal Density	1300 Gbit/in <sup>2</sup>	700 Gbit/in <sup>2</sup>	700 Gbit/in <sup>2</sup>
-- Minimum Feature	0.012 um	0.016 um	0.016 um
-- Chip Capacity	32 GB	24 GB	24 GB
-- SSD Capacity <sup>2</sup>	2 TB	1.2 TB	1.2 TB
-- Volumetric Density	480 GB/in <sup>3</sup>	300 GB/in <sup>3</sup> ✱	300 GB/in <sup>3</sup> ✱

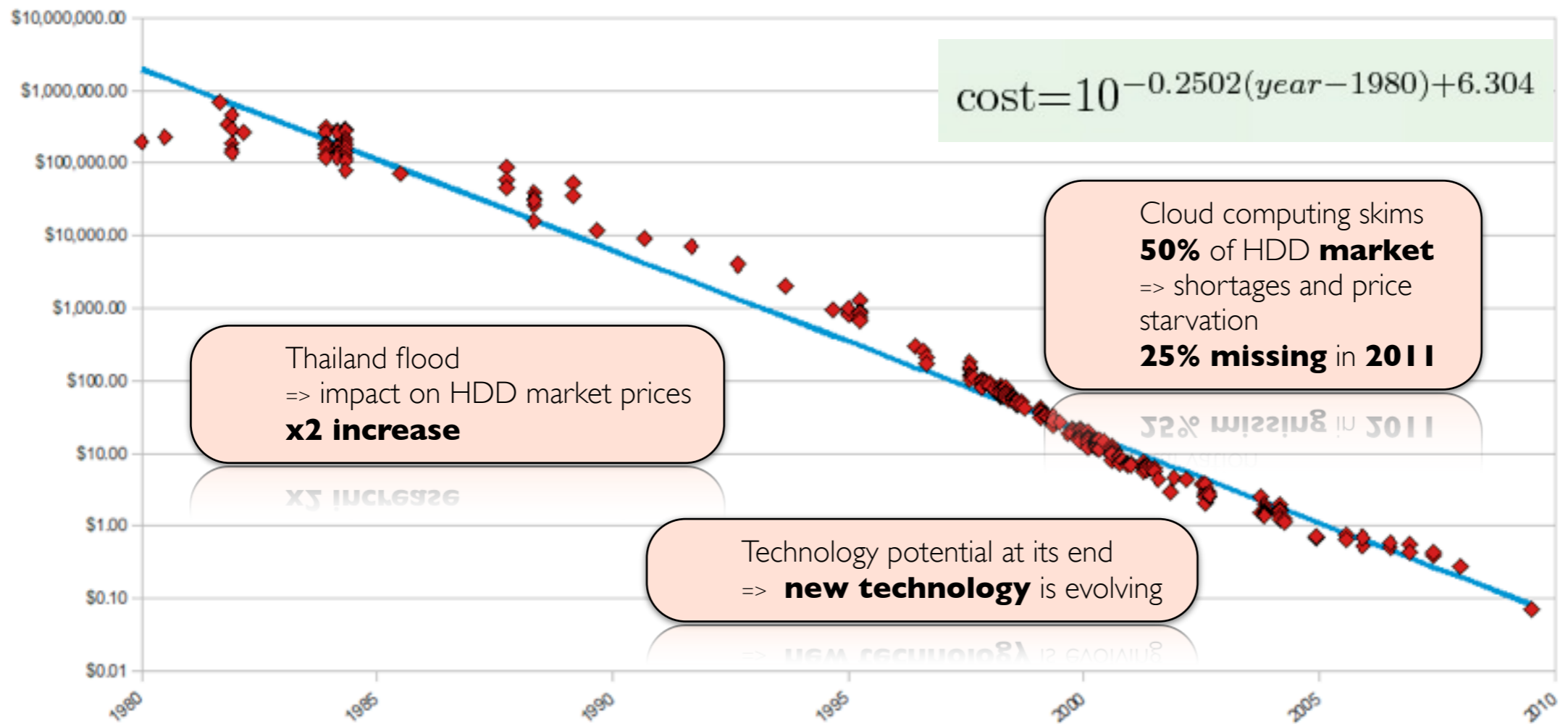
6.0 TB      6.0 TB      15.0 TB

**Tape** growth rate  
**2-4x** HDD & NAND



# Consumer HDD Prices

Hard Drive Cost per Gigabyte  
1980 - 2009



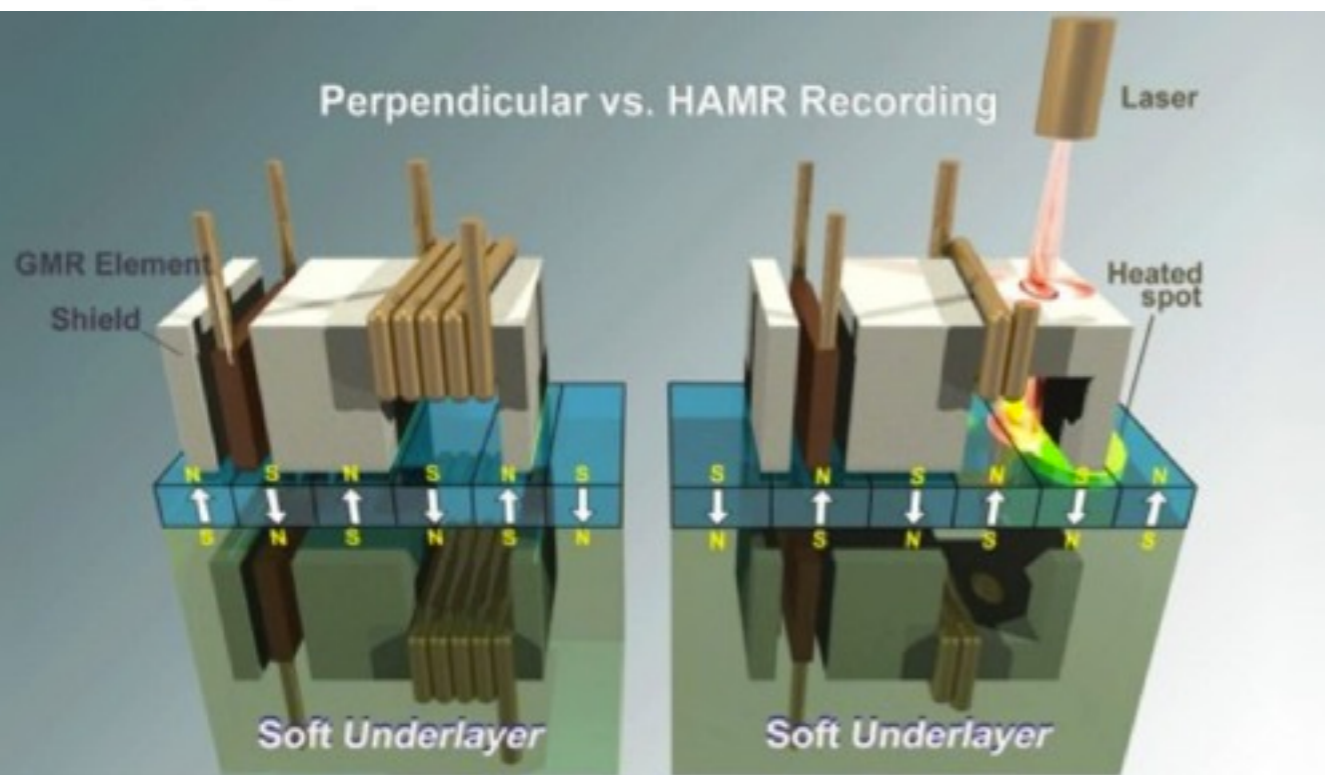
[mkomo.com]





# New HDD Technology

## *Heat Assisted Magnetic Recording*

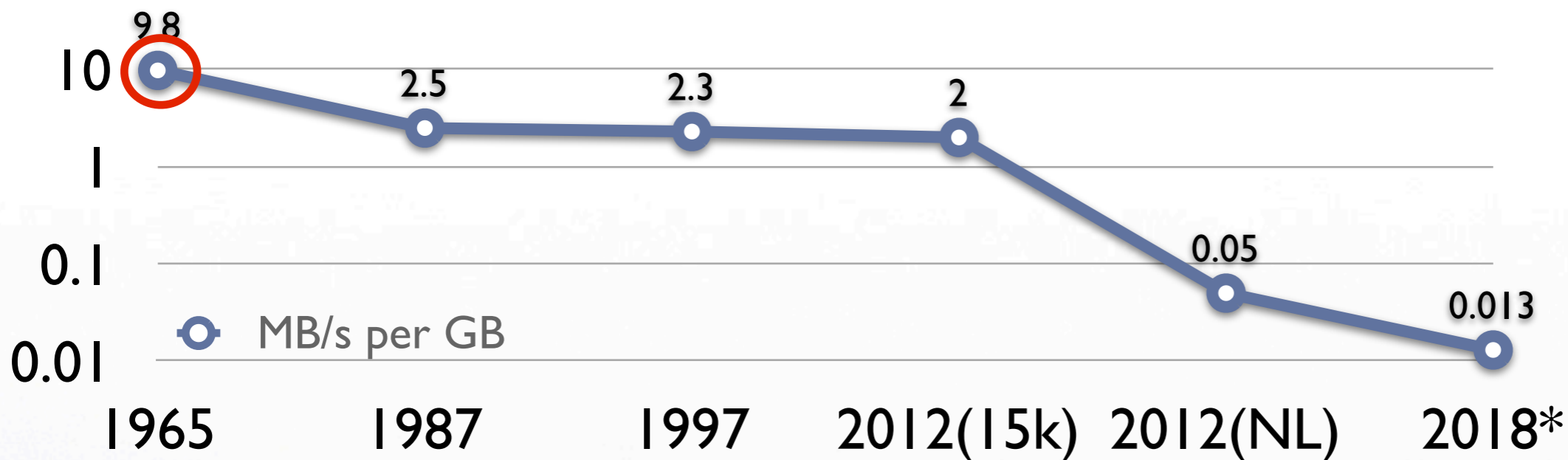


- HAMR limit is 5-10 Tbit/in<sup>2</sup>  
~**30-60 TB** 3.5" HDD
- large production level **2016-2017**
- expensive technology
- areal growth rate lower ~**20-25%**

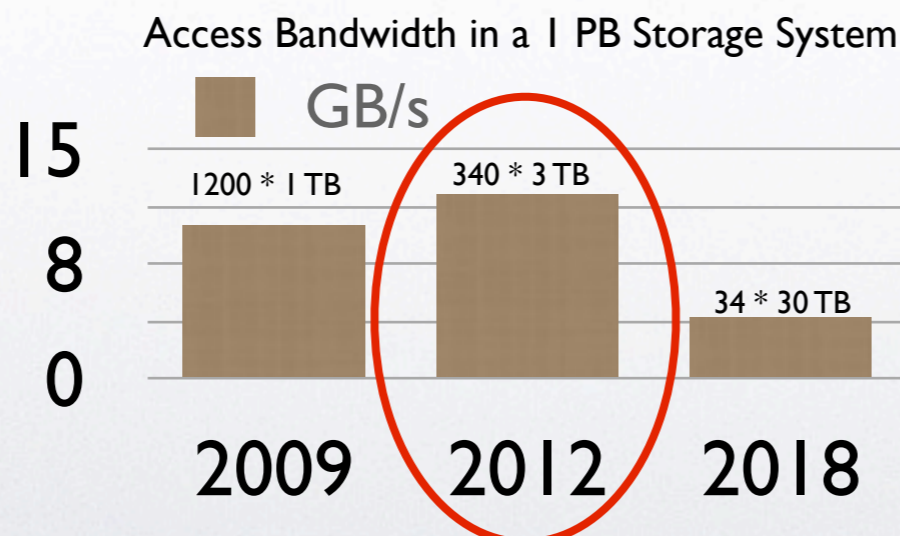
Seagate has achieved a milestone 1 terabit per square inch storage density using heat-assisted magnetic recording (HAMR) technology



# HDD Performance Capacity Ratio



Assumption:  
2018 30 TB HDD @ 380 MB/s



\* Gary Grider, Exa-Scale FSIO, 07/2010, LANL.

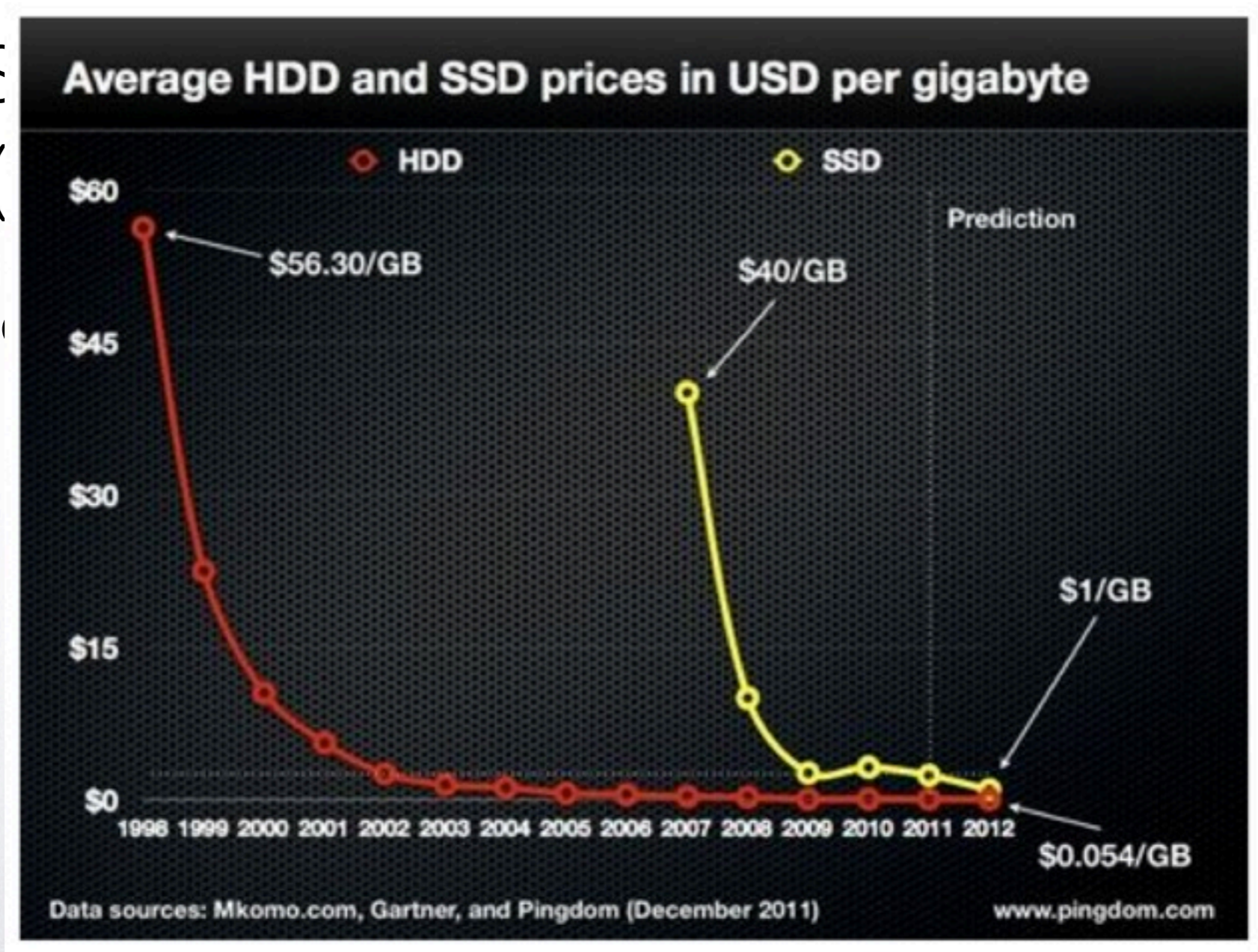
Source: Raymond L. Paden, Ph.D. - IBM Deep Computing - 28th IEEE Conference on Massive Data Storage

Increasing **capacity**  
**reduces bandwidth** to data!  
Critical crossover point?



# SSD

- S
- (
- I
- 
- 





# SSD

## Examples

Type	Capacity	Streaming Rate r/w	IOPS r/w
<b>Controller</b> (Rack of RamSan-820)	1 PB	168 GB/s	18.9M
<b>Block Device</b>	1 TB	500/380 MB/s	15-35k
<b>PCIe</b>	1.2 TB (ioDrive2) 12/16 TB (Z-Drive)	1.5/1.3 GB/s 7.2 GB/s	500k/140k 2.5M

## 5.1.2012 “Fusion-io Breaks One Billion IOPS Barrier”

**SSDs can deliver IOPS & bandwidth en masse**

Interesting for **LSS** to handle meta data work loads

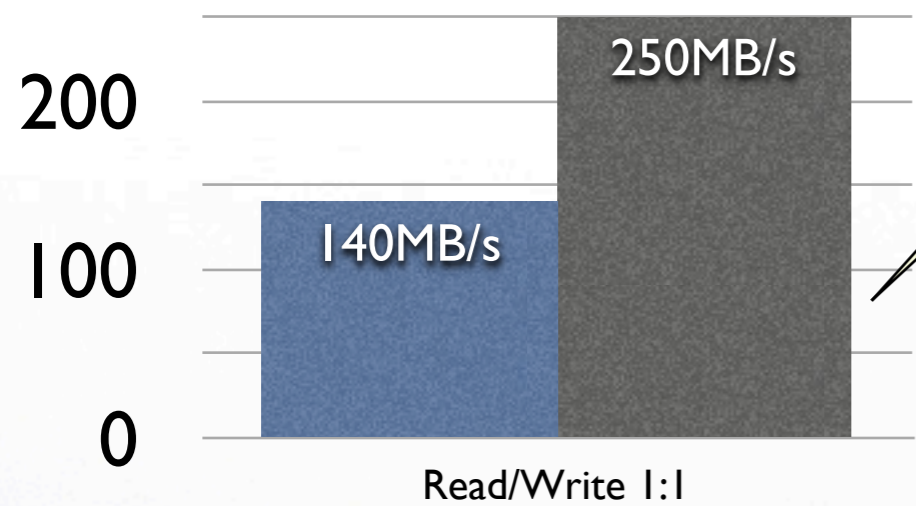


... **alive** ... market is growing!

# Tape

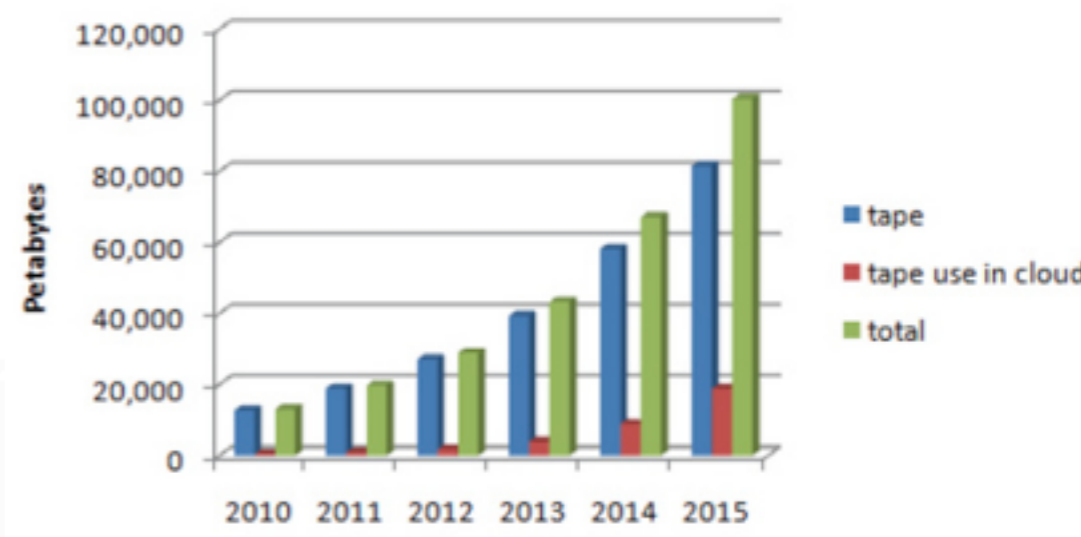
## Streaming Performance

- LTO5 1.5 TB
- StorageTek 5 TB



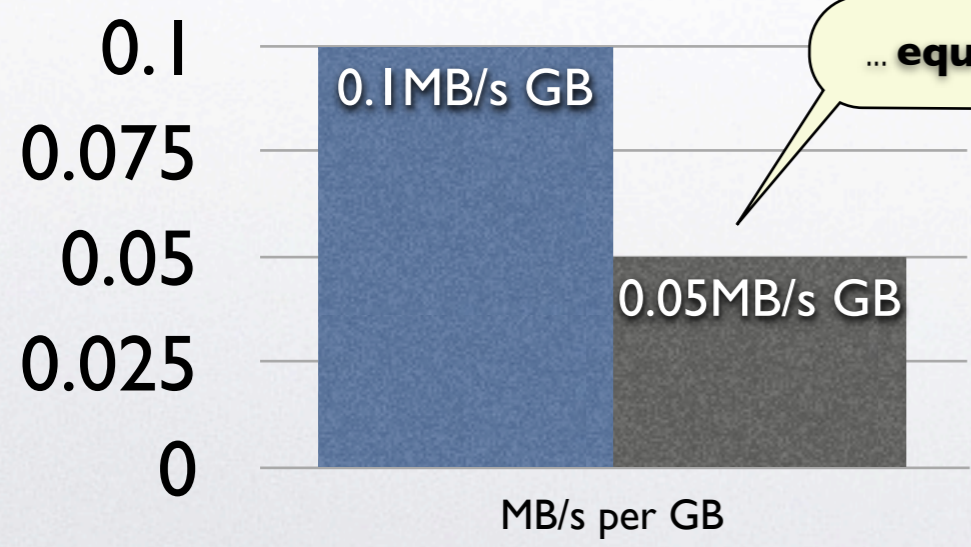
... **faster** than disks ...

## Archive Data in PB



Source: Enterprise Strategy Group

## Performance:Capacity Ratio



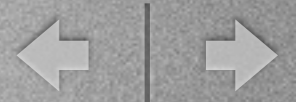
... **equal/better** than NL disks ...

## ULTRIUM LTO Eight-Generation Roadmap

	Generation 1	Generation 2	Generation 3	Generation 4	Generation 5	Generation 6	Generation 7	Generation 8
Compressed Capacity	200 GB	400 GB	800 GB	1.6 TB	3 TB	6 TB	16 TB	32 TB
Native Capacity	100 GB	200 GB	400 GB	800 GB	1.5 TB	3 TB	6.4 TB	12.8 TB
Compressed Data Rate	up to 40 MB/s	up to 80 MB/s	up to 160 MB/s	up to 240 MB/s	up to 280 MB/s	up to 525 MB/s	up to 788 MB/s	up to 1180 MB/s
Native Data Rate	up to 20 MB/s	up to 40 MB/s	up to 80 MB/s	up to 120 MB/s	up to 140 MB/s	up to 210 MB/s	up to 315 MB/s	up to 472 MB/s

Note: Compressed capacities for generations 1-5 assume 2:1 compression. Compressed capacities for generations 6-8 assume 3:1 compression (achieved with larger compression history buffer).  
Source: The LTO Program. The LTO Ultrium roadmap is subject to change without notice and represents goals and objectives only.  
Linear Tape-Open, LTO, the LTO logo, Ultrium, and the Ultrium logo are registered trademarks of HP, IBM and Quantum in the US and other countries.

**Tape most reliable and cheapest** volume storage for archiving and backup **in future projections**



# Storage Connectivity

## Infiniband

	SDR	DDR	QDR	FDR-10	FDR	EDR
Year	1999	2004	2008		2011	
<b>4X</b>	<b>8 Gbit/s</b>	<b>16 Gbit/s</b>	<b>32 Gbit/s</b>	<b>41.2 Gbit/s</b>	<b>54.54 Gbit/s</b>	<b>100 Gbit/s</b>
<b>12X</b>	24 Gbit/s	48 Gbit/s	96 Gbit/s	123.6 Gbit/s	163.64 Gbit/s	200 Gbit/s
<b>1X</b>	2 Gbit/s	4 Gbit/s	8 Gbit/s	10.3 Gbit/s	13.64 Gbit/s	25 Gbit/s

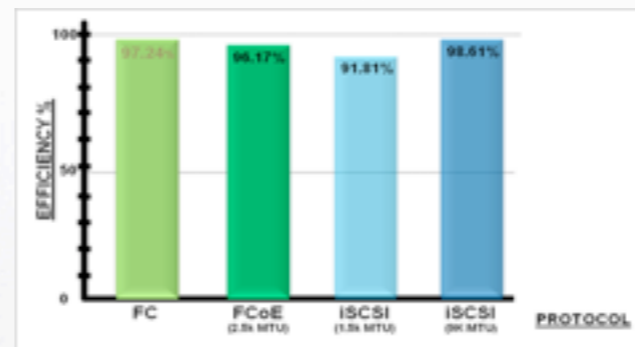
## SAS

Year	Gbit/s
2004	3
2009	6
2013 (?)	<b>12</b>

## Fibre Channel

ME	Gbit/s full duplex	Availability
1GFC	1.6	1997
2GFC	3.2	2001
4GFC	6.4	2004
8GFC	12.8	2005
10GFC Serial	20.4	2008
16GFC	25.6	2011
20GFC	<b>40.8</b>	20??

## iSCSI/iSCSIoE



## Ethernet

1 Gbit	1999
10 Gbit	2002
40/100 Gbit	2009/?

## PCIe

2003 1.0a	2 Gbit/s
2005 1.1	2 Gbit/s
2007 2.0	4 Gbit/s
2010 3.0	8 Gbit/s
? 4.0	16 Gbit/s

Standard: x1, x4, x8, 16, x32 lanes  
Common: x1, x4, x8 cards



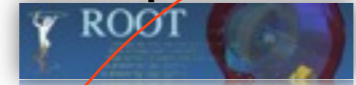
# “Storage Biosphere”

What can **we** learn from Internet Storage?

Environment



Computation



Storage



- look at large storage in its context
- storage is serving applications, not driving them
- LHC storage is used in a distributed environment implementing a lot of functionality on a higher level

Web Revolution:  
**Redesign** storage **around** their **applications**

**Don't stick** to standards and/or products!

**Set** or choose your **standard** according to **requirements** and use it consistently everywhere!

- What we should think about:
- 1 Unification of Storage, Computing and Network
  - 2 LHC appliances for storage and computing
  - 3 Rethink tier/replication model
  - 4 Can we move redundancy into the application layer instead of storage systems?
  - 5 Can we base LHC storage on object stores?

What do we learn?



# Future

Importance of storage systems will increase in the future

- internet + emersion of mobile devices drives **unseen growth of storage** needs  
=> huge market implications + technology push
- chance to **profit** from and **contribute** to large community projects
- commercial solutions follow market demands => options not only for HPC
- extreme large scale systems based on **elastic object store** in combination with **elastic databases** providing meta data views
  - LHC storage approach is compatible - useful to adopt big data technology
  - over time LHC storage might leave 'comfort zone' where things still scale easily with used technology
- **Exabyte** storage for big data mining will become a **new norm** within few years

Thank you for your attention!