# Tutorial

## Storage Lessons from HPC:
## Extreme Scale Computing Driving
## High Performance Data Storage

**Gary Grider**
HPC Division Leader
LANL/US DOE

# About the Instructor

- Gary Grider is the Leader of the High Performance Computing (HPC) Division at Los Alamos National Laboratory. As Division Leader, Gary is responsible for all aspects of High Performance Computing technologies and deployment at Los Alamos. Additionally, Gary is responsible for managing the R&D portfolio for keeping the new technology pipeline full to provide solutions to problems in the Lab's HPC environment, through funding of university and industry partners.

- Gary is also the US Department of Energy Exascale Storage, IO, and Data Management National Co-Coordinator. In this role, Gary helps managed the US government investments in Data Management, Mass Storage, and IO. Gary has 30 active patents/applications in the data storage area and has been working in HPC and HPC related storage since 1984.

# SNIA Legal Notice

# Eight Decades of Production Weapons Computing to Keep the Nation Safe

Maniac     IBM Stretch     CDC     Cray 1     Cray X/Y     CM-2
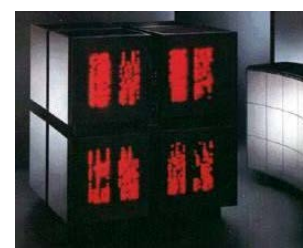
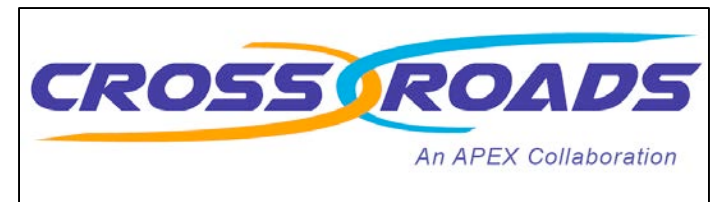CM-5     SGI Blue Mountain     DEC/HP Q     IBM Cell Roadrunner     Cray XE Cielo
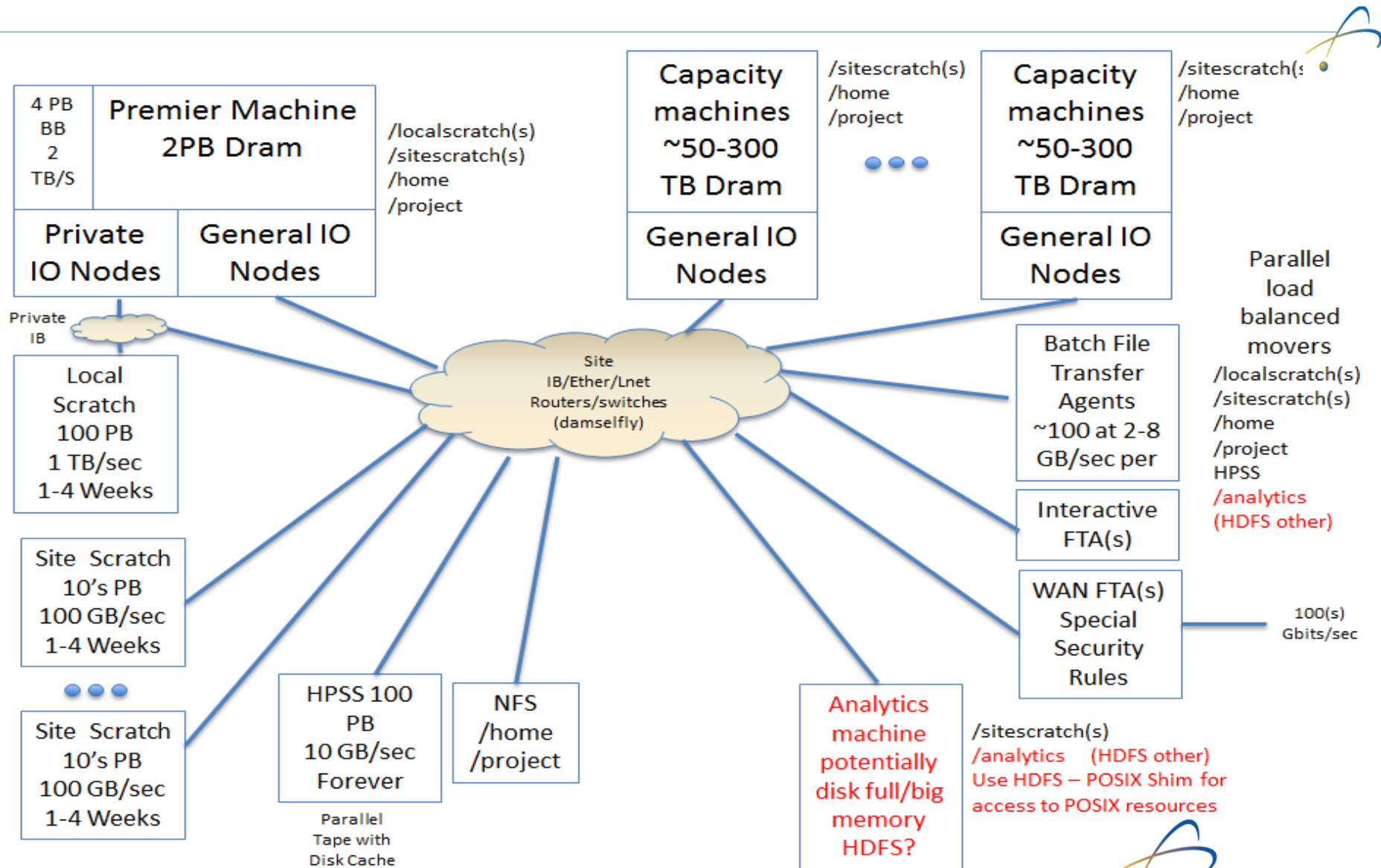
Cray Intel KNL Trinity     Ising DWave     Cross Roads

CROSS ROADS
An APEX Collaboration

Los Alamos
NATIONAL LABORATORY
EST. 1943

# Simple View of our Computing Environment



4 PB BB 2 TB/S

**Premier Machine 2PB Dram**

/localscratch(s)
/sitescratch(s)
/home
/project

**Private IO Nodes**

**General IO Nodes**

Private IB

Local Scratch
100 PB
1 TB/sec
1-4 Weeks

Site Scratch
10's PB
100 GB/sec
1-4 Weeks

Site Scratch
10's PB
100 GB/sec
1-4 Weeks

HPSS 100 PB
10 GB/sec
Forever

Parallel Tape with Disk Cache

NFS
/home
/project

**Capacity machines ~50-300 TB Dram**

/sitescratch(s)
/home
/project

**General IO Nodes**

**Capacity machines ~50-300 TB Dram**

/sitescratch(s
/home
/project

**General IO Nodes**

Site IB/Ether/Lnet Routers/switches (damselfly)

**Batch File Transfer Agents ~100 at 2-8 GB/sec per**

**Interactive FTA(s)**

**WAN FTA(s) Special Security Rules**

Parallel load balanced movers
/localscratch(s)
/sitescratch(s)
/home
/project
HPSS
/analytics
(HDFS other)

100(s) Gbits/sec

**Analytics machine potentially disk full/big memory HDFS?**

/sitescratch(s)
/analytics    (HDFS other)
Use HDFS – POSIX Shim for access to POSIX resources

**SNIA** Tutorial

**Los Alamos** NATIONAL LABORATORY
EST.1943

# Large Machines and Infrastructure

## Trinity

- Haswell and KNL
- 20,000 Nodes
- Few Million Cores
- 2 PByte DRAM
- 4 PByte NAND Burst Buffer ~ 4 Tbyte/sec
- 100 Pbyte Scratch PMR Disk File system ~1.2 Tbyte/sec
- 30PByte/year Sitewide SMR Disk Campaign Store ~ 1 Gbyte/sec/Pbyte (30 Gbyte/sec currently)
- 60 PByte Sitewide Parallel Tape Archive ~ 3 Gbyte/sec



Pipes for Trinity Cooling

- 30-60MW
- Single machines in the 10k nodes and > 18 MW
- Single jobs that run across 1M cores for months
- Soccer fields of gear in 3 buildings
- 20 Semi's of gear this summer alone

# HPC Driving Industry

### Flash Burst Buffers

**DDN STORAGE**

**INFINITE MEMORY ENGINE™**

**IBM GPFS**

**EMC²**

**HPSS** High Performance Storage System

**THE SUPERCOMPUTER COMPANY**
*Data Warp*

### Parallel File Systems

**lustre**

**panasas**

**pNFS**

**IBM GPFS**

**ceph**

### Parallel Archives

**HPSS** High Performance Storage System

**UniTree software**

**DataTree CFS**

### Other

**HDF**

**Tokutek**

**SNIA** Tutorial

# Ceph begins

UNIVERSITY OF CALIFORNIA, SANTA CRUZ

BERKELEY • DAVIS • IRVINE • LOS ANGELES • RIVERSIDE • SAN DIEGO • SAN FRANCISCO     SANTA BARBARA • SANTA

DEPARTMENT OF COMPUTER SCIENCE                                    SANTA CRUZ, CALIFORNIA 95064

April 8, 2001

Lawrence Livermore National Laboratory
Attention: Barbara Larson, L-550
P.O. Box 808
Livermore, CA 94551

Dear Ms. Larson,

We are pleased to submit this white paper to the ASCI ASAP Request for Expressions of Interest in Level 2 Strategic Investigations. We are responding to the *Scalable and Parallel I/O and File Systems* area of interest.

We propose to improve both file system performance and functionality by building a storage system from object-based storage devices (OBSDs) connected by high-speed networks. The key advantage of OBSDs in a high-performance environment is the ability to delegate low-level block allocation and synchronization for a given segment of data to the device on which it is stored, leaving the file system to decide only on which OBSD a given segment should be placed. Since this decision is quite simple and allows massive parallelism, each OBSD need only manage concurrency locally, allowing a file system built from thousands of OBSDs to achieve massively parallel data transfers. Additionally, OBSDs can each manage their own storage consistency, removing the need to run a system-wide consistency check that could take days on a petabyte-scale traditional file system.
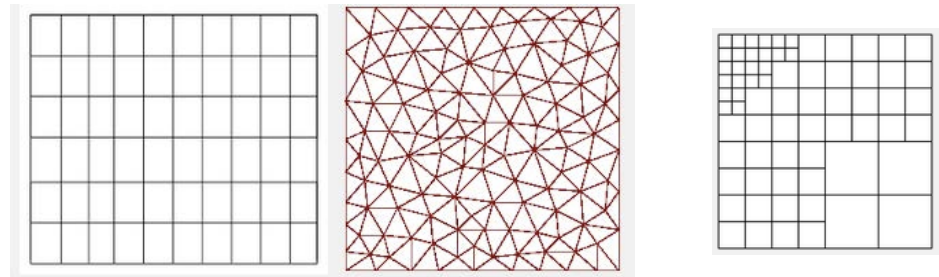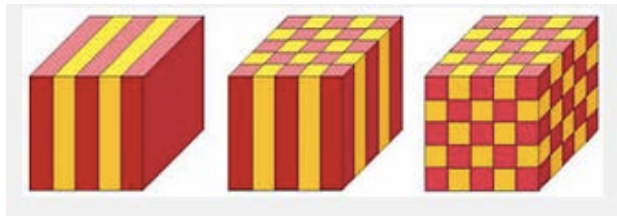
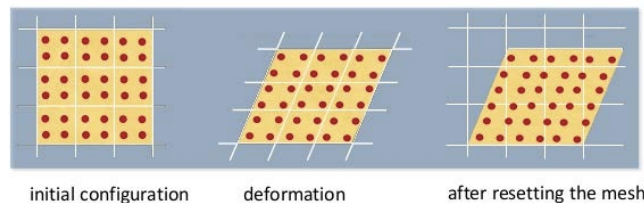# HPC Simulation Background

## Link scales

## Meshes

1D 2D 3D Structured Unstructured Resolutions



## Methods Lagrangian Eulerian ALE AMR



in each calculation step :

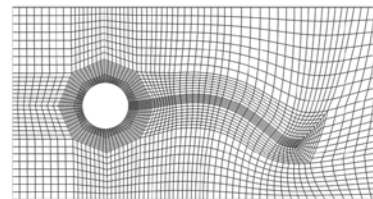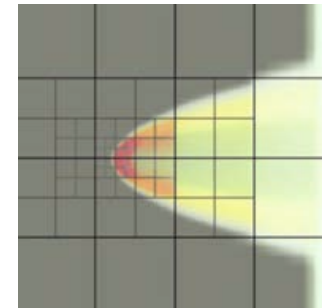initial configuration    deformation    after resetting the mesh

Lagrangian    Eulerian

UNIVERSITY OF CAMBRIDGE    PIAP-GA-2012-324522

http://media.archnumsoft.org/10305/

ALE

Eulerian AMR

SNIA. Tutorial

# Scaling and Programming Models

- The first is *strong scaling*, which is defined as how the solution time varies with the number of processors for a fixed *total* problem size.
- The second is *weak scaling*, which is defined as how the solution time varies with the number of processors for a fixed problem size *per processor*.
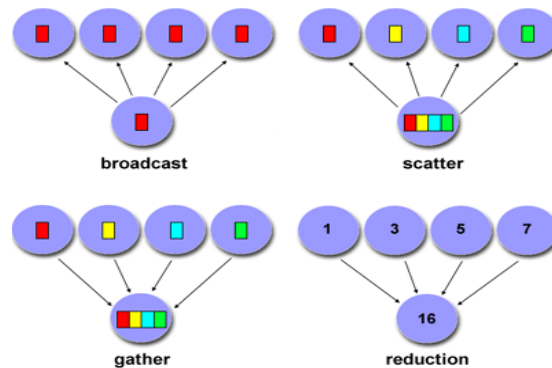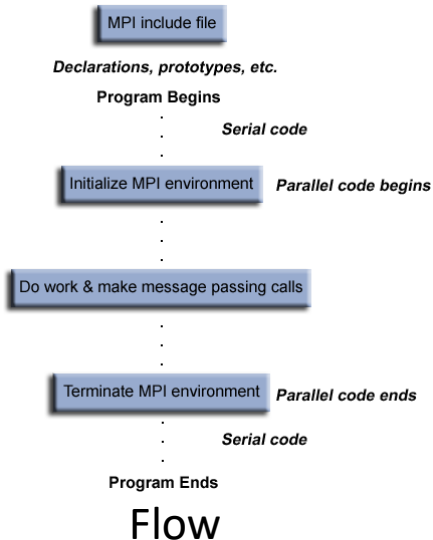
http://media.archnum
soft.org/10305/



Amdahl's law:
Parallel speedup vs. Sequential fraction

## Scaling

## Process Parallel Dominates the Past



Flow

Collectives
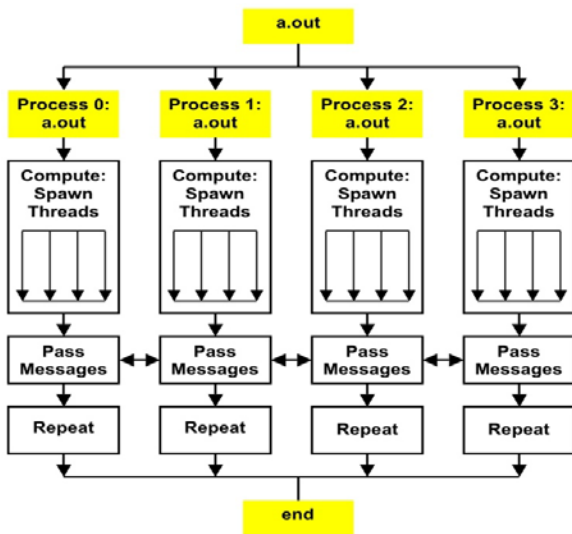
Path of a message buffered at the receiving process
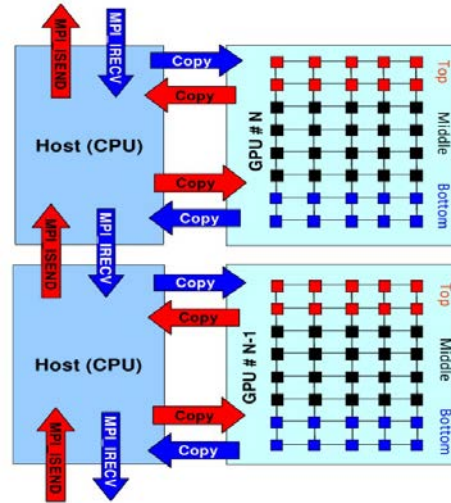
Point to Point

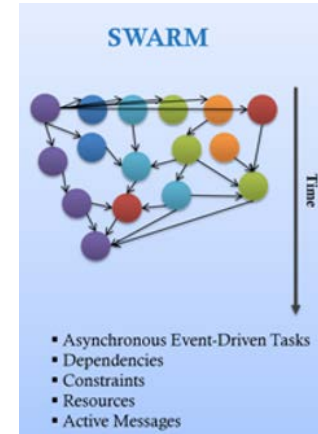# Current Programming Model Directions Combining Process Parallel, Data Parallel, Async Tasking, and Threading

## MPI+ Threads



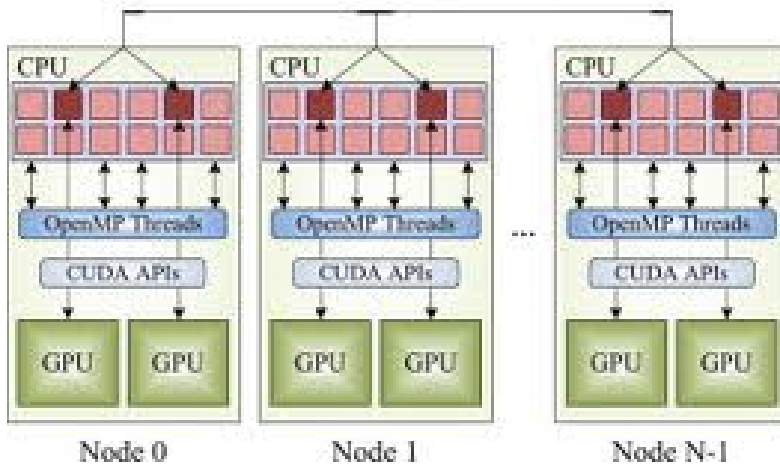## MPI+Data Parallel



## Async Tasking



SWARM

- Asynchronous Event-Driven Tasks
- Dependencies
- Constraints
- Resources
- Active Messages



MPI+ Threads+ Data Parallel

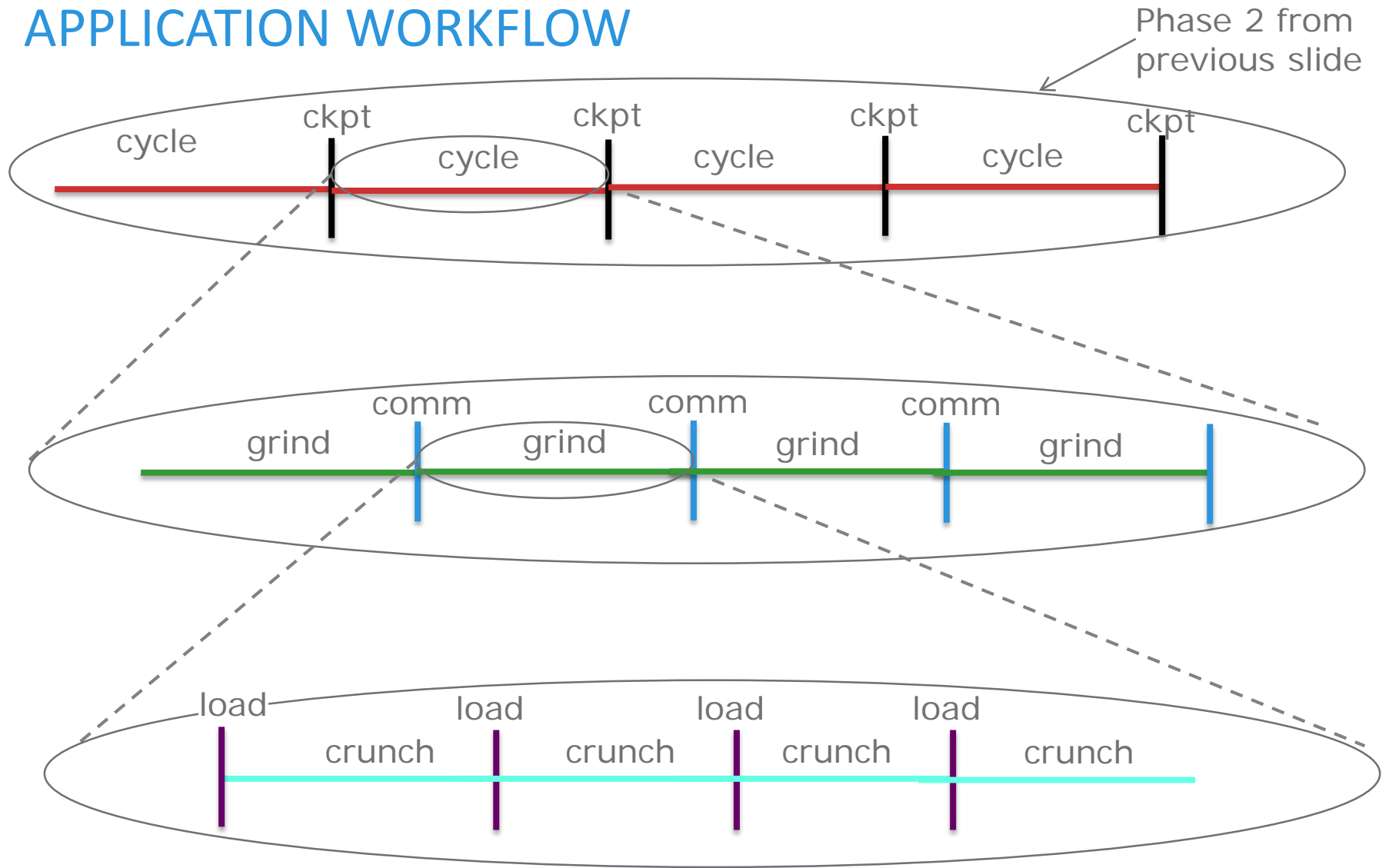## Domain Specific Languages



Engine

Model **DSL** Syntax

API

Others: OpenShmem, CAF, Chapel, X10 are possible as well

# Workflow Taxonomy from APEX Procurement
# A Simulation Pipeline



Figure 1: An example of an APEX simulation science workflow.

# APPLICATION WORKFLOW

Key observation: a grind (for strong-scaling apps) traverses all of memory.
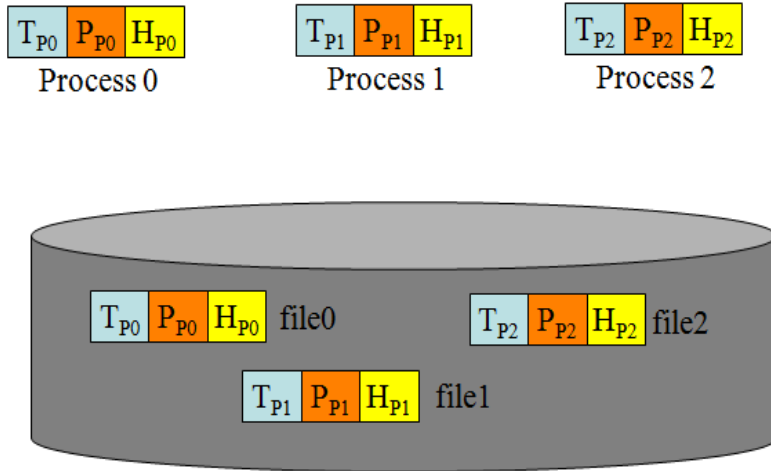
# Not Just Computation

failure
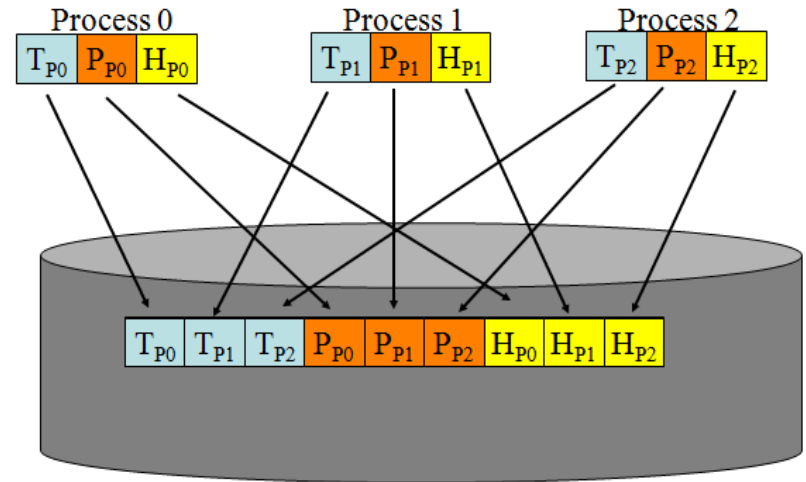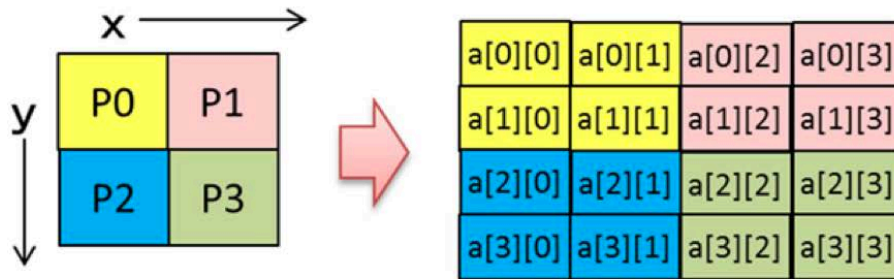
# HPC IO Patterns



- Million files inserted into a single directory at the same time
- Millions of writers into the same file at the same time
- Jobs from 1 core to N-Million cores
- Files from 0 bytes to N-Pbytes
- Workflows from hours to a year (yes a year on a million cores using a PB DRAM)

# A Simple collective 2D layout as an example
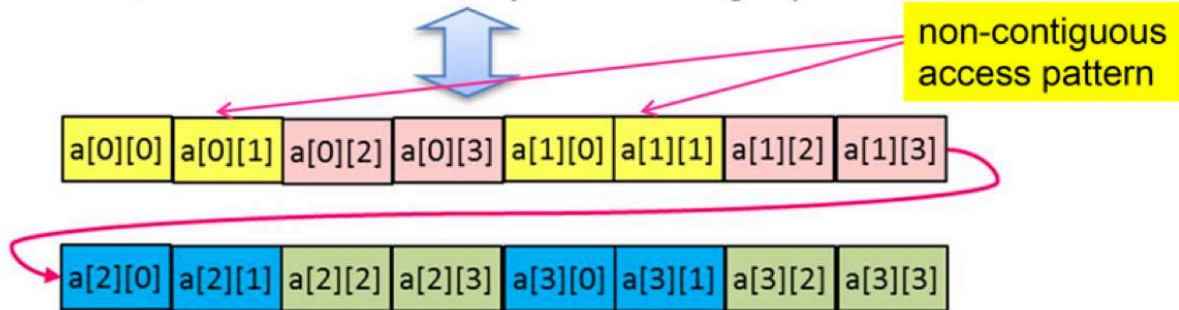
## Collective I/O for 2-Dimensional Data

- 2-Dimensional data accesses by 4 processes



Global view of a 2-D array data among 4 processes

non-contiguous access pattern

File view of a 2-D array data among 4 processes

Apps can map their 1,2,3,N dimensional data onto files in any way they think will help them
This leads to formatting middleware like HDF5, NetCDF, MPI-IO

# Why N->1 strided can be problematic

Process 1　　Process 2　　Process 3　　Process 4

| 11 | 12 | 13 | 14 | | 21 | 22 | 23 | 24 | | 31 | 32 | 33 | 34 | | 41 | 42 | 43 | 44 |

RAID Group 1　　RAID Group 2　　RAID Group 3

PSC
Lustre

LANL
GPFS

# PLFS: A Checkpoint Filesystem for Parallel Applications

John Bent[*], Garth Gibson[‡], Gary Grider; Ben McClelland;
Paul Nowoczynski[§], James Nunez; Milo Polte[†], Meghan Wingate[*]

SC 2009

## ABSTRACT

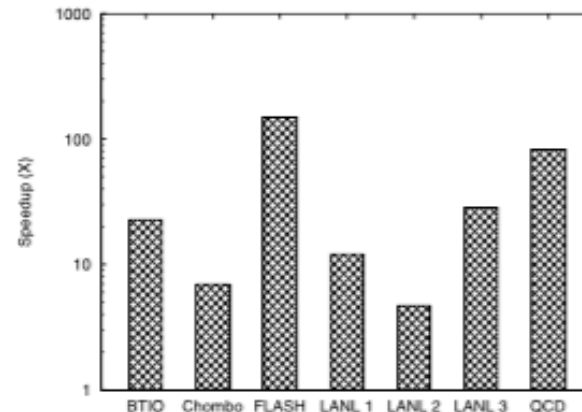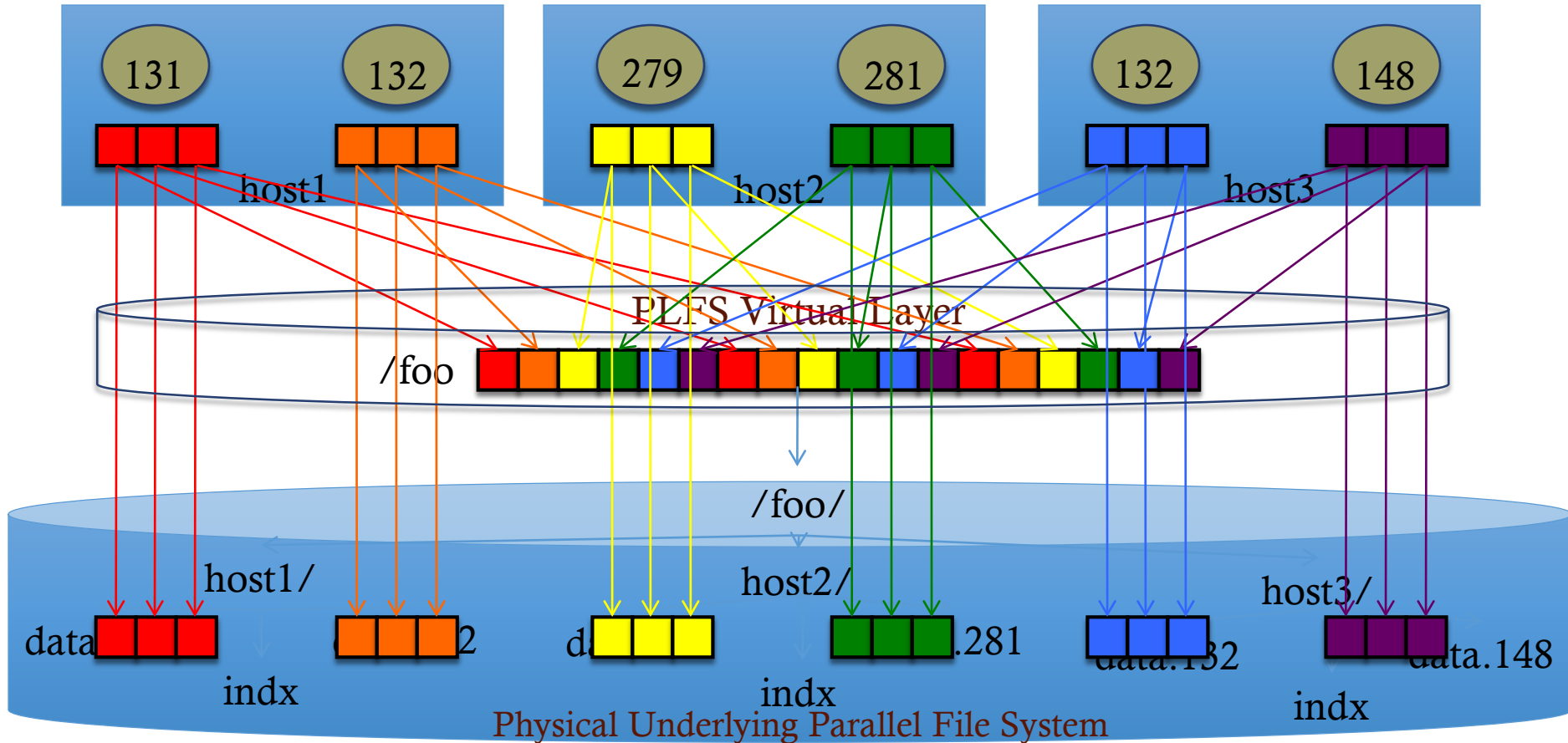Parallel applications running across thousands of processors must protect themselves from inevitable system failures. Many applications insulate themselves from failures by checkpointing. For many applications, checkpointing into a shared single file is most convenient. With such an approach, the size of writes are often small and not aligned with file system boundaries. Unfortunately for these applications, this preferred data layout results in pathologically poor performance from the underlying file system which is optimized for large, aligned writes to non-shared files. To address this fundamental mismatch, we have developed a virtual parallel log structured file system, PLFS. PLFS remaps an application's preferred data layout into one which is optimized for the underlying file system. Through testing on PanFS, Lustre, and GPFS, we have seen that this layer of indirection and reorganization can reduce checkpoint time by an order of magnitude for several important benchmarks and real applications without any application modification.

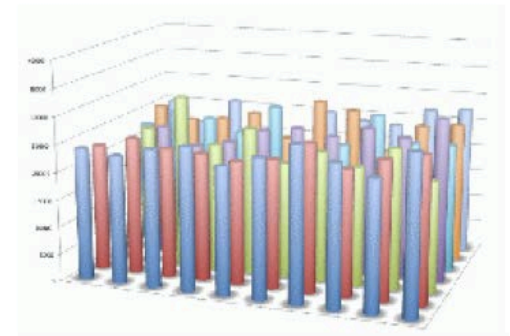1: **Summary of our results.** *This graph summarizes our results which will be explained in detail in Section 4. The key observation here is that our technique has improved checkpoint bandwidths for all seven studied benchmarks and applications by up to several orders of magnitude.*

SNIA. Tutorial

# Decouples Logical from Physical



- N to 1 is fixed by PLFS
- But won't scale to exascale (N-N unscalable (billions of files)

**Tightly Coupled Parallel Application**

Concurrent, unaligned, interspersed IO

**Parallel File System**

Concurrent, aligned, interleaved IO

**Tape Archive**

HPC Storage Stack, 2002-2015

# Economics have shaped our world
## The beginning of storage layer proliferation  2009
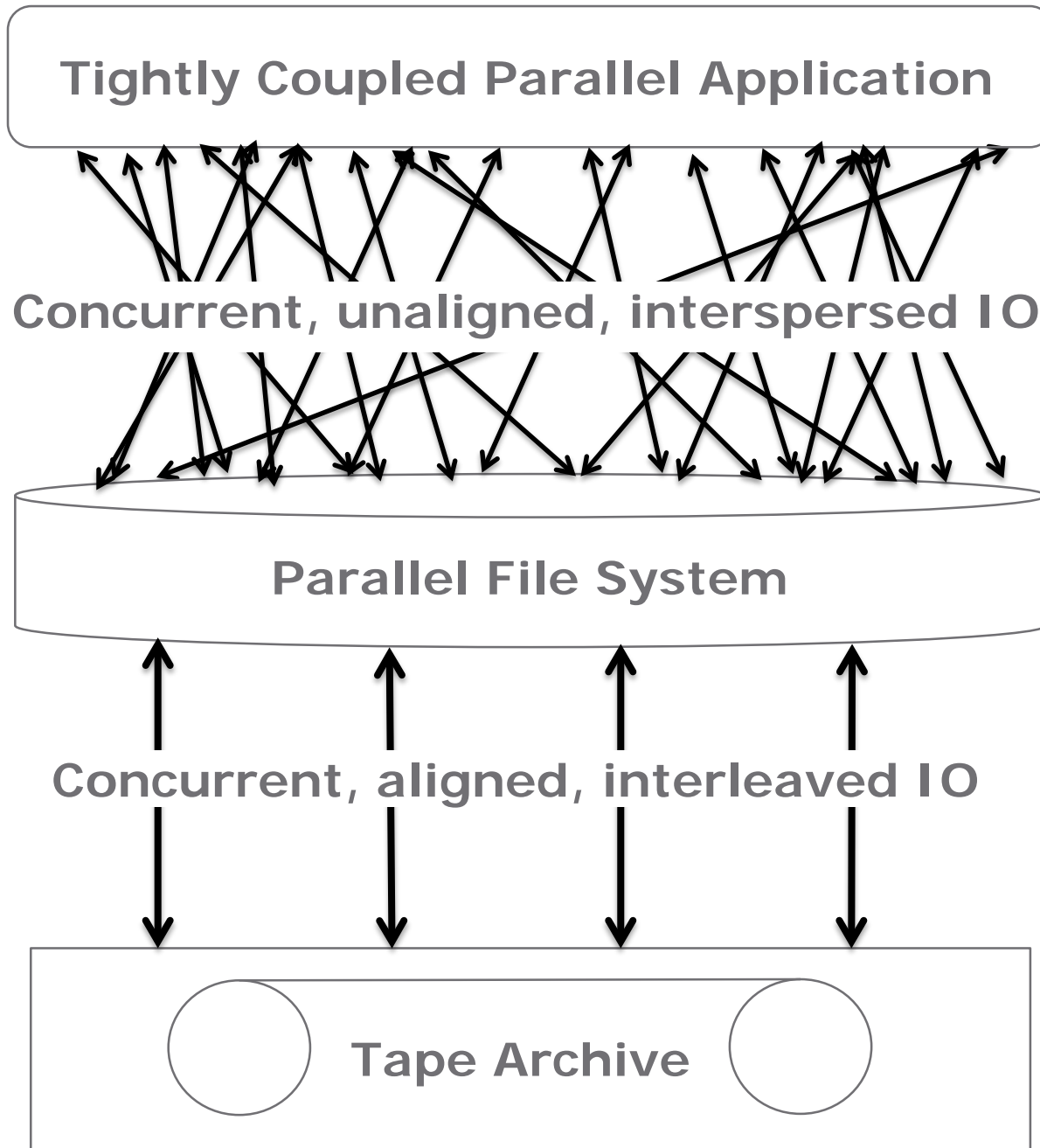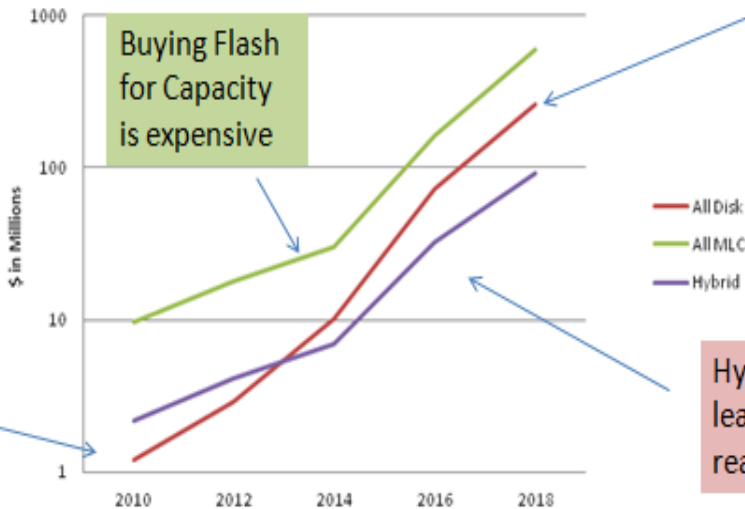


$ in Millions for IO Subsystem for Machine Progression

Buying Flash for Capacity is expensive

Buying disk for BW is expensive

Disk buy for capacity, get BW for Free

Hybrid is at least within reason

- All Disk
- All MLC
- Hybrid

- Economic modeling for large burst of data from memory shows bandwidth / capacity better matched for solid state storage near the compute nodes

■ Economic modeling for archive shows bandwidth / capacity better matched for disk

### Hdwr/media cost 3 mem/mo 10% FS



- new servers
- new disk
- new cartridges
- new drives
- new robots

HPC Storage Stack, 2015-2016

**Tightly Coupled Parallel Application**

**Burst Buffer**

**Parallel File System**

**Tape Archive**

# HPC environment

- Big difference from cloud: parallel, tightly coupled, extremely simple nodes to lower jitter and job failure due to tightly coupled behavior ( one code syncs between all neighbors every 1 millisecond

Burst Buffers

Low-latency interconnect
(IB, vendor proprietary torus)

SAN

Parallel File System
(Lustre, GPFS)

Routing nodes
(IO nodes)

HPC Storage Stack, 2016-2020

Tightly Coupled Parallel Application

Burst Buffer

Parallel File System

Object Store

Tape Archive

**Tightly Coupled Parallel Application**

**Burst Buffer**

**Object Store**

HPC Storage Stack, 2020-

# What about the Capacity Tier:  Won't cloud technology provide the capacity solution?

- Erasure to utilize low cost hardware
- Object to enable massive scale
- Simple minded interface, get put delete

- Problem solved    --$\rightarrow$ NOT

- Works great for apps that are newly written to use this interface
- Doesn't work well for people, people need folders and rename and …
- Doesn't work for the $trillions of apps out there that expect some modest name space capability (parts of POSIX)

# How about a Scalable Near-POSIX Name Space over Cloud style Object Erasure: MarFS

- Best of both worlds
  - Objects Systems
    - Provide massive scaling and efficient erasure techniques
    - Friendly to applications, not to people. People need a name space.
    - Huge Economic appeal (erasure enables use of inexpensive storage)
  - POSIX name space is powerful but has issues scaling
- The challenges
  - Mismatch of POSIX an Object metadata, security, read/write semantics, efficient object/file sizes.
  - No update in place with Objects
  - How do we scale POSIX name space to trillions of files/directories

# MarFS

## What it is

- <span style="color:red">100-1000 GB/sec, Exabytes, Billion files in a directory, Trillions of files total</span>
- Near-POSIX global scalable name space over many POSIX and non POSIX data repositories (Scalable object systems - CDMI, S3, etc.)
  - <span style="color:red">(Scality, EMC ECS, all the way to simple erasure over ZFS's)</span>
- It is small amount of code (C/C++/Scripts)
  - A small Linux Fuse
  - A pretty small parallel batch copy/sync/compare/ utility
  - A moderate sized library both FUSE and the batch utilities call
- Data movement scales just like many scalable object systems
- Metadata scales like NxM POSIX name spaces both across the tree and within a single directory
- It is friendly to object systems by
  - <span style="color:red">Spreading very large files across many objects</span>
  - <span style="color:red">Packing many small files into one large data object</span>

## What it isnt

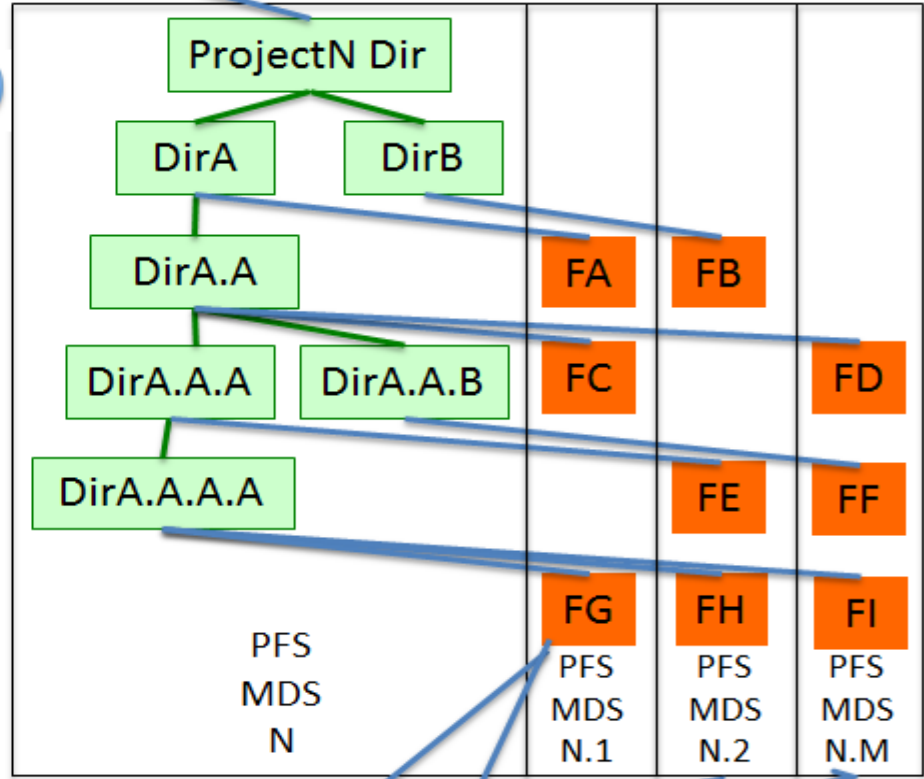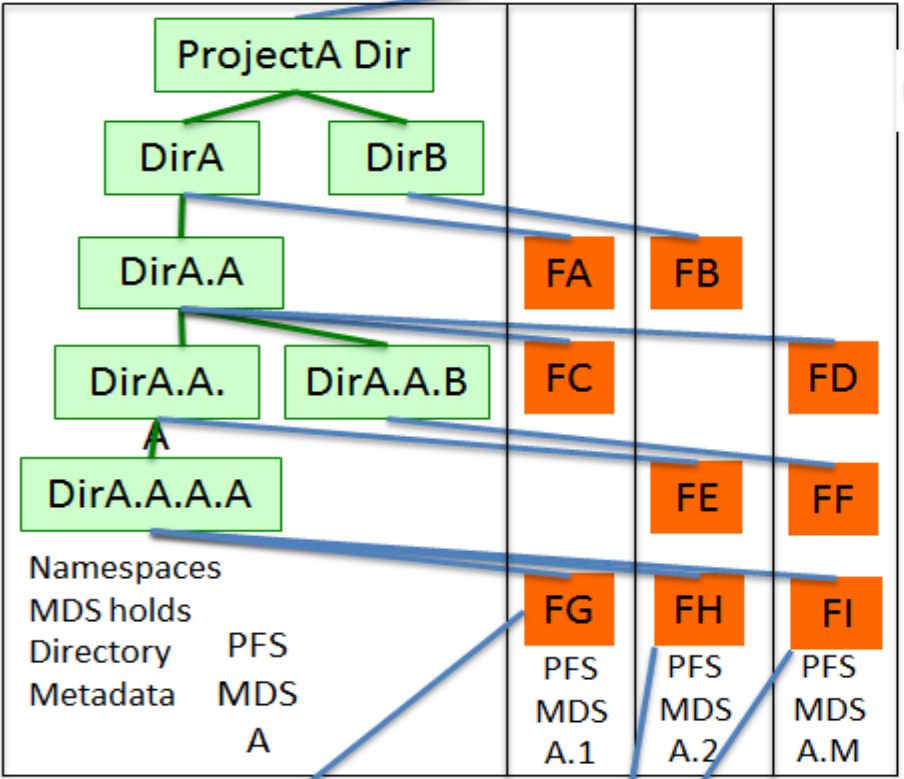- <span style="color:red">No Update in place!  Its not a pure file system, Overwrites are fine but no seeking and writing.</span>

Los Alamos
NATIONAL LABORATORY
EST.1943

# MarFS Scaling



Namespace Project A

Namespace Project N

MarFS

ProjectA Dir
- DirA
  - DirA.A
    - DirA.A.
    - DirA.A.B
      - DirA.A.A.A
- DirB

FA  FB  FC  FD  FE  FF  FG  FH  FI

Namespaces MDS holds Directory Metadata

PFS MDS A

PFS MDS A.1  PFS MDS A.2  PFS MDS A.M

N Name Spaces

ProjectN Dir
- DirA
  - DirA.A
    - DirA.A.A
    - DirA.A.B
      - DirA.A.A.A
- DirB

FA  FB  FC  FD  FE  FF  FG  FH  FI

PFS MDS N

PFS MDS N.1  PFS MDS N.2  PFS MDS N.M

**N X M MDS File Systems (for metadata only)**

File Metadata is hashed over M multiple MDS

Uni Object File    Packed Object File    **Object Repo A**    Multi Object File    **Object Repo X**
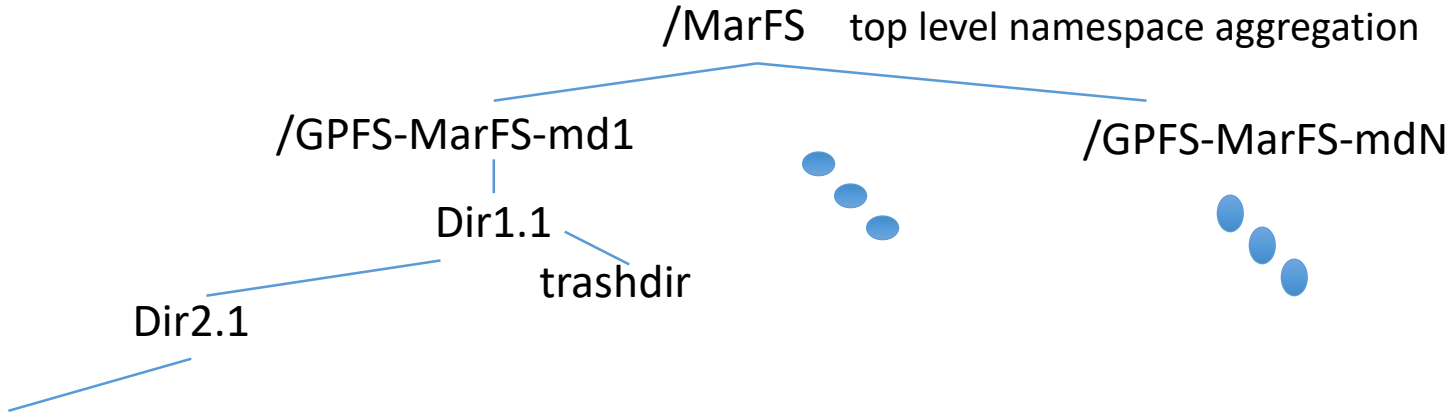
Scaling test on our retired Cielo machine:
835M File Inserts/sec Stat single file < 1 millisecond
> 1 trillion files in the same director
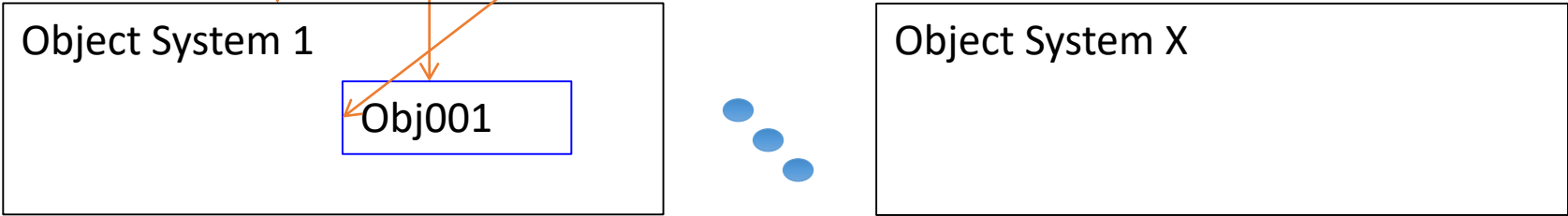
Striping across 1 to X Object Repos

# MarFS Internals Overview Uni-File

/MarFS    top level namespace aggregation

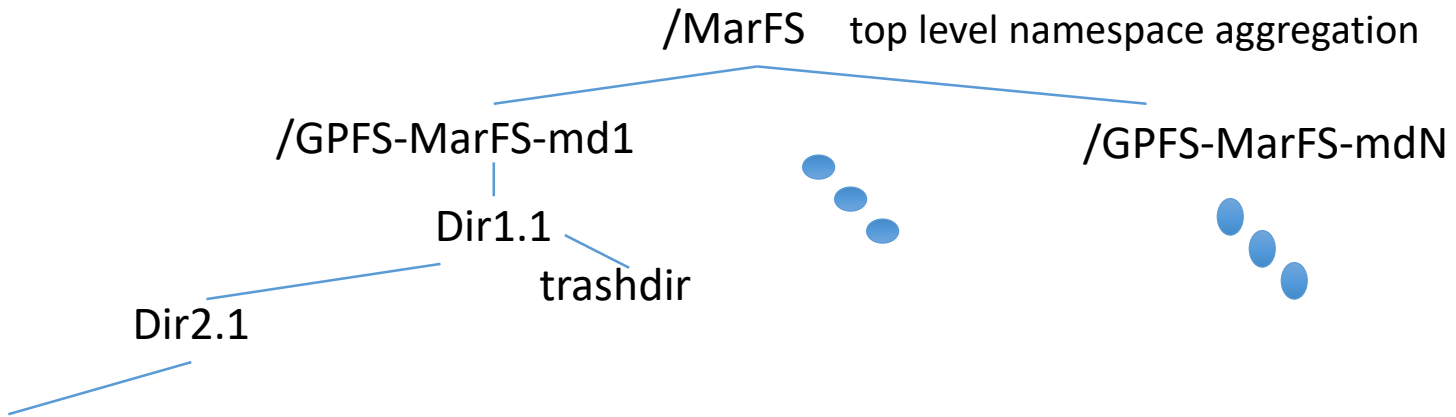/GPFS-MarFS-md1                                    /GPFS-MarFS-mdN

Dir1.1

trashdir

Dir2.1

UniFile  - Attrs: uid, gid, mode, size, dates, etc.
Xattrs - objid repo=1, id=Obj001,  objoffs=0,  chunksize=256M, Objtype=Uni, NumObj=1, etc.

M
e
t
a
d
a
t
a

D
a
t
a

| Object System 1 | | Object System X |
| --- | --- | --- |
| Obj001 | | |

# MarFS Internals Overview Multi-File (striped Object Systems)

/MarFS    top level namespace aggregation

M
e
t
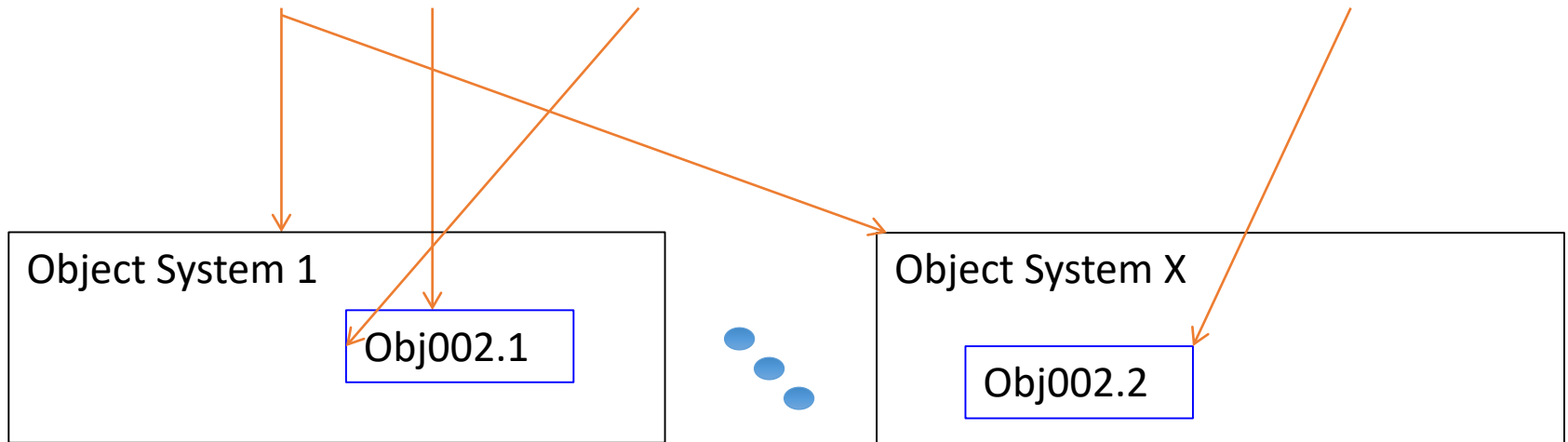a
d
a
t
a

/GPFS-MarFS-md1

/GPFS-MarFS-mdN

Dir1.1

trashdir

Dir2.1

MultiFile  - Attrs: uid, gid, mode, size, dates, etc.
Xattrs - objid repo=S, id=Obj002.,  objoffs=0,  chunksize=256M, ObjType=Multi, NumObj=2, etc.
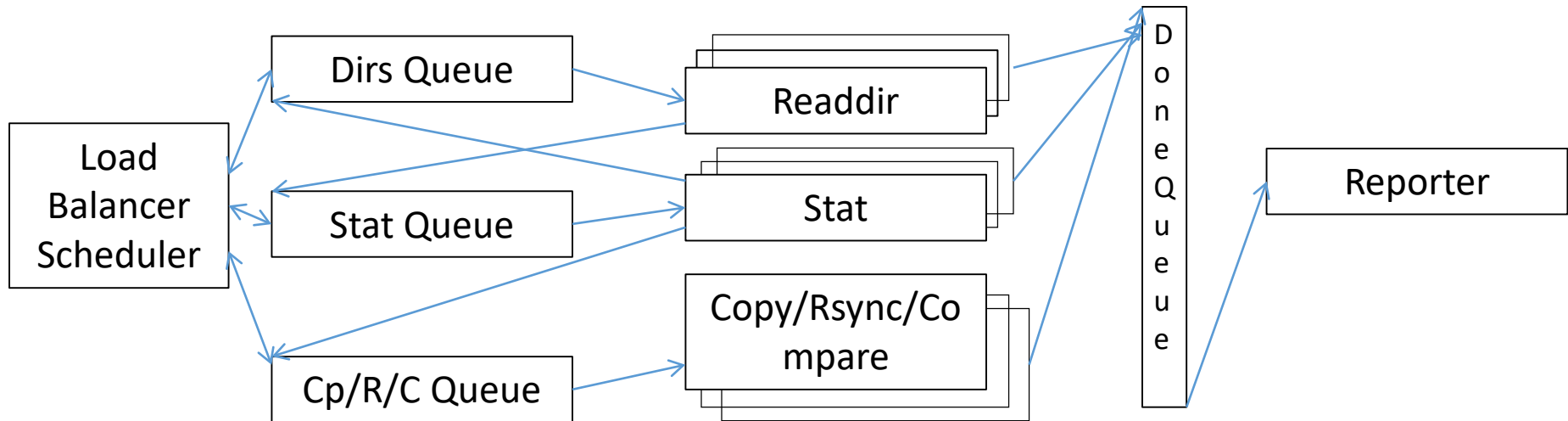
D
a
t
a

Object System 1

Obj002.1

Object System X

Obj002.2

# MarFS Internals Overview Packed-File

/MarFS    top level namespace aggregation

M
e
t
a
d
a
t
a

/GPFS-MarFS-md1

/GPFS-MarFS-mdN

Dir1.1

trashdir

Dir2.1

UniFile  - Attrs: uid, gid, mode, size, dates, etc.
Xattrs - objid repo=1, id=Obj003,  objoffs=4096,  chunksize=256M, Objtype=Packed, NumObj=1, Ojb=4 of 5, etc.

D
a
t
a

Object System 1

Obj003

Object System X

# Pftool – parallel copy/rsync/compare/list tool

- Walks tree in parallel, copy/rsync/compare in parallel.
  - Parallel Readdir's, stat's, and copy/rsinc/compare
- Dynamic load balancing
- Restart-ability for large trees or even very large files
- Repackage: breaks up big files, coalesces small files
- To/From NFS/POSIX/parallel FS/MarFS

# How does it fit into our environment in FY16

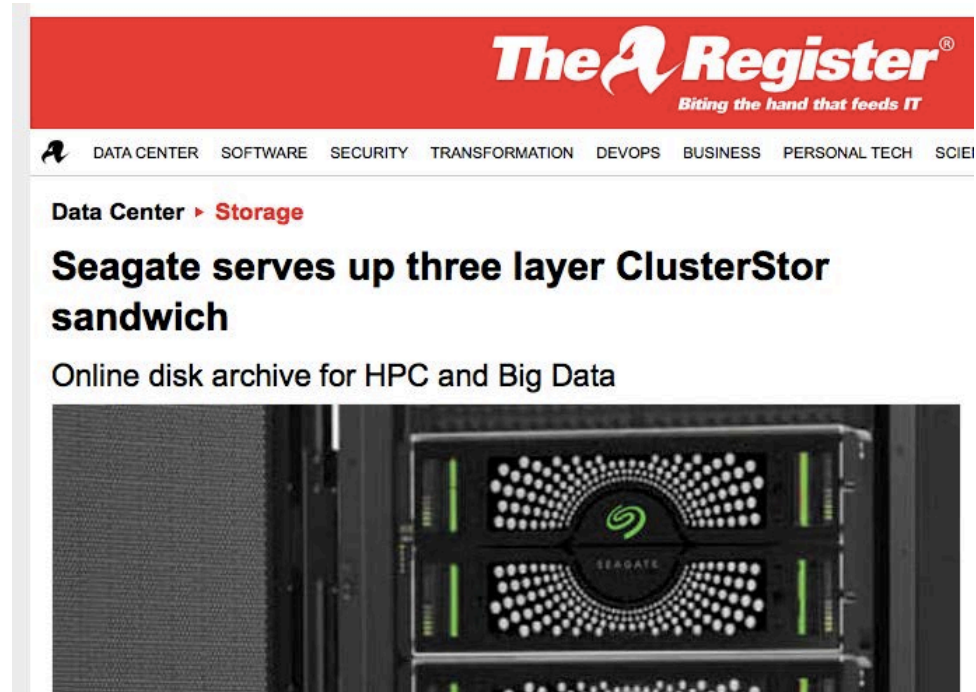# Not Just LANL Developing This Tier



Spectra Logic Campaign Storage LLC



Seagate ClusterStor A200

# From 2 to 4 and back again



| 2002-2015 | 2015-2016 | 2016-2020 | 2020- |



**THE USENIX MAGAZINE**

## Serving Data to the Lunatic Fringe: The Evolution of HPC Storage

John Bent, Brad Settlemyer, and Gary Grider

**Article Section:** STORAGE

Before the advent of Big Data, the largest storage systems in the world were found almost exclusively within high performance computing centers such as those found at US Department of Energy national laboratories. However, these systems are now dwarfed by large datacenters such as those run by Google and Amazon. Although HPC storage systems are no longer the largest in terms of total capacity, they do exhibit the largest degree of concurrent write access to shared data. In this article, we will explain why HPC applications must necessarily exhibit this degree of concurrency and the unique HPC storage architectures required to support them.
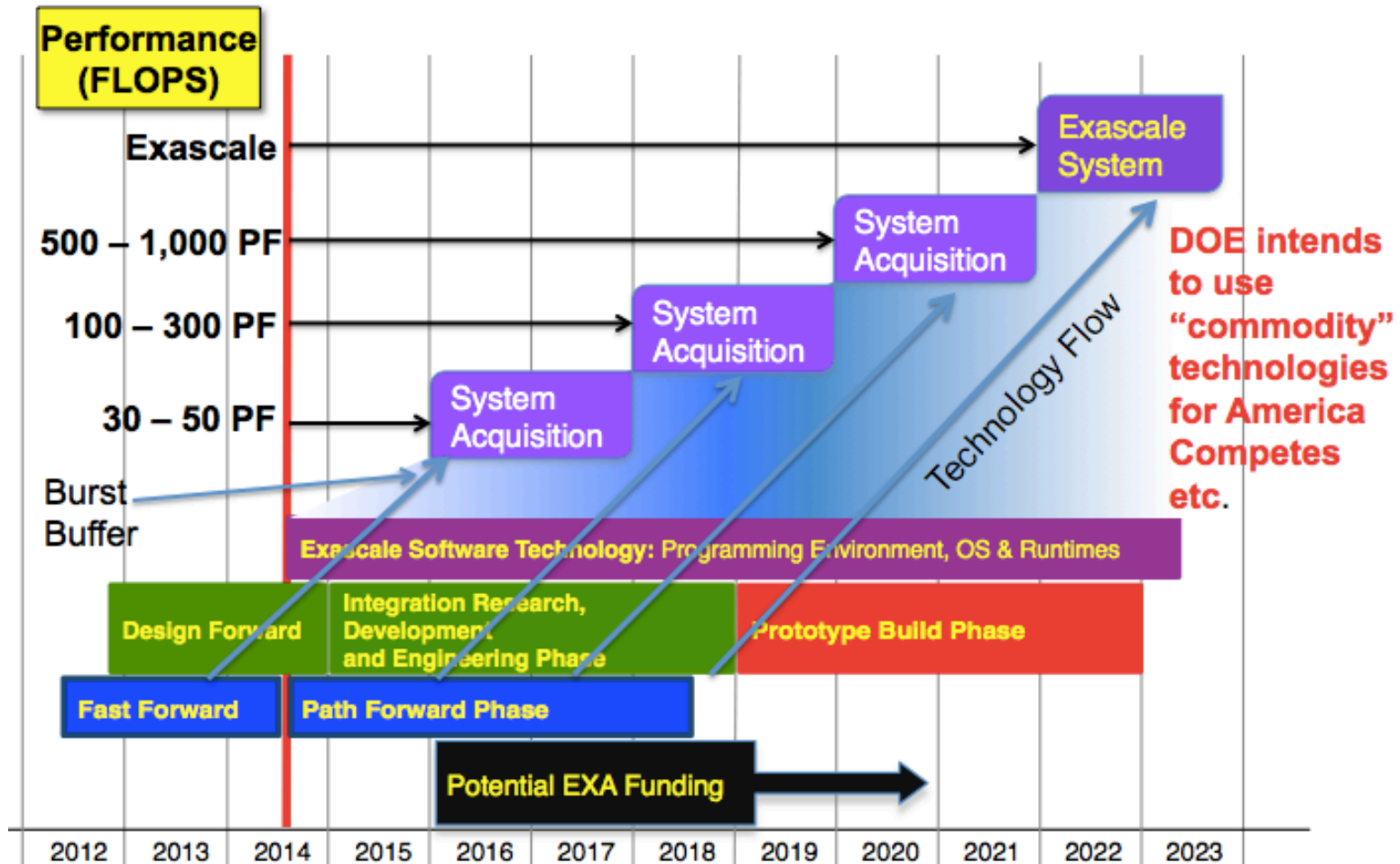
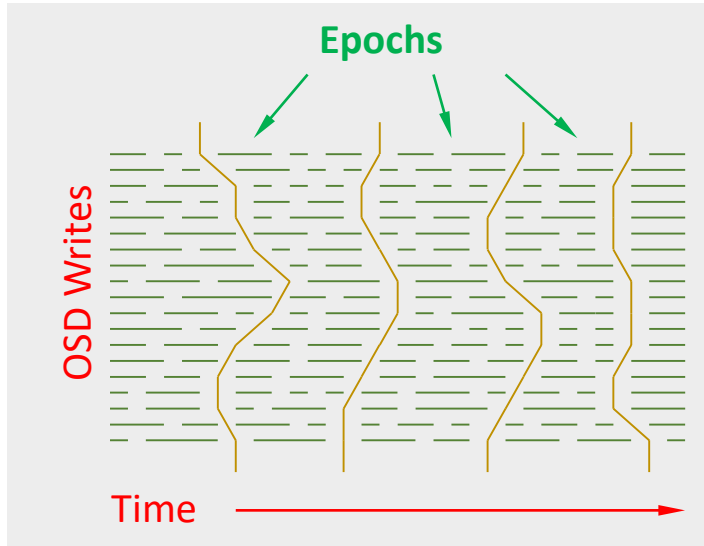# DOE Exascale Computing and Future Considerations

- R&D and integration required to deploy Applications on Exascale computers in 2023+

- Partnership involving: Government, Computer industry, DOE laboratorie, Academia

- Target System Characteristics
  - **1-10 Billion** degrees of concurrency
  - **20-30 MW** Power requirement for one machine (machine only)
  - **<300** cabinets
  - **Development and execution time productivity improvements**
  - **100 PB working sets**
  - **Checkpoint times < 1 minute ( constant failure)**
  - **Storage systems need to be very reliable as the machine wont be**
  - **Leverage->Exploit new technology, dense flash/SCM/low latency high bandwidth byte addressable networks**

# Exascale Computing Timeline
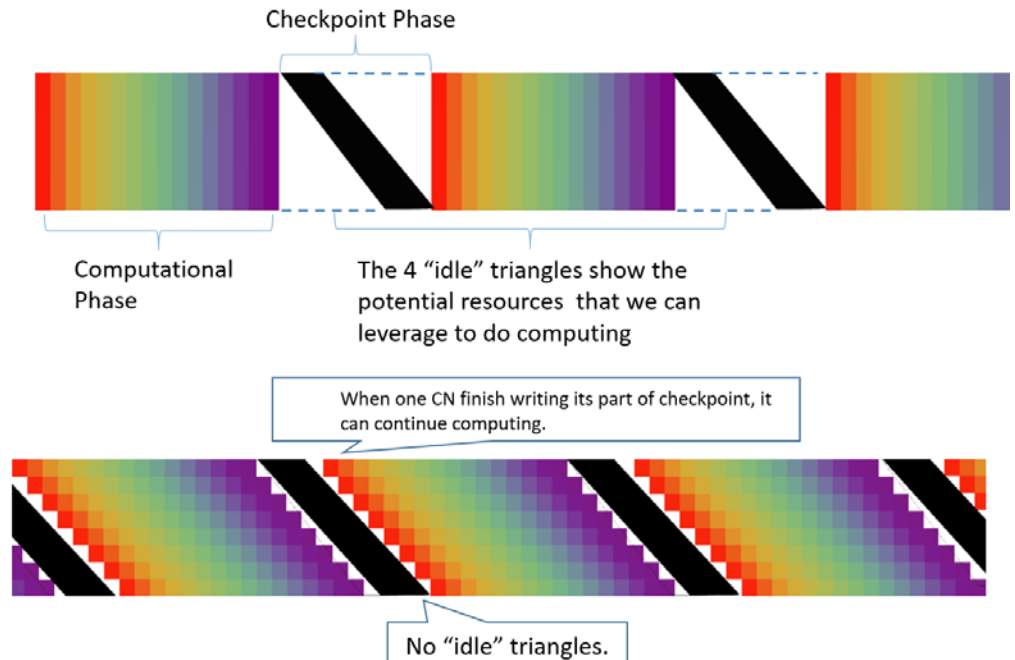## System Perspective

# Transactional/Versioning Coupled with Async Programming Models



## Save every microsecond
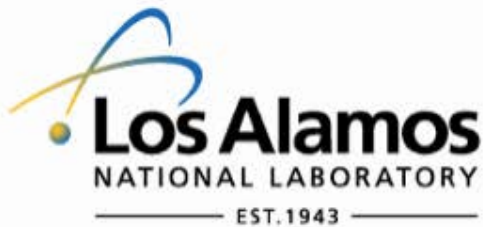
Spread the load over as much time as the app will allow

# Exploiting New Techology (dense flash, SCM, one sided networks)

- Todays storage stacks do not allow full exploitation of even flash latencies never the less SCM – trapped IOP's in software

- Must move away from thin client to heavy server to block storage (Lustre, Ceph, MongoDB).

- These stacks are incredibly thick.  Client KV/Object -> MDS/Access Server->KV/Object Server->Kernel File System->block or network block.

- Head towards embedding more of this in the client and provide light weight one side enabled servers.

- Applications use Middleware (HDF5 for example) to form Objects that go directly to light weight object servers, no real server just light weight kvs/object on the hardware.
    - HDF5 (HDF Group) Vol interface is an example
    - MDHIM (LANL) MultiDimensional Hierarchical Middleware  is a user space distributed KVS middleware
    - DeltaFS (CMU/LANL) is a user space file system Middleware

https://github.com/mar-file-system/marfs
https://github.com/pftool/pftool)

# Thank You For Your Attention

# Attribution & Feedback

The SNIA Education Committee thanks the following Individuals for their contributions to this Tutorial.

| Authorship History | Additional Contributors |
|---|---|
| Gary Grider July 2017 | |

*Please send any questions or comments regarding this SNIA Tutorial to **tracktutorials@snia.org***

SNIA. Tutorial