

Clustered and Parallel Storage System Technologies FAST10

Joel Jacobson, Marc Unangst
{jjacobson,mju}@panasas.com
Panasas, Inc.

1



About Us

Marc Unangst (mju@panasas.com)

- Software Architect, Panasas
- CMU NASD object storage & distributed filesystem
- Panasas ActiveScale Filesystem

Joel Jacobson (jjacobson@panasas.com)

- Senior Software Engineer, Panasas
- Lead of Panasas OSDFS development group
- SNIA OSD technical working group

Thanks also to:

- Rob Latham (robl@mcs.anl.gov)
- Rob Ross (rross@mcs.anl.gov)
- Brent Welch (welch@panasas.com)

2



Outline of the Day

Part 1

- Introduction
- Storage System Models
- Parallel File Systems
 - GPFS
 - PVFS
 - Panasas
 - Lustre

Part 2

- Benchmarking
- MPI-IO
- Future Technologies

Outline of the Day

Part 1

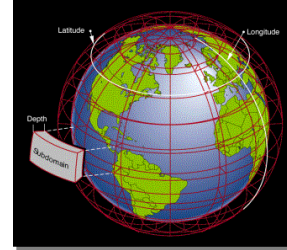
- Introduction**
- Storage System Models
- Parallel File Systems
 - GPFS
 - PVFS
 - Panasas
 - Lustre

Part 2

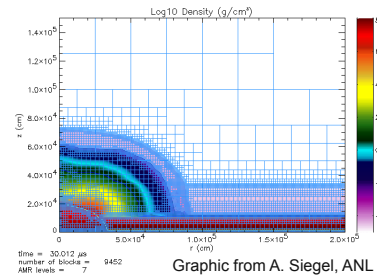
- Benchmarking
- MPI-IO
- Future Technologies

Applications, Data Models, and I/O

- Applications have data models appropriate to domain
 - Multidimensional typed arrays, images composed of scan lines, variable length records
 - Headers, attributes on data
- I/O systems have very simple data models
 - Tree-based hierarchy of containers
 - Some containers have streams of bytes (files)
 - Others hold collections of other containers (directories or folders)
- Someone has to map from one to the other!



Graphic from J. Tannahill, LLNL



Graphic from A. Siegel, ANL

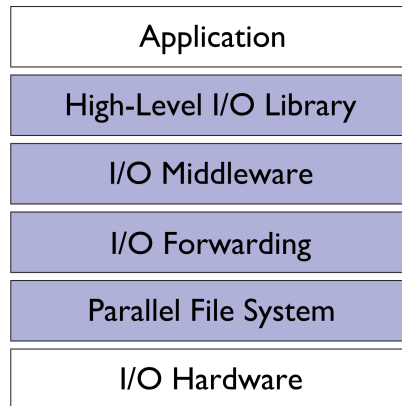
I/O for Computational Science

High-Level I/O Library
maps application abstractions onto storage abstractions and provides data portability.

HDF5, Parallel netCDF, ADIOS

I/O Forwarding
bridges between app. tasks and storage system and provides aggregation for uncoordinated I/O.

IBM cioid, Cray DVS



I/O Middleware
organizes accesses from many processes, especially those using collective I/O.

MPI-IO

Parallel File System
maintains logical space and provides efficient access to data.

PVFS, PanFS, GPFS, Lustre

Additional I/O software provides improved performance and usability over directly accessing the parallel file system. Reduces or (ideally) eliminates need for optimization in application codes.

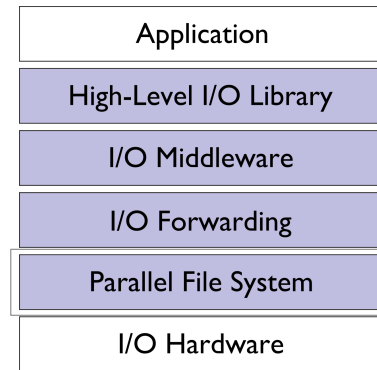
Parallel File System

■ Manage storage hardware

- Present single view
- Stripe files for performance

■ In the I/O software stack

- Focus on concurrent, independent access
- Publish an interface that middleware can use effectively
 - Rich I/O language
 - Relaxed but sufficient semantics



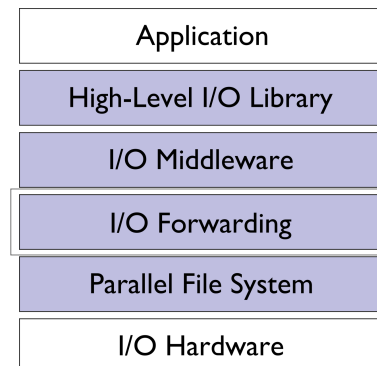
I/O Forwarding

■ Newest layer in the stack

- Present in some of the largest systems
- Provides bridge between system and storage in machines such as the Blue Gene/P

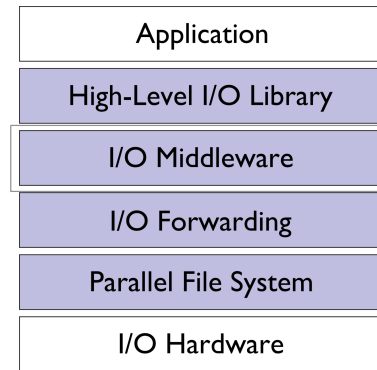
■ Allows for a point of aggregation, hiding true number of clients from underlying file system

■ Poor implementations can lead to unnecessary serialization, hindering performance



I/O Middleware

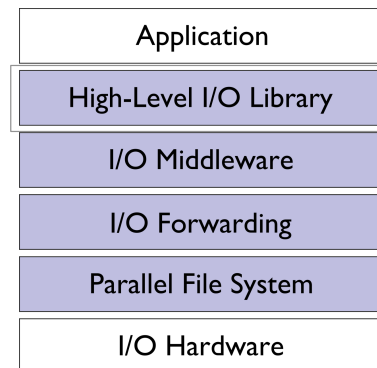
- Match the programming model (e.g. MPI)
- Facilitate concurrent access by groups of processes
 - Collective I/O
 - Atomicity rules
- Expose a generic interface
 - Good building block for high-level libraries
- Efficiently map middleware operations into PFS ones
 - Leverage any rich PFS access constructs, such as:
 - Scalable file name resolution
 - Rich I/O descriptions



9

High Level Libraries

- Match storage abstraction to domain
 - Multidimensional datasets
 - Typed variables
 - Attributes
- Provide self-describing, structured files
- Map to middleware interface
 - Encourage collective I/O
- Implement optimizations that middleware cannot, such as
 - Caching attributes of variables
 - Chunking of datasets



10

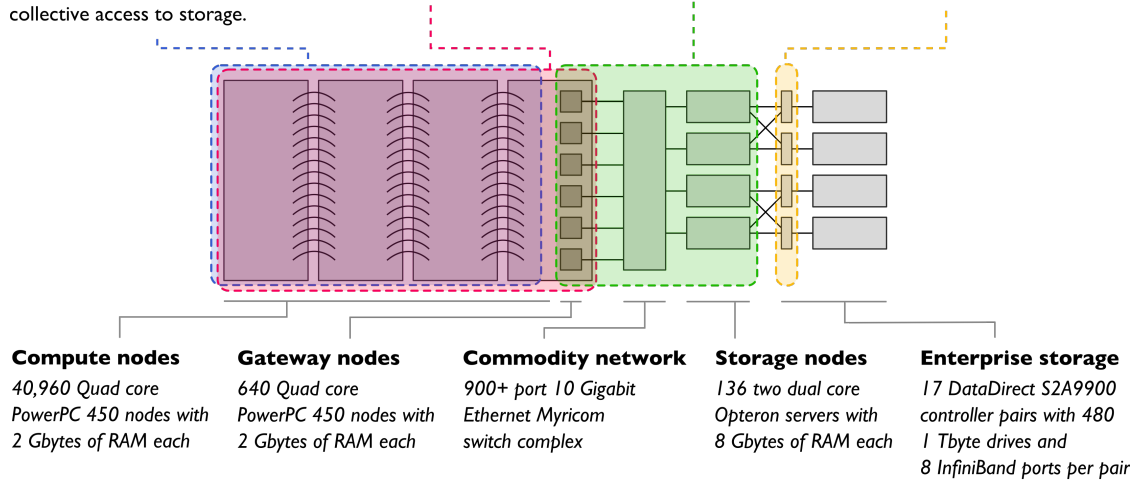
I/O Hardware and Software on Blue Gene/P

High-level I/O libraries execute on compute nodes, mapping application abstractions into flat files, and encoding data in portable formats.
I/O middleware manages collective access to storage.

I/O forwarding software runs on compute and gateway nodes, bridges networks, and provides aggregation of independent I/O.

Parallel file system code runs on gateway and storage nodes, maintains logical storage space and enables efficient access to data.

Drive management software or firmware executes on storage controllers, organizes individual drives, detects drive failures, and reconstructs lost data.



Architectural diagram of the 557 TFlop IBM Blue Gene/P system at the Argonne Leadership Computing Facility.

Outline of the Day

Part 1

Introduction

Storage System Models

Parallel File Systems

- GPFS
- PVFS
- Panasas
- Lustre

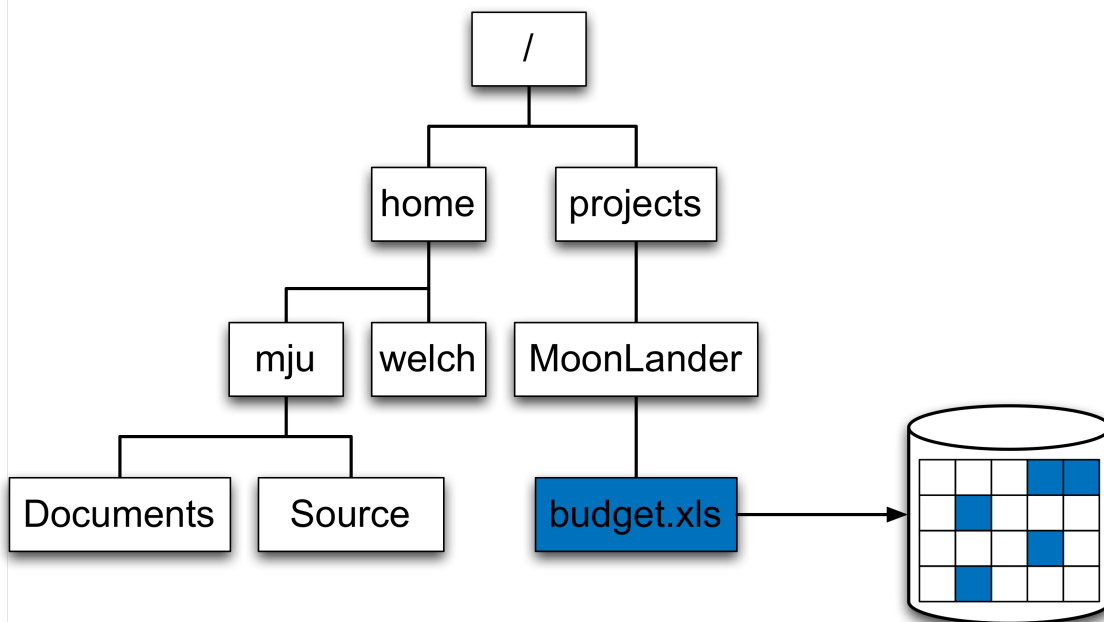
Part 2

Benchmarking

MPI-IO

Future Technologies

Role of the File System



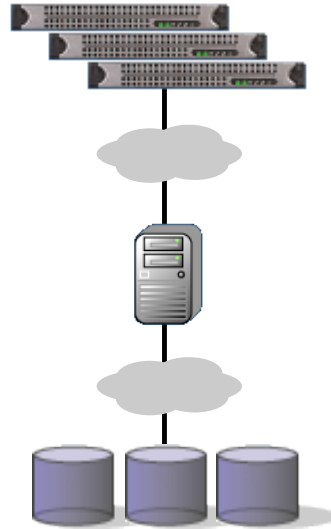
13

Parallel File System Design Issues

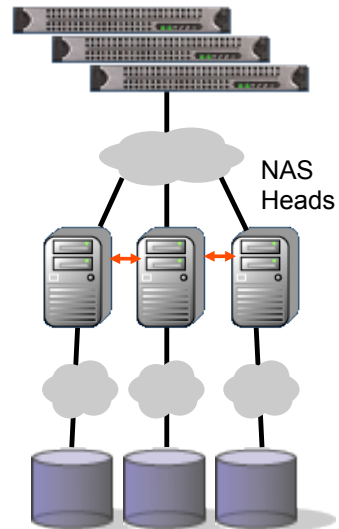
- Same problems as local filesystem
 - Block allocation
 - Metadata management
 - Data reliability and error correction
- Additional requirements
 - Cache coherency
 - High availability
 - Scalable capacity & performance

14

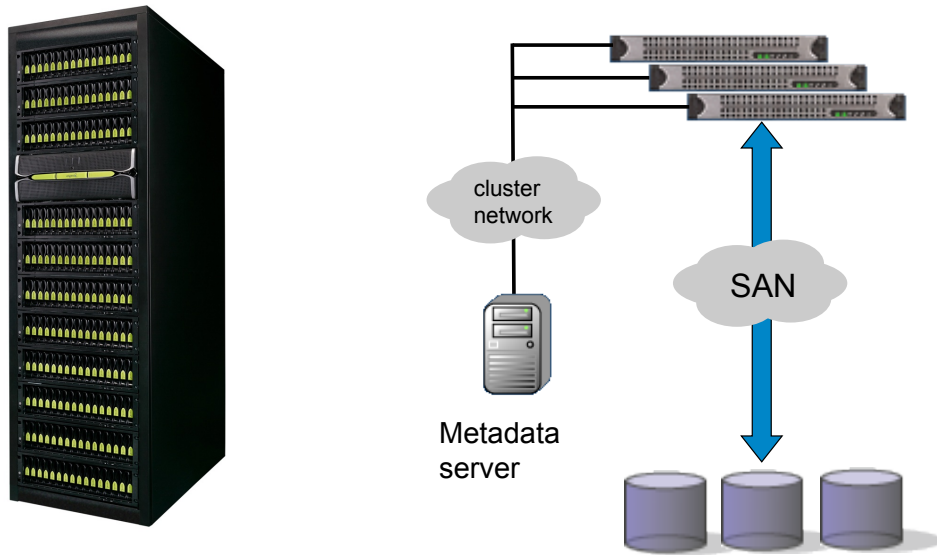
Network Attached Storage (NAS)



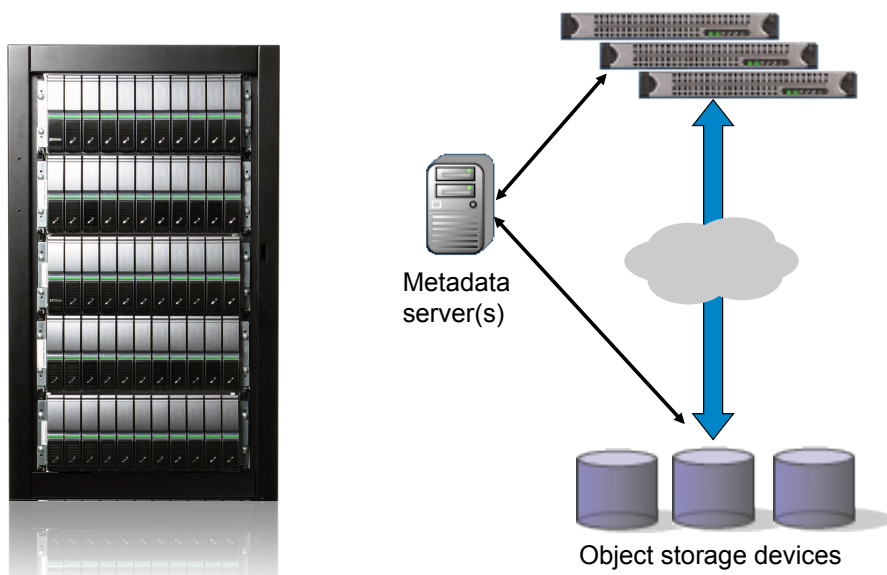
Clustered NAS



SAN Shared Disk File Systems



Object-based Storage Clusters



Object Storage Architecture

Operations

Read block
Write block

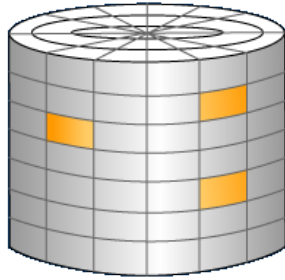
Addressing

Block range

Allocation

External

Block Based Disk



Operations

Create object
Delete object
Read object
Write object
Get Attribute
Set Attribute

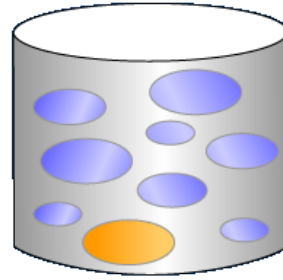
Addressing

[object, byte range]

Allocation

Internal

Object Based Disk



What's in an OSD?



+



Lustre OSS
PVFS storage node

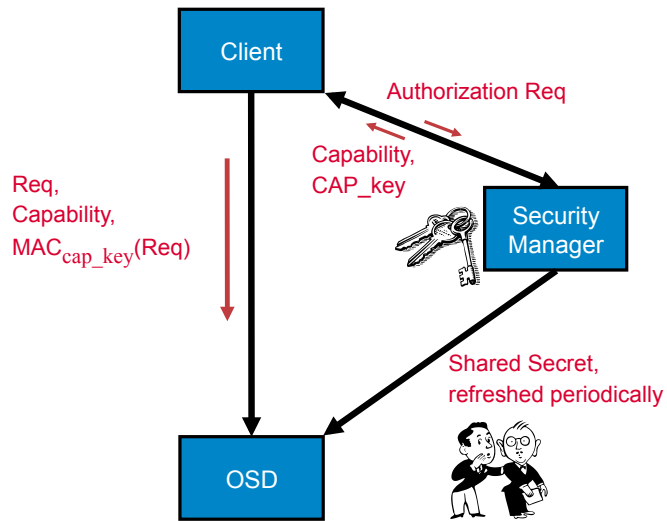


Panasas
StorageBlade



Seagate
prototype

SCSI T10 OSD Security Model



21

Strengths of Object Storage

- Scalable block allocation
- Data relationships exposed to OSD
- Extensible metadata
- Fine-grained security
- Command set friendly to embedded devices

22

Outline of the Day

Part 1

Introduction
Storage System Models
Parallel File Systems

- GPFS
- PVFS
- Panasas
- Lustre

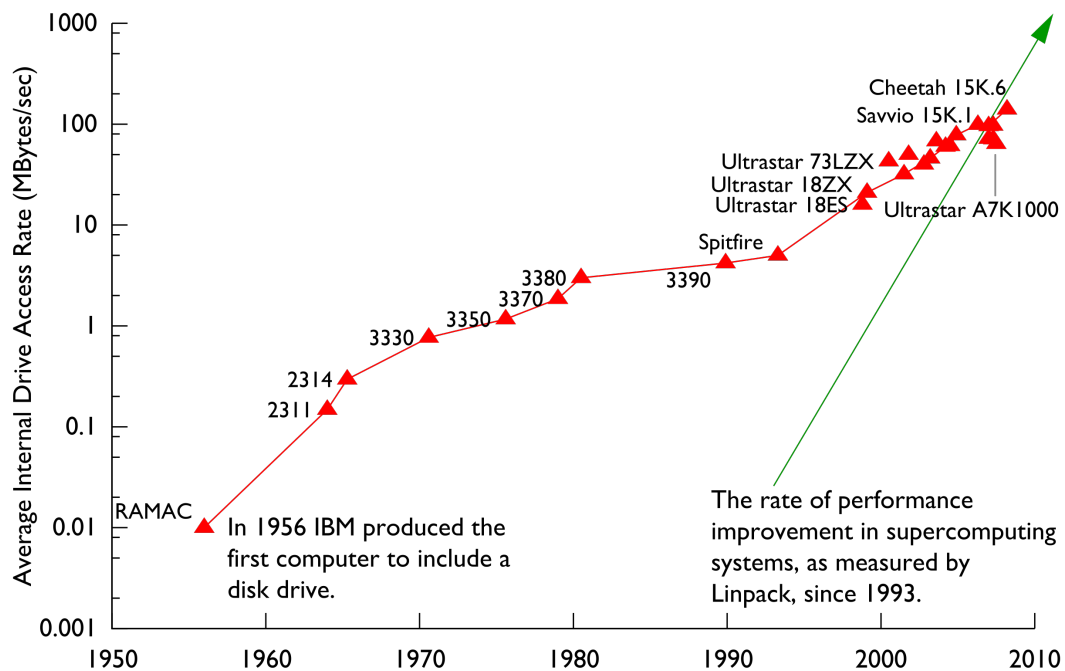
Part 2

Benchmarking
MPI-IO
Future Technologies

23



Disk Access Rates over Time

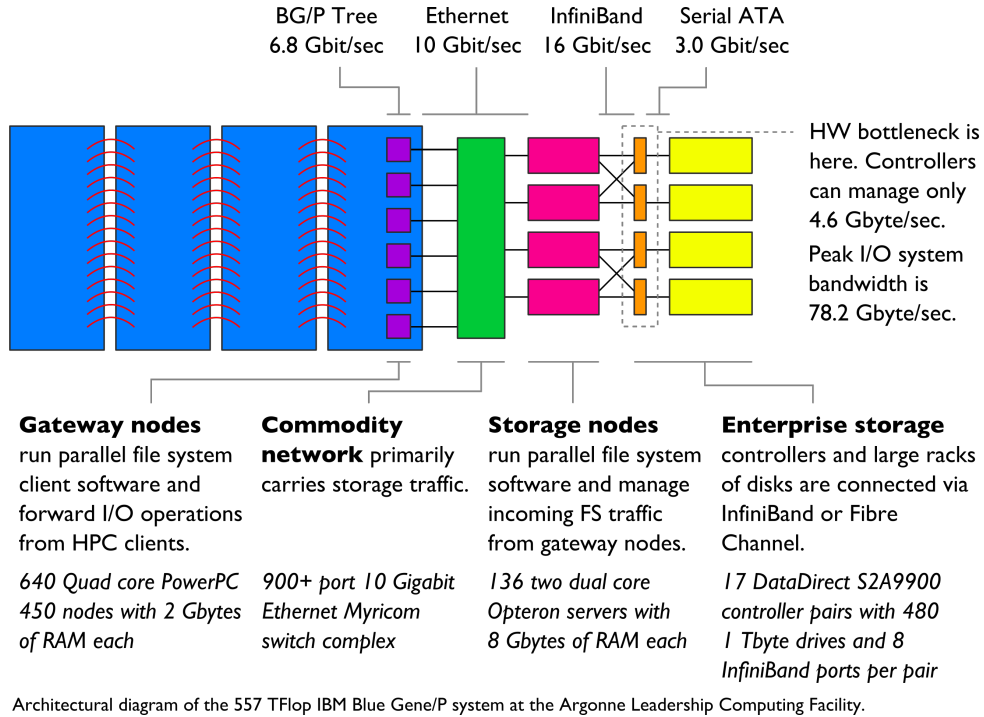


Thanks to R. Freitas of IBM Almaden Research Center for providing much of the data for this graph.

24



Blue Gene/P Parallel Storage System



25

panasas

Production Parallel File Systems

- All four systems scale to support the very largest compute clusters
 - LLNL Purple, LANL RoadRunner, Sandia Red Storm, etc.
- All but GPFS delegate block management to “object-like” data servers or OSDs
- Approaches to metadata vary
- Approaches to fault tolerance vary
- Emphasis on features & “turn-key” deployment vary

GPFS

panasas

PVFS

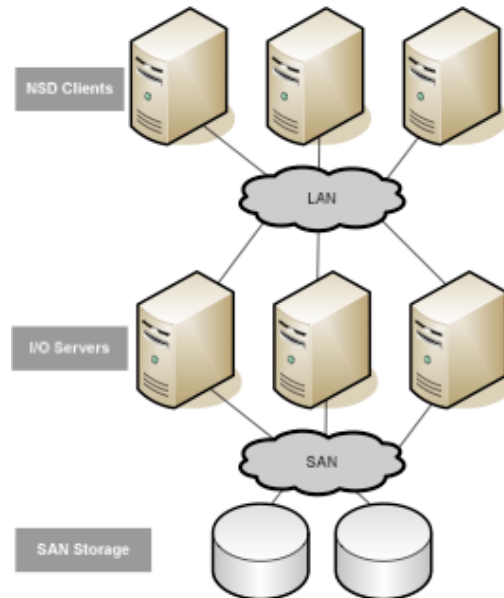
lustre

26

panasas

IBM GPFS

- General Parallel File System
- Legacy: IBM Tiger multimedia filesystem
- Commercial product
- Lots of configuration flexibility
 - AIX, Linux, Windows
 - Direct storage, Virtual Shared Disk, Network Shared Disk
 - Clustered NFS re-export
- Block interface to storage nodes
- Distributed locking



27

panasas

GPFS: Block Allocation

- I/O server exports exports local disk via block-oriented protocol
- Block allocation map shared by all nodes
 - Block map split into N regions
 - Each region has 1/Nth of each I/O server's blocks
- Writing node performs block allocation
 - Locks a region of the block map to find free blocks
 - Updates inode & indirect blocks
 - If # regions \approx # client nodes, block map sharing reduced or eliminated
- Stripe each file across multiple I/O servers (RAID-0)
- Large block size (1-4 MB) typically used
 - Increases transfer size per I/O server
 - Match block size to RAID stripe width
 - Minimizes block allocation overhead
 - Not great for small files – min allocation is 1/32 of a block

28

panasas

GPFS: Metadata Management

- Symmetric model with distributed locking
- Each node acquires locks and updates metadata structures itself
- Global token server manages locking assignments
 - Client accessing a shared resource contacts token server
 - Token server gives token to client, or tells client current holder of token
 - Token owner manages locking, etc. for that resource
 - Client acquires read/write lock from token owner before accessing resource
- inode updates optimized for multiple writers
 - Shared write lock on inode
 - “Metanode token” for file controls which client updates inode
 - Other clients send inode updates to metanode

GPFS: Caching

- Clients cache reads and writes
- Strong coherency, based on distributed locking
- Client acquires R/W lock before accessing data
- Optimistic locking algorithm
 - First node accesses 0-1M, locks 0...EOF
 - Second node accesses 8M-9M
 - First node reduces its lock to 0...8191K
 - Second node locks 8192K...EOF
 - Lock splitting assumes client will continue accessing in current pattern (forward or backward sequential)
- Client cache (“page pool”) pinned and separate from OS page/buffer cache

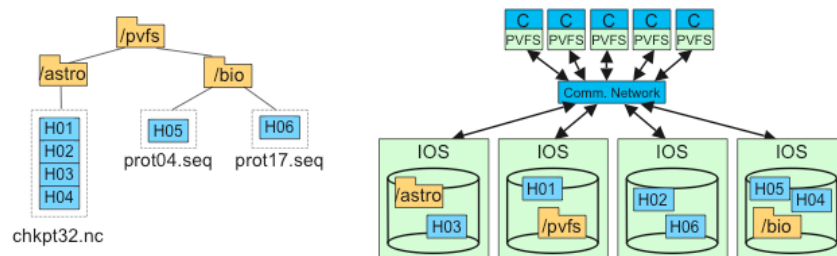
GPFS: Reliability

- RAID underneath I/O server to handle disk failures & sector errors
- Replication across I/O servers supported, but typically only used for metadata
- I/O server failure handled via dual-attached RAID or SAN
 - Backup I/O server takes over primary's disks if it fails
- Nodes journal metadata updates before modifying FS structures
 - Journal is per-node, so no sharing/locking issues
 - Journal kept in shared storage (i.e., on the I/O servers)
 - If node crashes, another node replays its journal to make FS consistent
- Quorum/consensus protocol to determine set of “online” nodes

31

PVFS

- Parallel Virtual Filesystem
- Open source
- Linux based
- Community development
 - Led by Argonne National Lab
- Asymmetric architecture (data servers & clients)
- Data servers use object-like API
- Focus on needs of HPC applications
 - Interface optimized for MPI-IO semantics, not POSIX



32

PVFS: Block Allocation

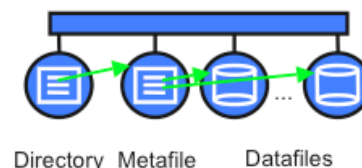
- I/O server exports file/object oriented API
 - Storage object (“dataspace”) on an I/O server addressed by numeric handle
 - Dataspace can be stream of bytes or key/value pairs
 - Create dataspace, delete dataspace, read/write
- Files & directories mapped onto dataspaces
 - File may be single dataspace, or chunked/striped over several
- Each I/O server manages block allocation for its local storage
- I/O server uses local filesystem to store dataspaces
- Key/value dataspace stored using Berkeley DB table

33



PVFS: Metadata Management

- Directory dataspace contains list of names & metafile handles
- Metafile dataspace contains
 - Attributes (permissions, owner, xattrs)
 - Distribution function parameters
 - Datafile handles
- Datafile(s) store file data
 - Distribution function determines pattern
 - Default is 64 KB chunk size and round-robin placement
- Directory and metadata updates are atomic
 - Eliminates need for locking
 - May require “losing” node in race to do significant cleanup
- System configuration (I/O server list, etc.) stored in static file on all I/O servers



34



PVFS: Caching

- Client only caches immutable metadata and read-only files
- All other I/O (reads, writes) go through to I/O node
- Strong coherency (writes are immediately visible to other nodes)
- Flows from PVFS2 design choices
 - No locking
 - No cache coherency protocol
- I/O server can cache data & metadata for local dataspaces
- All prefetching must happen on I/O server
- Reads & writes limited by client's interconnect

PVFS: Reliability

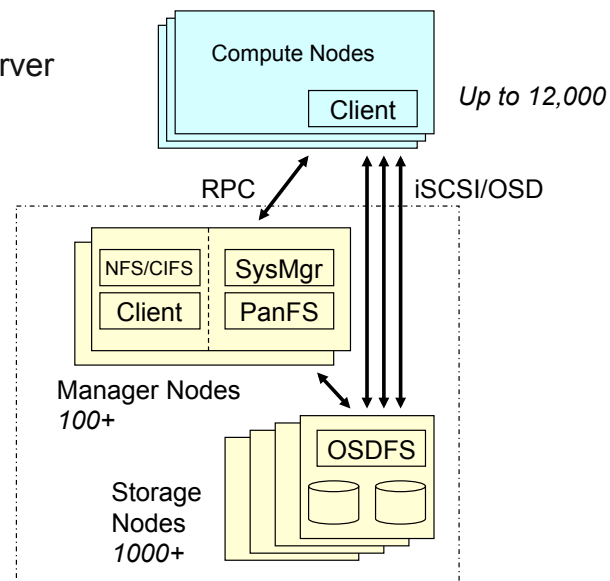
- Similar to GPFS
 - RAID underneath I/O server to handle disk failures & sector errors
 - Dual attached RAID to primary/backup I/O server to handle I/O server failures
- Linux HA used for generic failover support
- Sequenced operations provide well-defined crash behavior
 - Example: Creating a new file
 - Create datafiles
 - Create metafile that points to datafiles
 - Link metafile into directory (atomic)
 - Crash can result in orphans, but no other inconsistencies

Panasas ActiveScale (PanFS)

- Commercial product based on CMU NASD research
- Complete “appliance” solution (HW + SW), blade server form factor

- DirectorBlade = metadata server
- StorageBlade = OSD

- Coarse grained metadata clustering
- Linux native client for parallel I/O
- NFS & CIFS re-export
- Integrated battery/UPS
- Integrated 10GE switch
- Global namespace



37

PanFS: Block Allocation

- OSD exports object-oriented API based on T10 OSD
 - Objects have a number (object ID), data, and attributes
 - CREATE OBJECT, REMOVE OBJECT, READ, WRITE, GET ATTRIBUTE, SET ATTRIBUTE, etc.
 - Commands address object ID and data range in object
 - Capabilities provide fine-grained revocable access control
- OSD manages private local storage
 - Two SATA drives, 500-2000 GB each, 1-4 TB total capacity
 - 10-40 TB per shelf enclosure
- Specialized filesystem (OSDFS) stores objects
 - Delayed floating block allocation
 - Efficient copy-on-write support
- Files and directories stored as “virtual objects”
 - Virtual object striped across multiple container objects on multiple OSDs

38

PanFS: Metadata Management

- Directory is a list of names & object IDs in a RAID-1 virtual object
- Filesystem metadata stored as object attributes
 - Owner, ACL, timestamps, etc.
 - Layout map describing RAID type & OSDs that hold the file
- Metadata server (DirectorBlade)
 - Checks client permissions & provides map/capabilities
 - Performs namespace updates & directory modifications
 - Performs most metadata updates
- Client modifies some metadata directly (length, timestamps)
- Coarse-grained metadata clustering based on directory hierarchy

PanFS: Caching

- Clients cache reads & writes
- Strong coherency, based on callbacks
 - Client registers callback with metadata server
 - Callback type identifies sharing state (unshared, read-only, read-write)
 - Server notifies client when file or sharing state changes
- Sharing state determines caching allowed
 - Unshared: client can cache reads & writes
 - Read-only shared: client can cache reads
 - Read-write shared: no client caching
 - Specialized “concurrent write” mode for cooperating apps (e.g. MPI-IO)
- Client cache shared with OS page/buffer cache

PanFS: Reliability

- RAID-1 & RAID-5 across OSDs to handle disk failures
 - Any failure in StorageBlade is handled via rebuild
 - Declustered parity allows scalable rebuild
- “Vertical parity” inside OSD to handle sector errors
- Integrated shelf battery makes all RAM in blades into NVRAM
 - Metadata server journals updates to in-memory log
 - Failover config replicates log to 2nd blade’s memory
 - Log contents saved to DirectorBlade’s local disk on panic or power failure
 - OSDFS commits updates (data+metadata) to in-memory log
 - Log contents committed to filesystem on panic or power failure
 - Disk writes well ordered to maintain consistency
- System configuration in replicated database on subset of DirectorBlades

41

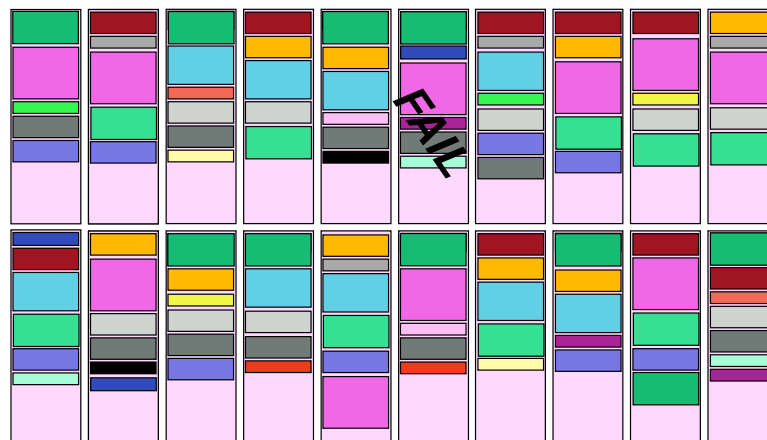


PanFS: Declustered RAID

- Each file striped across different combination of StorageBlades
- Component objects include file data and file parity
- File attributes replicated on first two component objects
- Components grow & new components created as data written
- Declustered, randomized placement distributes RAID workload

20 OSD
Storage Pool

Mirrored
or 9-OSD
Parity
Stripes



Read
about
half of
each
surviving
OSD

Write a
little
to each
OSD

Scales up
in larger
Storage
Pools

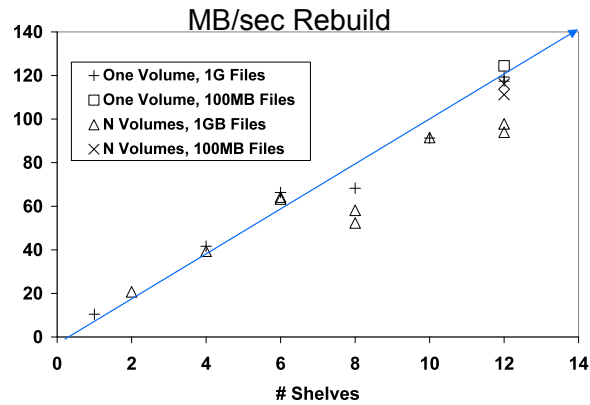
42



Panasas Scalable Rebuild

- Two main causes of RAID failures
 - 1) 2nd drive failure in same RAID set during reconstruction of 1st failed drive
 - Risk of two failures depends on time-to-repair
 - 2) Media failure in same RAID set during reconstruction of 1st failed drive

- Shorter repair time in larger storage pools
 - From 13 hours to 30 minutes
- Four techniques to reduce MTTR
 - Use multiple “RAID engines” (DirectorBlades) in parallel
 - Spread disk I/O over more disk arms (StorageBlades)
 - Reconstruct data blocks only, not unused space
 - Proactively remove failing blades (SMART trips, other heuristics)

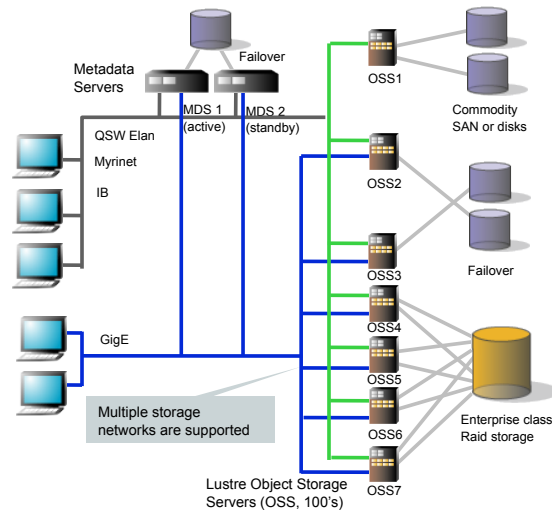


43



Lustre

- Open source object-based parallel file system
 - Based on CMU NASD architecture
 - Lots of file system ideas from Coda and InterMezzo
 - ClusterFS acquired by Sun, 9/2007
- Originally Linux-based, Sun now porting to Solaris
- Asymmetric design with separate metadata server
- Proprietary RPC network protocol between client & MDS/OSS
- Distributed locking with client-driven lock recovery



Lustre material from www.lustre.org and various talks

44



Lustre: Block Allocation

- Each OSS (object storage server) manages one or more OSTs (object storage target)
 - Typically 2-25 OSTs per OSS (max OST size 8 TB)
 - Client communicates with OSS via proprietary RPC protocol
 - RPC built on LNET message-passing facility (based on Sandia Portals)
 - LNET supports RDMA over IB, Myrinet, and Quadrics Elan
- OST stores data in modified ext3 file system
- Currently porting OST to ZFS
 - User-level ZFS via FUSE on Linux
 - In-kernel ZFS on Solaris
- RAID-0 striping across OSTs
 - No dynamic space management among OSTs (i.e., no object migration to balance capacity)
- Snapshots and quota done independently in each OST

Lustre: Metadata

- Metadata server (MDS) hosts metadata target (MDT), which stores namespace tree and file metadata
- MDT uses a modified ext3 filesystem to store Lustre metadata
 - Directory tree of “stub” files that represents Lustre namespace
 - Lustre metadata stored in stub file’s extended attributes
 - Regular filesystem attributes (owner, group, permissions, size, etc.)
 - List of object/OST pairs that contain file’s data (storage map)
 - Single MDS and single MDT per Lustre filesystem
 - Clustered MDS with multiple MDTs is on roadmap (Lustre 2.0)
- Distributed lock protocol among MDS, OSS, and clients
 - “Intents” convey hints about the high-level file operations so the right locks can be taken and server round-trips avoided
 - If a failure occurs (MDS or OSS), clients do lock recovery after failover

Lustre: Caching

- Clients can cache reads, writes, and some metadata operations
- Locking protocol used to protect cached data and serialize access
 - OSS manages locks for objects on its OSTs
 - MDS manages locks on directories & inodes
 - Client caches locks and can reuse them across multiple I/Os
 - MDS/OSS recalls locks when conflict occurs
 - Lock on logical file range may span several objects/OSTs
- Directory locks allow client to do CREATE without round-trip to MDS
 - Only for unshared directory
 - Create not “durable” until file is written & closed
 - Non-POSIX semantic but helpful for many applications
- Client cache shared with OS page/buffer cache

Lustre: Reliability

- Block-based RAID underneath OST/MDT
- Failover managed by external software (Linux-HA)
- OSS failover (active/active or clustered)
 - OSTs on dual-ported RAID controller
 - OSTs on SAN with connectivity to all OSS nodes
- MDS failover (active/passive)
 - MDT on dual-ported RAID controller
 - Typically use dedicated RAID for MDT due to different workload
- Crash recovery based on logs and transactions
 - MDS logs operation (e.g., file delete)
 - Later response from OSS cancels log entry
 - Some client crashes cause MDS log rollback
 - MDT & OST use journaling filesystem to avoid fsck
- LNET supports redundant networks and link failover

Design Comparison

	GPFS	PVFS	Panasas	Lustre
Block mgmt	Shared block map	Object based	Object based	Object based
Metadata location	With data	With data	With data	Separate
Metadata written by	Client	Client	Client, server	Server
Cache coherency & protocol	Coherent; distributed locking	Cache immutable/ RO data only	Coherent; callbacks	Coherent; distributed locking
Reliability	Block RAID	Block RAID	Object RAID	Block RAID

Other File Systems

■ GFS (Google)

- Single metadata server + 100s of chunk servers
- Specialized semantics (not POSIX)
- Design for failures; all files replicated 3+ times
- Geared towards colocated processing (MapReduce)

■ Ceph (UCSC)

- OSD-based parallel filesystem
- Dynamic metadata partitioning between MDSs
- OSD-directed replication based on CRUSH distribution function (no explicit storage map)

■ Clustered NAS

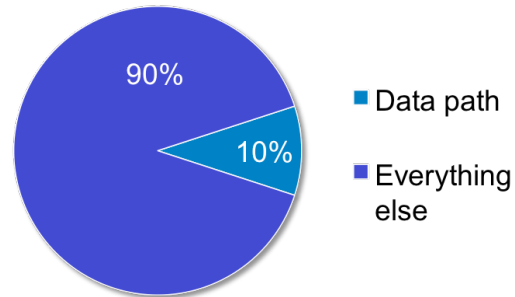
- NetApp GX, Isilon, BlueArc, etc.

Other Issues

What about...

- Monitoring & troubleshooting?
- Backups?
- Snapshots?
- Disaster recovery & replication?
- Capacity management?
- System expansion?
- Retiring old equipment?

Development Effort



51

Themes

"A supercomputer is a device for turning compute-bound problems into I/O-bound problems."

- Ken Batcher

- Scalable clusters need scalable storage
- Avoid centralized/single anything
- File/object storage API superior to blocks
- Reliability is important

52

Outline of the Day

Part 1

- Introduction
- Storage System Models
- Parallel File Systems
 - GPFS
 - PVFS
 - Panasas
 - Lustre

Part 2

- Benchmarking**
- MPI-IO
- Future Technologies

Performance Measurement

- Lots of different performance metrics
 - Sequential bandwidth, random I/Os, metadata operations
 - Single-threaded vs. multi-threaded
 - Single-client vs. multi-client
 - N-to-N (file per process) vs. N-to-1 (single shared file)
- Ultimately a method to try to estimate what you really care about
 - “Time to results”, aka “How long does my app take?”
- **Benchmarks are best if they model your real application**
 - Need to know what kind of I/O your app does in order to choose appropriate benchmark
 - Similar to CPU benchmarking – e.g., LINPACK performance may not predict how fast your codes run

What is a benchmark?

- Standardized way to compare performance of different systems
- Properties of a good benchmark
 - Relevant: captures essential attributes of real application workload
 - Simple: Provides an understandable metric
 - Portable & scalable
 - Consistent & repeatable results (on same HW)
 - Accepted by users & vendors
- Types of benchmark
 - Microbenchmark
 - Application-based benchmark
 - Synthetic workload

Microbenchmarks

- Measures one fundamental operation in isolation
 - Read throughput, write throughput, creates/sec, etc.
- Good for:
 - Tuning a specific operation
 - Post-install system validation
 - Publishing a big number in a press release
- Not as good for:
 - Modeling & predicting application performance
 - Measuring broad system performance characteristics
- Examples:
 - IOzone
 - IOR
 - Bonnie++
 - mdtest
 - metarates

Application Benchmarks

- Run real application on real data set, measure time
- Best predictor of application performance on your cluster
- Requires additional resources (compute nodes, etc.)
 - Difficult to acquire when evaluating new gear
 - Vendor may not have same resources as their customers
- Can be hard to isolate I/O vs. other parts of application
 - Performance may depend on compute node speed, memory size, interconnect, etc.
 - Difficult to compare runs on different clusters
- Time consuming – realistic job may run for days, weeks
- May require large or proprietary dataset
 - Hard to standardize and distribute

Synthetic Benchmarks

- Selected combination of operations (fractional mix)
 - Operations selected at random or using random model (e.g., Hidden Markov Model)
 - Operations and mix based on traces or sampling real workload
- Can provide better model for application performance
 - However, inherently domain-specific
 - Need different mixes for different applications & workloads
 - The more generic the benchmark, the less useful it is for predicting app performance
 - Difficult to model a combination of applications
- Examples:
 - SPEC SFS
 - TPC-C, TPC-D
 - FLASH I/O

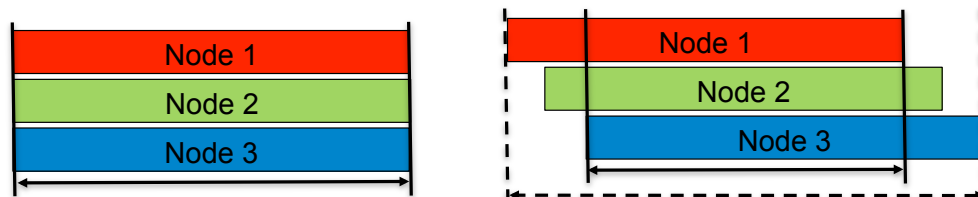
Benchmarks for HPC

- Unfortunately, there are few synthetic HPC benchmarks that stress I/O
- HPC Challenge
 - Seven sub-benchmarks, all “kernel” benchmarks (LINPACK, matrix transpose, FFT, message ping-pong, etc.)
 - Measures compute speed, memory bandwidth, cluster interconnect
 - No I/O measurements
- SPEC HPC2002
 - Three sub-benchmarks (CHEM, ENV, SEIS), all based on real apps
 - Only SEIS has a dataset of any size, and even it is tiny
 - 2 GB for Medium, 93 GB for X-Large
- NAS Parallel Benchmarks
 - Mix of kernel and mini-application benchmarks, all CFD-focused
 - One benchmark (BTIO) does significant I/O (135 GB N-to-1/collective write)
- FLASH I/O Benchmark
 - Simulates I/O performed by FLASH (nuclear/astrophysics application, Net-CDF/HDF5)
- Most HPC I/O benchmarking still done with microbenchmarks
 - IOzone, IOR (LLNL), LANL MPI-IO Test, mdtest, etc.

59

Benchmarking Pitfalls

- Not measuring what you think you are measuring
 - Most common with microbenchmarks
 - For example, measuring write or read from cache rather than to storage
 - Watch for “faster than the speed of light” results
- Multi-client benchmarks without synchronization across nodes
 - Measure aggregate throughput only when all nodes are transferring data
 - Application with I/O barrier may care more about when last node finishes



- Benchmark that does not model application workload
 - Different I/O size & pattern, different file size, etc.

60

Analyzing Results

- Sanity-checking results is important
- Figure out the “speed of light” in your system
- Large sequential accesses
 - Readahead can hide latency
 - 7200 RPM SATA 60-140 MB/sec/spindle
 - 15000 RPM FC 100-170 MB/sec/spindle
- Small random access
 - Seek + rotate limited
 - Readahead rarely helps (and sometimes hurts)
 - 7200 RPM SATA avg access 15 ms, 75-100 ops/sec/spindle
 - 15000 RPM FC avg access 6 ms, 150-200 ops/sec/spindle

Some Examples

- Let's look at some results from our example filesystems
- We'll look at:
 - File create/stat/delete (mdtest)
 - Read/write bandwidth (IOR)
- Results measured on:
 - Lustre: LLNL Thunder
 - PVFS: OSC Opteron cluster
 - GPFS: ASC Purple (LLNL)
 - Panasas: Pittsburgh Lab

The Fine Print

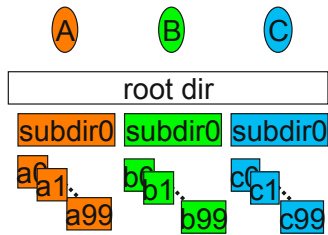
- Results are not directly comparable due to wildly different system specs
 - Lustre: 64 OSSs, 170 TB, ~765 disks
 - PVFS: 16 I/O nodes, 120 TB, ?? disks
 - GPFS: 125 I/O nodes, 300 TB, ~1,500 disks
 - Panasas: 30 OSDs, 30 TB, 60 disks
- Therefore, focus on trends instead of absolute numbers
 - Shape of curves
 - Differences between workloads on same filesystem

Metadata Performance

- Storage is more than reading & writing
- Metadata operations change the namespace or file attributes
 - Creating, opening, closing, and removing files
 - Creating, traversing, and removing directories
 - “Stat”ing files (obtaining the attributes of the file, such as permissions and file size)
- Several users exercise metadata subsystems:
 - Interactive use (e.g. “ls -l”)
 - File-per-process POSIX workloads
 - Collectively accessing files through MPI-IO (directly or indirectly)

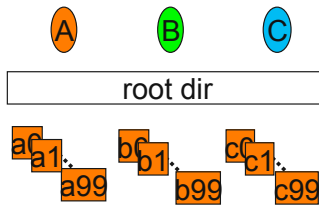
mdtest: Parallel Metadata Performance

Unique Directory



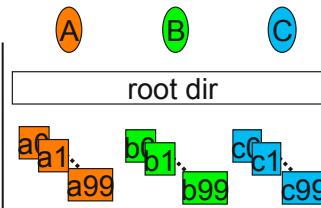
- 1) Each process (A, B, C) creates own subdir in root directory, then chdirs into it.
- 2) A, B, and C create, stat, and remove their own files in the unique subdirectories.

Single Process



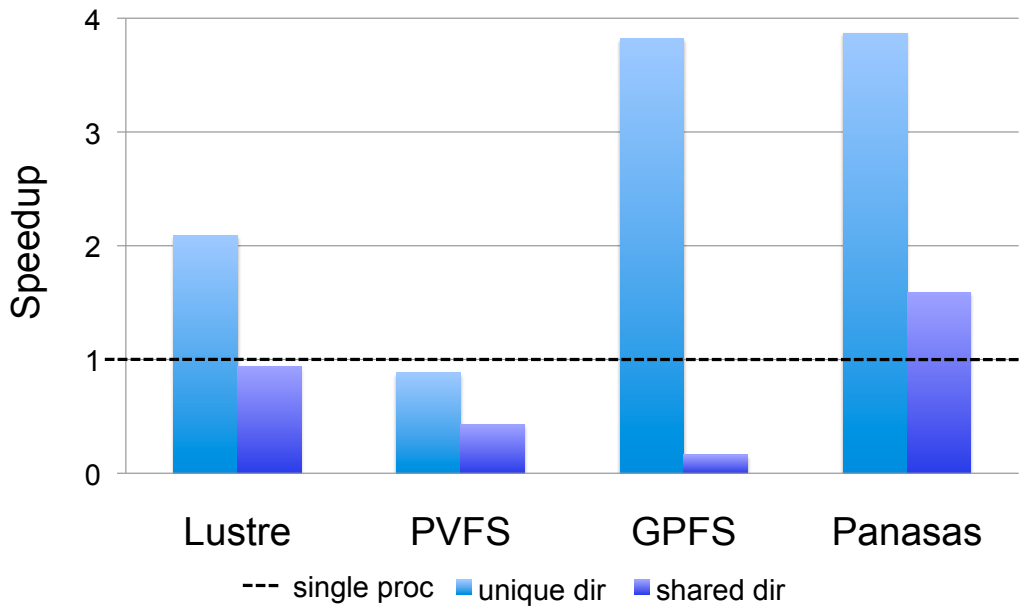
- 1) Process A creates files for all processes in root directory.
- 2) Processes A, B, and C open, stat, and close their own files.
- 3) Process A removes files for all processes.

Shared Directory

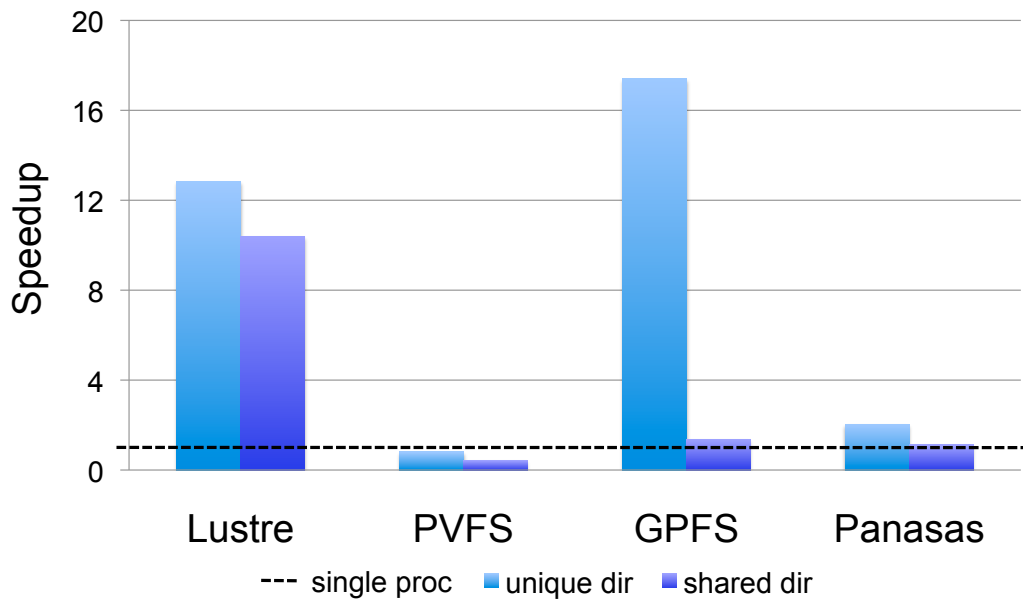


- 1) Each process (A, B, C) creates, stats, and removes its own files in the root directory.

mdtest: Create File



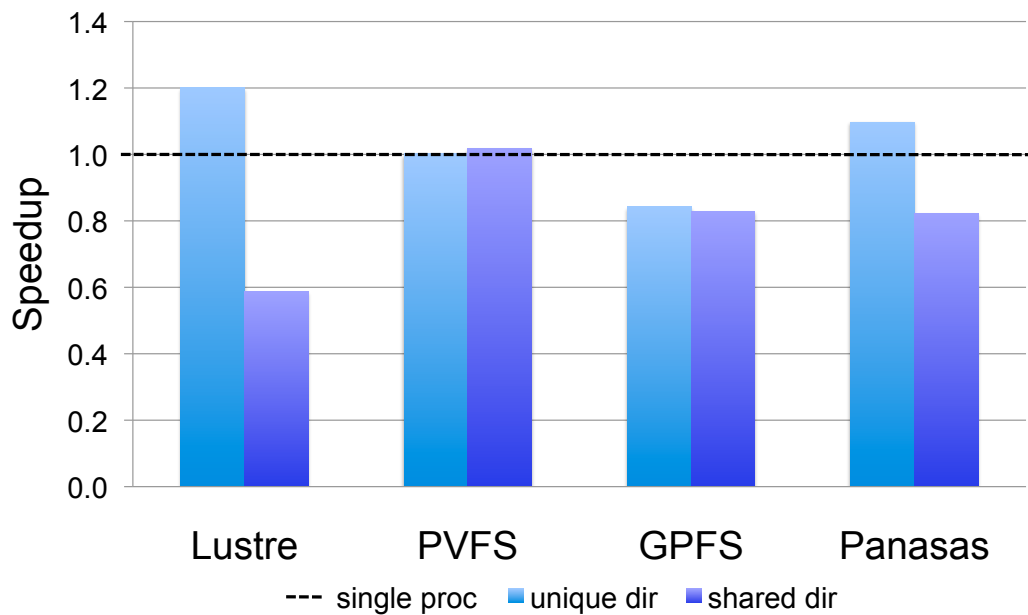
mdtest: Remove File



67

panasas

mdtest: Stat File



68

panasas

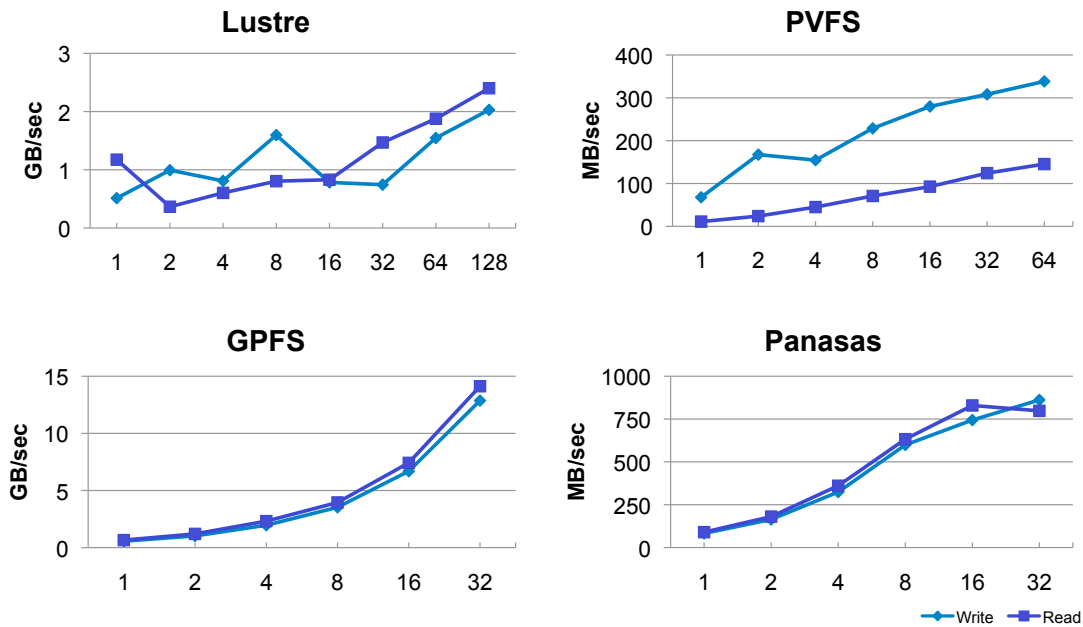
mdtest Analysis

- Relative analysis can hide some issues
- PVFS
 - Different dirs managed by different storage nodes
 - No client caching limits stat performance
- GPFS
 - Very high cost to operating in the same directory due to token interactions
- Lustre
 - Single MDS, directory lock limits shared dir case
- Panasas
 - Coarse-grained metadata clustering not active
 - Directory lock on metadata server limits parallelism

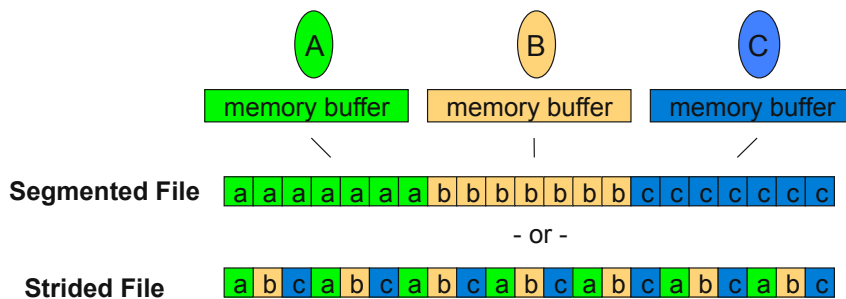
IOR: File System Bandwidth

- Written at LLNL
- Named for the acronym 'interleaved or random'
- POSIX, MPI-IO, HDF5, Parallel-NetCDF APIs
 - Shared or independent file access
 - Collective or independent I/O (when available)
 - Configurable interleaving patterns
- Employs MPI for process synchronization
- We show POSIX API results

IOR: File per proc, read vs. write

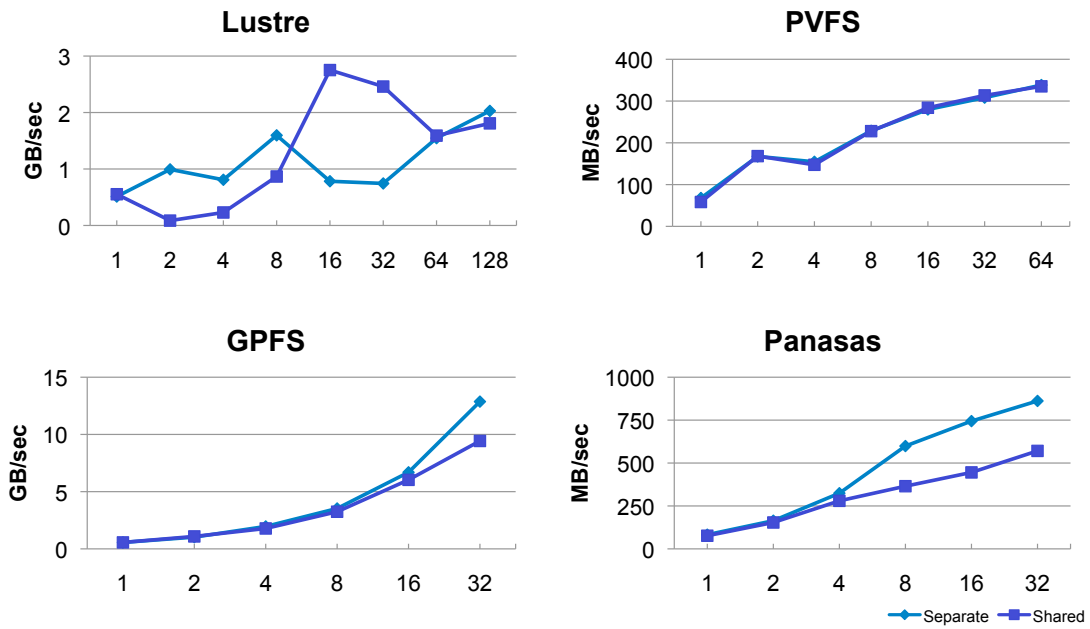


IOR Access Patterns for Shared Files



- Primary distinction between the two major shared-file patterns is whether each task's data is contiguous or noncontiguous
- Contiguous (segmented) easier for the filesystem
- Noncontiguous (strided) often easier for application
- We show segmented pattern results

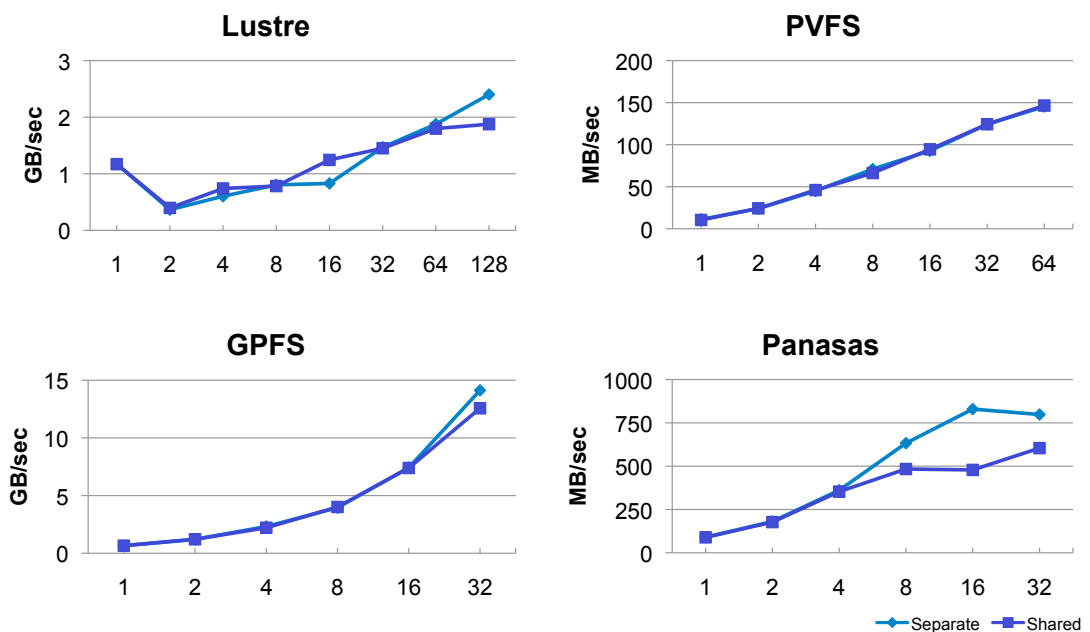
IOR: Write, file per proc vs. shared file



73



IOR: Read, file per proc vs. shared file



74



IOR Analysis

- Shared file is usually slower than file-per-proc
 - However, shared file is often much more convenient for app developers
 - File-per-proc also unwieldy at large node counts
- Segmented access pattern is a compromise between FS and app
 - Maybe makes neither one happy
- Much of slowdown comes from POSIX API
 - MPI-IO coming up...

Outline of the Day

Part 1

- Introduction
- Storage System Models
- Parallel File Systems
 - GPFS
 - PVFS
 - Panasas
 - Lustre

Part 2

- Benchmarking
- MPI-IO**
- Future Technologies

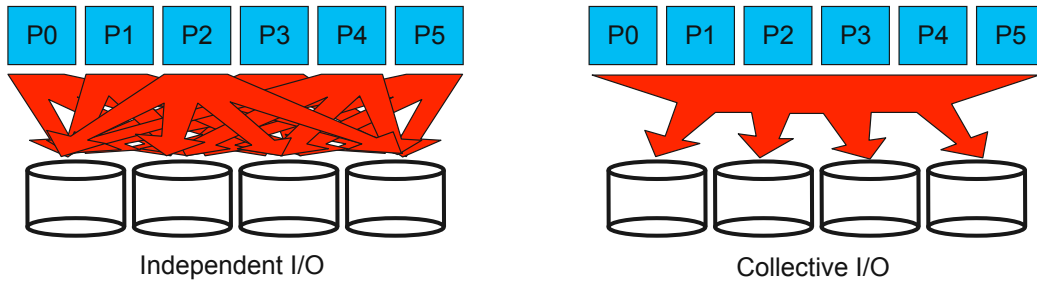
What's wrong with POSIX?

- It's a useful, ubiquitous interface for basic I/O
- It lacks constructs useful for parallel I/O
 - Cluster application is really one program running on N nodes, but looks like N programs to the filesystem
 - No support for noncontiguous I/O
 - No hinting/prefetching
- Its rules hurt performance for parallel apps
 - Atomic writes, read-after-write consistency
 - Attribute freshness
- POSIX should not be used (directly) in parallel applications that want good performance
 - But developers use it anyway

MPI-IO

- I/O interface [specification](#) for use in MPI apps
- Data model is same as POSIX
 - Stream of bytes in a file
- Features:
 - Collective I/O
 - Noncontiguous I/O with MPI datatypes and file views
 - Nonblocking I/O
 - Fortran bindings (and additional languages)
 - System for encoding files in a portable format (external32)
 - Not self-describing - just a well-defined encoding of types
- Implementations available on most platforms (more later)

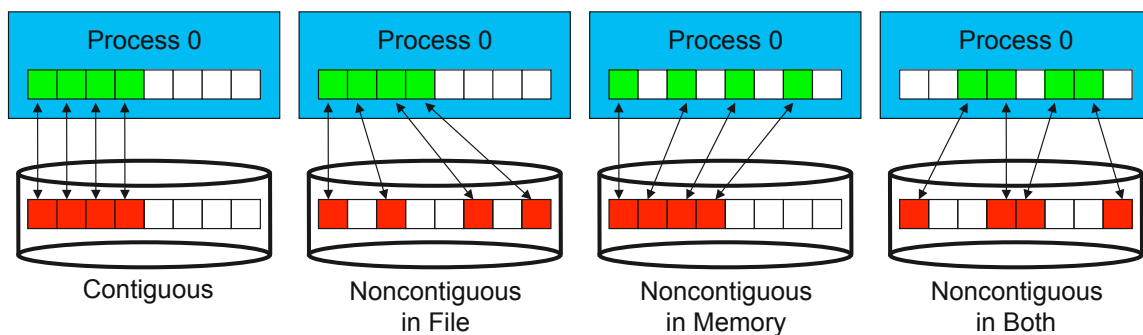
Independent and Collective I/O



- **Independent I/O** operations specify only what a single process will do
 - Independent I/O calls do not pass on relationships between I/O on other processes
- Many applications have phases of computation and I/O
 - During I/O phases, all processes read/write data
 - We can say they are **collectively** accessing storage
- **Collective I/O** is coordinated access to storage by a group of processes
 - Collective I/O functions are called by all processes participating in I/O
 - **Allows I/O layers to know more about access as a whole, more opportunities for optimization in lower software layers, better performance**

79

Contiguous and Noncontiguous I/O



- **Contiguous I/O** moves data from a single memory block into a single file region
- **Noncontiguous I/O** has three forms:
 - Noncontiguous in memory, noncontiguous in file, or noncontiguous in both
- Structured data leads naturally to noncontiguous I/O (e.g. block decomposition)
- Describing noncontiguous accesses with a single operation passes more knowledge to I/O system

80

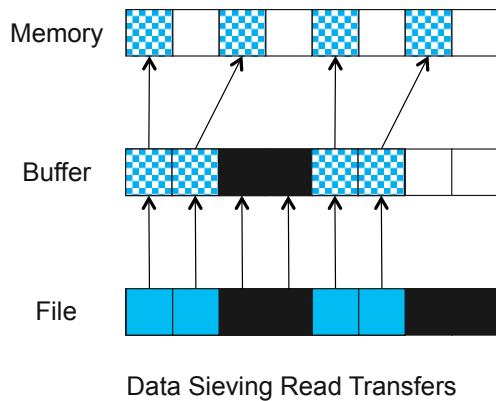
Nonblocking and Asynchronous I/O

- Blocking/synchronous I/O operations return when buffer may be reused
 - Data in system buffers or on disk
- Some applications like to overlap I/O and computation
 - Hiding writes, prefetching, pipelining
- A **nonblocking** interface allows for submitting I/O operations and testing for completion later
- If the system also supports **asynchronous I/O**, progress on operations can occur in the background
 - Depends on implementation
- Otherwise progress is made at start, test, wait calls

Under the Covers of MPI-IO

- MPI-IO implementation gets a lot of information
 - Collection of processes reading data
 - Structured description of the regions
- Implementation has some options for how to perform the data reads
 - Noncontiguous data access optimizations
 - Collective I/O optimizations

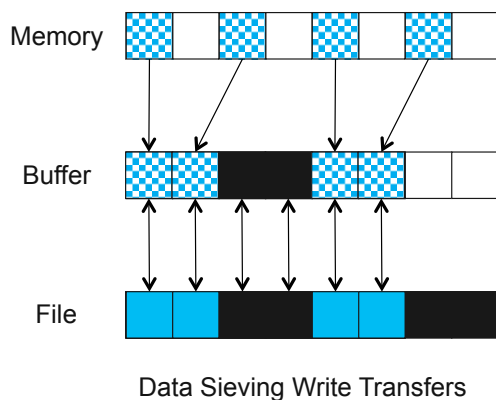
Noncontiguous I/O: Data Sieving



- Data sieving is used to combine lots of small accesses into a single larger one
 - Remote file systems (parallel or not) tend to have high latencies
 - Reducing # of operations important
- Similar to how a block-based file system interacts with storage
- Generally very effective, but not as good as having a PFS that supports noncontiguous access

83

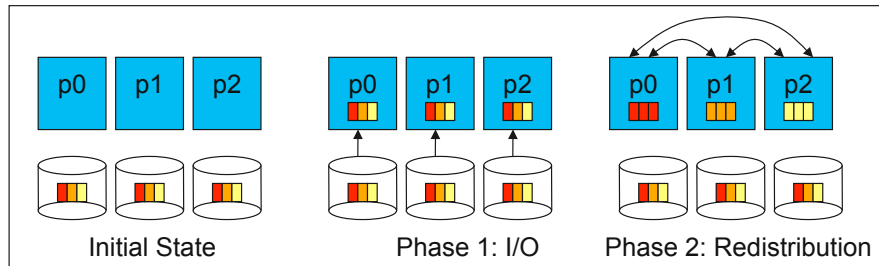
Data Sieving Write Operations



- Data sieving for writes is more complicated
 - Must read the entire region first
 - Then make changes in buffer
 - Then write the block back
- Requires locking in the file system
 - Can result in false sharing (interleaved access)
- PFS supporting noncontiguous writes is preferred

84

Collective I/O and Two-Phase I/O

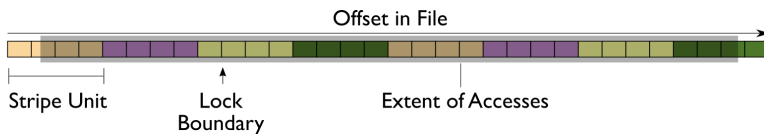


Two-Phase Read Algorithm

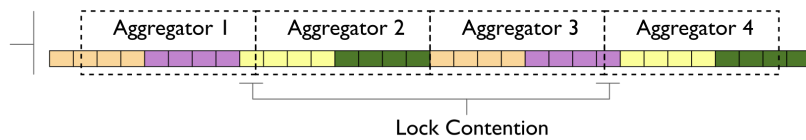
- Problems with independent, noncontiguous access
 - Lots of small accesses
 - Independent data sieving reads lots of extra data, can exhibit false sharing
- Idea: Reorganize access to match layout on disks
 - Single processes use data sieving to get data for many
 - Often reduces total I/O through sharing of common blocks
- Second “phase” redistributes data to final destinations
- Two-phase writes operate in reverse (redistribute then I/O)
 - Typically read/modify/write (like data sieving)
 - Overhead is lower than independent access because there is little or no false sharing
- Note that two-phase is usually applied to file regions, not to actual blocks

Two-Phase I/O Algorithms

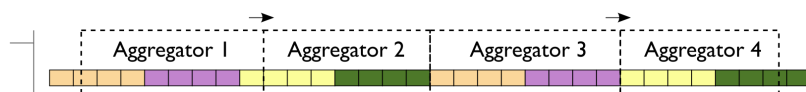
Imagine a collective I/O access using four aggregators to a file striped over four file servers (indicated by colors):



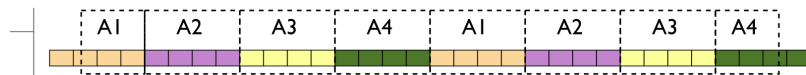
One approach is to evenly divide the region accessed across aggregators.



Aligning regions with lock boundaries eliminates lock contention.



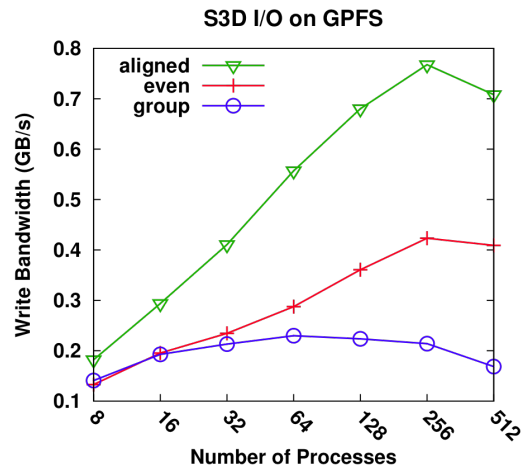
Mapping aggregators to servers reduces the number of concurrent operations on a single server and can be helpful when locks are handed out on a per-server basis (e.g., Lustre).



For more information, see W.K. Liao and A. Choudhary, “Dynamically Adapting File Domain Partitioning Methods for Collective I/O Based on Underlying Parallel File System Locking Protocols,” SC2008, November, 2008.

Impact of Two-Phase I/O Algorithms

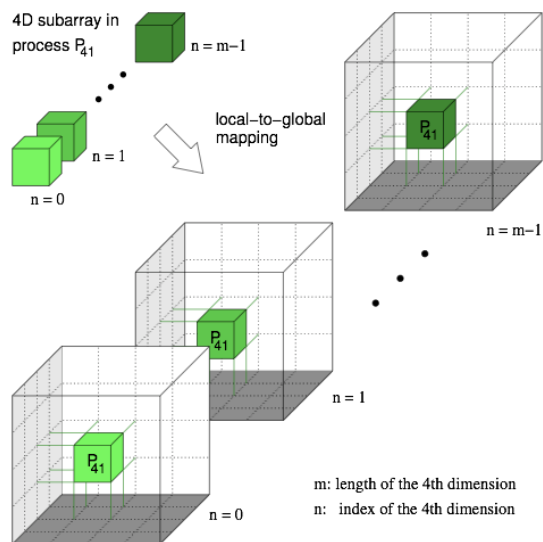
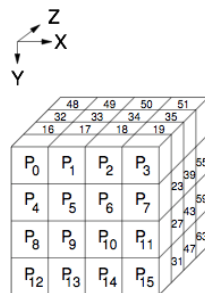
- This graph shows the performance for the S3D combustion code, writing to a single file.
- Aligning with lock boundaries doubles performance over default “even” algorithm.
- “Group” algorithm similar to server-aligned algorithm on last slide.
- Testing on Mercury, an IBM IA64 system at NCSA, with 54 servers and 512KB stripe size.



W.K. Liao and A. Choudhary, “Dynamically Adapting File Domain Partitioning Methods for Collective I/O Based on Underlying Parallel File System Locking Protocols,” SC2008, November, 2008.

S3D Turbulent Combustion Code

- S3D is a turbulent combustion application using a direct numerical simulation solver from Sandia National Laboratory
- Checkpoints consist of four global arrays
 - 2 3-dimensional
 - 2 4-dimensional
 - 50x50x50 fixed subarrays



Thanks to Jackie Chen (SNL), Ray Grout (SNL), and Wei-Keng Liao (NWU) for providing the S3D I/O benchmark, Wei-Keng Liao for providing this diagram.

Impact of Optimizations on S3D I/O

- Testing with PnetCDF output to single file, three configurations, 16 processes
 - All MPI-IO optimizations (collective buffering and data sieving) disabled
 - Independent I/O optimization (data sieving) enabled
 - Collective I/O optimization (collective buffering, a.k.a. two-phase I/O) enabled

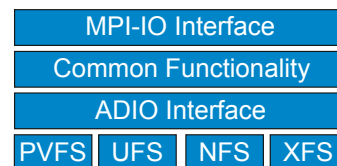
	Coll. Buffering and Data Sieving Disabled	Data Sieving Enabled	Coll. Buffering Enabled (incl. Aggregation)
POSIX writes	102,401	81	5
POSIX reads	0	80	0
MPI-IO writes	64	64	64
Unaligned in file	102,399	80	4
Total written (MB)	6.25	87.11	6.25
Runtime (sec)	1443	11	6.0
Avg. MPI-IO time per proc (sec)	1426.47	4.82	0.60

89



MPI-IO Implementations

- Different MPI-IO implementations exist
- Three better-known ones are:
 - ROMIO from Argonne National Laboratory
 - Leverages MPI-1 communication
 - Supports local file systems, network file systems, parallel file systems
 - UFS module works GPFS, Lustre, and others
 - Includes data sieving and two-phase optimizations
 - MPI-IO/GPFS from IBM (for AIX only)
 - Includes two special optimizations
 - **Data shipping** -- mechanism for coordinating access to a file to alleviate lock contention (type of aggregation)
 - **Controlled prefetching** -- using MPI file views and access patterns to predict regions to be accessed in future
 - MPI from NEC
 - For NEC SX platform and PC clusters with Myrinet, Quadrics, IB, or TCP/IP
 - Includes listless I/O optimization -- fast handling of noncontiguous I/O accesses in MPI layer



ROMIO's layered architecture.

90



MPI-IO Wrap-Up

- MPI-IO provides a rich interface allowing us to describe
 - Noncontiguous accesses in memory, file, or both
 - Collective I/O
- This allows implementations to perform many transformations that result in better I/O performance
- Also forms solid basis for high-level I/O libraries
 - But they must take advantage of these features!

Outline of the Day

Part 1

- Introduction
- Storage System Models
- Parallel File Systems
 - GPFS
 - PVFS
 - Panasas
 - Lustre

Part 2

- Benchmarking
- MPI-IO
- Future Technologies**

Storage Futures

■ pNFS

- An extension to the NFSv4 file system protocol standard that allows direct, parallel I/O between clients and storage devices
- Eliminates the scaling bottleneck found in today's NAS systems
- Supports multiple types of back-end storage systems, including traditional block storage, other file servers, and object storage systems

■ FLASH and other non-volatile devices

- New level in storage hierarchy

Why a Standard for Parallel I/O?

■ NFS is the only network file system standard

- Proprietary file systems have unique advantages, but can cause lock-in

■ NFS widens the playing field

- Panasas, IBM, EMC want to bring their experience in large scale, high-performance file systems into the NFS community
- Sun and NetApp want a standard HPC solution
- Broader market benefits vendors
- More competition benefits customers

■ What about open source

- NFSv4 Linux client is very important for NFSv4 adoption, and therefore pNFS
- Still need vendors that are willing to do the heavy lifting required in quality assurance for mission critical storage

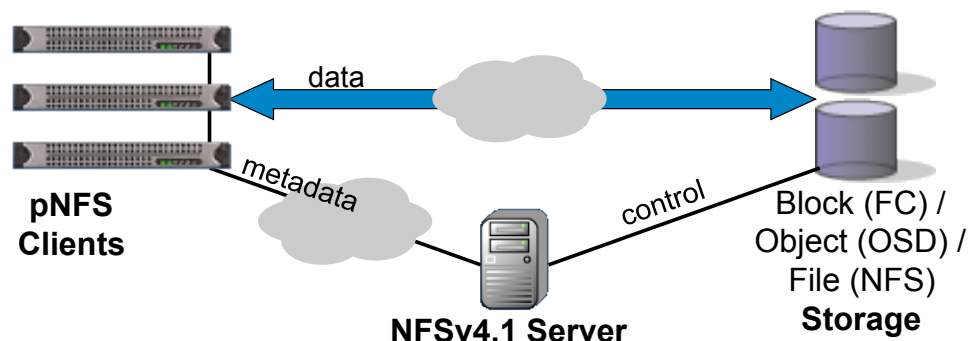
NFSv4 and pNFS

- NFS created in '80s to share data among engineering workstations
- NFSv3 widely deployed
- NFSv4 several years in the making, lots of new stuff
 - Integrated Kerberos (or PKI) user authentication
 - Integrated File Locking and Open Delegations (stateful server!)
 - ACLs (hybrid of Windows and POSIX models)
 - Official path to add (optional) extensions
- NFSv4.1 adds even more
 - pNFS for parallel IO
 - Directory Delegations for efficiency
 - RPC Sessions for robustness, better RDMA support

95

pNFS: Standard Storage Clusters

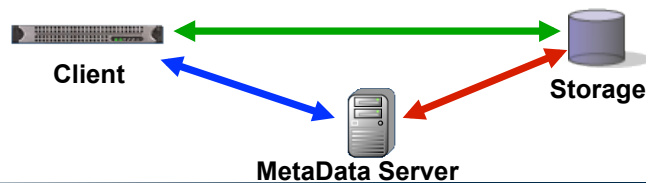
- pNFS is an extension to the Network File System v4 protocol standard
- Allows for parallel and direct access
 - From Parallel Network File System clients
 - To Storage Devices over multiple storage protocols
 - Moves the NFS (metadata) server out of the data path



96

The pNFS Standard

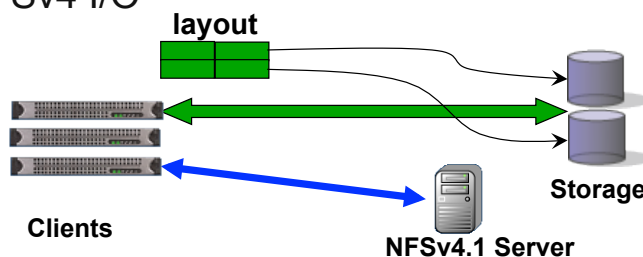
- The **pNFS** standard defines the NFSv4.1 protocol extensions between the **server and client**
- The **I/O** protocol between the **client and storage** is specified elsewhere, for example:
 - SCSI **Block** Commands (**SBC**) over Fibre Channel (**FC**)
 - SCSI **Object**-based Storage Device (**OSD**) over iSCSI
 - Network **File** System (**NFS**)
- The **control** protocol between the **server and storage** devices is also specified elsewhere, for example:
 - SCSI **Object**-based Storage Device (**OSD**) over iSCSI



97

pNFS Layouts

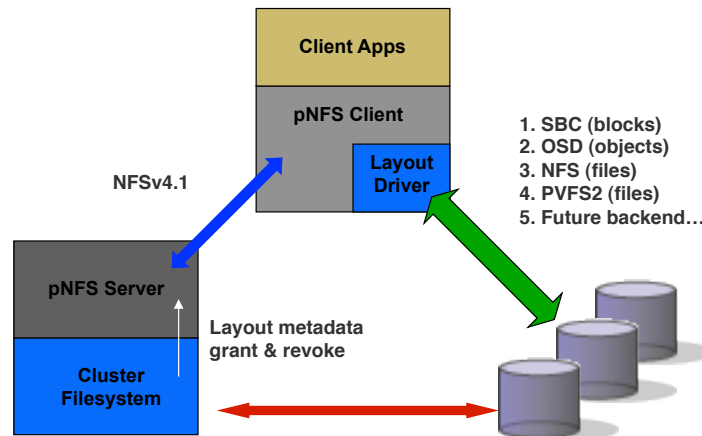
- Client gets a *layout* from the NFS Server
- The layout maps the file onto storage devices and addresses
- The client uses the layout to perform direct I/O to storage
- At any time the server can recall the layout
- Client commits changes and returns the layout when it's done
- pNFS is optional, the client can always use regular NFSv4 I/O



98

pNFS Client

- Common client for different storage back ends
- Wider availability across operating systems
- Fewer support issues for storage vendors



99

panasas

pNFS is not...

- Improved cache consistency
 - NFS has open-to-close consistency enforced by client polling of attributes
 - NFSv4.1 directory delegations can reduce polling overhead
- Perfect POSIX semantics in a distributed file system
 - NFS semantics are good enough (or, all we'll give you)
 - But note also the POSIX High End Computing Extensions Working Group
 - <http://www.opengroup.org/platform/hecewg/>
- Clustered metadata
 - Not a server-to-server protocol for scaling metadata
 - But, it doesn't preclude such a mechanism

100

panasas

Is pNFS Enough?

- Standard for out-of-band metadata
 - Great start to avoid classic server bottleneck
 - NFS has already relaxed some semantics to favor performance
 - But there are certainly some workloads that will still hurt
- Standard framework for clients of different storage backends
 - Files
 - Objects
 - Blocks
 - PVFS
 - Your project... (e.g., dcache.org)

Key pNFS Participants



- Panasas (Objects)
- ORNL and ESSC/DoD funding Linux pNFS development
- Network Appliance (Files over NFSv4)
- IBM (Files, based on GPFS)
- EMC (Blocks, HighRoad MPFSi)
- Sun (Files over NFSv4)
- U of Michigan/CITI (Linux maintainers, EMC and Microsoft contracts)
- DESY – Java-based implementation

pNFS Standard Status

- IETF approved Internet Drafts in December 2008
- RFCs for NFSv4.1, pNFS-objects, and pNFS-blocks published January 2010
 - RFC 5661 - Network File System (NFS) Version 4 Minor Version 1 Protocol
 - RFC 5662 - Network File System (NFS) Version 4 Minor Version 1 External Data Representation Standard (XDR) Description
 - RFC 5663 - Parallel NFS (pNFS) Block/Volume Layout
 - RFC 5664 - Object-Based Parallel NFS (pNFS) Operations

pNFS Implementation Status

- Implementation interoperability continues
 - San Jose Connect-a-thon March '06, February '07, May '08, June '09
 - Ann Arbor NFS Bake-a-thon September '06, October '07
 - Dallas pNFS inter-op, June '07, Austin February '08, Sept '08, Oct '09
- Server vendors waiting for Linux client
 - Sun, NetApp, EMC, IBM, Panasas, ...
 - 2.6.30
 - exofs object storage file system (local) and iSCSI/OSDv2
 - 2.6.31 has been released.
 - most of nfsv4.1: sessions, 4.1 as an option, no pnfs yet
 - server back channel is absent
 - 2.6.32 in stabilization, finish nfsv4.1 including server callbacks
 - Adds server back-channel support.
 - Goal to complete patch adoption by Q3 2010

How to use pNFS today

- Up-to-date GIT tree from Linux developers
 - bhalevy@panasas.com manages the source trees
- Red Hat/Fedora RPMs that include pNFS
 - steved@redhat.com builds experimental packages
- pNFS mailing list, pnfs@linux-nfs.org
- <http://open-osd.org>
 - Useful to get to OSD target, the user level program
 - Exofs uses kernel initiator, need the target

The problem with rotating media

- Areal density increases by 40% per year
 - Per drive capacity increases by 50% to 100% per year
 - 2008: **1 TB**
 - 2009: **2 TB**
 - 2010: **3 TB** (maybe 4 TB late in the year)
 - Drive vendors prepared to continue like this for years to come
- Drive interface speed increases by 15-20% per year
 - 2008: 500 GB disk (WD RE2): **98 MB/sec**
 - 2009: 1 TB disk (WD RE3): **113 MB/sec** (+15%)
 - 2010: 2 TB disk (WD RE4): **138 MB/sec** (+22%)
- Takes longer and longer to completely read each new generation of drive
- Seek times and rotational speeds not increasing
 - 15,000 RPM and 2.5 ms/sec still the norm for high end

FLASH is...

■ Non-volatile

- Each bit is stored in a "floating gate" that holds value without power
- Electrons can leak, so shelf life and write count is limited

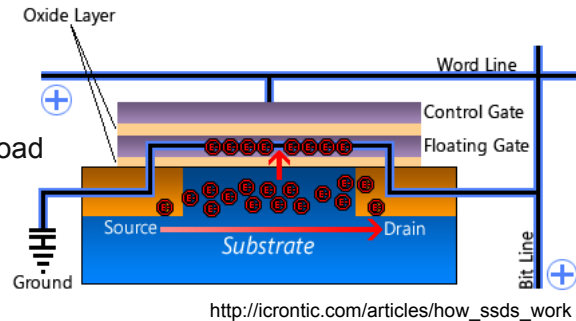
■ Page-oriented

- Read, write, and erase operations apply to large chunks
- Smaller (e.g., 4K) read/write block based on addressing logic
- Larger (e.g., 256K) erase block to amortize the time it takes to erase

■ Medium speed

- Slower than DRAM
- Faster than disks for reading
- Write speed dependent on workload

■ Relatively cheap



http://icrantic.com/articles/how_ssds_work

107

panasas

FLASH Reliability

■ SLC – Single Level Cell

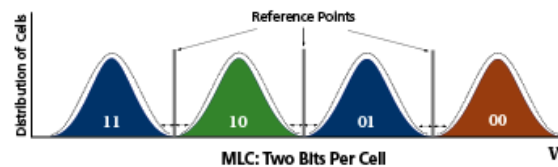
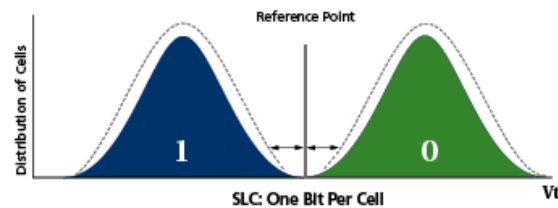
- One threshold, one bit
- 10^5 to 10^6 write cycles per page

■ MLC – Multi Level Cell

- Multiple thresholds, multiple bits (2 bits now, 3 & 4 soon)
- N bits requires 2^N Vt levels
- 10^4 write cycles per page
- Denser and cheaper, but slower and less reliable

■ Wear leveling is critical

- Pre-erase blocks before writing is required
- Page map indirection allows shuffling of pages to do wear leveling



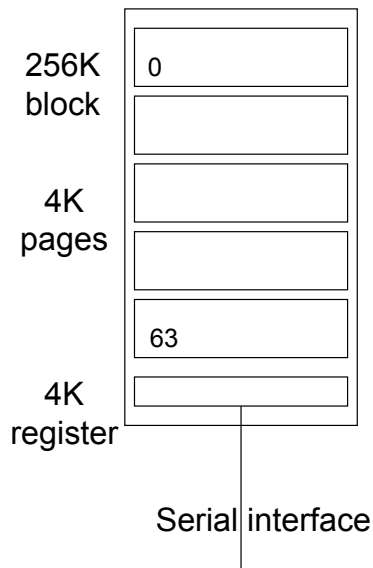
<http://www.micron.com/handcom/>

108

panasas

FLASH Speeds

Samsung 4 GB Device



100 usec	Transfer 4K over serial interface	40 MB/sec
25 usec	Load 4K register from Flash	160 MB/sec
125 usec	Read latency	32 MB/sec
200 usec	Store 4K register to FLASH	20 MB/sec
225 usec	Write latency	16 MB/sec
1.5 msec	Erase 256K block	170 MB/sec
1.725 msec	Worse case write	2.3 MB/sec

- Write performance heavily dependent on workload and wear leveling algorithms
- Writes are slower with less free space

109

panasas 

FLASH in the Storage Hierarchy

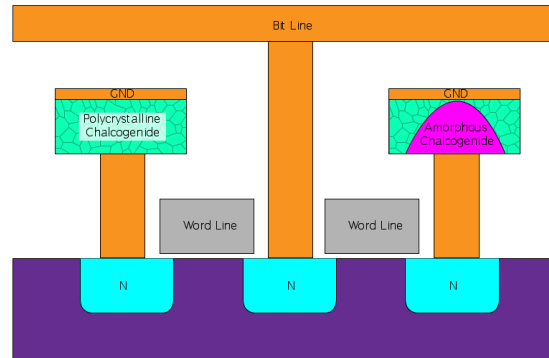
- On the compute nodes
 - High reliability local storage for OS partition
 - Local cache for memory checkpoints?
 - Device write speeds vary widely
 - 4 MB/sec for a cheap USB
 - 80 or 100 MB/sec for MTron or Zeus
 - 600 MB/sec for Fusion-io ioDrive
 - One Fusion-io board could double cost of node
- On the storage server
 - Metadata storage
 - Low latency log device
 - Replacement for NVRAM? Probably not enough write bandwidth to absorb all the write data

110

panasas 

Phase Change Memory

- Newer technology that may replace FLASH
- Based on state change instead of stored electrons
 - Crystalline vs. amorphous
 - Think ice vs. water
 - Except it's Germanium-Antimony-Tellurium (GST) chalcogenide glass
- Change state by heating to 650°C and then cooling
 - Cool quickly ⇒ amorphous
 - Cool slowly ⇒ crystalline
- Samsung, Numonyx shipping devices now
- Maybe DRAM replacement in 2015?



Courtesy <http://en.wikipedia.org/wiki/User:Cyferz>

111

panasas

Wrapping Up

- We've covered a lot of ground in a short time
 - Disk drives & filesystems
 - Benchmarking
 - Programming middleware
 - pNFS and FLASH
- There is no magic in high performance I/O
 - Under the covers it looks a lot like shared memory or message passing
 - Knowing how things work will lead you to better performance
- Things will continue to get more complicated, but hopefully easier too!
 - Remote access to data
 - More layers to I/O stack
 - Domain-specific application interfaces

112

panasas

Thank you!

Joel Jacobson, Marc Unangst
{jjacobson,mju}@panasas.com
Panasas, Inc.

Printed References

- John May, Parallel I/O for High Performance Computing, Morgan Kaufmann, October 9, 2000.
 - Good coverage of basic concepts, some MPI-IO, HDF5, and serial netCDF
- William Gropp, Ewing Lusk, and Rajeev Thakur, Using MPI-2: Advanced Features of the Message Passing Interface, MIT Press, November 26, 1999.
 - In-depth coverage of MPI-IO API, including a very detailed description of the MPI-IO consistency semantics

Online References: Filesystems

- ROMIO MPI-IO
 - <http://www.mcs.anl.gov/romio/>
- POSIX I/O Extensions
 - <http://www.opengroup.org/platform/hecewg/>
- PVFS
 - <http://www.pvfs.org/>
- Panasas
 - <http://www.panasas.com/>
- Lustre
 - <http://www.lustre.org/>
- GPFS
 - http://www.almaden.ibm.com/storagesystems/file_systems/GPFS/

Online References: Benchmarks

- LLNL I/O tests (IOR, fdtree, mdtest)
 - <http://www.llnl.gov/icc/lc/siop/downloads/download.html>
- Parallel I/O Benchmarking Consortium (noncontig, mpi-tile-io, mpi-md-test)
 - <http://www.mcs.anl.gov/pio-benchmark/>
- FLASH I/O benchmark
 - <http://www.mcs.anl.gov/pio-benchmark/>
 - http://flash.uchicago.edu/~jbgallag/io_bench/ (original version)
- b_eff_io test
 - http://www.hlr.de/organization/par/services/models/mpi/b_eff_io/
- mpiBLAST
 - <http://www.mpiblast.org>
- HPC Challenge
 - <http://icl.cs.utk.edu/hpcc/>
- SPEC HPC2002
 - <http://www.spec.org/hpc2002/>
- NAS Parallel Benchmarks
 - <http://www.nas.nasa.gov/Resources/Software/npb.html>

Online References: pNFS

■ NFS Version 4.1

- RFC 5661 - Network File System (NFS) Version 4 Minor Version 1 Protocol
- RFC 5662 - Network File System (NFS) Version 4 Minor Version 1 External Data Representation Standard (XDR) Description
- RFC 5663 - Parallel NFS (pNFS) Block/Volume Layout
- RFC 5664 - Object-Based Parallel NFS (pNFS) Operations
- <http://tools.ietf.org/html/>

■ pNFS Problem Statement

- Garth Gibson (Panasas), Peter Corbett (Netapp), Internet-draft, July 2004
- <http://www.pdl.cmu.edu/pNFS/archive/gibson-pnfs-problem-statement.html>

■ Linux pNFS Kernel Development

- <http://www.citi.umich.edu/projects/asci/pnfs/linux>

How to use pNFS today

■ Benny's git tree:

`git://linux-nfs.org/~bhalevy/linux-pnfs.git`

■ The kernel rpms can be found at:

<http://fedorapeople.org/~steved/repos/pnfs/i686>

http://fedorapeople.org/~steved/repos/pnfs/x86_64

■ The source rpm can be found at:

<http://fedorapeople.org/~steved/repos/pnfs/source/>

■ Bug database

<https://bugzilla.linux-nfs.org/index.cgi>

■ OSD target

<http://open-osd.org/>

Benchmarking References

119

panasas 

PVFS Test Platform: OSC Opteron Cluster

- 338 nodes, each with
 - 4 AMD Opteron CPUs at 2.6 GHz, 8 GB memory
- Gigabit Ethernet network
 - Switch Hierarchy with multiple GBit uplinks
- 16 I/O servers (also serving metadata)
 - 2 2-core Xeon CPU at 2.4 GHz, 3 GB memory
- 120 TB parallel file system
 - Each server has Fibre Channel interconnect to back-end RAID



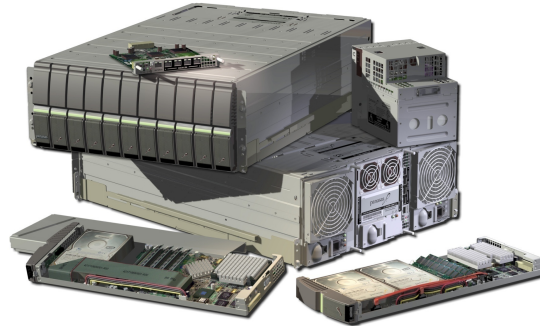
Ohio Supercomputer Center

120

panasas 

Panasas Test Platform: Pittsburgh Lab

- Small test system from our Pittsburgh development lab
- 3 Panasas Shelves, each with
 - 10 SB-1000a-XC StorageBlades
 - (1.5GHz Celeron, 2GB RAM, 2x500GB SATA, 1GE)
 - 1 DB-100a DirectorBlade
 - (1.8GHz 475, 4GB RAM, 1GE)
 - 18-port switch with 10GE uplink
- 48 client nodes
 - 2.8 GHz Xeon, 8GB, 1GE
- GE Backbone
 - 40 GB/s between clients and shelves



121

panasas

GPFS Test Platform: ASC Purple

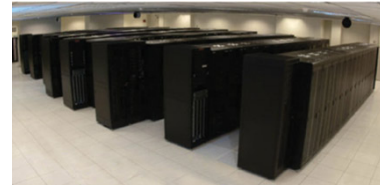
- 1536 nodes, each with
 - 8 64-bit Power5 CPUs at 1.9 GHz
 - 32 GB memory
- Federation high-speed interconnect
 - 4Gbyte/sec theoretical bisection bandwidth per adapter
 - ~5.5 Gbyte/sec measured per I/O server w/dual adapters
- 125 I/O servers, 3 metadata servers
 - 8 64-bit Power5 CPUs at 1.9 GHz
 - 32 GB memory
- 300 TB parallel file system
 - HW RAID5 (4+P, 250 GB SATA Drives)
 - 24 RAID5 per I/O server



122

panasas

Lustre Test Platform: LLNL Thunder



- 1024 nodes each with
 - 4 64-bit Itanium2 CPUs at 1.4 GHz
 - 8 GB memory
- Quadrics high-speed interconnect
 - ~900 MB/s of bidirectional bandwidth
 - 16 Gateway nodes with 4 GigE connections to the Lustre network
- 64 object storage servers, 1 metadata server
 - I/O server - dual 2.4 Ghz Xeons, 2GBs ram
 - Metadata Server - dual 3.2 Ghz Xeons, 4 GBs ram
- 170 TB parallel file system
 - HW RAID5 (8+P, 250 GB SATA Drives)
 - 108 RAIDs per rack
 - 8 racks of data disk

123



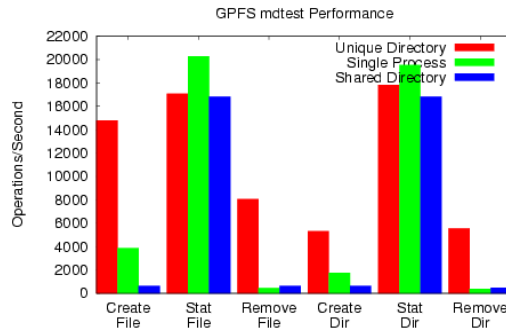
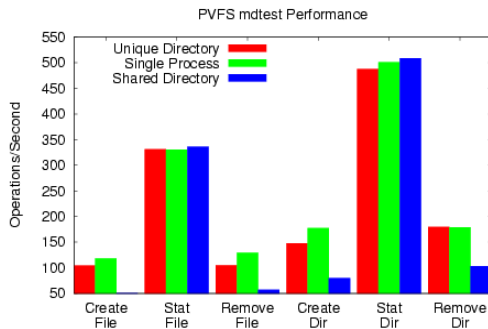
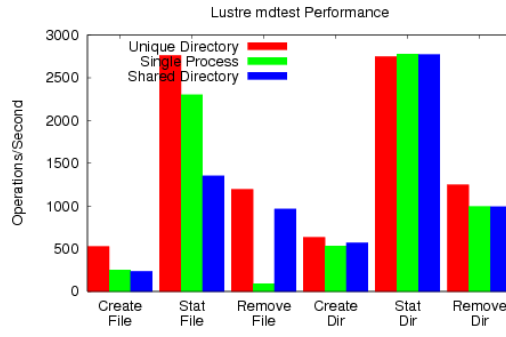
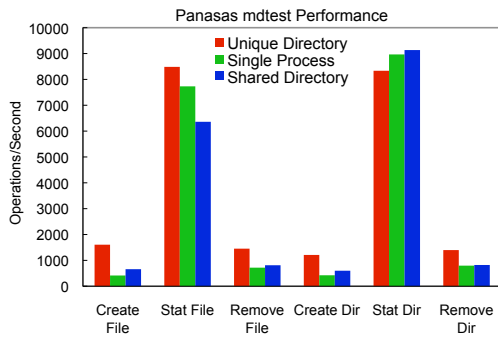
mdtest Details

- We ran three variations, each with 64 processes:
 - mdtest -d \$DIR -n 100 -i 3 -N 1 -v -u
 - Each task creates 100 files in a unique subdirectory
 - mdtest -d \$DIR -n 100 -i 3 -N 1 -v -c
 - One task creates 6400 files in one directory
 - Each task opens, removes its own
 - mdtest -d \$DIR -n 100 -i 3 -N 1 -v
 - Each task creates 100 files in a single shared directory
- GPFS tests use 16 tasks with 4 tasks on each node
- Panasas tests use 48 tasks on 48 nodes

124



mdtest Results



IOR: File System Bandwidth

■ IOR parameters used:

- POSIX API
- Segmented (contiguous) I/O mode
- 4 MB transfer size
- 4 GB segment size per task

■ Command lines:

- Single shared file:
 - `./IOR -a POSIX -C -i 3 -t 4M -b 4G -e -v -v -o $FILE`
- File per process:
 - `./IOR -a POSIX -C -i 3 -t 4M -b 4G -e -v -v -F -o $FILE`

IOR POSIX Segmented Results

