

# Addressing Scalable I/O Challenges for Exascale

Approved for Public Release SAND2011-3588C

27<sup>th</sup> IEEE Symposium on Massive Storage  
Systems and Technologies

May 24, 2011

*Ron Oldfield*  
*Sandia National Laboratories*



# I/O Challenges at Exascale

---

- Power is the largest hurdle for exascale computing
  - It drives nearly every aspect of design
  - Data movement is a big problem
    - Memory, network, storage... each layer costs ~10x in power
- Current usage model and programming models are inadequate
  - Checkpoints are a **HUGE** concern
  - App workflow uses storage as a communication conduit
    - Simulate, store, analyze, store, refine, store, ... most of the data is transient
- Current file systems do not handle scale or handle faults very well
  - Expect millions of clients, faults are the norm
- We are attacking the problem from two directions
  - Approaches that reduce I/O (smarter resilience, integrated analysis)
  - File systems designed for extreme scale that expect faults

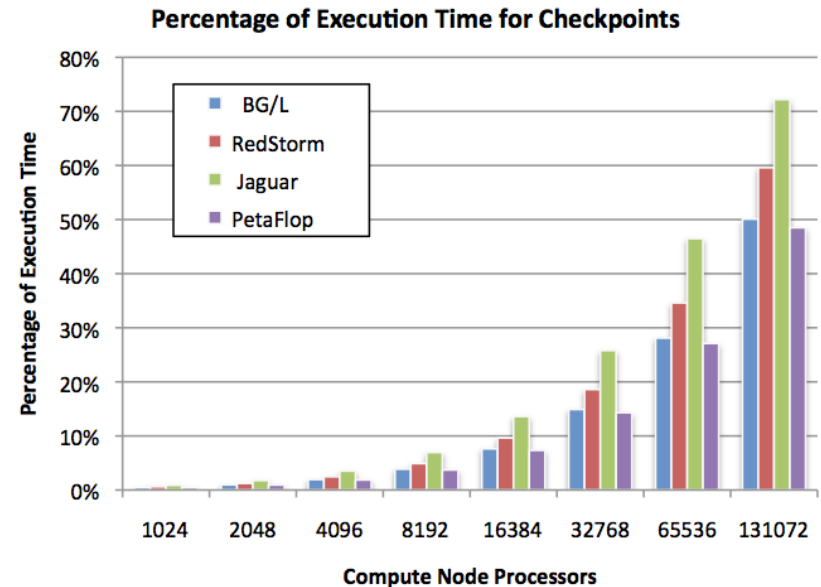
# Resilience... Our Biggest I/O Challenge

## Most of our I/O is for resilience

- Application-directed checkpoints are the primary protection against faults
- Application characteristics
  - Require large fractions of system
  - Resource constrained
  - **Cannot survive a failure**
- Probability of failure is based on application size.
- Frequency of checkpoint is based on probability of failure

## Our resilience efforts reduce I/O

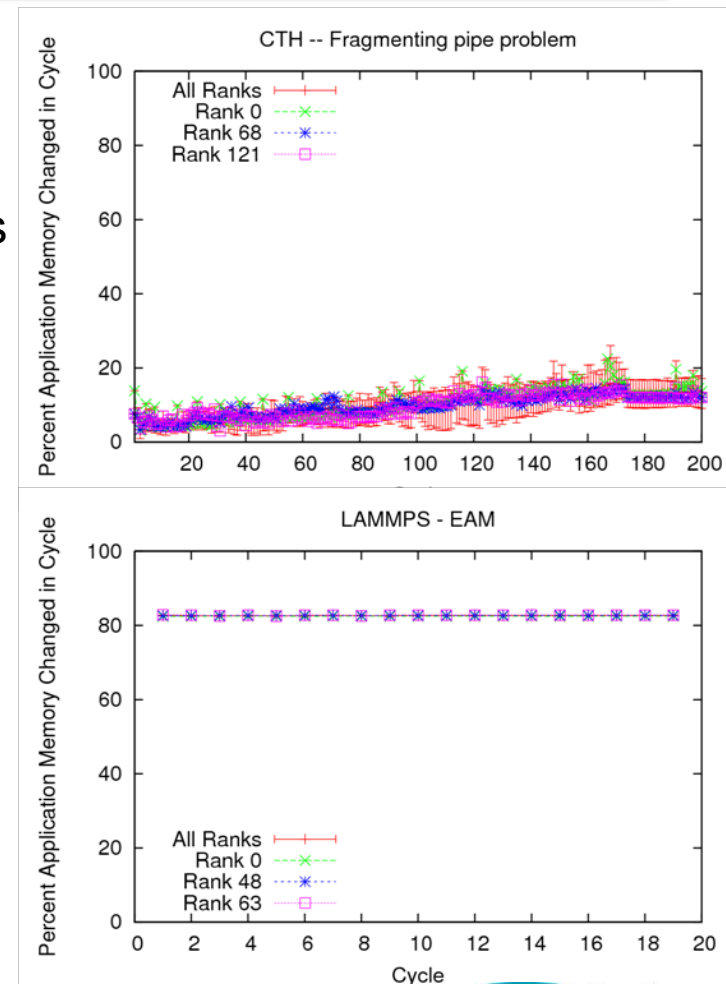
- System-influence on how/when to chkpt
- **Viability of incremental checkpoints,**
- Diskless checkpoints,
- **Redundant computation**



Oldfield et al. Modeling the impact of checkpoints on next-generation systems. In *Proceedings of the 24th IEEE MSST*, Sept. 2007

# The Case For/Against Incremental Checkpoints

- Lightweight lib to identify modified memory
  - Page-table trickery identifies modified pages
  - Crypto-hash (MD5) identifies modified blocks
  - No app changes required; user/system specifies interval to collect memory statistic
  - User & kernel-space version for Catamount. CNL user-space version in testing.
- Results
  - Runtime overhead < 10% (~free with GPU)
  - CTH: modified memory within 8% of app
  - LAMMPS: modified memory 4x larger than checkpoint .



*Submitted to EuroMPI 2011*

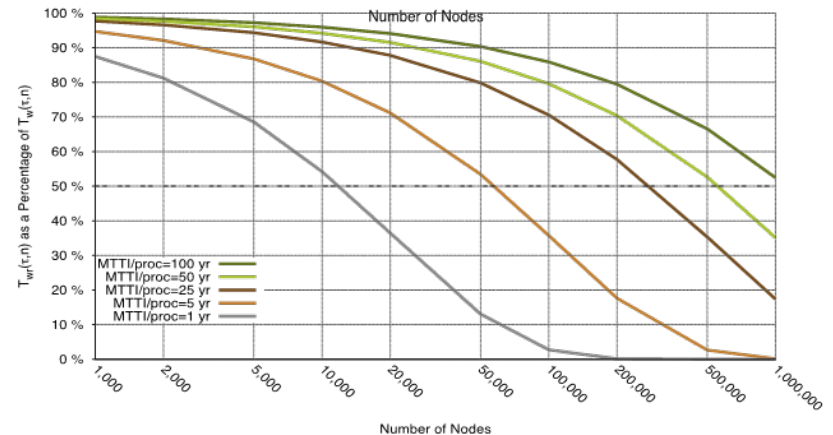
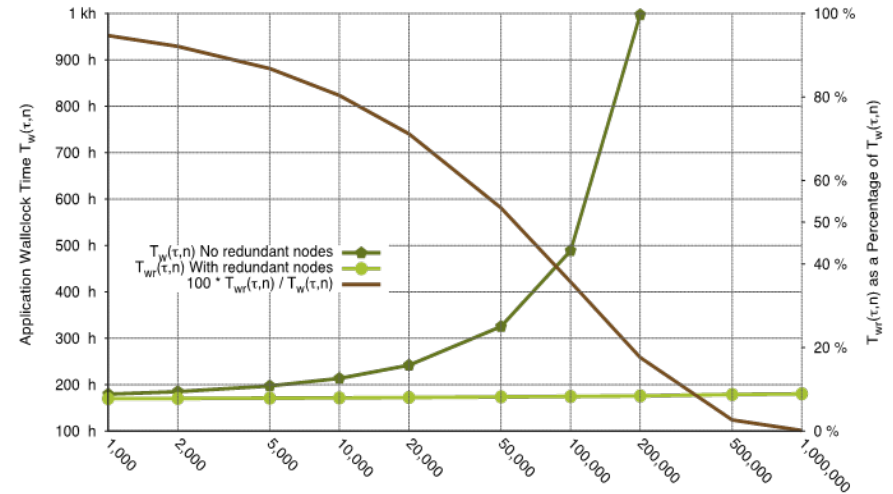
# Exploring Redundant Computation

- Motivation

- Overhead of checkpoint unacceptable
- Increase MTTI
- Reduce defensive I/O
- Hypothesis: at large scale, overhead of redundant computation is less than checkpoint/restart

- rMPI library

- Between application and MPI
- Replicates ranks 0..n
- Checkpoint still required (just not as often)
- rMPI almost a full MPI implementation
  - MPI\_Wtime, MPI\_Probe, ... need to return same answer for both nodes
  - Message order and other MPI semantics must be preserved



Submitted to SC2011

# Scalable I/O Services

Even our I/O research is about reducing I/O

## Purpose

- Leverage available compute/service node resources for I/O caching and data processing

## Application-Level I/O Services

- Lightweight File System (authr, authn, storage)
- **Shock physics particle tracking**
- **PnetCDF caching service**
- SQL Proxy (for NGC)
- Sparse-matrix visualization (for NGC)

## Other Plans

- PnetCDF caching
- Investigate placement issues
- ADIOS I/O services for fusion, climate, combustion apps on Jaguar

Client Application  
(compute nodes)



I/O Service  
(compute/service nodes)



Processed Data



Cache/aggregate /process



Lustre File System



Visualization Client

NETEZZA

LexisNexis\*



Network-Scalable Service Interface

# Scalable I/O Services

We did this for Salvo Seismic Imaging (circa 1996)

## Salvo's I/O Partition

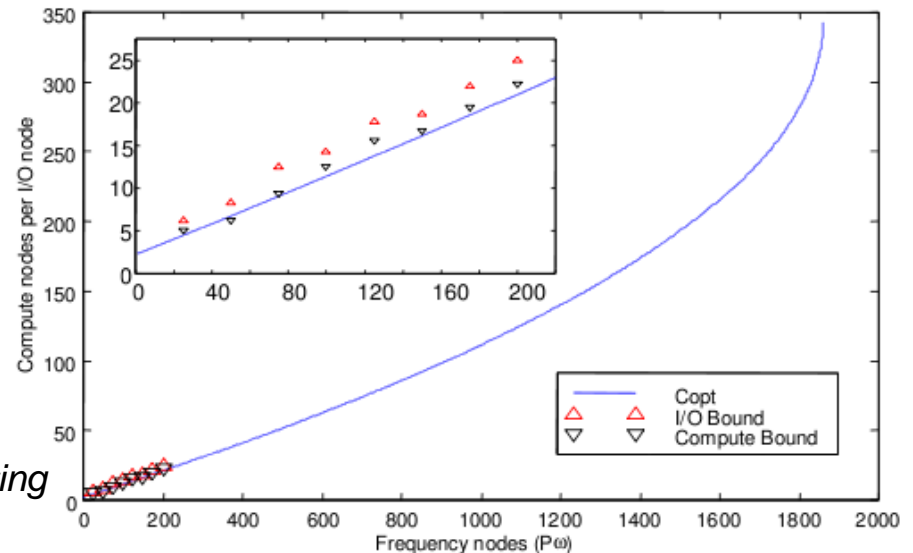
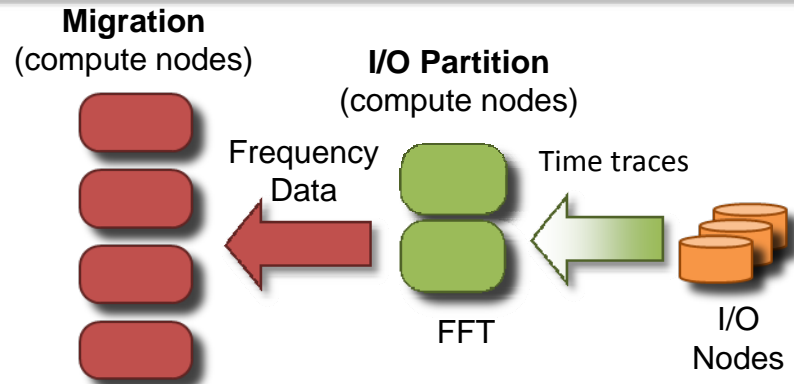
- Partition of application processors (used separate MPI Communicator for I/O)
- Used for FFT, I/O cache, and interpolation
- Async I/O allowed overlap of I/O and computation (pre-process next step)

## Results

- +10% nodes led to +30% in performance
- Modeling I/O and compute costs helped find the right balance of compute and I/O nodes

**Contacts:** Ron Oldfield, Curtis Ober  
{raoldfi,ccoer}@sandia.gov

Oldfield, et al. Efficient parallel I/O in seismic imaging.  
*The International Journal of High Performance Computing Applications*, 12(3), Fall 1998



# Scalable I/O Services

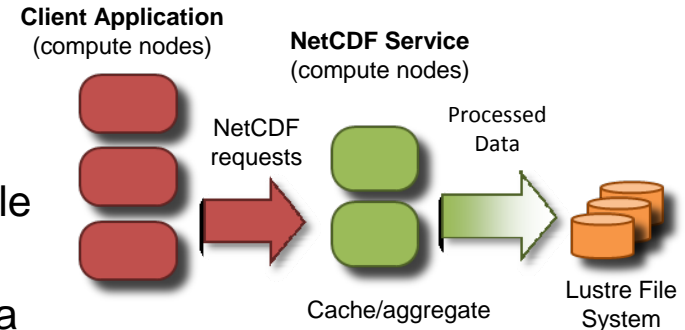
## NetCDF I/O Cache

### Motivation

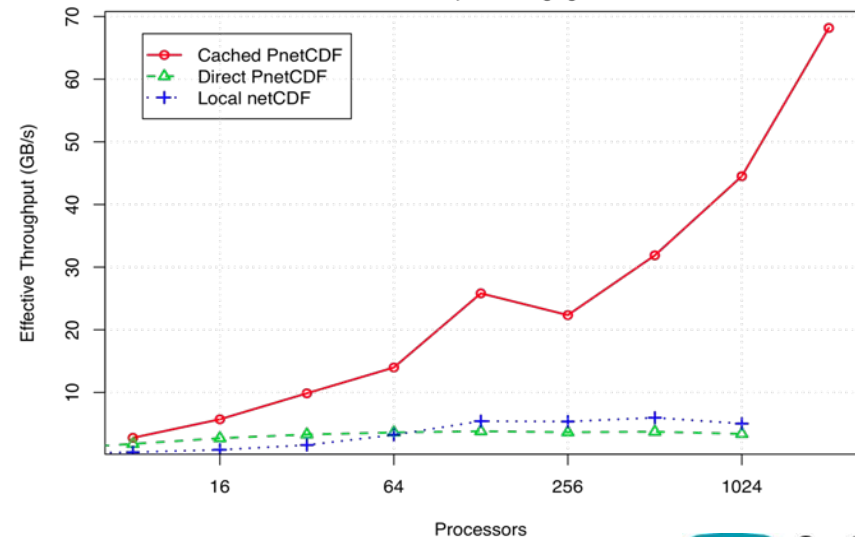
- Synchronous I/O libraries require app to wait until data is on storage device
- Not enough cache on compute nodes to handle “I/O bursts”
- NetCDF is basis of important I/O libs at Sandia (Exodus)

### NetCDF Caching Service

- Service aggregates/caches data and pushes data to storage
- Async I/O allows overlap of I/O and computation



IOR Performance on Red Storm  
4:1 ratio of compute to staging nodes





# Scalable I/O Services

## CTH Fragment Detection

### Motivation

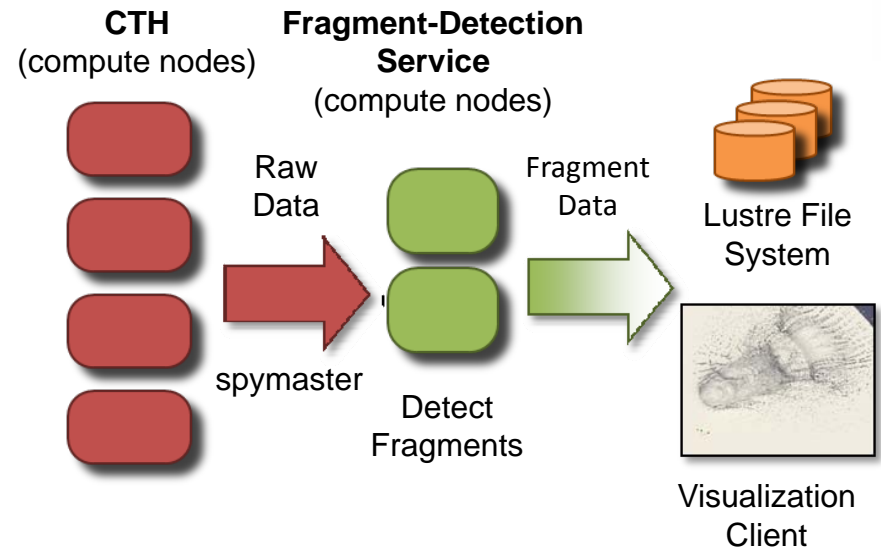
- Fragment detection process takes 30% of time-step calculation
- Fragment tracking requires data from every time step (too data intensive for post processing)
- Integrating detection software with CTH is intrusive on developer

### CTH fragment detection service

- Extra compute nodes provide in-line processing (overlap fragment detection with time step calculation)
- Only output fragments to storage (reduce I/O)
- Non-intrusive
  - Looks like normal I/O (pvspy interface)
  - Can be configured out-of-band

### Status

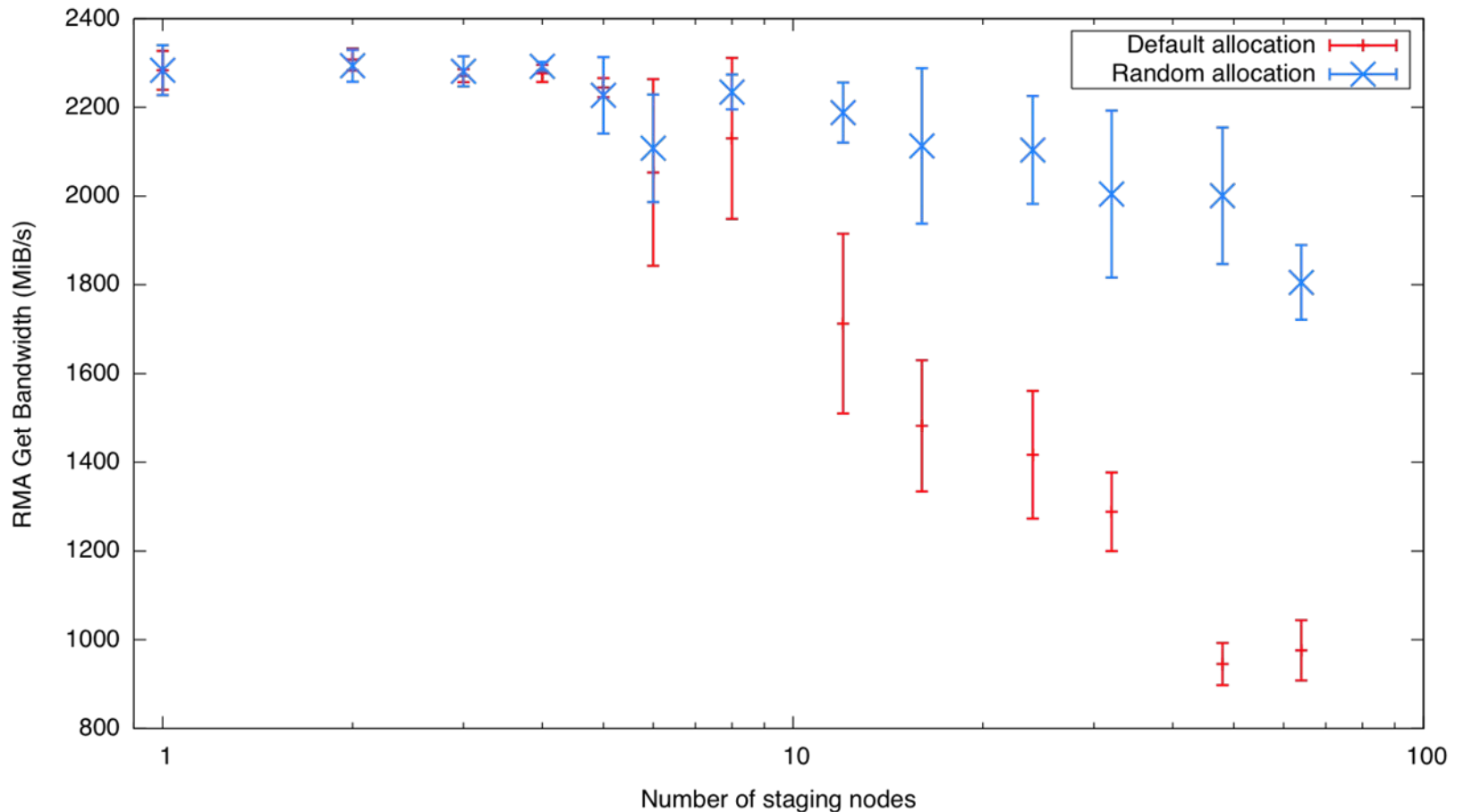
- Developing client/server stubs for pvspy
- Developing Paraview frag detect service



*Fragment detection service provides on-the-fly data analysis with no modifications to CTH.*

# Placement Issues for I/O Services

8:1 ratio of application to staging nodes



# Peer-to-Peer File System (New)

## Lee Ward

---

- Motivation
  - Current FS designs use centralized management to coordinate access to shared devices
  - Current FS require level of predictable performance and reliability not practical at exascale
    - Device failures, performance variability, device contention, all have a big impact on application performance and system uptime
- Features
  - Decentralized management of devices
  - Support for heterogeneity in a system of inherently unreliable networks and storage devices
- “Smart” servers are the key
  - Pervasive in the computing system (they’re everywhere!)
  - Support for a variety of local and remote media (disk, tape, memory, NVRAM)
  - Directly handle I/O reqs, initiate 3<sup>rd</sup> party transfers, or replicate data as needed

# Other Gaps to Fill

---

- **System software**
  - Support for dynamic allocation and reconfiguration
    - Data services: balanced workflow, reduce data movement, dynamic deployment
    - Smart placement (topologically aware scheduling)
    - Resilience: failed node replacement (reduce I/O for checkpoint)
  - Integrated support for NVRAM as a memory device
  - System support for application-driven RDMA
- **Programming models**
  - Standard approaches for integrating sim and analysis
  - Standard approaches for programming services (CPU, GPU, FPGA)
- **Resilience**
  - Storage-efficient app resilience is still a problem after 20+ years of research
  - Data service resilience: services use memory for transient data, how do we ensure resilience in such a model? We are working on this... let's talk again next year ;0)

# Summary

---

## Scalable I/O Research

- Sandia is Involved in Leading-Edge R&D for SIO
  - Peer-to-peer File System (expect faults, handle extremely large scale)
  - Storage-efficient resilience
  - Scalable I/O Services
  - I/O Characterizations, Tracing, and Simulation\*
  - I/O System use of Accelerators (GPGPU RAID6 – Submitted for R&D 100)\*
- Scalable I/O Services (Nessie)
  - Integrated simulation, analysis, caching
  - Already demonstrated value for Seismic (Salvo)
  - New functionality for HPC systems
    - To manage bursts of I/O: netCDF cache
    - In-transit fragment detection/tracking (to reduce I/O)

*Smart data movement and I/O reduction is critical for exascale*

\* Not discussed in this talk

# Addressing Scalable I/O Challenges for Exascale

Approved for Public Release SAND2011-3588C

27<sup>th</sup> IEEE Symposium on Massive Storage  
Systems and Technologies

May 24, 2011

*Ron Oldfield*  
*Sandia National Laboratories*

