# LANL SciDAC Petascale Data Storage Institute Operational Data Releases

**James Nunez ( jnunez@lanl.gov ),**

**Gary Grider, John Bent,**

**HB Chen, Meghan Quist, Alfred Torrez**

**Los Alamos National Lab**

**February 2008**

LA-UR-07-7449

UNCLASSIFIED

# Publicly Available Data from LANL Supercomputers

- ## Available Data
  - Nine years of computer operational failure data, over 23,000 records for several thousand machines
  - Several million usage records (job size, processors/machines used, duration, time, etc.)
  - Disk failure data for a single Supercomputer
  - Machine Layout Information (Building, room, Rack location in room, node location in rack, hot/cold rows, etc.)

- ## Recently Released Data
  - I/O Traces of MPI-IO Based Synthetic Application

- ## In-Progress
  - Disk failure, node and scratch file system, data for several Supercomputer
  - Hundreds of workstation File Systems Statistics Information

- ## Soon to be Released Data
  - Physics Application I/O Traces

UNCLASSIFIED

# Failure/Usage/Event Data Sets From LANL Supercomputers

| description | size | records | name | |
|---|---|---|---|---|
| all systems failure/interrupt data 1996-2005 | 2963538 | 23741 | LA-UR-05-7318-failure-data-1996-2005.csv | |
| system 20 usage with domain info | 51675641 | 489376 | LA-UR-06-0803-MX20_NODES_0_TO_255_DOM.TXT | |
| system 20 usage with node info nodes number from zero | 43926669 | 489376 | LA-UR-06-0803-MX20_NODES_0_TO_255_NODE-Z.TXT | |
| system 20 event info nodes number from zero | 33120015 | 433490 | LA-UR-06-0803-MX20_NODES_0_TO_255_EVENTS.csv | |
| system 15 usage with node info nodes number from zero | 2416139 | 17823 | LA-UR-06-0999-MX15-NODE-Z.TXT | |
| system 16 usage with node info nodes number from one | 321293488 | 1630479 | LA-UR-06-1446-MX16-NODE-NOZ.TXT | |
| system 23 usage with node info nodes number from one | 60674531 | 654927 | LA-UR-06-1447-MX23-NODE-NOZ.TXT | |
| system 8 usage with node info nodes number from one | 67291020 | 763293 | LA-UR-06-3194-MX8-NODE-NOZ.TXT | |

# Machine Information:
# 23 LANL Supercomputers

| system CMU paper number | system data machine number | system type | number nodes | number cpus | cpus/node | install date | production date | decommision date | fru | mem per node | cpu type | number of interconnects | use type | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 7 | smp | 1 | 8 | 8 | before tracking | before tracking | Dec-99 | part | 16 | 3 | 0 | compute | |
| 2 | 24 | smp | 1 | 32 | 32 | before tracking | before tracking | Dec-03 | part | 8 | 7 | 1 | compute | |
| 3 | 22 | smp | 1 | 4 | 4 | before tracking | before tracking | Apr-03 | part | 1 | 6 | 0 | compute | |
| 4 | 8 | cluster | 164 | 328 | 2 | Mar-01 | Apr-01 | current | part | 1 | 4 | 1 | compute | GET Machine Layout |
| 5 | 20 | cluster | 512 | 2048 | 4 | Oct-01 | Dec-01 | current | part | 16 | 2 | 2 | compute | GET Machine Layout |
| 6 | 21 | cluster | 128 | 512 | 4 | Aug-01 | Sep-01 | Jan-02 | part | 16 | 2 | 2 | compute | |
| 7 | 18 | cluster | 1024 | 4096 | 4 | Mar-02 | May-02 | current | part | 16 | 2 | 2 | compute | GET Machine Layout |
| 8 | 19 | cluster | 1024 | 4096 | 4 | Aug-02 | Oct-02 | current | part | 16 | 2 | 2 | compute | GET Machine Layout |
| 9 | 3 | cluster | 128 | 512 | 4 | Aug-03 | Sep-03 | current | part | 4 | 2 | 1 | compute | GET Machine Layout |

# Sample Machine Layout Information

Machine 8, bldg 1, room 1

,RackPosition,RackPosition,Position in Rack,Row facing

NODE NUM,East/West,North/South,Vertical Position,Single row cluster

N#,1 to 26,28 to 35,"1 to 37, top to bottom",

1,23,28,1,rear to N/Hot

2,23,28,2,rear to N/Hot

3,23,28,3,rear to N/Hot

4,23,28,4,rear to N/Hot

5,23,28,5,rear to N/Hot

6,23,28,6,rear to N/Hot

7,23,28,7,rear to N/Hot

8,23,28,8,rear to N/Hot

…

159,23,34,19,rear to N/Hot

160,23,34,20,rear to N/Hot

161,23,34,21,rear to N/Hot

162,23,34,22,rear to N/Hot

163,23,34,23,rear to N/Hot

164,23,34,24,rear to N/Hot

Los Alamos
NATIONAL LABORATORY
EST.1943

pdsi

# File Systems Statistics Survey

- **Based on CMU/Panasas File System Statistics Survey (fsstats)**

- **Purpose**
  - Develop better understanding of file system components
  - Gather and Build large db of static file tree statistics

- **Usage**
  - Parallel statistics collection with LANL's MPI-File Tree Walk
  - DB query interface

- **LANL Data**
  - Production file system statistics from LANL Supercomputers & Testbeds
  - Hundreds of statistics from backups of workstations Lab-wide

- **Output: Histogram and statistics on**
  - File Size
  - Capacity Used
  - Directory Size
  - File Name Size

UNCLASSIFIED

# Sample File Systems Statistics Survey

histogram,file size
count,2400,items
average,195.763750,KB
min,0,KB
max,18629,KB
bucket min,bucket max,count,percent,cumulative pct
0,2,14,0.005833,0.005833
2,4,20,0.008333,0.014167

…

histogram,filename length
count,48175,items
average,19.143830,chars
min,0,chars
max,164,chars
bucket min,bucket max,count,percent,cumulative pct
0,7,4150,0.086144,0.086144
8,15,24112,0.500509,0.586653
16,23,10512,0.218204,0.804857

…

# Tracing Mechanisms and Trace Data

- ## Desired Attributes in a Trace Mechanism
  - Minimum overhead
  - Bandwidth preserving

- ## Methods being Reviewed
  - Linux ltrace/strace
  - Tracefs from SUNY Stony Brook

- ## Information Collected
  - Time Stamps to detect node clock skew and drift
  - Library and system calls and summery information
  - Listing of directories

- ## Traces Availability
  - Available Now - Traces of Open Source Benchmark: MPI-IO based synthetic
  - Real Physics Application I/O Traces

UNCLASSIFIED

# Tracing Mechanisms and Trace Data

## N-to-N

|  | 64 KB | 256 KB | 448 KB | 512 KB | 1024 KB | 4096 KB | 8192 KB | 16386 KB | 32772 KB | 65544 KB |
|---|---|---|---|---|---|---|---|---|---|---|
| 32 Procs | TGZ | TGZ | TGZ | TGZ | TGZ | TGZ | TGZ | TGZ | TGZ | TGZ |
| 96 Procs | TGZ | TGZ | TGZ | TGZ | TGZ | TGZ | TGZ | TGZ | TGZ | TGZ |

## N-to-1 nonstrided

|  | 64 KB | 256 KB | 448 KB | 512 KB | 1024 KB | 4096 KB | 8192 KB | 16386 KB | 32772 KB | 65544 KB |
|---|---|---|---|---|---|---|---|---|---|---|
| 32 Procs | TGZ | TGZ | TGZ | TGZ | TGZ | TGZ | TGZ | TGZ | TGZ | TGZ |
| 96 Procs | TGZ | TGZ | TGZ | TGZ | TGZ | TGZ | TGZ | TGZ | TGZ | TGZ |

## N-to-1 strided

|  | 64 KB | 256 KB | 448 KB | 512 KB | 1024 KB | 4096 KB | 8192 KB | 16386 KB | 32772 KB | 65544 KB |
|---|---|---|---|---|---|---|---|---|---|---|
| 32 Procs | TGZ | TGZ | TGZ | TGZ | TGZ | TGZ | TGZ | TGZ | TGZ | TGZ |
| 96 Procs | TGZ | TGZ | TGZ | TGZ | TGZ | TGZ | TGZ | TGZ | TGZ | TGZ |

# Sample MPI-Based Synthetic I/O Trace

- ## Timing Information

#Barrier before benchmark_call

7: testbedX_node113 (10378) Entered barrier at 1159808385.170918

7: testbedX_node113 (10378) Exited barrier at 1159808385.173167

3: testbedX_node117 (11335) Entered barrier at 1159808385.166396

3: testbedX_node117 (11335) Exited barrier at 1159808385.168893

5: testbedX_node115 (10373) Entered barrier at 1159808385.168842

5: testbedX_node115 (10373) Exited barrier at 1159808385.171370

- ## Traced Application Data

10:59:47.092996 MPI_File_open(92, 0x80675c0, 37, 0x80675a8, 0xbfdfe5e4 <unfinished ...>

10:59:47.093718 SYS_statfs64(0x80675c0, 84, 0xbfdfe410, 0xbfdfe410, 0xbd3ff4) = 0 <0.011131>

10:59:47.108352 SYS_open("/FS/scratch/usr1/O"..., 32832, 0600) = 3 <0.000745>

10:59:47.109189 SYS_close(3) = 0 <0.000063>

10:59:47.109310 SYS_open("/FS/scratch/usr1/O"..., -2147450814, 0600) = 3 <0.000564>

10:59:47.110912 <... MPI_File_open resumed> ) = 0 <0.017855>

10:59:47.110955 MPI_Wtime(0x8063830, 0xb7fee8dc, 0xbfdfe568, 0xb7f66ad7, 0x8059020 <unfinished ...>

# Sample MPI-Based Synthetic I/O Trace

- ## Summary Information

# SUMMARY COUNT OF TRACED CALL(S)

# Function Name Number of Calls Total time (s)

#=====================================================

MPIO_Wait 2 0.000118

MPI_Barrier 29 2.156431

MPI_File_close 2 0.108482

MPI_File_delete 1 0.102532 …

# SUMMARY COUNT OF CALLS WITHIN 1 MPI_File_delete CALL(S)

# Function Name Number of Calls Total time (s)

#=====================================================
    =

SYS_ipc 8 0.000140

SYS_statfs64 1 0.009227

SYS_unlink 1 0.092411

# Additional Sources of Data

- NERSC I/O Trace Data and Workload Profiles
  - http://pdsi.nersc.gov/benchmarks.htm

- USENIX Computer Failure Data Repository
  - The USENIX, the Advanced Computing Systems Association has begun a project to index all of these large scale operational data releases on the internet, at cfdr.usenix.org.  LANL and the PDSI are helping to organize and lead this effort.

UNCLASSIFIED

# Publications Based on LANL Data or Software

- Eric Lalonde. "A Characterization of LANL HPC Systems", Masters Thesis, University of California, Santa Cruz. 2007

- Olen Davis, Kari Macklin, BailyKelly "Parallel Search using multiple Google Desktops in Parallel", Colorado School of Mines, Technical report, 2007.

- Clay Baenziger, Bruce Bugbee, Ryan Ford, Charlie, Grammon. "LANL Supercomputing Data Analysis", Colorado School of Mines Technical Report, 2007.

- Bianca Schroeder, Garth Gibson. "The computer failure data repository." Invited contribution to the Workshop on Reliability Analysis of System Failure Data (RAF'07)to be held at MSR Cambridge, UK.

- Bianca Schroeder, Garth Gibson. "Disk failures in the real world: What does an MTTF of 1,000,000 hours mean too you?" 5th UsenixConference on File and Storage Technologies (FAST 2007). Winner of best paper award. This paper has also been featured in an article on slashdot, which so far has received more than 75,000 hits!

- Bianca Schroeder, Garth Gibson. "A large scale study of failures in high-performance-computing systems." International Symposium on Dependable Systems and Networks (DSN 2006). As one of the best DSN'06 papers invited to IEEE Transactions on Dependable and Secure Computing (TDSC).

- Michael Mesnier, Matthew Wachs, Raja R. Sambasivan, Alice Zheng, Gregory R. Ganger. Modeling the relative fitness of storage. International Conference on Measurement and Modeling of Computer Systems (SIGMETRICS 2007). San Diego, CA. June 12-14, 2007. ACM. Awarded Best Paper. Michael Mesnier, Matthew Wachs, Raja R. Sambasivan, Julio Lopez, James Hendricks, Gregory R. Ganger, David O'Hallaron //TRACE: Parallel Trace Replay with Approximate Causal Events.Fifth Conference on File and Storage Technologies (FAST'07). San Jose,CA. February 12-13, 2007. USENIX

# Links to Data & Codes

- **Machine Failure/Usage/Event/Location & Disk Failure Data Sets**
  - http://institutes.lanl.gov/data/fdata

- **Traces of MPI-IO Based Synthetic**
  - http://institutes.lanl.gov/data/tdata/

- **MPI-IO based synthetic, MPI-File Tree Walk, LANL Trace**
  - http://institutes.lanl.gov/data/software/

- **File Systems Statistics Survey (fsstats) Code**
  - http://www.pdsi-scidac.org/fsstats/

- **USENIX Computer Failure Data Repository**
  - http://cdfr.usenix.org

- **Contact E-mail:** jnunez@lanl.gov

UNCLASSIFIED