

MPI-IO Performance on Franklin Katie Antypas

Andrew Uselton, Harvey Wasserman, Shane Canon, Steve Luzmoor, Tom Davis

Office of Science







I/O Performance on Franklin

- Intended to talk about poor MPI-IO performance on Franklin
- But ... the system maintenance on March 17th drastically improved MPI-IO performance and I/O performance over all
- However, none of the March 17th changes were intended to improve I/O performance -- (so we thought)
- What caused the performance improvements?
- Knowing what caused the improvement is important so we can avoid regressions and assure this improvement is available on all XTs







Franklin's I/O System

- Lustre file system -- /scratch
- 20 OSSs, 80 OSTs
- 346 TB disk
- At this time measured peak performance ~11 GB/sec









MPI-IO shared file performance compared to file-per-processor performance Feb 2009

Franklin I/O Rates - IOR 64 processor









March 17th Maintenance

- OS Upgrade from CLE 2.1 to CLE 2.1UP01 with patch sets 01, 01A and 02
- Repair down hardware links
- Change PAM/LDAP settings
- Change network configuration from 42net to 41net
- Reduce max wall clock time for regular queues from 36 to 24 hours
- Insert new SIO modules (converted from compute nodes)









FLASH Checkpoint Read 2048 Procs









rrrr

BERKELEY

Franklin IO Plotfile Write 2048 cores





IO profile of MADBench, MPI-IO test.



Writes -blue, Reads- Red, x axis time, y axis processors





Speed up of over 4 times after the March 17th maintenance



Data from Andrew Uselton and Noel Keen



MADBench2 using POSIX Before and After Maintenance







Significant performance change?





S3D I/O Write Performance

- Not an MPI-IO code, file per processor
- 200 MB/core
- Average rate before: 955 MB/sec
- Average rate after: 4972 MB/sec









IOR MPI-IO performance before and after maintenance

MPI-IO IOR Performance Before and After March Upgrade ~2GB/proc 64 processors









Possible Explanations

- Repair down hardware links
 - If crucial links were down could this cause poor performance?
 - Cray staff did congestion analysis and found no differences before and after upgrade
- OS Upgrade from CLE 2.1 to CLE 2.1UP01 with patch sets 01, 01A and 02
 - Many changes in UP01
- Workload Changes?







U.S. DEPARTMENT OF ENERGY

CLE2.1 UP01

Async journal commit mod? – Reduces frequency of writes by 8





Data from Tom Davis

BERKELEY L



Workload Changes?

Aprun commands per day in 2009





Ending Thoughts

- Since March 17th completely reformatted I/O and moved I/O modules around
- Stability and performance of machine increased
- So can not do tests on Franklin system to recreate situation
- Priority is to keep I/O well performing, no regression
- Stresses the importance of performance monitoring
- Changes in one part of system can have unintended (but positive) consequences



