

# Extending Scalability of Collective IO Through Nessie and Staging

**Parallel Data Storage Workshop** 

November 13, 2011

Jay Lofstead (SNL) Ron Oldfield (SNL) Todd Kordenbrock (HP) Charles Reiss (UCB)









## Motivation

- Collective & Two-Phase IO proven beneficial
  Relatively modest data volumes, 1-D, 2-D
- But....
- Trade-off of inter-node communication for data reorganization to save IO not always beneficial
  - Large datasets
  - 3-D domain decompositions particularly bad





### **Motivation**

- Problem: technique is central to some IO APIs – netCDF4, HDF5
- Problem: Changing IO techniques/file format may not be an option for some applications
  - CESM climate model is committed to netCDF file format
- Problem: Continued scaling of problem sizes and resolution making things worse





# Solution

- Use efficient transport layer and data staging transparently in the IO stack
  - Re-implement native IO API (link-time compatible)
  - Ensure format on disk is identical
- Requirements
  - Efficient, portable transport layer
  - Staging area functionality to reduce time to completion





### **Solution Architecture**



www.lofstead.org gflofst@sandia.gov

PDSW 2011

Sandia National Laboratories

5



## **Nessie Transport Layer**

- Network Scalable Services Interface (Nessie)
  - Originally developed for the Lightweight File Systems project
  - RPC-like asynchronous API layer supporting RDMA
  - Physical layer support
    - InfiniBand
    - Portals
    - Cray Gemini
  - Server-directed for bulk data
    - Writes: pull from client
    - Reads: push to client







# **Staging Functionality**

- Collect data packets prior to writing to storage – Cache data chunks to afford other optimizations
- Perform data rearrangement
  - Perform partial data rearrangement like two-phase IO
- Use different techniques for writing to storage – Direct, Caching, Aggregation







Performance of xfer\_write\_rdma on Red Storm

Jay Lofstead www.lofstead.org gflofst@sandia.gov

PDSW 2011







Jay Lofstead www.lofstead.org gflofst@sandia.gov

PDSW 2011





# **NetCDF Staging Operation**

- 1. Initiate Request
- 2. Start data retrieval
- 3. Move data
- 4. Put completion/ result
- 5. Process in staging area





# **NetCDF Staging Functionality**

- NetCDF4 and PnetCDF API supported
- Direct synchronous with client calls
- Cache Independent asynch with client calls
- Aggregate Independent asynch with client calls, but aggregate data prior to writing (on node only)





# **NetCDF Staging Functionality**

#### Untested functionality for this paper

- Collective IO versions of cache and aggregate like independent versions, but a maximal number of collective IO calls made for writing
- Other data manipulation
  - Different implementation using Nessie for data analysis hosting





#### **NetCDF Staging Performance**



#### **Testing on JaguarPF using S3D IO kernel**

Jay Lofstead www.lofstead.org gflofst@sandia.gov

PDSW 2011





### **Current Status**

#### Nessie and NetCDF staging now part of Trilinos

- Trios capability area
- Port to BlueGene underway
- Integration of accelerators for staging processing





## **Future Work**

- Finish tests on RedSky to isolate Lustre issues
- Test collective IO routines
- Examine impact on reading performance
- For Nessie, other applications
  - 'In flight' data analysis routines
  - Transactions for resilience in data staging (see our poster!)
  - Hybrid, high level IO routine complications
    - Exodus + NetCDF

