# The Data Services Archive

Rena A. Haynes, Wilbur R. Johnson
Sandia National Laboratories
rahayne@sandia.gov, wrjohns@sandia.gov

NASA IEEE MSST 2004
12th NASA Goddard/21st IEEE Conference on
Mass Storage Systems & Technologies
The Inn and Conference Center
University of Maryland University College
Adelphi MD USA
April 13-16, 2004

# Archival Storage Configuration

High Performance Storage System (HPSS)
StorageTek and IBM tape libraries
Parallel tape
    4-way, 2-way, 1-way tape striping
Minimal HPSS managed disk storage
Access through FTP and PFTP

The Data Services Archive capability is designed to simplify and optimize the process of archiving large data sets.

# Massive Data and File Movement Issues

- User interface
    - File specification
    - Determining Progress
    - Error Detection
    - Recovery Procedures

- Resource requirements
    - Disk
    - Network
    - Archival

Sandia National Laboratories

Advanced Simulation & Computing™

# Movement across Distance Issues

- Data integrity and robustness
    End-to-end integrity
    Error recovery can exacerbate problems

- Archival device and media latencies can be magnified

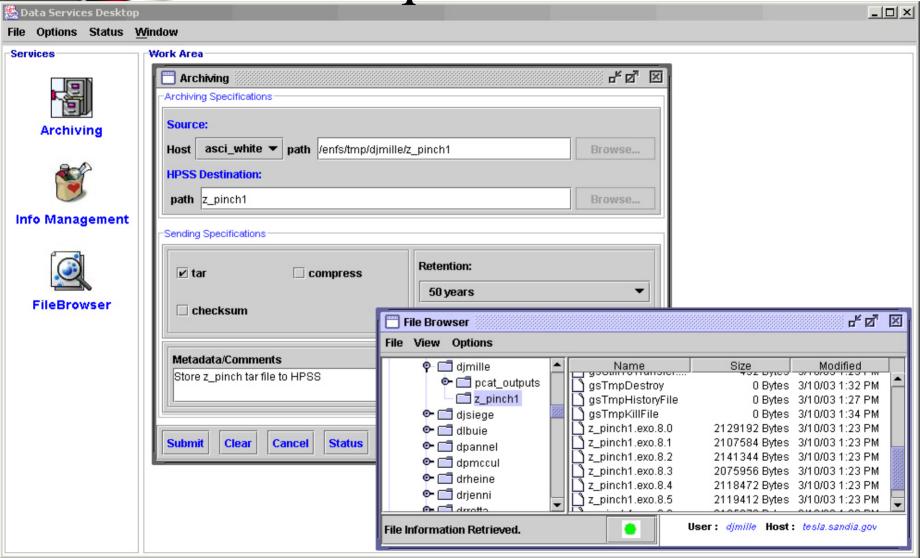- Networks and protocols not tuned for large archival data transfers

Sandia
National
Laboratories

Advanced
Simulation &
Computing™

# On-Demand Access to Parallel Archive Issues

- Inefficient use of resources
  Large transfer followed by small transfer followed by large transfer
  Multiple large transfers can block smaller transfers
  Use of parallel tapes to store small files

- Over-subscription of resources
  Can cause failures due to credential timeout and ultimately denial of service
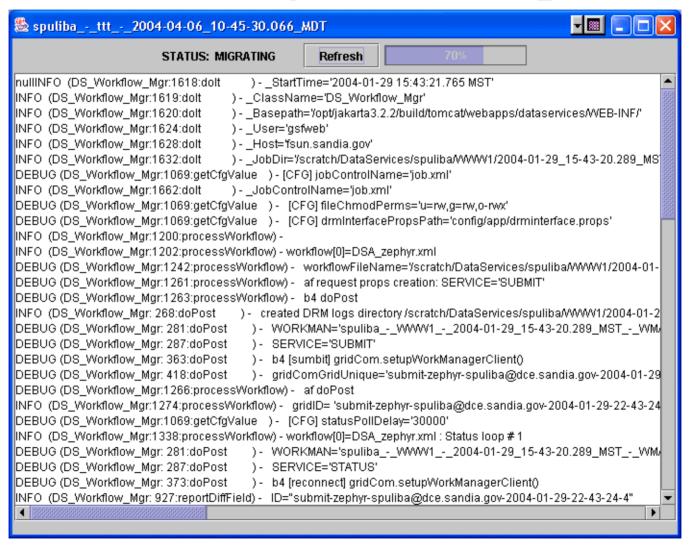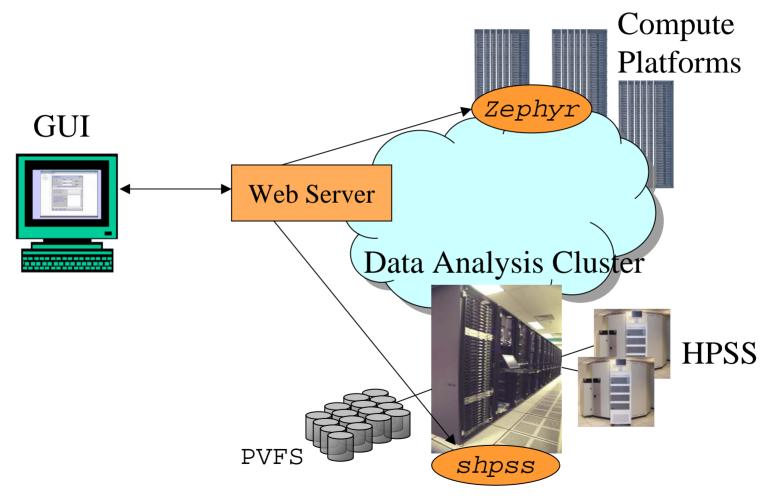  Cascades into additional resource requirements, e.g., disk

Sandia
National
Laboratories

Advanced
Simulation &
Computing™

# DSA Request Interface

# Monitoring DSA Request

# DSA is a Distributed Application

GUI

Compute Platforms

*Zephyr*

Web Server

Data Analysis Cluster

HPSS

PVFS

*shpss*

Sandia National Laboratories

Advanced Simulation & Computing™

# DSA use of Grid Technologies

- Grid workflow processor to sequence data transfers to intermediate disks then to tape

- Globus toolkit used for submitting partitions to queues

- Globus toolkit enhancements
  - Immediate feedback from components
  - Ability to request specific resources

# Resource Contention

- Staging files to intermediate disk cache
  - Reduces potential for retries/retransmissions caused by tape latencies
  - Allows grouping, or partitioning, of data to optimize use of archive resources
  - Simplifies interface to parallel archive
- Intermediate disk cache directly accessible by HPSS movers reduces external network traffic
- Resource management to minimize contention

Sandia
National
Laboratories

ADVANCED
SIMULATION &
COMPUTING™

# Dataset Partitioning

- Place dataset files into transfer groups based on file characteristics (currently, file size)

- Number of files in a partition are limited to facilitate recovery and prevent large transfers from starving smaller transfers

- Partitions are scheduled based on parallel tape resource requirements

# Scheduling Techniques

HPSS state and resources are modeled as nodes in a compute cluster

- – Holds ftp sessions that cannot be satisfied
- – Increase HPSS utilization
- – Requires knowledge of HPSS configuration
- – Requires knowledge of HPSS striping policies

Sandia National Laboratories

Advanced Simulation & Computing™

# Scheduling Techniques

DSA queuing system scheduler is configured to use backfill

- Allows pending transfers which can complete with available resources in the time to start the transfer at the head of the queue

- Requires estimating time to transfer

$$T_{job} = T_{login} + \sum_{i=1}^{i=N} \left[ T_{startup} + \left\lceil \left( X_i / R_{rate} \right) \right\rceil + \left( S_{width} \times T_{load} \right) \right]$$

# Performance Observations

- Staging
  - Allows for local file puts (up to 30MB/s per tape stripe)
  - Relieves source space more efficiently
- Partitioning
  - A tunable scheduling parameter for optimizing archive flow into HPSS
- Scheduling with backfill
  - Increases resource utilization and archive request throughput

# Conclusions

- Reducing resource contention, managing data partitions, and scheduling transfers have increased overall performance of the HPSS parallel tape archival system.

- More performance data should be collected and analyzed to adjust transfer time estimations and determine additional tuning parameters.

- Partitioning scheme should be optimized based on partition data size as well as file count.