# A Solution for Maintaining File Integrity within an Online Data Archive

Dan Scholes

PDS Geosciences Node

Washington University

# Presentation Will Discuss

- PDS Geosciences Node background
- Threats to online data archives
- Methods to identify corrupt files
- PDS Geosciences Node approach to ensuring data archive file integrity

# Planetary Data System (PDS)

- A NASA organization that archives science data from NASA's planetary missions.
- PDS responsibilities are:
  - To help NASA missions and other data providers to organize and document their digital planetary data
  - To collect complete, well-documented planetary data into archives that are peer-reviewed
  - To make the planetary data available and useful to the science community
  - **To ensure the long-term preservation and usability of the data.**

# PDS Geosciences Node's Data Holdings

- Planetary science data related to geoscience studies

  ◦ Surface and interior of the terrestrial planets and satellites (Moon, Mars, Mercury, Venus).

- Currently maintain:

  ◦ Archives from over 20 NASA missions

  ◦ Archive consists of over 40 TB of data

  ◦ Over 13 million files

# Access to Geosciences Node's Archive

- Direct Access
  - FTP and HTTP
- Web Interfaces
  - Providing search and retrieval capabilities
- Custom User Request
  - External hard drive

# Geosciences Node
# Data Storage Architecture

- Primary online data archive (SAN)

- Secondary online replication site
- Tape backups
- Deep archive at
  National Space Science Data Center (NSSDC)

# Threats to Online Data Archives

- Accidental change by staff
- Software error
- Hardware failure
- Malicious threats: Hacker or Virus
- Natural disaster

# Defenses

- Firewall settings
- Network security policies
- Proactive hardware maintenance
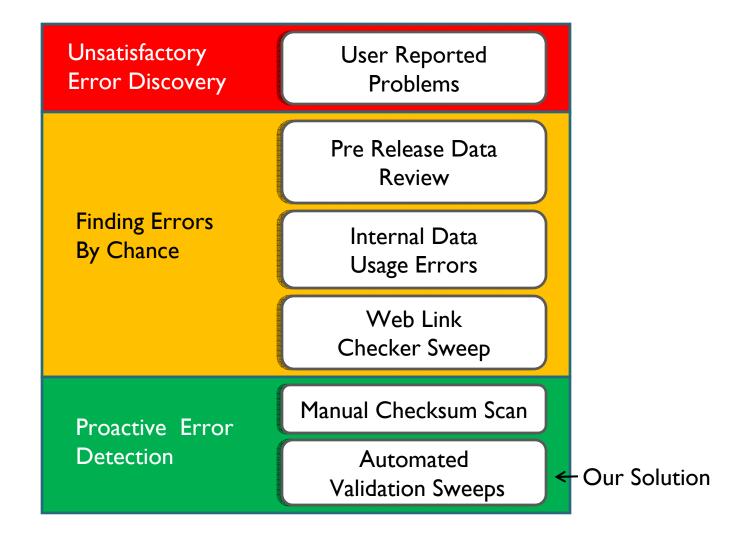- Multiple backup copies of the data

# Typical Recovery

- Restoration from offline backup
  - Tapes
  - External hard drive
  - DVD/CD
- Restoration from online secondary copy
  - Mirror site
  - Replication site

- **How do you know the recovered copy is not corrupt?**

# Bigger Question

**How do you know if a change or corruption has occurred in the data archive?**

# Identifying Corrupt Files

| | |
|---|---|
| **Unsatisfactory Error Discovery** | User Reported Problems |
| **Finding Errors By Chance** | Pre Release Data Review |
| | Internal Data Usage Errors |
| | Web Link Checker Sweep |
| **Proactive Error Detection** | Manual Checksum Scan |
| | Automated Validation Sweeps ← Our Solution |

# Checksum

- Checksum – a digital signature created by a hashing algorithm
  - File:  frt000027e2_01_if156l_trr2.img
  - MD5 Checksum: 5F393DAD7B36F6418045A9299E605E51
- The Geosciences Node uses MD5
  - Commonly used
    - Many client tools for data providers
  - Fast calculation

# Initial Data Integrity Study

- Manual Process
  - Create and compare checksum index files of data archive
- Advantages
  - Technically worked
  - Lessons learned
- Disadvantages
  - Time consuming
  - Difficult to manage
  - Difficult to update with new or replacement files
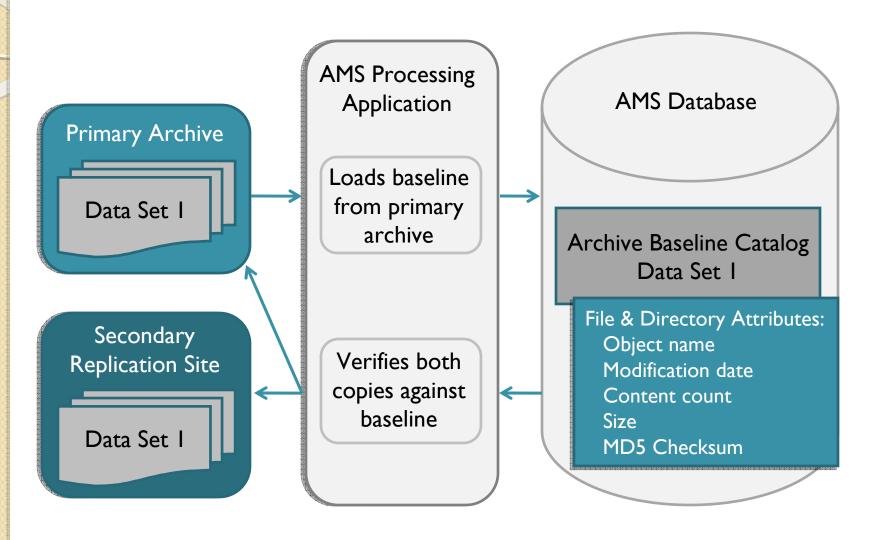
# Application System Requirements

- Create catalog of data archive contents
- Track multiple archive copies
- Update catalog as archive grows
- Verify archive against cataloged contents
- Provide processing speed for monthly archive validations
- Provide an easy to use application interface
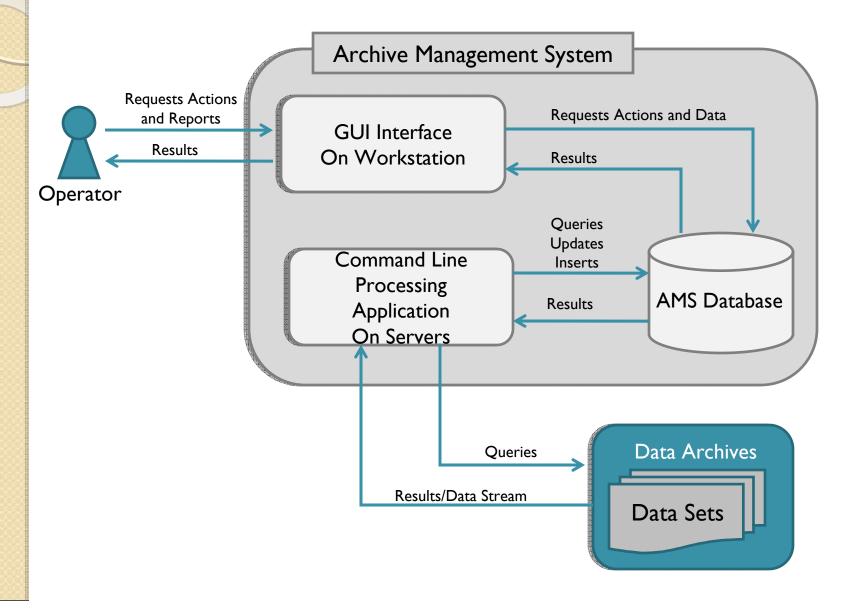
# Archive Management System (AMS)

- Custom application
- Components
  - Graphical user interface (GUI)
  - Command line processing application
  - Relational database
- Concept
  - Archive baseline catalogs
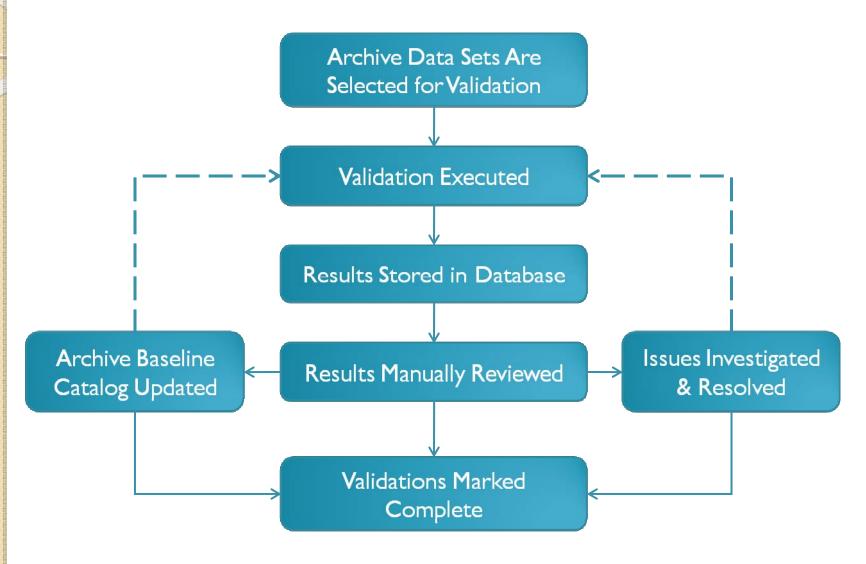
# Archive Baseline Catalog Concept

**Primary Archive**

Data Set 1

**Secondary Replication Site**

Data Set 1

**AMS Processing Application**

Loads baseline from primary archive

Verifies both copies against baseline

**AMS Database**

Archive Baseline Catalog Data Set 1

File & Directory Attributes:
Object name
Modification date
Content count
Size
MD5 Checksum

# AMS Overview

Archive Management System

Operator

Requests Actions and Reports

Results

GUI Interface On Workstation

Requests Actions and Data

Results

Command Line Processing Application On Servers

Queries Updates Inserts

Results

AMS Database

Queries

Results/Data Stream

Data Archives

Data Sets

# AMS Processing

- Create new archive baseline catalog
- Monthly validation scans
- Baseline is updated when new data is received

- Data recovery situations
  - Verify restored data against archive baseline catalog

# AMS Monthly Validation Scans

# Full Scan Validation

- File and Directory attributes scanned
  - Object name – case sensitive
  - Modification date
  - Content count (directory's file count)
  - Size
  - MD5 checksum (file validation only)
- Advantage
  - Thorough validation
- Disadvantages
  - Consumes more resources
  - Time consuming - entire archive up to 9 days

# Quick Scan – no checksum

- File and Directory attributes scanned
  - ◦ Object name – case sensitive
  - ◦ Modification date
  - ◦ Content count (directory's file count)
  - ◦ Size
- Advantages
  - ◦ Very fast processing speed - entire archive 28 hours
  - ◦ Identifies most accidental changes
- Disadvantage
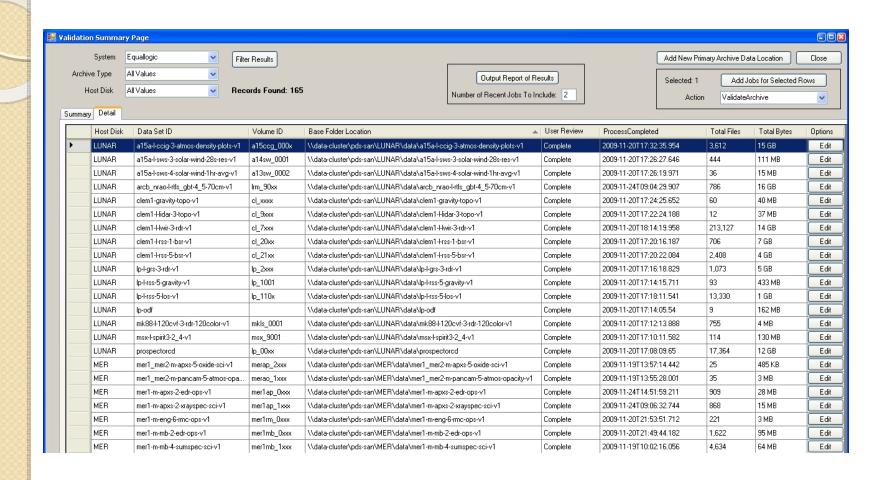  - ◦ Will not detect subtle file corruption

# Categories of Validation Results

- No differences are detected
- File/Directory attributes are different
- New archive content is discovered
- Archive content no longer exists
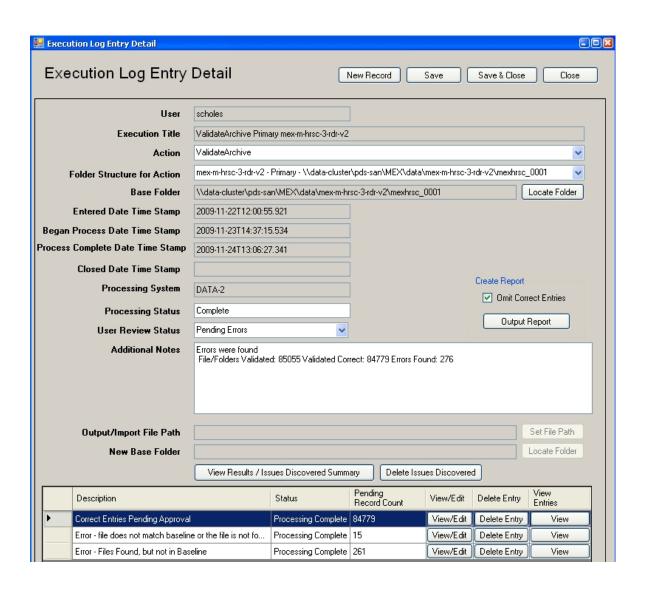
Differences require further review

# Interpreting Validation Results

- No differences are detected
  - Correct – no changes

- File/Directory attributes are different
  - Correct – revised data deployed to the data archive
  - Error – files were modified or corrupted

- New archive content is discovered
  - Correct – data added to the archive
  - Error - files accidently copied into archive

- Archive content no longer exists
  - Correct – items removed for archive revision
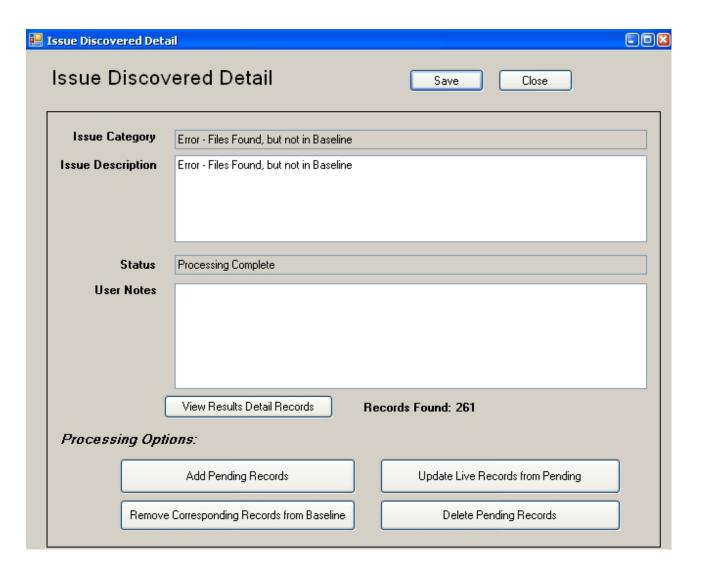  - Error - mistakenly or maliciously removed

# Archive Status List

# Validation Result Screen

# Validation Issue Resolution Screen

# AMS Results

- Geosciences Node has used the AMS for nearly a year.
  - Minimal personnel time to manage, monitor, and add new archives
  - Full scan of the entire archive 12 times
    - Can take up to 9 days of processing (full scan)
  - Two accidental archive changes
  - No file loss or corruptions
- Provides the Geosciences Node with a better degree of data integrity

# Future

- Geosciences Node's data archives continue to rapidly grow with current and future missions.
- Further performance review
  - Network switch configurations
  - Server Configurations
  - Disk Performance
  - Simultaneous processing streams
  - Possible code modifications

# Questions

- Contact Information
  - Dan Scholes
  - Applications Programmer
  - PDS Geosciences Node
  - Washington University in St. Louis
  - scholes@wunder.wustl.edu