

Overview

The Geosciences Node (<https://pds-geosciences.wustl.edu>) of NASA's Planetary Data System (PDS) stores and distributes ~250 terabytes of archive data. In addition to storing PDS archive files, hundreds of terabytes of disk space are used for archive preparation tasks, databases, virtual server images, etc. By migrating subsets of our systems to the cloud, we have eased the burden of managing time consuming IT-related tasks, while positioning the node to begin investigating more advanced cloud processing techniques.

Current Cloud Uses

The PDS Geosciences Node identified several IT operations to migrate to the cloud. The decision to migrate onsite systems are made based on their required effort to manage, cost to operate, and the possibility of improved functionality by moving to the cloud.

“Cold” Backup to Offline Cloud

For disaster recovery purposes, and to adhere to PDS data integrity policy, a tertiary offline copy of the PDS archive and supporting system files are maintained offsite.

- Historically, the disaster recovery backup was copied to tape, physically transported, and stored at an offsite facility. This process was time consuming and expensive.
- A more efficient process was developed in 2019 using AWS Glacier Deep Archive, which is intended to be a long-term offline cloud storage platform and serves a similar function to tape backups. (see Figure 1)
- Disaster recovery backups were migrated from magnetic tape to Glacier Deep Archive to address a number of IT challenges including effort required to manage backup jobs and perform hardware maintenance, high hardware support costs, and the requirement for an offsite storage location.
- Existing data protection software was used to redirect backups from the tape library to the cloud.
- The node is currently backing up over 450 terabytes of PDS archive data, databases, virtual server images, etc. to AWS.
- Once requested, data stored in Glacier are made available to download in less than 12 hours. The recovery time using this processes is comparable to tape backups.
- Close to a 50% cost saving was observed by migrating disaster recovery backups to the cloud.

“Warm” Backup to Online Cloud

In an effort to reduce costs, the PDS Geosciences Node migrated its warm (secondary) copy of PDS archive data to Azure's Blob tier of cloud storage instead of replacing onsite storage systems that were reaching the end of service life. Storing the secondary copy of PDS archive data in the Azure cloud not only saves on hardware and maintenance costs, but also positions the node to begin developing more advanced methods to process and analyze data in the cloud. (see Figure 1)

- The secondary copy of PDS archive data in Azure is an online mirror of the primary on-premises storage system.
- Files can be accessed at the secondary location in the event of data loss on the primary copy.
- Azure Blob cloud storage was selected in 2020 for this project due to contracts between the university and Microsoft that provide competitive rates for cloud offerings and egress fees.
- The PDS Geosciences Node uses a combination of an in-house file catalog and the Microsoft Azcopy tool to maintain synchronization of the secondary cloud copy with the local primary copy.
- The high bandwidth between on-premises systems and Azure is sufficient to adequately keep both copies of archive data synced within a few hours.

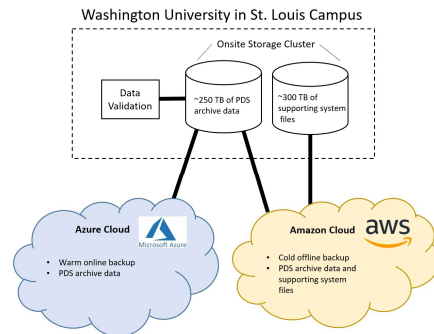


Figure 1. PDS Geosciences Node Cloud Backups

Pilot of CRISM Processing in the Cloud

The PDS Geosciences Node is also exploring in-cloud data analysis using the secondary warm archive backup. This is meant to allow users of the data to take advantage of cost benefits by performing data analysis in virtual machines co-located in the same zone as the cloud data. The lack of cloud egress when a user accesses data from such a VM reduces costs.

- Users benefit from this process because they do not need to purchase and maintain on-premises hardware for high-performance workloads; they might stand up a virtual machine that meets their computing requirements, run their programs, and then decommission the VM.
- The user does not need to download all the data to be analyzed to their on-premises computer over the internet. The transfer from our secondary backup to their cloud VM uses a high-speed in-cloud network. (see Figure 2).
- The pilot project in the cloud involved the analysis of CRISM spectrometer data using JCAT (Java CRISM Analysis Tool).
- JCAT was installed on an Azure VM and CRISM data was transferred from our Azure storage to the VM for analysis.
- After using the cloud VM to analyze the data, the small result files were sent back to the on-premises computer.

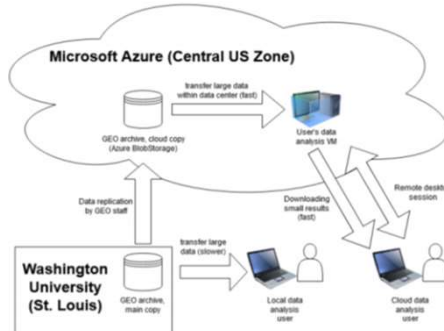


Figure 2. Diagram showing local versus cloud-based data processing

Lessons Learned

Data Replication

- Data replication to the cloud proved to be challenging because the on-premises and cloud storage systems use different metadata systems to track file modified and creation times.
- Traditional replication tools could not correctly mirror the two storage systems. The solution was to use Azure's Azcopy CLI sync tool in conjunction with in-house developed data validation software to detect differences.

How to Access Files in the Cloud

- Accessing files stored in the Azure Blob tier of cloud storage is not possible using traditional storage protocols such as SMB, NFS, or iSCSI. Custom APIs were written using the Azure SDK.
- There are options to use file servers in the cloud that use traditional storage protocols, but they were cost prohibitive.

Cloud Education

- There is a deep learning curve to understand cloud security, authentication, networking, etc.
- Sufficient time should be dedicated to researching the cloud to determine if the workload is a good candidate for migration.

Data Egress

- Understanding how data egress is charged and how to mitigate these costs is important.
- One way to address egress fees is to require users to access cloud data from within the same zone.
- There are also more advanced methods to reduce egress that governs the number of bytes that each user can transfer out of a cloud zone.

Next Steps

Additional cloud implementations to add value to the PDS Geosciences' archive data are being considered.

- In an effort to improve data availability for users, the node is considering methods in which PDS archives could be served directly from the cloud in the event of catastrophic or network failure to the primary on-premises copy.
- Developing new methods in which users can use cloud computing techniques against our data will be the primary focus for future cloud implementations. (see Figure 3)
- In addition to building on cloud pilots that have already been completed, the node plans to investigate cloud computing methods in which users are able to mount PDS archive data to their own compute instances.
- The node is also investigating a workflow that involves a user “checking out” a pre-built VM with data analytic tools already installed.

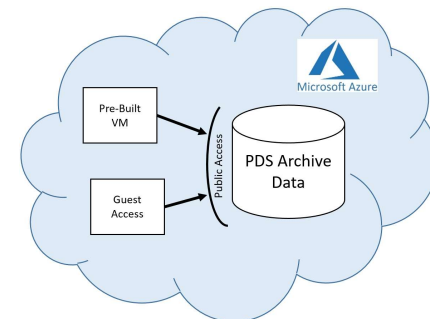


Figure 3. Public access to Azure Cloud Data