



## **USGS Looks to Distributed Multi-Cloud Storage to Accelerate the Understanding of our Planet**



# CONTENTS

- CONTENTS.....1**
- About United States Geological Survey .....2**
- About ScienceBase .....2**
- About Spectra Logic.....2**
- Data as Vast as the Planet .....2**
- Environment Overview.....5**
- Implementing a Hybrid Storage Solution .....8**
- How Data Moves from Site to Site.....10**
- Utilizing BlackPearl and Object Storage .....12**
- Utilizing Vail and Hybrid Cloud Storage .....13**
- USGS Exploratory Work.....14**
- Additional Resources.....14**

Copyright ©2021 Spectra Logic Corporation. All rights reserved worldwide. Spectra and Spectra Logic are registered trademarks of Spectra Logic. All other trademarks and registered trademarks are property of their respective owners. This document is provided for information purposes only, and the contents hereof are subject to change without notice. This document is not warranted to be error-free, nor subject to any other warranties or conditions, whether expressed orally or implied in law. We specifically disclaim any liability with respect to this document, and no contractual obligations are formed either directly or indirectly by this document. This document may not be reproduced or transmitted in any form or by any means, electronic or mechanical, for any purpose, without our prior written permission. All opinions in this white paper are those of Spectra Logic and are based on information from various industry reports, news reports and customer interviews.



## About United States Geological Survey

The United States Geological Survey (USGS) provides information and predictions about the natural hazards that threaten lives and livelihoods: the water, energy, minerals, and other natural resources we rely on; the health of our ecosystems and environment; and the impacts of climate and land-use change. USGS scientists develop new methods, tools, and data sets to supply timely, relevant, and useful information about the Earth and its processes.

## About ScienceBase

The United States Geological Survey is committed to enhancing and expanding information sharing and sound data management practices by developing [ScienceBase](#), a collaborative scientific data and information management platform used directly by science teams. ScienceBase provides access to aggregated data derived from many research and information domains, including feeds from existing data systems, metadata catalogs, and scientists contributing new and original content. ScienceBase architecture is designed to help science teams and data practitioners centralize their data and information resources to create a foundation needed for their work. ScienceBase, both original software and engineered components, is released as an open-source project to promote involvement from the larger scientific programming community both inside and outside the USGS.

## About Spectra Logic

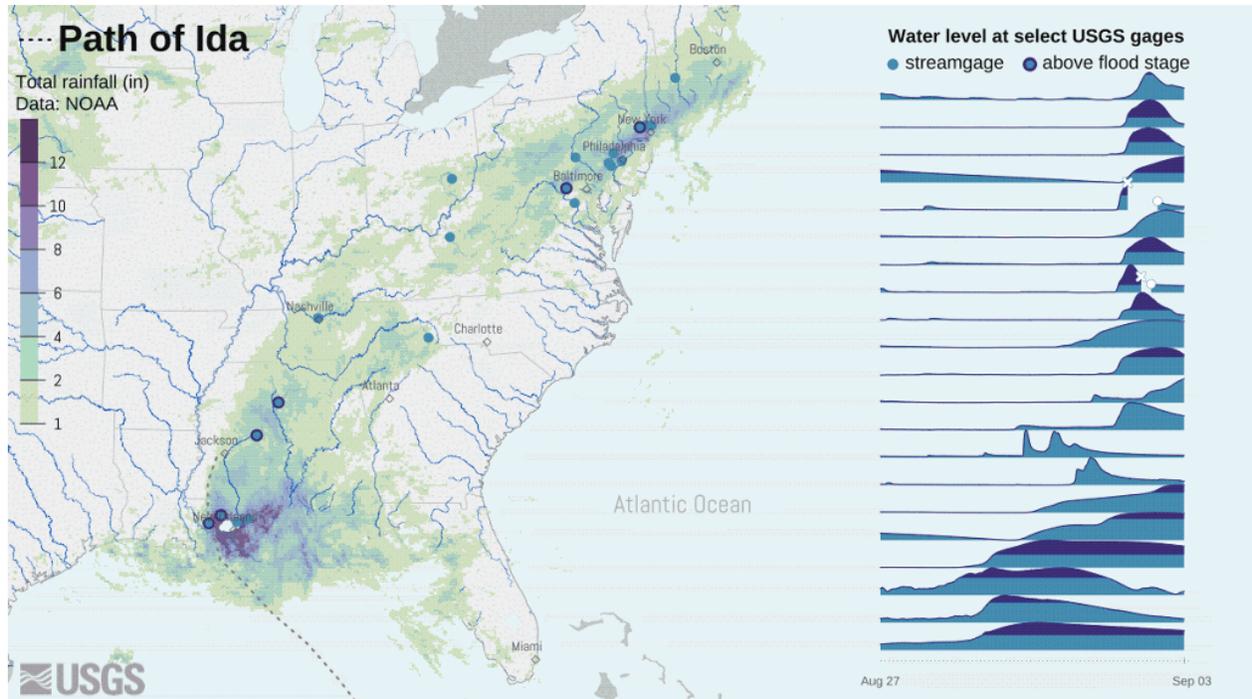
Spectra Logic develops a full range of Attack Hardened™ data management and data storage solutions for a multi-cloud world. Dedicated solely to data storage innovation for more than 40 years, Spectra Logic helps organizations modernize their IT infrastructures and protect and preserve their data with a broad portfolio of solutions that enable them to manage, migrate, store and preserve business data long-term, along with features to make them ransomware resilient, whether on-premises, in a single cloud, across multiple clouds, or in all locations at once. To learn more, visit [www.spectralogic.com](http://www.spectralogic.com).

## Data as Vast as the Planet

The USGS collects, transmits, stores and centralizes data from around the world for current and future scientific study. Disaster preparedness and response agencies depend on the timely access of data and results from the USGS and its supercomputers to help them plan and execute effective emergency services and responses to natural disasters.

Researchers use USGS data for long-term studies that affect how we live. The ScienceBase platform rivals current industry offerings in both the public and private sectors, making data available anywhere regardless of its physical location, while minimizing costs.





*Image: Hurricane Ida made landfall as a Category 4 hurricane in Louisiana and brought widespread precipitation and flooding along its path and up the northeastern coast of the U.S. in the following week. USGS streamgages provide critical information during storms to flood forecasters and emergency managers as they make decisions that contribute to protecting lives and property. Understanding river levels and locations of flooding can make a huge difference in these dangerous storms. The water footprint visualization shows patterns of precipitation and river discharge of 21 USGS streamgages in the path of Hurricane Ida.*

*Photographer: USGS Water Resources Mission Area*

To store, curate, and deliver data to other research teams, collaborators, and the public, the USGS has historically relied on various methods, ranging from the use of USB drives, FTP (file transfer protocol), to uploading through web forms - often accessing these resources required VPN (Virtual Protocol Network) connections. However, without VPN connections or sufficient VPN bandwidth (due to the nature of remote locations), remote collection sites were unable to successfully complete timely transfers, leaving critical and expensive data isolated and inaccessible, or worse, incomplete and possibly corrupt.



*Image: A telephoto image of fountaining from the western vent in Halema'uma'u crater, at the summit of Kilauea. Spatter from the fountain continues to build up a cone around the vent, which is almost entirely out of view from this angle. This photo was taken from the western crater rim on October 4, 2021. Photographer: USGS photo by M. Patrick.*

In an attempt to mitigate the problems caused by outlying, inaccessible data, the USGS evaluated the public cloud<sup>1</sup> as a means to supplement the existing data transfers over VPN connections. Unfortunately, the costs associated with using the public cloud proved confusing and excessively difficult to predict. In an attempt to defray costs, the USGS then explored distributing cloud expenses amongst scientists; however, based on project needs and anticipated data use, scientists were unwilling to risk unpredictable cloud costs.

---

<sup>1</sup> A public cloud is a platform that uses the standard cloud computing model to make resources - such as virtual machines, applications or storage - available to users remotely. Public cloud services may be free or offered through a variety of subscription or on-demand pricing schemes, including a pay-per-usage model.

Sarah Neenan, Kathleen Casey, and Alan R. Earls, "What Is Public Cloud? Everything You Need to Know," SearchCloudComputing (TechTarget, August 20, 2020), <https://searchcloudcomputing.techtarget.com/definition/public-cloud>.



## USGS Vision

- *This agency collects tens of terabytes per day from remote locations, various instrument types, and data formats. two HPC sites, 50 branch offices, +100 remote research locations*
- *Varied hardware, software and connectivity created high cost/labor to coordinate the system*
- *Considering going all-cloud for convenience, cost, and standardization with projects pulled down for compute*

Many of the USGS's main challenges stemmed from the daily volume of data amassed from their collection points (one site could produce upward of 10TB/day) paired with their diverse collection methods, at times across remote locations from many collection sites. To mitigate these challenges, a workflow was pursued that could ease the friction associated with collecting and delivering data within the community, as well as staging and preparing that data for analysis. Additionally, there was a need to curate this data for long-term retention in a central repository with multiple storage technologies. With a goal of supporting the transfer and collection of large data sets into primary staging areas or directly into a central repository, the USGS needed to establish a standards-based workflow that could capitalize on the USGS's cloud capacities in the [USGS Cloud Hosting Solutions](#) (CHS) team's AWS (Amazon Web Services) infrastructure. The new workflow needed to enable them to contain costs, and provide options for USGS science teams to keep "golden" copies of their data on USGS storage. Regardless of whether scientists process data in AWS or on USGS supercomputers, the goal was to establish common tools to initiate and manage data flows and to ease the burden associated with data movement.

## Environment Overview

The USGS infrastructure consists of two main data centers, remote offices throughout the United States, and multiple remote collection sites – where data is collected and then sent back to the main offices to be analyzed and stored. The main data centers are located in Denver, Colorado and Sioux Falls, South Dakota.

The Denver, Colorado site houses the USGS Yeti Supercomputer, whose sole purpose is running scientific workloads. It is a 3,728 core, 143 node system with approximately 1PB of high-speed storage. The USGS installed Spectra Vail®, a distributed multi-cloud data management software, along with a Spectra BlackPearl® system that stores data via object storage to a Spectra T950 library. Spectra Vail utilizes Shared Library Services to partition the T950 tape library, sharing it with other workloads, including traditional enterprise backup applications. The entire Spectra solution consists of approximately 6PB of total storage.

The second main site located in Sioux Falls, South Dakota is home to the USGS EROS Data Center and its two supercomputers. This location serves as the primary LandSat receiving station and primary location for LandSat data storage. The first supercomputer, Tallgrass, is a GPU-heavy system targeted at machine learning and deep learning workloads. The second supercomputer, Denali, is a CPU-only system designed for large-scale, multi-node simulation and data analytics with approximately 2.75PB of high-speed Lustre storage.



*Image: Photo of the USGS Denali supercomputer in Sioux Falls*

Another instance of the Spectra solution is installed at the Sioux Falls site, including a Spectra Vail, Spectra BlackPearl and a Spectra TFinity tape library. The entire Spectra solution provides ~5PB of storage. As with the Denver site, the Spectra Vail instance at Sioux Falls also uses Shared Library Services, which allows other applications to share the TFinity tape library.





*Image: The USGS's Spectra TFinity tape library at their Sioux Falls location.*

Operating in all 50 states and U.S. territories, the USGS has remote offices (often multiple) in every state. Many sites are located in remote areas without infrastructure available. Most have minimal internet connection and limited bandwidth. Additionally, most do not have data centers or event storage servers. In an attempt to create targeted strategic sites, the USGS implemented Spectra Vail virtual machines (VM) to act as “data concentrators” for other sites in the regions. For example, offices and sites on the west coast of the U.S. connect to a Vail node in Sacramento or Mountain View, CA. Either of those sites will then have a backbone connection to one of the main sites in Denver, Sioux Falls or Reston, Virginia (USGS Headquarters). Each of the data concentrator sites will consist of one Vail node in the form of a single VM instance, and disk storage with local VM storage (via the hypervisor).

Remote collection sites are located in the field - often in locations where natural disaster events have taken place (e.g., near a volcano, glacier, river, mountain). Data is captured via different sensors, such as unmanned aerial vehicles (UAV), cameras, fixed wing and rotary aircraft with cameras, LIDAR equipment, EM sensors, etc. After collection, researchers often return to their hotel room or to a location with a basic internet connection to upload data and remit it to the main offices to be analyzed and stored. Previously, the upload of data was achieved by using a VPN connection. The overhead of the VPN connection on a less than optimal network connection often led to incomplete or corrupted data transfers. With the new solution implementation, preliminary tests using AWS S3 storage demonstrated a significant performance improvement over the previous method - performing transfers via VPN.





## Challenges

- *Unpredictable cloud costs*
- *Decisions about which departments should pay for storage and access charges*
- *Amount of data being generated and stored growing exponentially*
- *Separate onsite and cloud workflows creating high costs and complexity*
- *Data retention requirements*

## Implementing a Hybrid Storage Solution

With the evolution of new computing methods that include employing on-premises private cloud<sup>2</sup> platforms as well as expanding public cloud offerings, the timing was perfect for the USGS to consider one or both of these methods to enhance their storage ecosystem and establish a solid data workflow. The USGS sought a solution that could combine a private cloud storage infrastructure with public cloud storage services to form a single efficient storage solution with the flexibility to leverage both cloud services and on-premises vendor-agnostic infrastructure options.

In researching options, Spectra Vail was identified as a technology that could help with some of the specific needs of the use case. As mentioned, Spectra Vail is a distributed multi-cloud data management software that provides a single global namespace for any combination of public cloud and on-premises storage. Vail uses public cloud services for fully integrated command and control. Vail consists of an S3-compliant interface that can unite any combination of on-premises storage with public cloud storage, and includes unlimited sites, objects, capacities, users and speeds. In addition, it incorporates a single global namespace with temporal and location policies. Vail provides cloud-based command and control operations, consisting of a cloud-based management portal across all storage, as well as automatic data management, which includes replication, data migration, data tiering (to disk, tape, etc.), encryption, compression, control and protection.

---

<sup>2</sup> Private cloud is a type of cloud computing that delivers similar advantages to public cloud, including scalability and self-service, but through a proprietary architecture. A private cloud, also known as internal or corporate cloud, is dedicated to the needs and goals of a single organization whereas public clouds deliver services to multiple organizations. <https://searchcloudcomputing.techtarget.com/definition/private-cloud>

Ben Lutkevich, "What Is Private Cloud?," SearchCloudComputing (TechTarget, December 6, 2019), <https://searchcloudcomputing.techtarget.com/definition/private-cloud>.



Vail consists of two main components: 1) a Vail node, which acts in several ways - as an access point or as an S3 endpoint; as a data mover; as a local access to the main controller; and as local storage, and 2) the Vail Sphere, a management server that acts as the command and control operator.

The USGS understands that service offerings from public cloud vendors (e.g., AWS, Microsoft, Google, etc.) provide flexibility and agility; however, it is also understood that using a public cloud service provider to serve petabytes of data per month to the public can come with significant cloud storage and access charges. The USGS struggled with balancing those costs with the benefits provided by the cloud. Spectra was able to demonstrate to the USGS that a deployment of Vail would enable USGS to garner the benefits of public cloud agility with the cost effectiveness of using existing facilities and storage equipment to offer scientific data delivery as a public service to the world.

The Spectra Vail solution fit the USGS's desired workflow and provided it with a hybrid cloud solution that paired all the strengths of the public cloud with the benefits of an on-premises solution. Additionally, Vail allows them to utilize the public cloud services and distribution model with on-premise data. Vail brings data that has entered the sphere of influence at the USGS into a single management view, where it can be managed and accessed regardless of its physical location, providing the USGS with a centralized data storage solution. By connecting all sites into a distinct storage sphere, including to the AWS public cloud, Vail creates an easy-to-manage hybrid cloud workflow for the collection, distribution and storage of USGS data.

With Vail, the USGS researchers can leverage any combination of public cloud and on-premises datastores and match the value of data at any point in time, while managing it based on the elasticity and reliability of AWS cloud services. This capability builds in a layer of flexibility for the overall USGS strategy for scientific data storage by supporting more diverse workflows and decreasing the risk of vendor lock-in for any single cloud or on-premise storage technology.

The USGS was able to implement an end-to-end data storage and delivery service that:

- Allows data collection from instruments and contributors from all over the world using a combination of Amazon Web Services (AWS) S3 and USGS servers depending on conditions at any collection point.
- Stores data in AWS and on-premises, according to a policy that meets service level agreements or science requirements.
- Allows other agencies, collaborators, and the public to access data via a combination of AWS and on-premises storage according to access policies, while managing access costs - to meet the needs of the data consumers and their budget resources.
- Keeps historical data in a multi-site, durable, and affordable datastore that can be accessed in minutes for no incremental cost beyond storing it.



## How Data Moves from Site to Site

Data is the primary asset of the USGS and virtually everything the organization does revolves around it. Data is collected from remote sites across the globe by different platforms, ranging from satellites to action cameras. After capture, data is further analyzed at various remote sites and developed into scientific findings and ultimately released as research papers and data sets for current response and future policy decisions.

The data lifecycle begins in the field where scientists gather and collect information, such as observing lava flows from an erupting volcano, exploring the deepest crevices in a moving glacier, or collecting electro-magnetic data in a bid to see below the Earth's surface. At the end of each collection, the scientist connects to a wireless or cellular internet connection, and transfers data to the primary data centers to be staged for analysis. This is accomplished by using the Amazon public cloud where scientists upload current data collections to a bucket in the AWS public cloud. When this happens, the upload triggers a Vail lifecycle rule such that all data that is uploaded into the AWS bucket will be synched to a local onsite storage bucket in one of the main USGS data centers, then replicated to the second main data center for disaster recovery (DR) purposes. With this workflow, scientists are able to focus on research during the day and transfer all data collected into an AWS bucket with one of the hundreds of S3 compatible tools available. The USGS then uses any of several S3 compliant applications to manually transfer all of the day's data into another AWS bucket that will be synched to local storage. Data that is stored in or downloaded from an AWS bucket is subject to storage and egress fees, but once the data is at the USGS offices, no further egress fees are incurred. This approach provides a hybrid cloud storage workflow that ties in directly with the USGS private cloud.

Remote offices with a secure and reliable connection can configure a Vail node to transfer data to primary offices based on the lifecycle policy associated with each bucket at those remote offices. This method bypasses the public cloud and eliminates the traditional egress charges incurred when the remote collection sites download data out of the public cloud by centralizing data at primary data centers.

Metadata is logged into the cloud management server to track and identify objects in all locations, and the data is stored in multiple on-prem and cloud locations. By using metadata to know where the data is located, data can be stored and transferred from any storage location, minimizing or eliminating storage and egress fees.

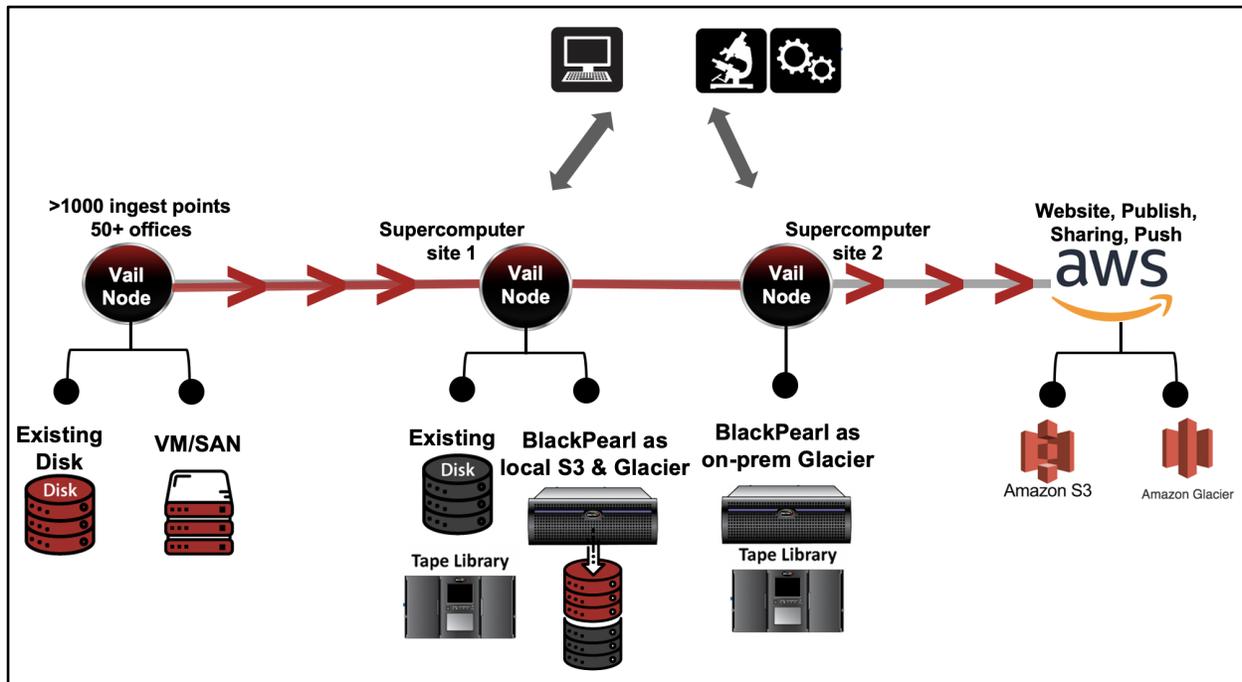


Image: USGS workflow

At the science centers or primary data centers, data is collected and consolidated. At that point analyses can be performed against the data. The USGS leverages their existing supercomputers to run analytics and interrogate data that has been collected and staged on high-speed storage systems such as their Lustre file systems. This produces detailed findings and results including, but not limited to, predicting future volcanic eruptions, hurricanes, earthquakes, and other natural disasters.

After data and analyses are completed, most data must be stored for long-term retention and preservation. The USGS has a no-delete policy for finalized datasets. Vail provides this by implementing a local, on-prem glacier\*, or cold storage repository, where infrequently accessed data can be kept, accessed when needed and protected for the long-term, on the USGS's prior installations of a Spectra Logic storage infrastructure, including a BlackPearl platform and Spectra tape library. This tier of storage acts as a local glacier but with the free retrieval of data. In addition to cost savings, local glacier tiers are accessed immediately, meaning restoration can begin within minutes (not hours as with cold public cloud storage).

With this data flow, the USGS is able to effectively create its own private cloud that acts as a single storage platform for scientific data, while leveraging the value and flexibility of the public cloud. Additionally, maintaining multiple AWS accounts and buckets across different cost centers allows storage usage to be compartmentalized and billed accordingly.



## Solution Benefits

- Vail brings additional flexibility of storage lifecycle management, including two synchronized copies retained forever
- Combines on-premises and AWS storage into a single platform managed in a highly available, elastic cloud like experience
- Automates storage policies to ensure data is placed geographically and in tiers according to plan
- Leverages free ingress offered by AWS and the agency's On-Ramp access for populating on-premises storage
- Match SLAs to economic value of data without compromising access by using Spectra Vail, BlackPearl and TFinity library
- Simple publishing of final data to publicly available buckets in the cloud
- Researchers can ingest data from any location – even directly to the cloud

After USGS researchers have analyzed their data, and when results are ready to be published and distributed to the scientific community around the world, many scientists publish their data on ScienceBase, a data platform developed and maintained by the bureau. As a trusted USGS data repository, ScienceBase serves a role in the USGS helping store, curate, and serve published data sets (in addition to in-progress and to-be-released data sets).

Vail acts as the data management and transfer application that moves the finished data from the research data center to the ScienceBase distribution platform. As data ages and is accessed less often, it is moved out of the distribution platform and kept for long-term retention in the local glacier archive repository.

## Utilizing BlackPearl and Object Storage

The USGS research computing group needed an avenue for moving data to and from their Lustre file systems. The system needed to integrate with applications that catalogued and managed delivery of the data; be accessible remotely via standard protocols; and be capable of sustained high transfer rates to the Lustre file system. By leveraging standard protocols, application development along with user training could be streamlined (with knowledge of AWS S3, the transition to any object platform is straightforward).

The BlackPearl is an object storage platform that stores data on tape, allowing the USGS to take advantage of its existing tape library, drives, and media. This provides a media migration path as new media becomes available and expands storage capacity cost-effectively as needed (often without increasing power draw or data center footprint). BlackPearl provides high-speed network connections

(bonded 40Gb/s ethernet) and NVMe-based cache to deliver the necessary transfer speeds to move data between the high-speed Lustre storage system and BlackPearl.

BlackPearl provides the foundation for location-agnostic high-capacity storage that meshes with cloud-based methods, workflows and technologies. It provides options for pre-built clients and software development kits in several common languages (or raw RESTful API access) that gives the USGS scientist and developers access to a cost-effective storage platform.

## Utilizing Vail and Hybrid Cloud Storage

With the current computing environment evolving to include cloud computing services, it is fundamental that the USGS develop methods and processes that allow seamless, frictionless movement of data between cloud and on-premises storage platforms – whether processing data on HPC equipment, using cloud services for analytics, or delivering scientific data and products to users around the world.

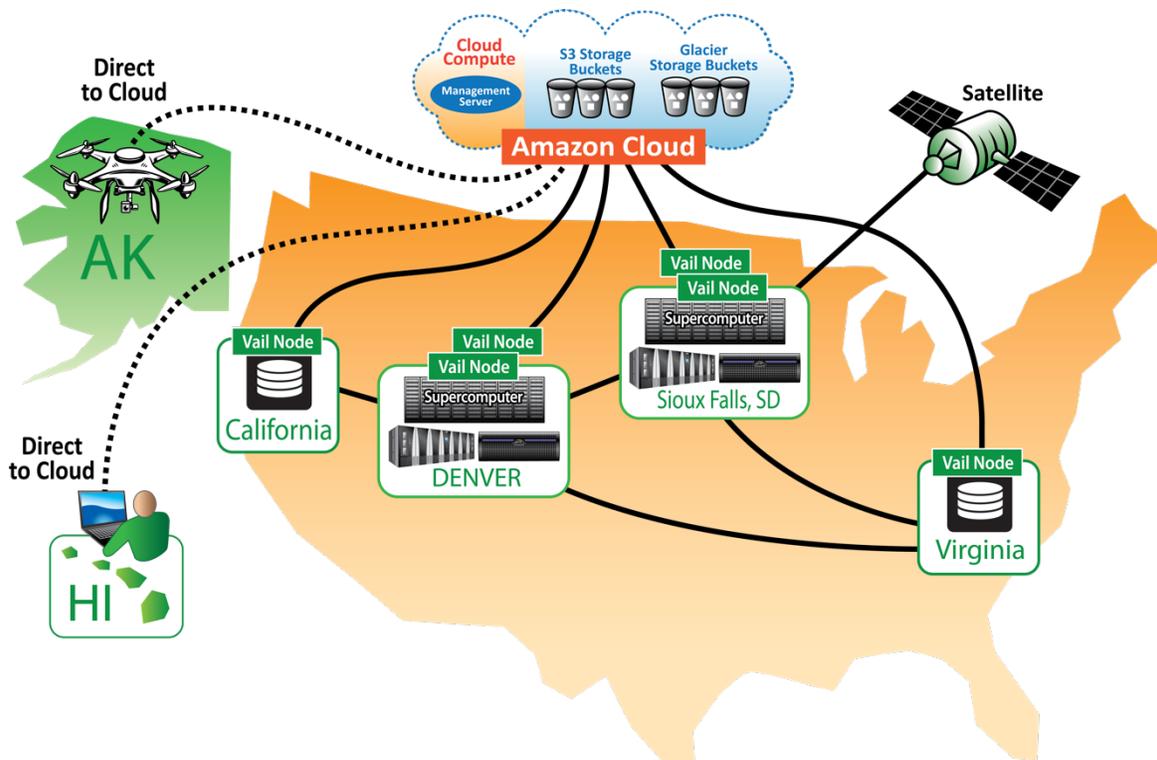


Image: Graphical representation highlighting USGS's Spectra Vail use case



## USGS Exploratory Work

While early in its implementation of Vail, the USGS is exploring numerous tools and techniques to access data in Vail.

Notable wins include:

- Scripts written using the AWS-CLI (<https://aws.amazon.com/cli/>) have worked without any code changes - only configuration items changed
- Using Goofys (<https://github.com/kahing/goofys>) to mount Vail buckets as FUSE mounts on data movers
- A single command to restore from on-premises glacier and download files
- Using Cyberduck as a generic data movement interface.

Additionally, the USGS is in the process of testing several tape libraries to offer additional programmatic interfaces to Vail.

## Additional Resources

- [Spectra Vail](#)
- [Spectra TFinity](#)
- [Spectra BlackPearl](#)



## About Spectra Logic Corporation

Spectra Logic develops a full range of Attack Hardened™ data management and data storage solutions for a multi-cloud world. Dedicated solely to data storage innovation for more than 40 years, Spectra Logic helps organizations modernize their IT infrastructures and protect and preserve their data with a broad portfolio of solutions that enable them to manage, migrate, store and preserve business data long-term, along with features to make them ransomware resilient, whether on-premises, in a single cloud, across multiple clouds, or in all locations at once.

To learn more, visit [www.Spectralogic.com](http://www.Spectralogic.com).

Copyright ©2021 Spectra Logic Corporation. All rights reserved worldwide. Spectra and Spectra Logic are registered trademarks of Spectra Logic. \*Amazon Glacier® is a registered trademark of Amazon Technologies, Inc. All other trademarks and registered trademarks are property of their respective owners. This document is provided for information purposes only, and the contents hereof are subject to change without notice. This document is not warranted to be error-free, nor subject to any other warranties or conditions, whether expressed orally or implied in law. We specifically disclaim any liability with respect to this document, and no contractual obligations are formed either directly or indirectly by this document.

303-449-6400 • 800-833-1132 • 6285 Lookout Road • Boulder, CO 80301 USA • [spectralogic.com](http://spectralogic.com)

V1-10142021