

# Research Computing and Cloud Resources at GWU: Use Cases and Solutions

CLARK GAYLORD

DIRECTOR, RESEARCH TECHNOLOGY SERVICES

[CGAYLORD@GWU.EDU](mailto:CGAYLORD@GWU.EDU)

# Cloud consultancy and deployments

- NUMEROUS ENGINEERED SOLUTIONS USING CLOUD RESOURCES.
- RETAIN EACH ENVIRONMENT'S CONFIGURATION FOR REUSABILITY, RECONFIGURATION, REDEPLOYMENT
- ENGAGE RESEARCHERS IN CONSULTANCY AND RESEARCH SOLUTIONS
- EXAMPLES:
  - CONTROLLED UNCLASSIFIED INFORMATION ENVIRONMENT
  - DEPARTMENTAL RESEARCH COMPUTING AND DATA PROCESSING ENVIRONMENT
  - DYNAMIC PROVISIONING SYSTEM FOR MORE COMMON USE CASES
  - PROJECT-SPECIFIC DATA MANAGEMENT, DATABASE HOSTING

## Protected data environment

- CONTROLLED UNCLASSIFIED INFORMATION (SIMILAR FOR HIPAA, ETC.)
- BIOMEDICAL RESEARCH FOR DEFENSE APPLICATION
  - “SECURE WORKSTATIONS” IN CLOUD
  - CONNECT WITH LABORATORY INSTRUMENTS
- PROJECT REQUIREMENT FOR DATA SHARING WITH OTHER INSTITUTIONS
  - FEDERATED IDENTITY: GLOBUS
  - SEPARATE AWS VPC FROM THE WORKSTATION ENVIRONMENT

## Department/lab research computing

- SEVERAL APPLICATIONS, SERVERS, WORKSTATIONS
- SUPPORTING VARIOUS DATABASES, TEACHING
- APPLICATIONS AND WORKFLOWS: R STUDIO SERVER, SHINY, GALAXY
- INSTEAD OF SEVERAL PHYSICAL SERVERS OR VIRTUAL MACHINES IN DATA CENTER, HOSTED IN AWS CLOUD, SOME ON-DEMAND

## “Dynamic cloud platform”

- BASED ON AMAZON’S SERVICE WORKBENCH
  - USEFUL FOR SPECIFIED TEMPLATES TO DEPLOY
  - SOME DATA SHARING BETWEEN DEPLOYMENTS CAN BE SUPPORTED
- STAND-ALONE VMs
  - LONG RUNNING INSTANCES, PERHAPS NEEDING DESKTOP EXPERIENCE (NOT BATCH FRIENDLY)
  - STUDENT WEB SITES
- TEMPLATED SOLUTIONS ARE THE KEY
  - COULD INCLUDE HPC OR GALAXY

## Project specific database hosting

- MANY PROJECTS FOR A SINGLE RESEARCH PROGRAM
- PROJECT-SPECIFIC REQUIREMENTS IN TERMS OF PERFORMANCE, COST, RETENTION
- OFTEN SECURITY CONTROLS ARE REQUIRED
  - CONFIGURABLE NETWORK ACCESS
- LEVERAGE *AWS RDS* WITH ANCILLARY COMPUTE AND STORAGE RESOURCES AS NEEDED

## Various other example research platforms

NOTE: ALL THESE ARE "ON PREMISE" SUPPORT EFFORTS

- MULTI-TB MYSQL SERVER FOR DATA MINING OF LARGE DATASET (VISUAL COGNITION)
- SPECIALIZED STORAGE AND VIRTUAL DESKTOP ENVIRONMENT (SCHOOL OF PUBLIC HEALTH)
- DATABASE AND COMPUTE ENVIRONMENT FOR ORCHESTRATED PIPELINES (BIOCOMPUTE)
- LAB SPECIFIC ENGINEERING SUPPORT (NANOFABRICATION AND IMAGING CENTER)
- NUMEROUS "BOUTIQUE" CLUSTERS DEDICATED TO RESEARCH GROUPS OR LABS (ENGINEERING, ARTS & SCIENCES, MEDICINE)



# SURA WORKSHOP TULANE UNIVERSITY



Tulane



# SAAS BROKER

## Researcher Need

- NIH recognizes the value of conducting studies across **multiple centers and** mine research data holistically
- Tulane Primate Center was awarded to act as a **data coordination center for a multi-site/center study** working on data collection and analysis platform for COVID non-human subject research
- Task Order needed a **highly secured platform** stood up within **2-3 months (FEDRAMP/FISMA)**

## Approach

- Tulane IT needs a **highly secured full stack Electronic Data Capture platform** with little lead time.
- **SaaS vendor** is the only solution with the **rich requirements in features and cybersecurity and short timeline.**
- IT will provide resources to administer the platform going forward after the 1st study as a **full managed service.**

## Technical Delivery

- Tulane had already completed an in-depth review of different EDC vendor choices.
- After selecting a vendor, it brokered a relationship with the SaaS provider to host a unified data collection platform with the required security to allow sharing research data across the institutions.
- Professional services collaborated with Tulane Primate Center PI & IT and delivered the solution in less than 2 months

## Business Value

- This managed clinical research SaaS service allows a faster start up for future clinical research opportunities with predictable cost structure.
- Tulane can provision new studies in weeks instead of months because of the existing approved master services agreement & skillset.
- Becomes a solid stepping-stone to create an RFP to extend this solution to become a Tulane wide Life Science Clinical Full-Service Platform





Tulane



# CMMC In The Microsoft Cloud

---

Secure Enclave in Microsoft Azure



THE UNIVERSITY OF  
SOUTHERN  
MISSISSIPPI®



# USM Research Requirements

---

The University of Southern Mississippi in preparation for Cybersecurity Maturity Model Certification (CMMC) requirements including those that will be required per the DoD DFARS 252.204-7012 Interim Rule has initiated this project to enable researchers to continue supporting DoD and other initiatives that will require this compliance.





# USM Research Requirements

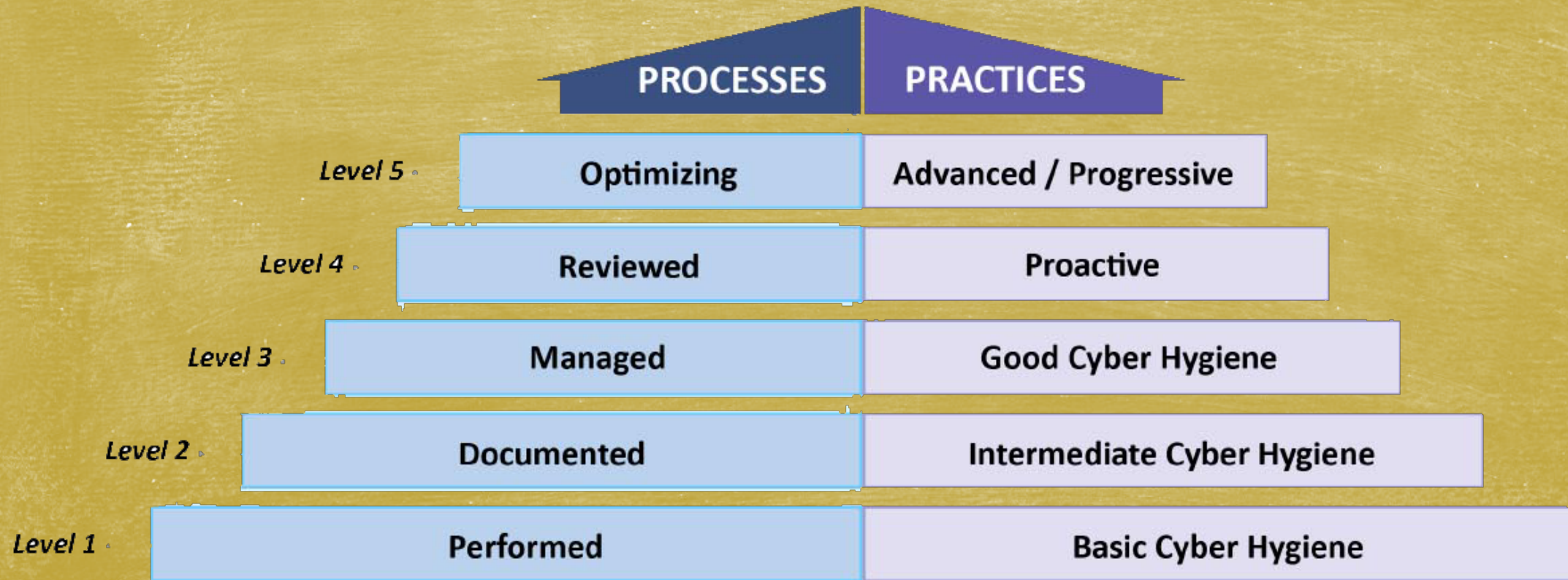
---

The University of Southern Mississippi is seeking a scalable Microsoft Office 365 GCC High messaging and collaboration platform for researchers requiring CMMC and NIST 800-171 compliance, enabling them to accept federally funded grant awards subject to these requirements. The platform **MUST** be DFARS -7012 and CMMC Level 3 compliant.





# CMMC Levels



<https://www.thecoresolution.com>



THE UNIVERSITY OF  
**SOUTHERN**  
**MISSISSIPPI**



# CMMC Level 3 Chosen

---

Processes: Managed

Level 3 requires that an organization establish, maintain and resource a plan demonstrating the management of activities for practice implementation.

The plan may include information on missions, goals, project plans, resourcing, required training, and involvement of relevant stakeholders.

<https://www.preveil.com/blog/what-are-the-5-levels-of-cmmc/>



THE UNIVERSITY OF  
**SOUTHERN**  
**MISSISSIPPI**



# CMMC Level 3 Chosen

---

## Practices: Good Cyber Hygiene

Level 3 focuses on the protection of CUI and encompasses all of the security requirements specified in NIST SP 800-171 as well as 20 additional practices to mitigate threats. Any contractor with a DFARS clause in their contract will need to at least meet Level 3 requirements. Note that DFARS clause 252.204-7012 applies and specifies additional requirements beyond NIST SP 800-171 security requirements such as incident reporting.

<https://www.preveil.com/blog/what-are-the-5-levels-of-cmmc/>



THE UNIVERSITY OF  
**SOUTHERN**  
**MISSISSIPPI**



# USM Research Requirements

---

DoD contracts currently require that projects implement controls to safeguard government owned data. USM is not easily able to comply with the government requirements in our current network environment.

Two possible approaches for addressing the information security requirements around DoD contracts are as follows:

- USM can physically deploy a separate network on campus
- USM can setup a secure enclave on the internet





# USM Bid Process

---

The University of Colorado, Boulder is seeking qualified quotes from Bidders for the CMMC MS O365 GCC-High Tenant Project

Microsoft O365 GCC High Platform-Office 365 Govt G5 License Procurement & Fast Track Implementation



University of Colorado

Boulder | Colorado Springs | Denver | Anschutz Medical Campus



THE UNIVERSITY OF  
SOUTHERN  
MISSISSIPPI



# USM Bid Process

---

<https://www.coloradobids.us/colorado-bids/bids-NBD15675706752464107.htm>

Search for "University of Colorado" "PSC-U-1411"

<https://oit.colorado.edu/oit-projects>



University of Colorado

Boulder | Colorado Springs | Denver | Anschutz Medical Campus



THE UNIVERSITY OF  
SOUTHERN  
MISSISSIPPI



# USM Vendor Selection

---



US Partner Award 2020



2 Parade St NW  
Huntsville, AL 35806  
256.585.6868  
info@summit7.us  
cmmc@summit7.us

<https://www.summit7.us>



THE UNIVERSITY OF  
**SOUTHERN**  
**MISSISSIPPI**



# What Is Included In The Solution

---

- Office 365 GCC High Licensing – recurring annually
- Azure Government Subscription – recurring monthly
- Office 365 GCC High Backup – recurring annually
- CMMC Level 3 Windows 10 MFA – setup cost, then recurring annually
- CMMC Level 3 Office 365 GCC High – one time cost to configure tenet
- CMMC Level 3 Secure Enclave – one time cost to architect enclave





# First Things First

---

- Apply for Office 365 Government GCC High
- This process can take 15 business days to receive a Category 3 eligibility from Microsoft, once received it is valid for 3 years
- There are 3 types of documentation you can submit for eligibility:
  - Submit a contract that calls out ITAR, Export Controlled, CUI or DFARS 7012 as requirements
  - Submit a Department of State DS-2032 form
  - Sponsorship letter signed by a government official





# Project Kickoff

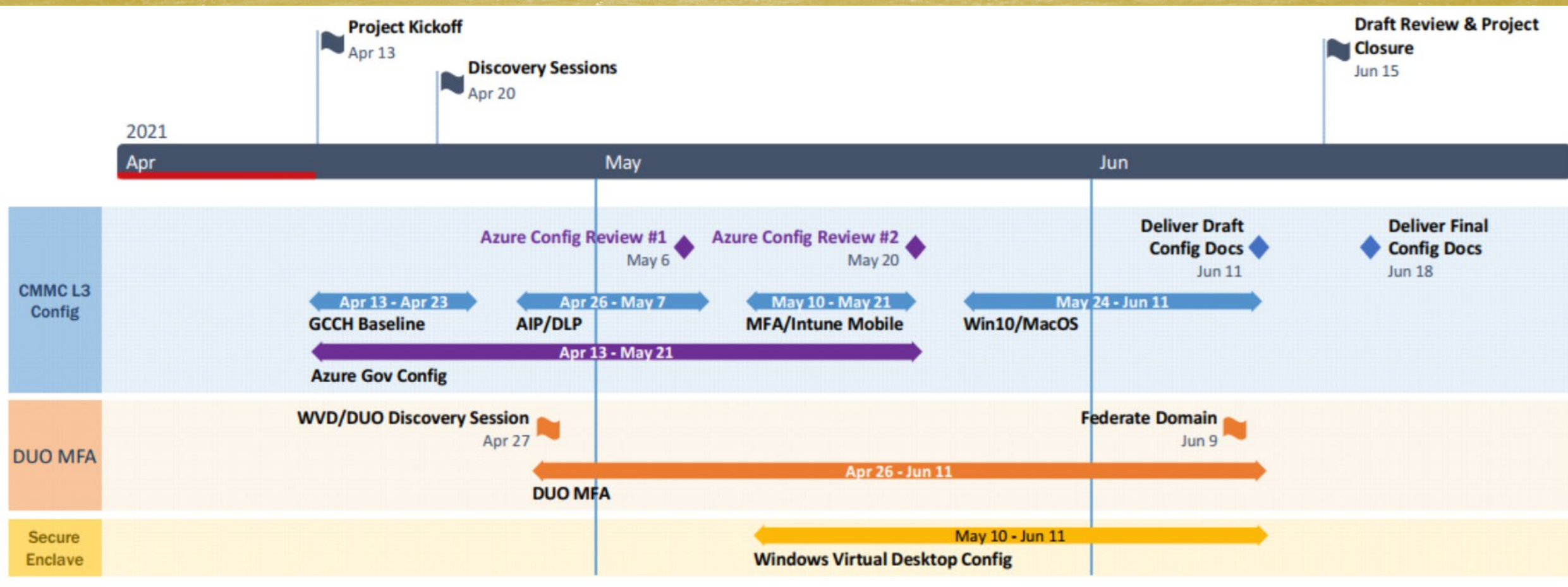
---

- Select and create a name for your GCC High tenant
- Multi-Factor Authentication Setup
  - DUO Security Federal Configuration
  - DUO ADFS & Federation
- Cloud Backup for Office 365
  - AvePoint Federal Cloud Service





# Project Timeline





# Office 365

---

- Office 365 Tenant CMMC Level 3 Configuration
  - Office 365 Tenant
  - Azure Identity Management
  - SharePoint Online
  - Exchange Online
  - Advanced Threat Protection (ATP)
  - Intune Enterprise Mobility (Mobile Devices: Phones, Tablets)
  - Intune Windows 10 Devices (Computers: Laptops, Desktops)
  - Azure Information Protection
  - Label Creation (CUI, Company Proprietary, [External] Proprietary, Public)





# Azure

---

- Azure Tenant CMMC Level 3 Configuration
  - Deploy Azure CMMC L3 Baseline Configuration
  - Deploy and Configure Next Generation Cloud Firewall - Azure Firewall
  - Deploy and Configure 3 Virtual Networks, 4 Network Security Groups
  - Build 2 Domain Controllers to CMMC L3 Standard
  - Build AD Connect Server and Sync with Azure Active Directory
  - Deploy and Configure Office 365 GCC High Backup (AvePoint)
  - Implement the Approved AD Design to CMMC Level 3





# Contact Information

---

David J. Sliman

david.sliman@usm.edu

601-266-4190

<https://www.usm.edu/itech/>



THE UNIVERSITY OF  
SOUTHERN  
MISSISSIPPI



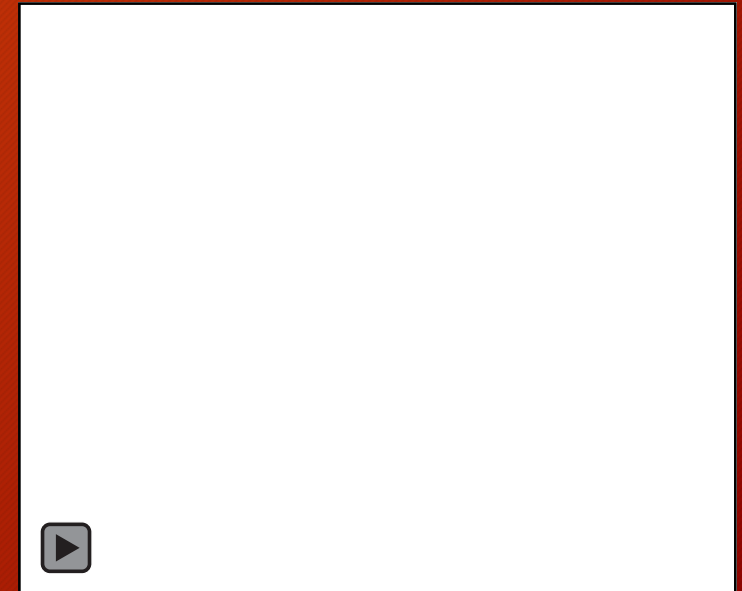
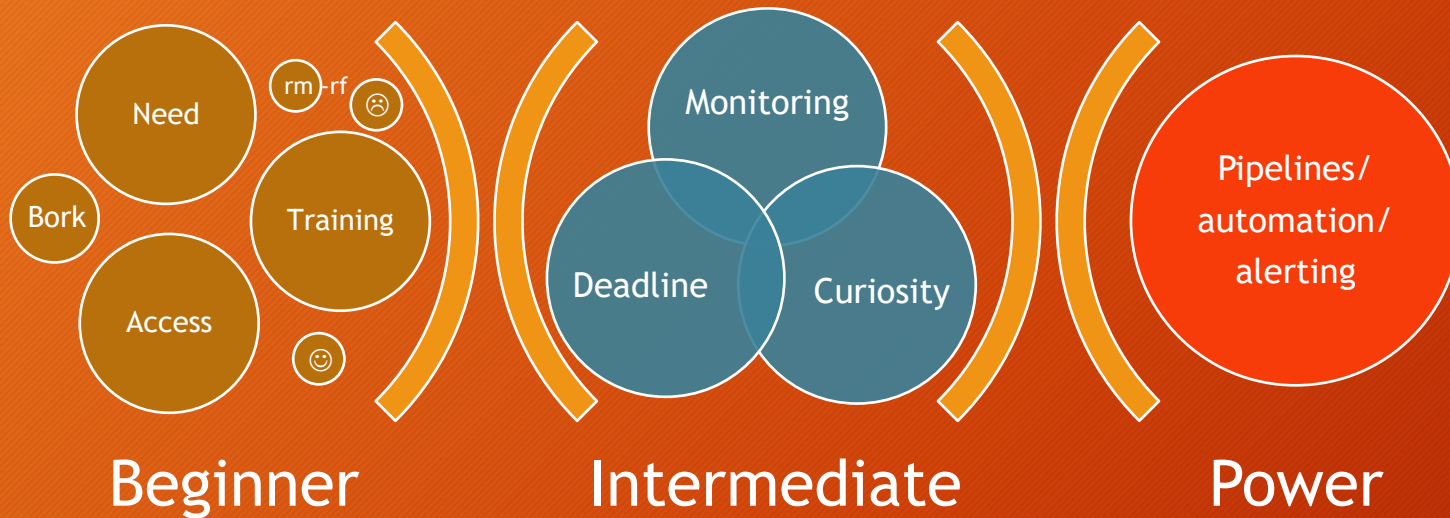


# Open OnDemand as a single pane of glass

Robert Settlege  
Advanced Research  
Computing @Virginia Tech  
April 2021

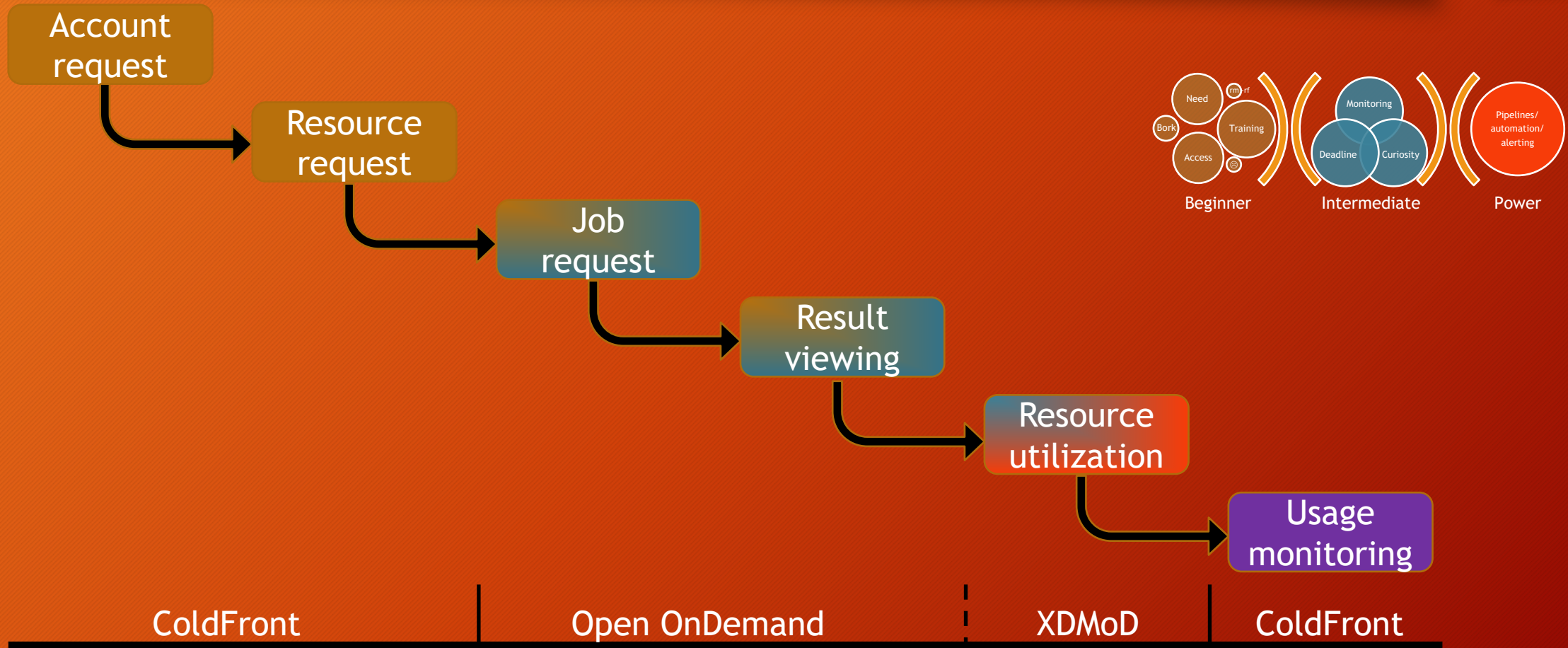


# User spectrum and experience





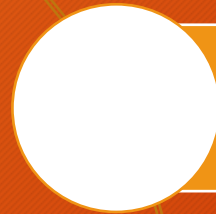
# GOAL: Facilitate and speed the progression



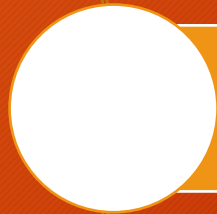


# Unify the user experience

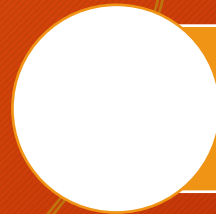
Open  
OnDemand



HPC



Kubernetes

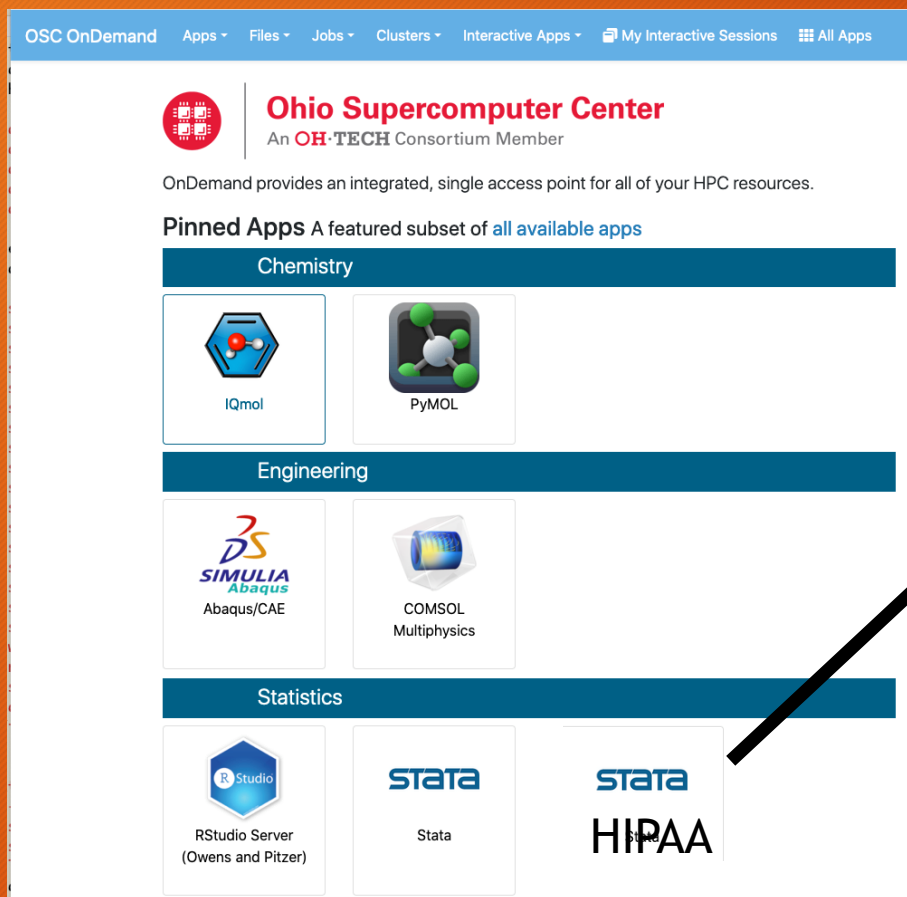


Cloud

CUI/HIPAA...

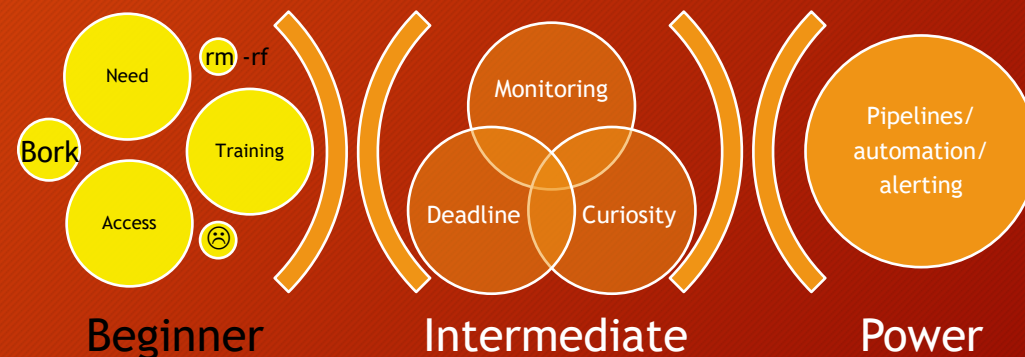


# Open OnDemand 2.0 - user perspective



What and how matter. Where is a distraction.

- Browser based HPC access
- Zero user install
- Site customizable
- ACL controls visibility/access to apps - CUI/HIPAA/Class/....





# Where is a (sometimes) optional choice, with consequences

**GAMS Studio**  
 This app will launch [GAMS Studio](#) on [TinkerCliffs](#) and [Infer](#) -- coming soon.  
 The underlying desktop container is based off the awesome Xfce image by [accetto](#).

**Cluster**    **ADD \$\$ projections**

tinkercliffs

**use container**

Container: accetto-ubuntu-vnc-xfce-firefox-g3-Feb2021.sif

This defines the remote desktop version to run.

**Account**

openondemand2

- The allocation you would like to use for SLURM.

**Reservation**

**Partition**

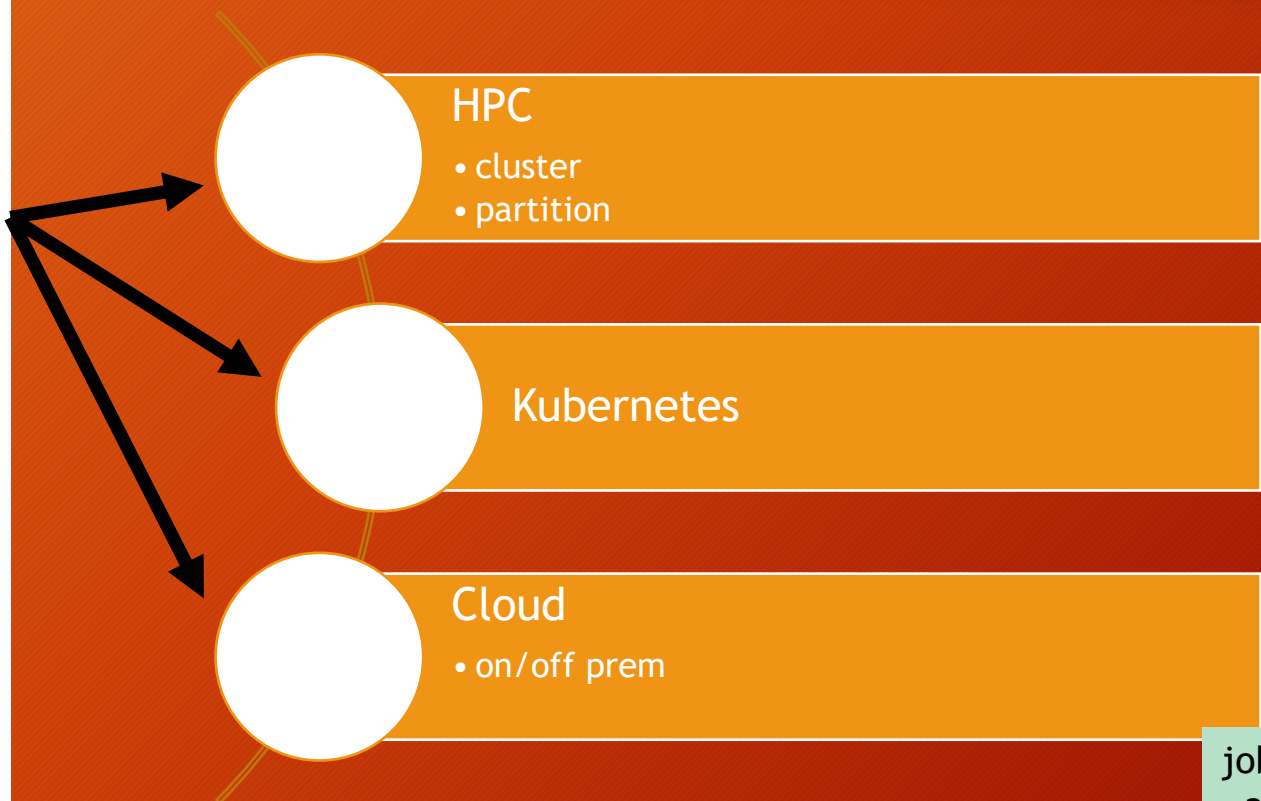
normal\_q

- To request a GPU enabled queue, preface it with v100\_. Example: v100\_normal\_q, t4\_normal\_q

**Number of hours (min-1, max-48)**

1

**Number of nodes (currently confined to 1)**



```
job:
  adapter: "kubernetes"
  cluster: "ood-dev"
  bin: "/usr/local/bin/kubectl"
```



# Open OnDemand: running jobs

**MATLAB (11421211.owens-batch.ten.osc.edu)** 1 node | 1 core | Running

Host: [>\\_o0467.ten.osc.edu](#) Delete

Created at: 2020-10-16 18:11:57 EDT

Time Remaining: 1 hour

Session ID: [5d9be97c-6480-4e0d-b27d-248aebec1e97](#)

noVNC Connection Native Instructions

**Compression** Image Quality

0 (low) to 9 (high) 0 (low) to 9 (high)

[Launch MATLAB](#)

**GAMS Studio (140342)** 1 node | 1 core | Running

Host: [>\\_tc056](#) Delete

Created at: 2021-04-22 07:29:37 EDT

Time Remaining: 59 minutes

Session ID: [2a99745c-601c-4ea9-a29c-dd33ea654668](#)

["Connect to VNC - /home in Documents"](#)

**Adding resource utilization metrics:**

- \$\$ equivalents
- Efficiency metrics
- This job is \$0.5/hour
- Your job is running at 50% efficiency

**Matlab (139321)** Completed

Created at: 2021-04-20 14:52:13 EDT

Session ID: [5ebcfc0f-992b-4914-9521-67ad95d80225](#) Delete

- This job cost \$1.00
- Your job was 50% efficient (link)

For debugging purposes, this card will be retained for 5 more days

- Reverse proxy setup automated
- NoVNC
  - Image/connection settings
- Direct SSH to node
- View only sharable link

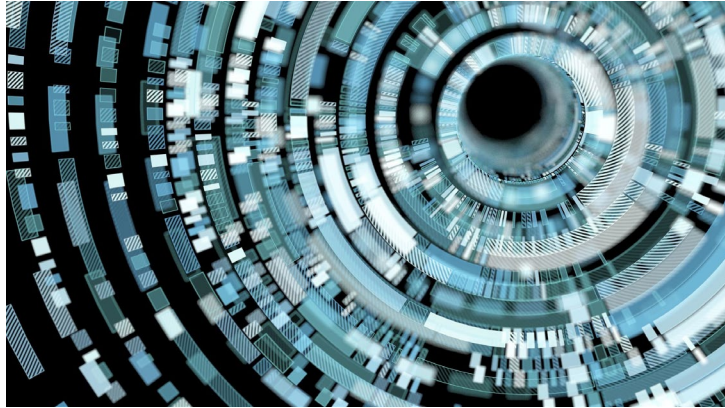


# OPEN OnDemand

Questions?

Robert Settlage  
Advanced Research  
Computing @Virginia Tech  
Oct 2020





# RESEARCH COMPUTING IN THE CLOUD

EXAMPLES OF CLOUD USAGE FOR RESEARCH AT USC



UNIVERSITY OF

**South Carolina**



# Research Computing using Cloud Services

- Campus IT uses AWS and Azure
  - Mostly storage for AWS, Azure for Microsoft Services
- UofSC Research Computing usually Cloud agnostic (AWS, GCP, Azure, OCI)
- Challenges in billing end users
- As a state institution, required to use resellers (Onyx, Enquizit)
- Challenges in building cost models for charging grants
  
- Free resources for teaching AWS  
<https://aws.amazon.com/education/awseducate/>



**Cloud Compute**



# SAS Viya in GCP for HIV ML Analysis



## Need:

Researcher to explore and identify new patterns/combinations of predictors associated with HIV medical care utilization.



## Approach:

Ingest and analyze data sets by combining both Statistical and Machine learning approaches.

## Technical Delivery:

Delivered SAS Viya package SaaS on GCP by building core infrastructure utilizing automation and custom machine types to minimize cost.



## Business Value:

HIV researcher can simply sign up, log in and get to work and tune the solution to their requirements and focus on solving analytic challenges and quickly realize value. Overall solution cuts infrastructure costs by providing on-demand allocations and services.



# Delivering Secure Healthcare Insight in AWS



## Need:

Provide analysis, core data, and visualizations to healthcare providers, insurance companies, and state government agencies to guide health care decisions, optimize cost. Includes massive healthcare provider data sets, like treatments and outcomes, to inform and optimize resourcing in specific hospital and provider locations. Must be highly secure and scalable.



## Approach:

Design and build PaaS to support dynamic and scalable data volume, access, and usage requirements while building framework for modular application structure.

## Technical Delivery:

Delivered fully automated PaaS on AWS by building automation, secured core framework and infrastructure, automated data and development pipelines, and provided training and knowledge transfer to use and manage platform.



## Business Value:

Enabled research organization efficiently collect data and use provided pipeline to provide insight through analysis and data visualizations, support dynamic scaling to as many users as needed, and trained developers in AWS framework. The flexible pipeline allows for the import of any dataset. The researchers Used pipeline infrastructure to provide COVID data to the state of South Carolina.



# Metagenomics on GCP



## Need:

Researcher needed to perform bioinformatics analysis on enormous metagenomics environmental datasets to investigate areas like climate change on ecosystems over time.



## Approach:

Build pipeline and software tools to parallelize and distribute work on massive scale.

## Technical Delivery:

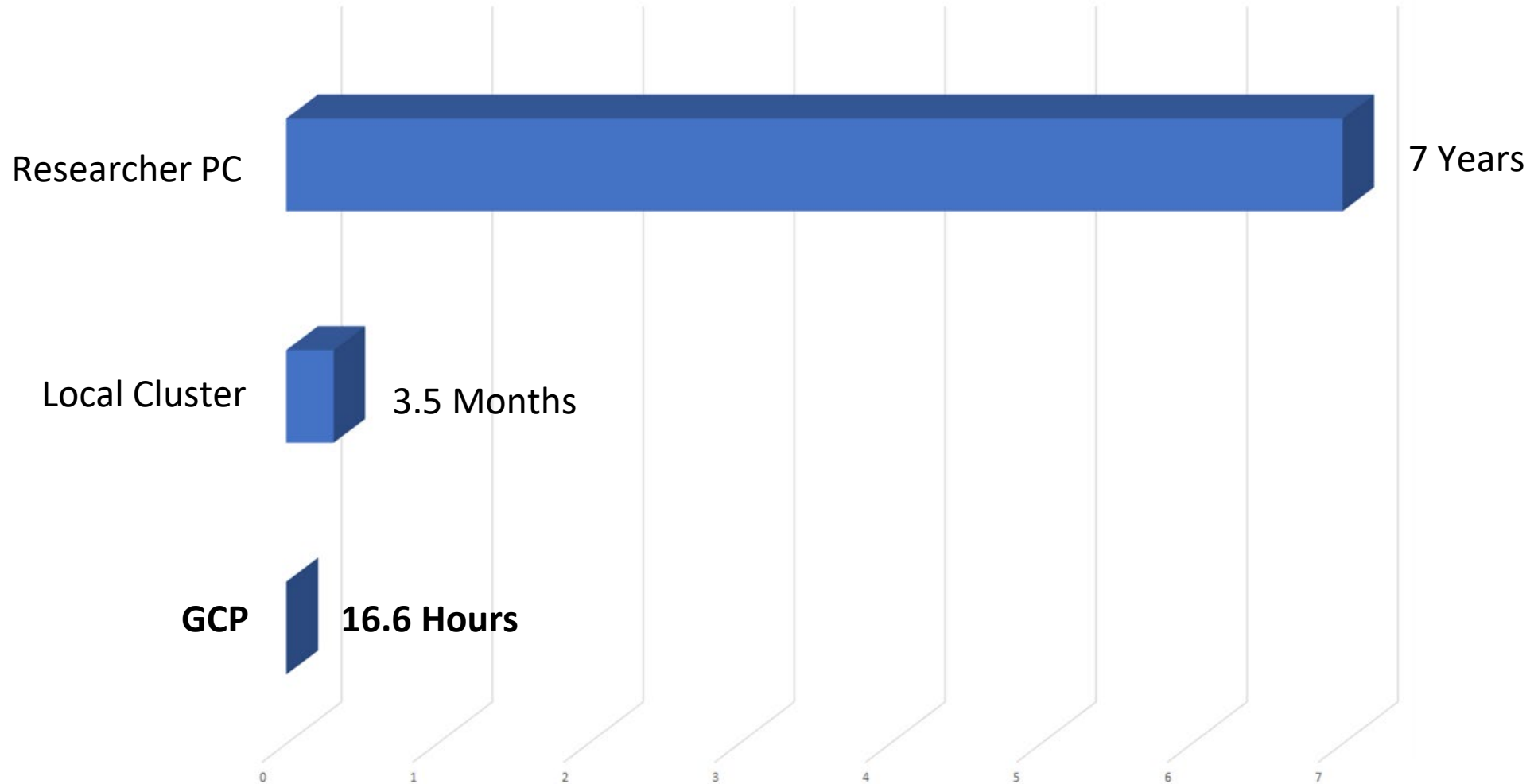
Developed automated HPC tools within GCP to automatically scale to size of data set and efficiently leverage persistent storage features. Minimize spin-up and spin-down to cut costs.



## Business Value:

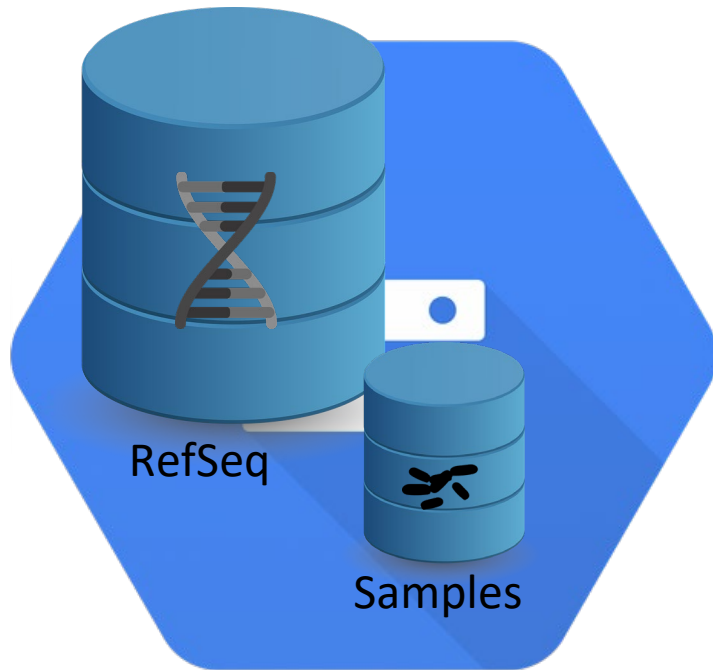
Scaled HPC style-job to massive scale, not available on any local resources. Cut computation time down from months to hours.

# Time to Compute

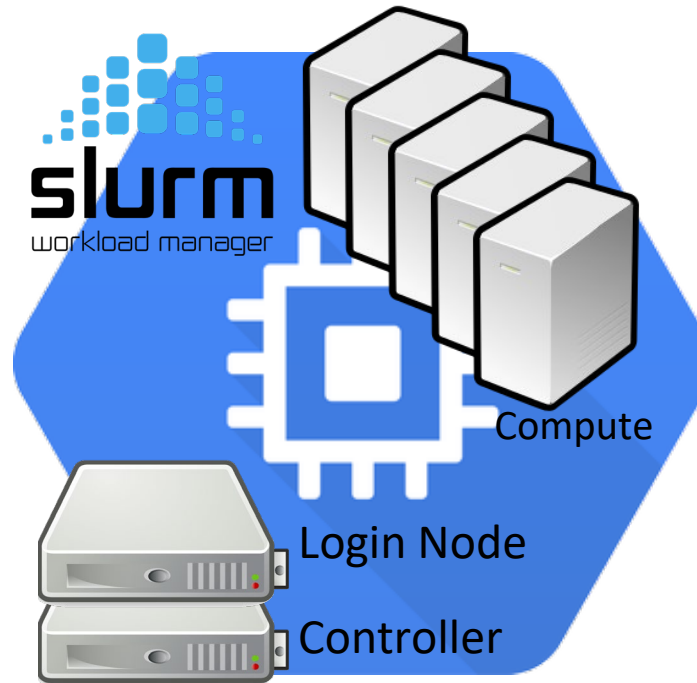




# Google Cloud Platform



Attached Persistent SSD  
(Read Only)



Dynamic HPC Cluster



Google Cloud Bucket

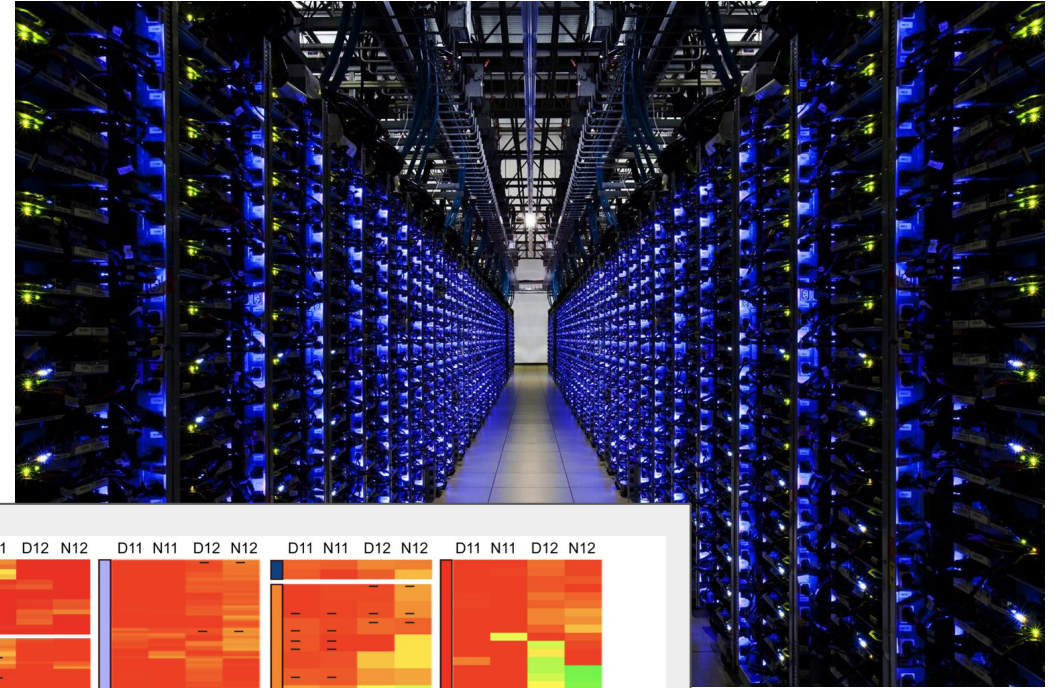
**124,352 CORES**

**3886 NODES**

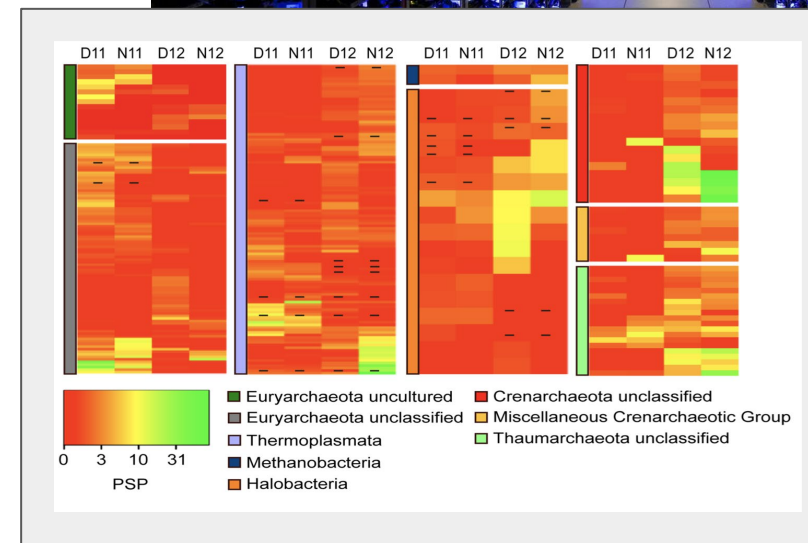
**16.6 HOURS**



# Challenges and Findings



- Enormously successful
- Encountered I/O bottlenecks during testing
  - Tried cluster NFS storage
  - Tried Elastifile storage
  - Tried copying to local disk
- Slurm FUN!
- Hyperthreading added run time



# Cloud Storage



# Research Computing Cloud Storage

## Qumulo Shift for AWS S3

- Replication of the qumulo based filesystem, holding our user home and department/lab shares to buckets on S3
- Allows RC to offer fast, seamless replication of user data to AWS for cloud processing or archive

## Globus Data Transfer Network

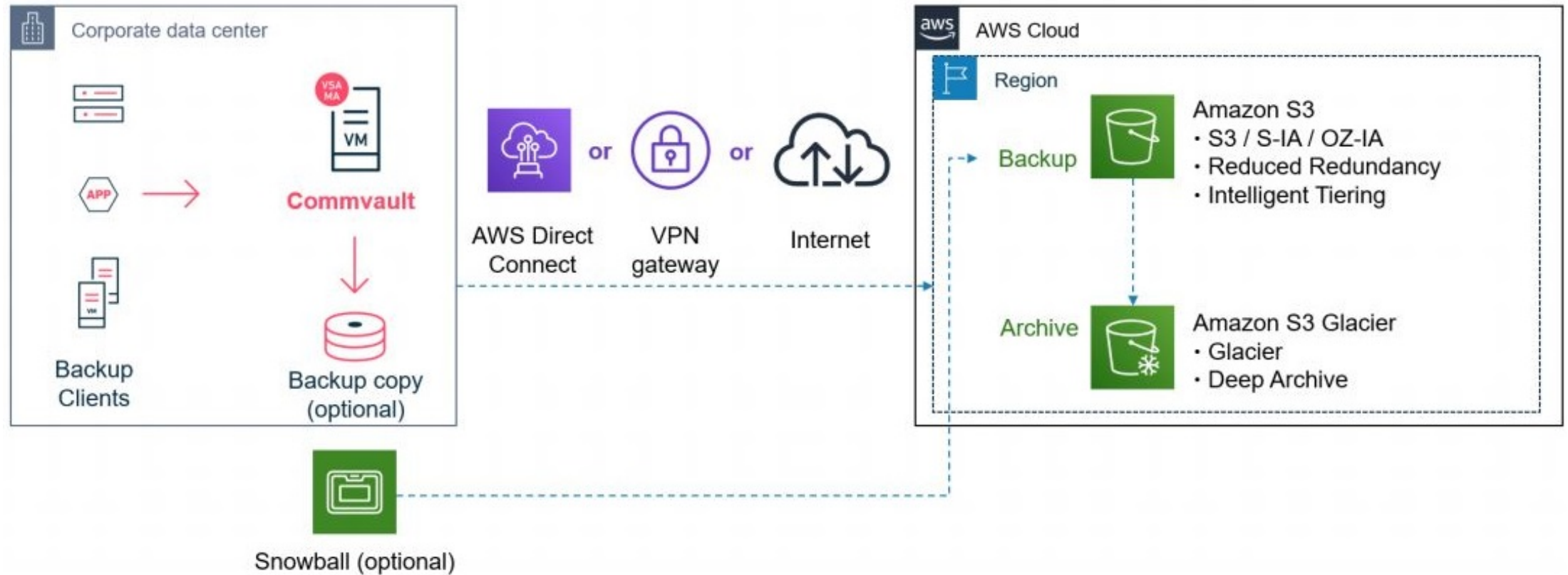
- Data transfer nodes connect UofSC research computing sites with research institutions globally.
- Allows high speed, secure data transfers between globus endpoints.
- Globus for Amazon S3 provides an instant link from S3 to Globus endpoint storage locations
- Log in and seamlessly move and share data stored in S3 across any storage resources using a unified, intuitive interface

# UofSC Restricted Data File Share in AWS

- AWS FSx for Windows
- Encrypted in Transit
- Restrict ingress through on campus firewalls, AWS security groups, isolated subnets for each FSx file system
- FSx access limited through on-campus VPN or virtual desktops with no internet access
- DFS Namespaces to create user friendly alias and flexibility to change source in the future
- Enable VSS (Volume Shadow Service) on FSx to enable user to perform data recovery operations



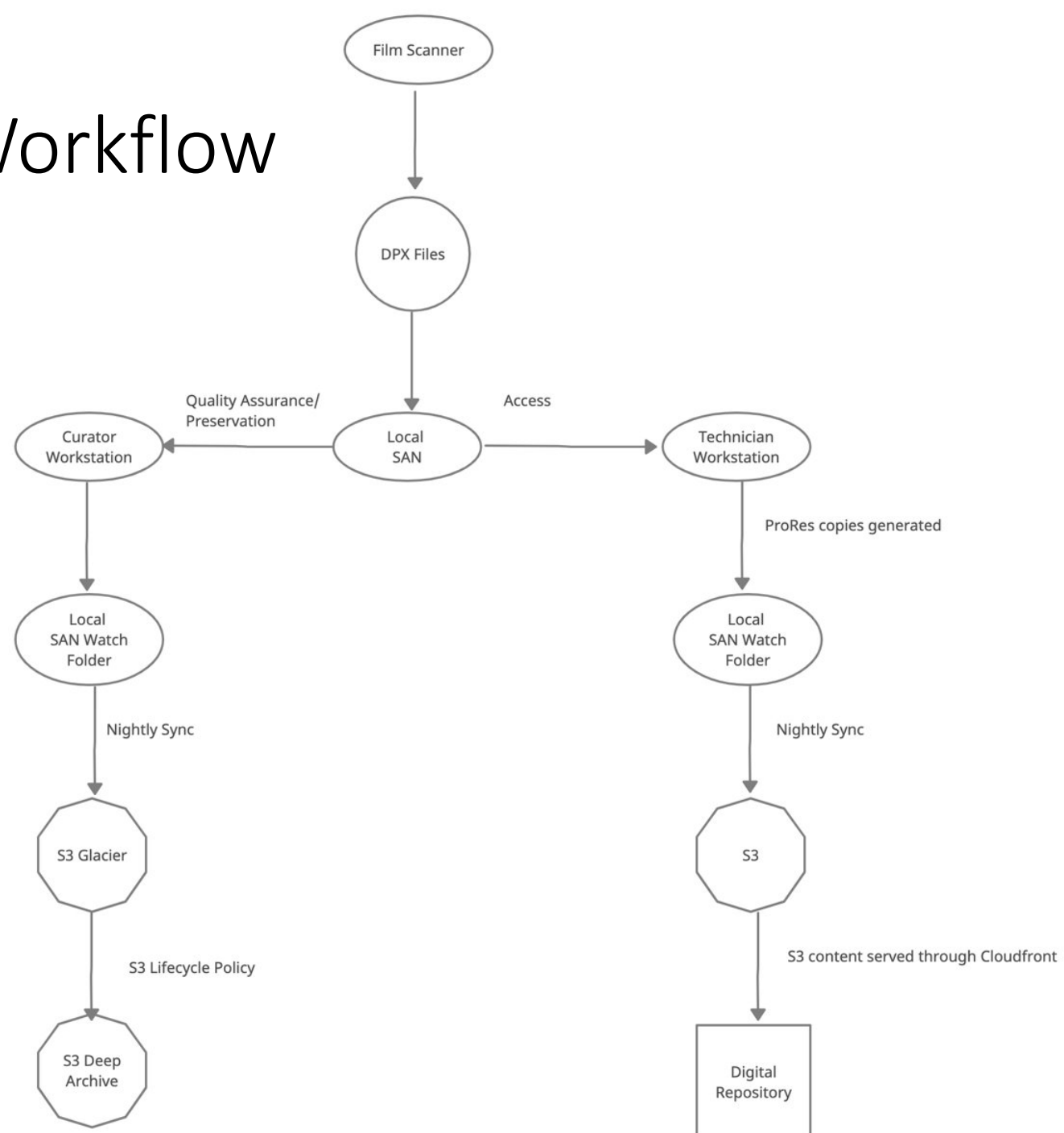
# UofSC Commvault Off-Site Backups



- Commvault direct to S3 backups bypassing need for on-site storage
- Use S3 Lifecycle Policies to archive into Glacier
- Use Commvault retention policies to retain data through end of life

# MIRC Film Digitization Workflow

- Physical media is digitized via on site film scanner
  - Preservation level DPX files
  - Stored on local SAN
- Preservation level files
  - File quality is verified and DPX files are [bagged](#)
  - Files are sent to sync watch folder in local SAN
  - Files are moved into S3 Glacier nightly
    - awscli in a cron job
  - Files are moved to Deep Archive after 14 days
    - Via a lifecycle policy
- Access level files
  - ProRes access copies are generated from DPX masters
  - Access copies are sent to sync watch folder in local SAN
  - moved into S3 nightly
    - awscli in a cron job
  - S3 bucket can be accessed by a CloudFront distribution
  - Streaming URL generated from CloudFront is integrated into digital repository





# Our future in the cloud...

- Federated HPC cluster bursting through Slurm (and Bright)
- Expand Science Gateways
- Connection to on-prem secure enclave to AWS, GCC High in Azure for collaboration tools like Teams
- Kubernetes
- Use cloud for hosting our workshops