

# Bringing Katahdin to the Masses

#### **Steve Cousins**

Interim Director of the Advanced Computing Group (ACG)

University of Maine System

https://acg.maine.edu



HTC23 Madison, Wisconsin July 11, 2023



#### Katahdin

- Maine's tallest Mountain at 5269 feet.
- Located in Baxter State Park
- Penobscot Indians named it Katahdin:
  The Greatest Mountain



- ACG HPC system: Katahdin
  - 122 nodes
    - OpenHPC with xCAT, stateless, SLURM
  - 4252 cores
    - Haswell, Broadwell, Skylake
    - Ерус3
  - Memory
    - 64GB, 128GB, 256GB, 512GB, 1TB
  - GPUs
    - 8 x Nvidia A100 in DGX system
    - 8 x Nvidia RTX 2080Ti in single system
  - Storage
    - 2.6 PB Ceph with CephFS for home
  - More soon!
    - 3 GPU nodes with 9 GPUs
    - Increasing storage to 4.8 PB







## University of Maine System

- Seven Universities plus Law School
  - Total of about 33,000 students
  - University of Maine (Umaine)
    - R1 Institution: majority of research for the System
  - University of Southern Maine (USM)
    - Portland and Gorham
  - University of Maine at Augusta (UMA)
  - University of Maine at Farmington (UMF)
  - University of Maine at Presque Isle (UMPI)
  - University of Maine at Fort Kent (UMFK)
  - University of Maine at Machias (UMM)
  - University of Maine School of Law (UMSL)







## The Advanced Computing Group :

- Formed in 2012 after over a decade as ad hoc group at UMaine
- Providing compute resources for researchers throughout the University System
- Services:
  - High Performance Computation
  - Cloud Computing Resources
  - A variety of data storage solutions
  - High performance visualization servers
  - Funding agency-compliant data management planning
  - Data Management Plan compliant archiving
  - Research Computing consulting







## Advanced Computing Group (continued)

#### • People:

- (Interim) Director:
- Outreach Specialist:
- Data Architect:
- Supercomputer Engineer/Administrator:
- Cloud Administrator:
- CyberInfrastructure Engineer\* :

Steve Cousins Ami Gaspar Open Position Steve Cousins Forrest Flagg Ongoing search

\* This position was created with funding from NSF CC\*DNI grant 1541346





## Sample Project

- ClimateReanalyzer.org
  - Sean Birkel, UMaine Climate Change Institute
  - Lots of press and traffic lately
  - Virtual Machines on ACG Cloud
  - Processing on ACG HPC system
  - NSF CC\* CyberTeams and REU students have contributed
  - Long term project dating back to 2007 or so during NSF DRL ITEST grant 0737583 (Segee)







- Kasey Legaard: UMaine Forestry Professor
  - Uses ACG combination of HPC resources and virtual machines to predict forest biomass, species, high resolution forest type mapping and other products (see next slide) from various sources of data: Satellite and LIDAR.
  - Long-term collaboration with ACG CI Engineer (NSF CC\*DNI Engineer: 1541346) to develop AI/ML methods and implementation to scale to high resolution.
  - Northeast Cyberteam (NSF CC\*Cyberteam: 1659377) project providing student experience in this process.

Quote: "The ACG has fundamentally altered the course of my research, and by extension those who use data I produce. And that's not a result of the hardware - that's the people."

Many thanks to the NSF CC\* program for helping us help Kasey and others.









#### Automated and accurate high-resolution forest mapping from Pareto optimized machine learning models



Kasey Legaard, Center for Research on Sustainable Forests, University of Maine

**Motivation:** Systematic error reduces the value of remote sensing data for forest management, monitoring, and modeling



Multispectral, multi-temporal satellite imagery



Terrain morphometry, visibility, hydrology



**Application:** Tree species and forest type mapping for the State of Maine High Resolution Land Cover project, 10-meter land cover and forest type data in partnership with the Maine GeoLibrary and NOAA C-CAP

Forest types from dominant/co-dominant species predicted by machine learning models



#### Additional projects:

**R&D:** Multi-objective (Pareto) optimization applied to machine

learning models controls or eliminates systematic error in forest maps

- Forest biomass and carbon mapping (State of Maine, NASA, NSF)
- Forest disturbance mapping (NASA)
- Forest modeling (NASA, USDA, NSF)
- Commercialization (private funding)

#### Student involvement:

- 2018 NSF Northeast Cyberteam and UMaine Advanced Computing Group, Pareto optimization of ML models and spatial prediction on HPC systems
- 2023 NSF CAREERS Cyberteam and UMaine Advanced Computing Group, seamless integration of multi-source remote sensing data on the cloud
- Monroe Community College, cloud and cloud shadow detection and cloud shadow detection



USFS field

plot data

Public Health Program, USM Muskie School of Public Service Cloud system for working with Compliant Data

- Transformed Medicaid Statistical Information System (T-MSIS Medicaid data)
  - Multiple HRSA-funded projects at the Maine Rural Health Research Center
  - Example study: Rural Health Clinic service use among pediatric and pregnancy-related populations enrolled in Medicaid in 20 states (Ahrens et al.)

- Maine Health Data Organization (inpatient, outpatient, and all payer claims)
  - Multiple Maine-based projects, various funders, often involve MPH students
  - Example study: Trends in maternal opioid use disorder and neonatal abstinence syndrome in Maine, 2016-2022 (Ahrens et al.)
  - 5 published papers from NIHfunded grant on maternal health







### External Collaborations

- NSF EPSCoR 1108153, NECC: Vermont, NH, Maine High Speed Networking
- CC\* Network Design OAC-1659142, Colby College (Maxwell): Low friction Science DMZ with DTNs at Colby, UMaine and Jackson Lab in Bar Harbor
- MGHPCC: Northeast CyberTeam OAC-1659377: Work with researchers and students with CI Engineer in between
- ERN: Eastern Regional Network (ERN): Prototype of Distributed compute clusters at institutions mainly in the northeast. Preparation for NSF Midscale pre-proposal. Turned into successful planning grant
- ERN: Ecosystem for Research ...: focus changed to remote instrument operation. Prototype CryoEM Remote Instrument. NSF OAC-2018927 and Grant OAC-1925482.





#### NSF 2020 CC\* Compute award

- Combination of high memory compute capabilities plus addition of storage to increase capacity and performance of the existing Ceph system.
- Goal of helping researchers who previously were unable to run jobs on existing systems due to lack of memory capabilities: 256 GB was too little.
- Grant required collaboration and 20% resource share with national systems. Proposed OSG, ERN and OSN.

- 14 High memory nodes:
  - AMD Epyc3: 2 x 48 core@ 2.3 Ghz
  - 512 GB RAM
  - 960GB NVMe drive for /scratch
  - HDR100 Infiniband
- 4 Highest memory nodes:
  - AMD Epyc3: 2 x 16 core @ 3.0 Ghz
  - 1 TB RAM
  - 960GB NVMe drive for /scratch
  - HDR100 Infiniband
- Addition of 1.2 PB of storage plus introduction of NVMe pool for fast shared scratch space for distributed jobs.







## Why OSG and Open Storage Network(OSN)?

- For compute share, first choice was ERN due to ongoing affiliation, however it is still in planning.
- OSG has a long history and it was a good opportunity to work with a national resource with a great reputation.
- Recommended by the solicitation as a strong option.

- OSN is composed of Ceph "Pods" distributed across the country, serving S3 storage.
- We wanted to extend existing Ceph system instead of adding a separate Pod.
- Worked with OSN with success, even using ZFS with MINIO until our Ceph storage is extended to provide capacity.







#### OSG Experience

- Tim Cartwright invited me to talk about our experience.
- I told him it will be a very boring talk: no conflict!
- Overall: extremely positive.
- Not one negative experience.
- Initial impression was that it would be simple.
- I dug deeper and thought there might be hurdles (conflict?):
  - Squid Server
  - Storage/caching system
  - Setting a CE server
- Tim assured me that we could start without these. End result: Simple





### OSN Experience (continued)

- We started simply: OSG hosted CE, No Squid Proxy, No caching storage.
- Just create "osg01" user and open up our firewall to a few hosts.
- many 1-core glideins (pilot jobs)
- Would it be more efficient to set up fewer glideins with more cores/memory?
- We decided on 12 cores per glidein, now can run parallel jobs.
- I asked about efficiency of jobs:
  - How much work is done per Glidein compared with SLURM allocation?
  - If we give each Glidein 80 GB of RAM, are jobs efficiently using that RAM?
  - Is there a way of showing these things?
  - Tim: We'll work on this. Always eager to extend the service.





#### **CPU Efficiency Observations**



#### • Checking SLURM seff on these jobs:

Job ID: 1138611 Cluster: katahdin User/Group: /osq01 State: COMPLETED (exit code 0) Nodes: 1 Cores per node: 12 CPU Utilized: 7-13:16:55 CPU Efficiency: 63.57% of 11-21:09:48 core-walltime Job Wall-clock time: 23:45:49 Memory Utilized: 2.64 GB Memory Efficiency: 5.50% of 48.00 GB

#### Overall for the month:

TOTAL CORE UTIL : 74097 TOTAL CORE HOURS: 112258.74732 Overall efficiency: 66%





#### OSG Experience: Preemption

- Prior to implementation, I believed this was a core principle:
  - Our cluster is fairly bursty :
    - Sometimes no queue of Pending jobs.
    - Sometimes the cluster is full and there are Pending jobs.
    - Between 50-60% utilization on average.
  - Let the OSG have access to some upper limit, say 50%.
  - OSG jobs start when the resources are available.
  - OSG jobs would run completely unless local jobs are queued.
  - Once there are local jobs Pending, preempt OSG jobs.
  - End result: near 100% utilization of the cluster and the OSG gets much more than 20%.
  - Good idea for clusters that are not full all of the time.
  - Not so for clusters that are completely full.
  - Need to keep an eye on this to make sure OSG gets at least 20%.





## OSG Experience: Preemption (continued)

- When it came time to implement preemption in SLURM I hoped that it would be well documented...
- I'm not a SLURM expert and it has **many** options.
- Tim to the rescue! He asked others, we met over Zoom and we put together a list of things to try:

| PreemptType=preempt/partition_prio | # Global                                 |
|------------------------------------|--|
| PreemptMode=CANCEL                 | # Global                                 |
| PreemptExemptTime=2:00:00          | # Global, not implemented on our version |
| PriorityTier=0 GraceTime=60        | # Partition                              |







## Take aways and Thankyous

- Highly recommend any institution with HPC resouces, with or without a CC\* grant, to collaborate with the OSG.
- For me, with what I have learned here, I will go back to Maine and promote use of OSG resources. I know we have researchers who can make use of those idle GPUs!
- To promote use, Campus Champions?
- Thanks Tim Cartwright for being gently persistent and always helpful.
- Thanks to the everyone at the OSG.
- Thanks to Kevin Thompson and the NSF for supporting our efforts to help our researchers, as well as connecting us to the OSG Consortium.





