

# The Cyberinfrastructure Landscape: Systems, Providers, Technologies

Mike Renfro<sup>1,2</sup>

<sup>1</sup>Tennessee Tech University

<sup>2</sup>Campus Champions Leadership Team

2025-06-23

# Who am I? (Who are any of us, really?)

## Back in the day

- ▶ ME student at a medium-sized public STEM-ish university who should have studied more instead of helping people do things in computer labs.
- ▶ Sysadmin/CAD/FEA co-op student at Oak Ridge National Lab before SGI Irix got its cameo in “Jurassic Park” (“It’s a Unix system: I know this!”).



Figure 1: Some skinny nerd, 1990

## Who am I? (Who are any of us, really?)

### Now

- ▶ Three ME degrees from the now-R2 university (1995, 1998, 2018)
- ▶ Mostly-solo practitioner of all things RCD at the same university (2000–2017, 2017–)
- ▶ Perpetually online member of multiple RCD organizations (2018–)
- ▶ Member of Campus Champions Leadership Team (2022–), CaRCC Emerging Centers Steering Committee (2024–)
- ▶ Compulsive advice-giver



Figure 2: Same nerd, not remotely skinny, 2023

# ACCESS Allocations Portal



<https://allocations.access-ci.org/>

## ACCESS Project Types

- Explore** for resource evaluation, grad student projects, small classes/training events, benchmarking, development/porting, other small-scale cases
- Discover** for grants with modest needs, Campus Champions, large classes/training events, NSF graduate fellowships, gateway development
- Accelerate** for experienced users with mid-scale needs, multi-grant programs, collaborative projects, growing gateways
- Maximize** largest-scale research activities



## ACCESS Project Types

Explore, Discover, Accelerate:

- ▶ run for grant duration or 12 months
- ▶ can be requested any time
- ▶ allow multiple projects

**Explore** 400k credits, only an overview required

**Discover** 1.5M credits, 1-page proposal required

**Accelerate** 3M credits, 3-page max proposal required, subject to merit review

**Maximize** usually only 1 project allowed, 10-page max proposal required, subject to merit review, requests accepted every 6 months



# ACCESS Resource Catalog

The screenshot shows a web browser window at `allocations.access-ci.org` displaying the ACCESS Resource Catalog. The page is titled "Compute & Storage Resources" and features a sidebar with filters. The main content area lists three resources: ACES, FASTER, and Neocortex. Each resource entry includes a thumbnail image, a title, a list of tags (e.g., Texas A&M University, GPU Compute), a brief description, and a link to learn more.

**Compute & Storage Resources**

**ACES**

Texas A&M University GPU Compute

ACCESS OnDemand Composable hardware

Accelerating Computing for Emerging Science (award number 2112356) that offers state of the art GPUs and

[Learn more about ACES »](#)

**FASTER**

Texas A&M University GPU Compute

CPU Compute Composable hardware fabric

Fostering Accelerated Scientific Transformation (award number 2019129) that offers state of the art GPUs and

MemoryExpress) based storage in a composable hardware fabric

[Learn more about FASTER »](#)

**Neocortex**

Pittsburgh Supercomputing Center Innovative / Novel Compute

Neocortex is a highly innovative resource that

**Resource Provider**

- Indiana University
- Institute for Advanced Computational Science at Stony Brook University
- National Center for Supercomputing Applications
- Northwestern University
- NSF National Center for Atmospheric Research
- Open Storage Network
- OSG Consortium
- Pittsburgh Supercomputing Center
- Purdue University
- Renaissance Computing Institute
- San Diego Supercomputer Center
- Science Gateways Center of Excellence
- Texas A&M University
- Texas Advanced Computing Center
- Texas Tech University
- University of Kentucky
- University of Texas at Austin

**Resource Type**

- Cloud
- CPU Compute
- GPU Compute
- Innovative / Novel Compute
- Service / Other



## New: DeltaAI at National Center for Supercomputing Applications (NCSA)

- ▶ 152 NVIDIA Grace-Hopper nodes, with 288 cores and 4 NVIDIA H100 GPUs (96 GB) each
- ▶ Superchip architecture
- ▶ 6 PB Lustre for \$HOME and \$SCRATCH
- ▶ 200 Gb HPE/Cray Slingshot networking





## Delta at National Center for Supercomputing Applications (NCSA)

- ▶ All using AMD 7763 CPUs
- ▶ 132 nodes with 128 cores and 256 GB RAM
- ▶ 100 nodes with 64 cores, 256 GB RAM, and 4 NVIDIA A40 GPUs
- ▶ 100 nodes with 64 cores, 256 GB RAM, and 4 NVIDIA A100 GPUs
- ▶ 6 nodes with 128 cores, 2048 GB RAM, and 8 NVIDIA A100 GPUs
- ▶ 1 node with 128 cores, 2048 GB RAM, and 8 AMD MI100 GPUs
- ▶ Same storage and networking as DeltaAI



## Stampede3 at Texas Advanced Computing Center (TACC)

- ▶ **New:** 24 Sapphire Rapids nodes with 1024 GB RAM, and 4 NVIDIA H100 (96 GB) GPUs
- ▶ 20 Intel Sapphire Rapids nodes each with 4 Intel GPUs, 128 GB HBM
- ▶ 560 Intel Sapphire Rapids nodes (no GPUs), 128 GB HBM
- ▶ 1060 Intel Skylake nodes, 192 GB RAM
- ▶ 224 Intel Ice Lake nodes, 256 GB RAM
- ▶ 10 PB VAST (\$SCRATCH) + 1 TB Lustre (\$WORK)
- ▶ 100 Gb Omni-Path networking
- ▶ Intended for:
  - ▶ parallel applications scalable to 10,000+ cores
  - ▶ general purpose computing
  - ▶ throughput computing



## Bridges-2 at Pittsburgh Supercomputing Center (PSC)

- ▶ 488 nodes with 128 AMD 7742 cores and 256 GB RAM
- ▶ 16 nodes with 128 AMD 7742 cores and 512 GB RAM
- ▶ 4 nodes with 96 Intel Cascade Lake cores and 4096 GB RAM
- ▶ **New:** 10 nodes with 104 Intel Sapphire Rapids cores, 2048 GB RAM, and 8 NVIDIA H100 GPUs (80 GB)
- ▶ 31 nodes with 40 Intel Cascade Lake cores, 102–512 GB RAM, and 8 NVIDIA V100 GPUs (16–32 GB)
- ▶ 1 node with 48 Intel Skylake cores, 1536 GB RAM, and 16 NVIDIA V100 GPUs (32 GB)
- ▶ 15 PB Lustre for \$PROJECT
- ▶ 200 Gb Infiniband networking



## Derecho at National Center for Atmospheric Research (NCAR)

- ▶ All using AMD 7763 CPUs
- ▶ 2488 nodes with 128 cores and 256 GB RAM
- ▶ 82 nodes with 64 cores, 512 GB RAM, and 4 NVIDIA A100 GPUs (40 GB)
- ▶ 200 Gb HPE/Cray Slingshot networking



## Anvil at Purdue University

- ▶ Supporting CPU, GPU simulation and large memory simulation
- ▶ “Composable Subsystem” offers Kubernetes support for science gateways and other workloads
- ▶ All using dual AMD 7763 CPUs (128 cores)
- ▶ 1000 CPU nodes with 256 GB RAM, 32 with 1024 GB RAM
- ▶ 16 GPU nodes with 512 GB RAM and 4 NVIDIA A100 GPUs



## Jetstream2 at Indiana University

- ▶ Hybrid cloud platform for flexible, on-demand, programmable cyberinfrastructure tools
- ▶ Interactive virtual machine services
- ▶ Infrastructure and orchestration services for research and education
- ▶ AMD Milan CPUs (128 per node)
- ▶ 360 NVIDIA A100 GPUs
- ▶ 512–1024 GB RAM
- ▶ 100 Gb Ethernet

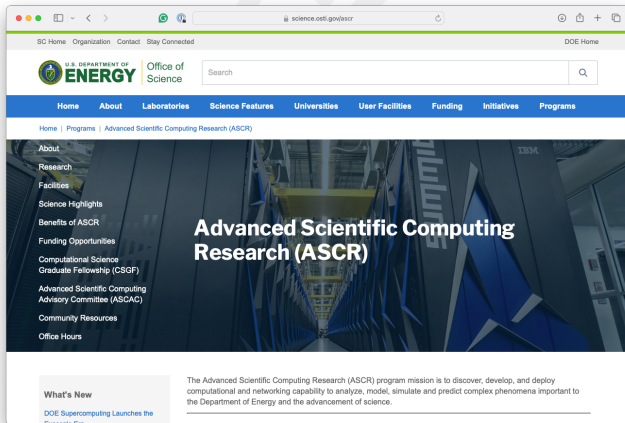


## ACES at Texas A&M University

- ▶ 130 nodes, 11888 cores
- ▶ Mostly Intel Sapphire Rapids, some Intel Ice Lake, Intel Cascade Lake, and AMD Rome
- ▶ Tons of mostly-composable accelerators:
  - ▶ GPUs: NVIDIA H100 and A30, Intel (coming soon)
  - ▶ FPGAs: Bittware Agilex, Intel D5005
  - ▶ Coprocessors: NextSilicon
  - ▶ Optane memory modules
- ▶ Non-composable accelerators:
  - ▶ Graphcore IPU: GC200, Bow-2000
  - ▶ NEC Vector Engine: Type 20B-P
- ▶ 2.3 PB Lustre
- ▶ 200 Gb Infiniband networking



# Advanced Scientific Computing Research (ASCR)





## Accessing ASCR Facilities

- ▶ **Innovative and Novel Computational Impact on Theory and Experiment (INCITE)**: multi-year awards for open science using majority of machine at Oak Ridge or Argonne
- ▶ **ASCR Leadership Computing Challenge (ALCC)**: 1-year awards for advancing DOE mission or broadening the community capable of using large computing resources at Oak Ridge, Argonne, or NERSC
- ▶ **Energy Research Computing Allocations Process (ERCAP)**: 1-year awards for advancing DOE Office of Science and SBIR/STTR mission at NERSC
- ▶ **Center Reserves**: 1-year awards for advancing science and engineering fields at **Oak Ridge**, **Argonne**, or **NERSC**



# National Artificial Intelligence Research Resource (NAIRR) Pilot

Early concept for a shared national research infrastructure connecting US researchers to:

- ▶ computational and AI,
- ▶ data,
- ▶ software,
- ▶ training, and
- ▶ educational

resources. Many of the federally-funded resources are available through ACCESS or other routes, but **NAIRR also facilitates access to commercial resources including AI models, inference services, and software as a service offerings.**



# Multi-tier Assistance, Training & Computational Help (MATCH) Plus

## MATCH Plus:

- ▶ takes requests from researchers with a support need,
- ▶ identifies a student and mentor that can provide that support,
- ▶ connects the researcher to the student and mentor with regular meetings and updates,
- ▶ for 5–10 student hours and 2–3 mentor hours per week for 3–6 months,
- ▶ at **no charge**.



# Multi-tier Assistance, Training & Computational Help (MATCH) Premier

## MATCH Premier:

- ▶ takes requests from already-funded projects,
- ▶ identifies an expert consultant and arranges payment,
- ▶ for a 6–12 month period.



## Engagement and Performance Operations Center (EPOC)

*EPOC provides researchers with a holistic set of tools and services needed to debug performance issues and enable reliable and robust data transfers. By considering the full end-to-end data movement pipeline, EPOC is uniquely able to support collaborative science, allowing researchers to make the most effective use of shared data, computing, and storage resources to accelerate the discovery process.*

– <https://epoc.global/>



## Science Gateways

*[Science g]ateways are online interfaces that give researchers, educators, and students easy access to shared resources that are otherwise inaccessible or unaffordable for a large segment of the scientific community.*

...

*The SGCI was founded to provide services and resources that advance the state of the art in science gateways, that help gateway creators use accepted practices in developing and operating gateways, and that catalyze the formation of a community that may be diverse in discipline but has a common need to advance science through gateways.*



## OpenHPC

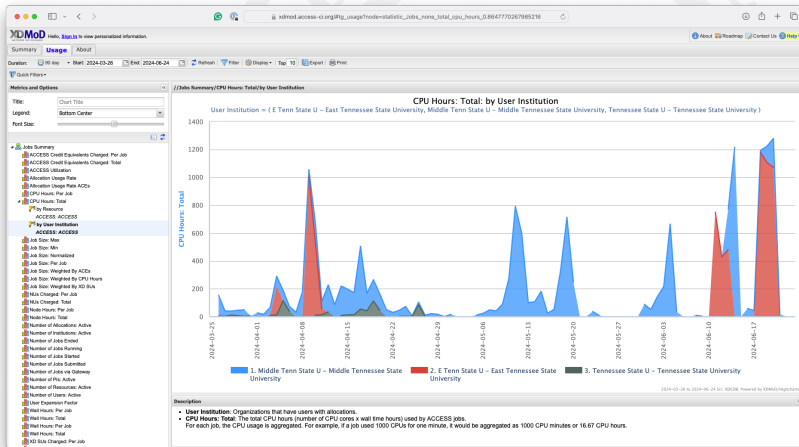
*OpenHPC is a Linux Foundation Collaborative Project whose mission is to provide a reference collection of open-source HPC software components and best practices, lowering barriers to deployment, advancement, and use of modern HPC methods and tools.*

*OpenHPC components and best practices will enable and accelerate innovation and discoveries by broadening access to state-of-the-art, open-source HPC methods and tools in a consistent environment, supported by a collaborative, worldwide community of HPC users, developers, researchers, administrators, and vendors.*

– <https://openhpc.community/about-us/>



# Open XDMoD





# Open OnDemand

The screenshot shows the Open OnDemand web interface. The top navigation bar includes links for Files, Jobs, Clusters, Interactive Apps, Develop, Help, and Log Out. The main content area features the Open OnDemand logo and a welcome message. On the left, the 'Active Jobs' section displays a table of running jobs. On the right, the 'Interactive Apps' sidebar lists available applications like Jupyter Server, RStudio Server, MATLAB, Spark Jupyter Server, and VMD. A preview of a Jupyter Notebook is shown on the far right, displaying a plot of a sine wave.

**Active Jobs**

ID	Name	User	Account	Time Used	Queue	Status
> 2904261	FW_job	petretto	lr_rnp	07:38:18	cd1	Completed
> 29054814	dlp	yanliru	ac_aome	00:22:57	id3	Completed
> 29005319	AIEn.7984.115	rporter	alice	44:37:07	alice	Completed
> 29005253	AIEn.7984.115	rporter	alice	45:56:49	alice	Completed
> 29007609	AIEn.7984.115	rporter	alice	43:08:40	alice	Completed
> 29041210	13_opt_SA0032_protobased.in.sh	jialang	mhg	03:50:10	mhg	Completed
> 28977391	start_all	khranki	nano	01:39:00	vulcan	Completed
> 29054279	swarm_56	sachlam	nano	00:07:16	extra-shared	Completed
> 29054280	swarm_57	sachlam	nano	00:07:15	extra-shared	Completed
> 29054281	swarm_58	sachlam	nano	00:07:12	extra-shared	Completed

**Interactive Apps [Sandbox]**

- Servers
- Jupyter Server
- RStudio Server
- MATLAB
- Spark Jupyter Server
- VMD

**Plot:**

Let's import the `numpy` module.

```
In [24]: import matplotlib
import matplotlib.pyplot as plt
import numpy as np

In [25]: # Data for plotting
x = np.arange(0.0, 2.0, 0.01)
y = 1 + np.sin(2 * np.pi * x)
fig, ax = plt.subplots()
ax.plot(x, y)

Out[25]: [<matplotlib.lines.Line2D at 0x2881cb9e490>]
```



# Spack

- ▶ Spack is a package management tool designed to support multiple versions and configurations of software on a wide variety of platforms and environments.
  - ▶ It was designed for large supercomputing centers, where many users and application teams share common installations of software on clusters with exotic architectures, using libraries that do not have a standard ABI.
  - ▶ Spack is non-destructive: installing a new version does not break existing installations, so many configurations can coexist on the same system.
- <https://spack.readthedocs.io/>

