# Agenda

- Intro to HPC
- Intro to Purdue clusters
- Clusters overview
- Storage overview
- How to log in
- Login vs. compute nodes
- Submitting a job
  - Monitoring a job
- Useful commands
- Open OnDemand
- Globus
- Application modules
- Engineering applications and licensing
- User support
- Questions
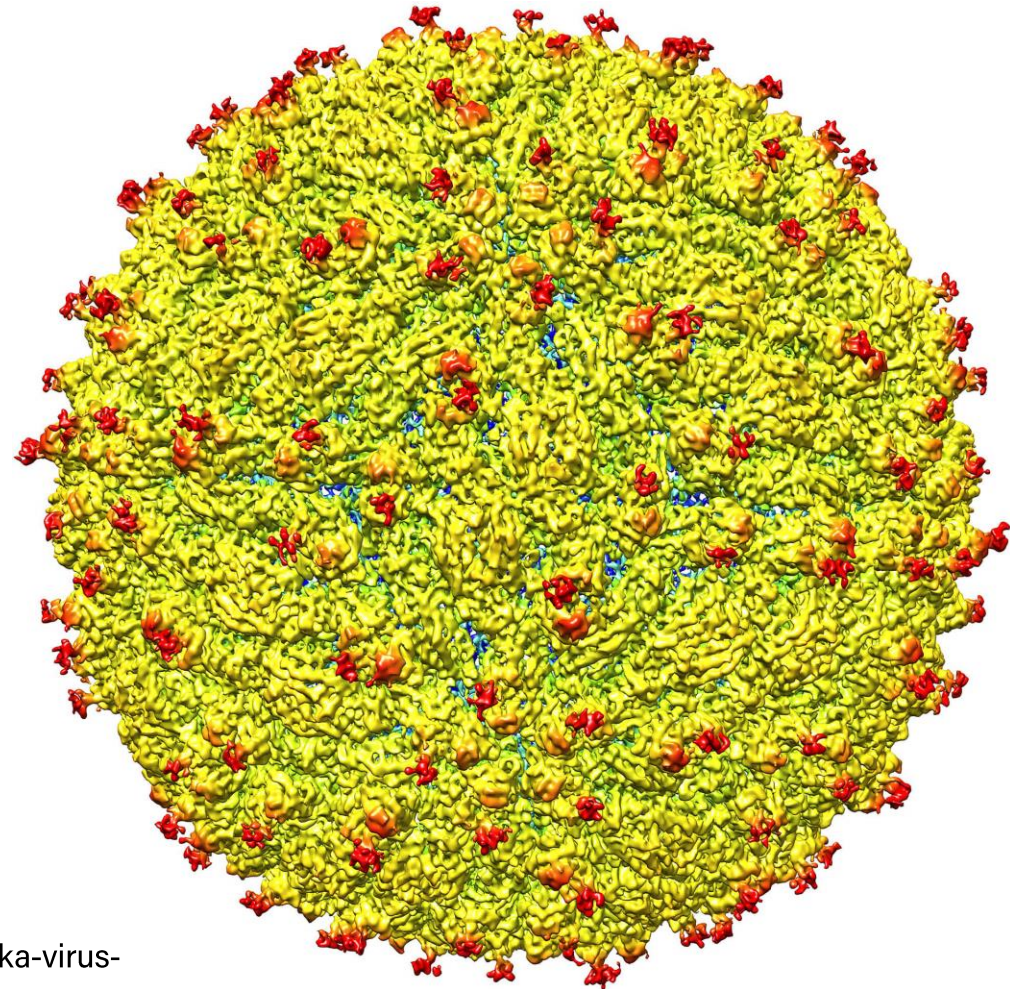
# Introduction to HPC

- High-performance computing (HPC) is technology that uses clusters of powerful processors that work in <u>parallel</u> to process <u>massive</u> multi-dimensional <u>data sets</u>, also known as big data, and solve complex problems at extremely <u>high speeds</u>.

  <u>www.ibm.com/topics/hpc</u>

- World's fastest supercomputer

  - Frontier (ORNL)

  - 1.102 exaflops

# *Structure of Zika Virus*

- Kuhn, Rossmann, *et al.*
- Combined Cryo-EM images of many Zika virus particles using RCAC clusters to create a 3-D structural map of the Zika virus.
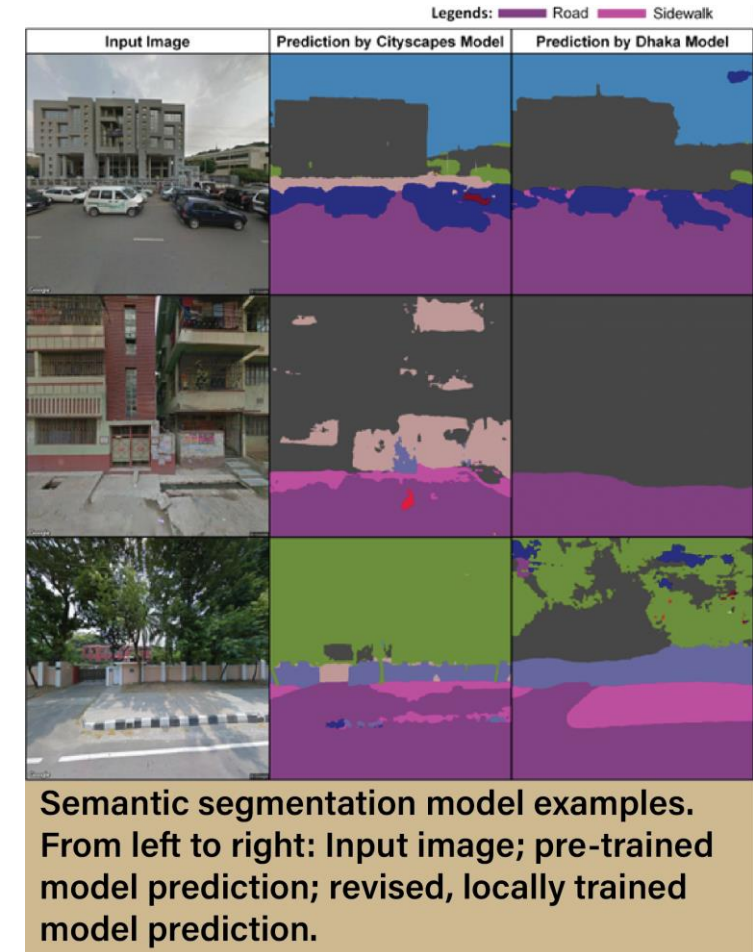- Work done on Snyder cluster.

https://www.purdue.edu/newsroom/releases/2016/Q1/researchers-reveal-zika-virus-structure,-a-critical-advance-in-the-development-of-treatments.html

# *Automated Sidewalk Mapping*

- Hamim, Kancharla, and Ukkusuri (Civil Engineering, Purdue)
- Used deep learning to create sidewalk maps from Google street view images.
- Work done on Anvil cluster.

https://www.rcac.purdue.edu/news/6424



Semantic segmentation model examples. From left to right: Input image; pre-trained model prediction; revised, locally trained model prediction.
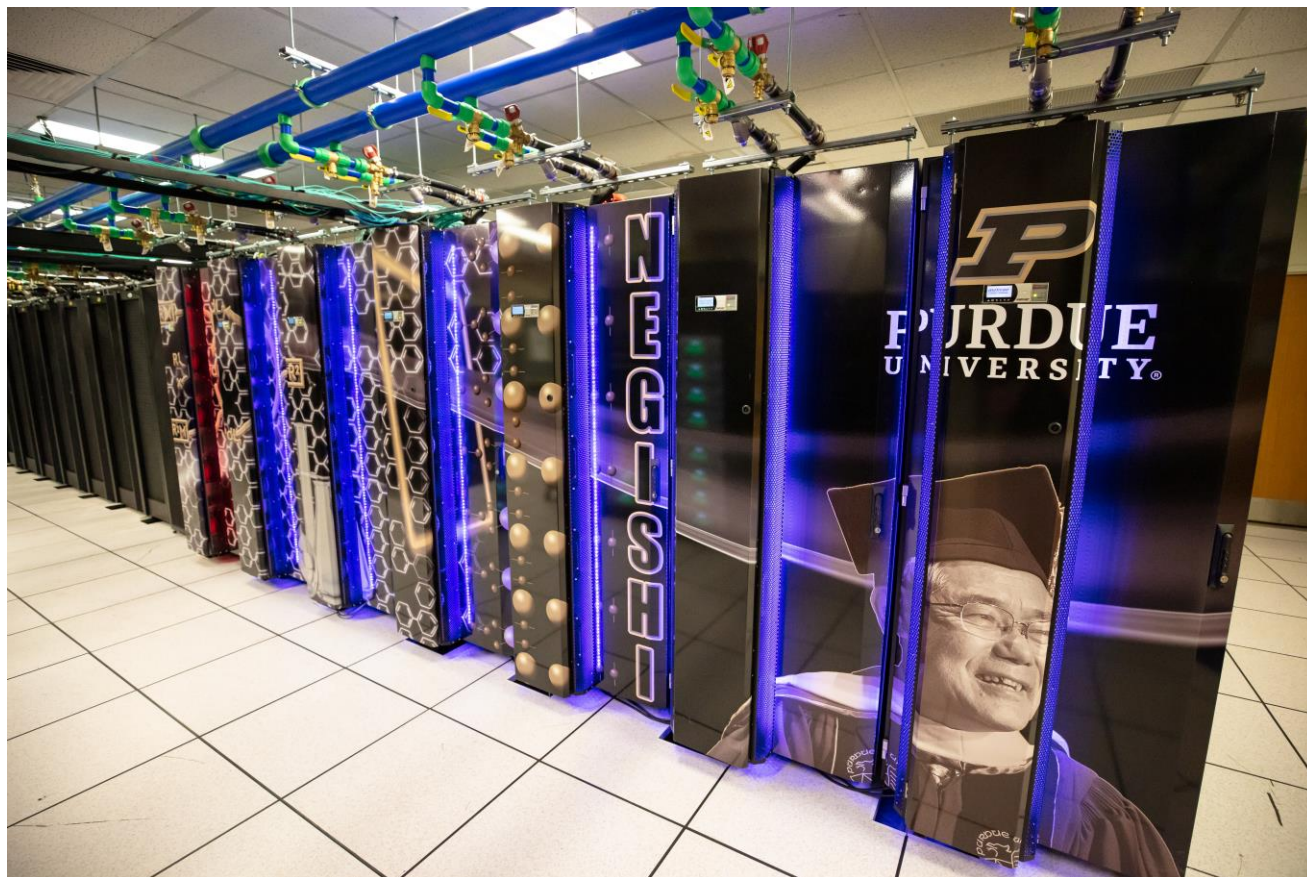
# RCAC Services

- Compute
- Storage
- Visualization
- Grant collaboration
- Training

# Community Cluster

- Faculty A needs 10 x 64-core nodes
- Faculty B needs 5 x 64-core nodes
- Faculty C needs 2 x 64-core nodes
- … … …
- Build a 100-node cluster for all the faculties

  - Ease of maintenance

  - Cost effective

  - Node failures do not lead to work stoppage

  - Use additional burst capacity when others are not using their nodes

  - Faculties buy "shares" on the cluster

# List of Community Clusters

| Name | Purpose | Hardware | Access |
|---|---|---|---|
| **Negishi** | CPU community cluster | CPU + AMD GPU | Community cluster purchase |
| **Gilbreth** | GPU community cluster | Nvidia GPU | Community cluster purchase |
| Anvil | NSF ACCESS resource | CPU + Nvidia GPU | NSF ACCESS allocations |
| Bell | CPU community cluster | CPU + AMD GPU | Community cluster purchase |
| Scholar | Teaching cluster | CPU + Nvidia GPU | Free |
| Weber | Export controlled research | CPU + Nvidia GPU | Community cluster purchase |
| Hammer | High-energy physics | CPU + Nvidia GPU | Community cluster purchase |

# Technical Specifications

**Negishi**

- 450+ nodes
  - 2 x 64-core AMD Milan processors
  - 256 GB memory
- 100 Gbps infiniband interconnect
- 6 x 1TB nodes
- 5 x 3 AMD MI210 GPUs
- 6.7 PB scratch storage

**Gilbreth**

- Heterogeneous cluster
- 100 Gbps infiniband interconnect
- 4.5 PB scratch storage
- Generations of Nvidia GPUs
  - V100
  - A100
  - H100
  - A10
  - A30

# List of Storage Resources

| Storage | Purpose | Capacity per user | Access | Access methods |
|---------|---------|-------------------|--------|----------------|
| Home | Persistent files, codes | 25 GB | With community cluster | Terminal, network drive, Globus |
| Scratch | Temporary files, data, results | 200 TB | With community cluster | Terminal, network drive, Globus |
| Data Depot | Persistent files, data, software (group shared) | On Demand | Purchased in units of 1 TB | Terminal, network drive, Globus |
| Fortress | Data archival | Unlimited | Free | SCP, HSI/HTAR, Globus |
| DBGAP | dbGaP-compliant storage | On Demand | With Negishi cluster | Terminal, Globus |

- 6.5 PB GPFS file system
- Shared workspace for research groups
- **Data owned by faculty/PI**
- **Fault tolerant**
  - All data duplicated at independent "sites"
- Regular snapshots for recovering old files
- Accessible from all clusters
- Use Globus for bulk data transfers
- Can be mounted as a network drive on laptop
- $70/TB per year
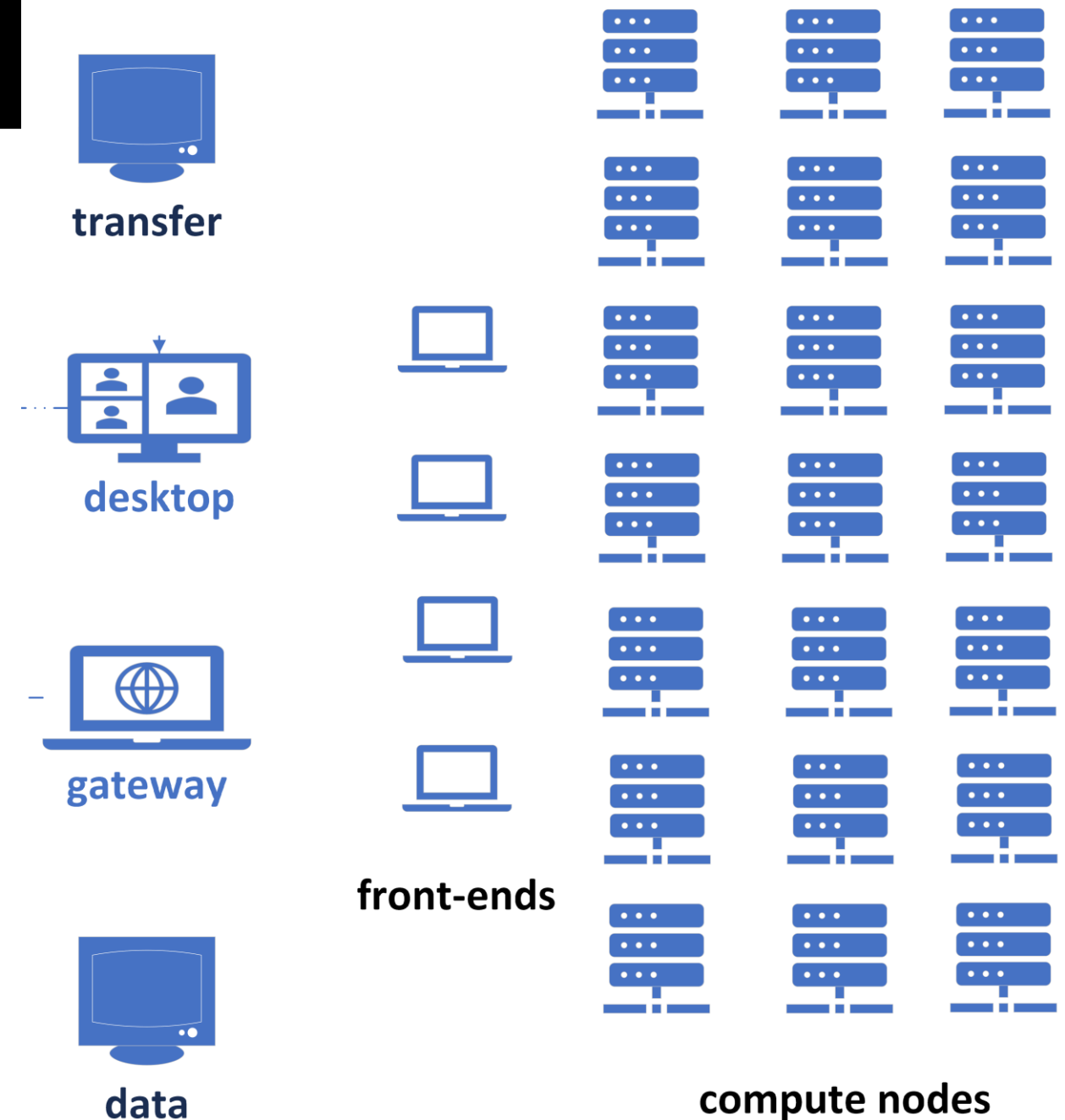
# Fortress Data Archive

- 200 PB HPSS tape archive
- Free for all Purdue researchers
- Practically unlimited storage
- Not for interactive work
- Data can be transferred using Globus or HSI/HTAR

# *How to Login to Purdue HPC Clusters*

- Thinlinc desktop
  - `desktop.clustername.rcac.purdue.edu`
- SSH client
  - `clustername.rcac.purdue.edu`
- Open OnDemand gateway
  - `gateway.clustername.rcac.purdue.edu`

# Quick Glance at a Cluster

- Login to the front-ends
  - Shared among all users
  - User for coding, data transfer, etc.
- Heavy computational work must be submitted to the back end nodes
- Submitting jobs
  - Command line and batch scripts
  - Graphical
  - Ondemand portal
- To find applications
  - Use the module command
  - Extensive listing on RCAC website

**transfer**

**desktop**

**gateway**

**data**

**front-ends**

**compute nodes**

# How to Submit a Job

- A "job" is a request for compute resources for a specific duration
- "Job request" is submitted using Slurm commands
  - Which queue
  - How many cores
  - How long
- Jobs are of two types
  - Batch: Submit a script
  - Interactive: Enter commands in the terminal

# Demo: Submitting a job

- **Batch job**
  - `sbatch myscript.sh`
- **Interactive job**
  - `sinteractive -N 1 -n 128 -A rcac -t 1:00:00`
  - Request an interactive job with 1 node and 128 cores for 1 hour under the queue "`rcac`"

# Useful Commands

- How do I find out which queues I can submit to?
  - `slist`
- How do I find out which jobs are currently running?
  - `squeue -u `myusername
- List details about a job
  - `jobinfo `jobid
- Show my storage usage
  - `myquota`

# Open OnDemand

- An easy web-based GUI for submitting jobs
  - `gateway.`<u>`clustername`</u>`.rcac.purdue.edu`
- Great for running GUI/interactive applications
  - Jupyterhub
  - Rstudio server
  - Matlab
  - VMD
  - Cryosparc/Relion
  - ...

# Data Transfer

- How do I transfer my files from my desktop to the HPC cluster?
  - SCP/SFTP
  - **Globus**
- Advantages of using **Globus**
  - Fast
  - Reliable
  - Intuitive web-based GUI
  - Can transfer to/from any location that supports Globus
- Login to `transfer.rcac.purdue.edu`

# Scientific Applications

- Compilers: GCC, Intel, AMD, Nvidia
- MPI libraries: OpenMPI, Intel, MVAPICH2
- Numerical libraries
- Data formats
- Popular applications
  - Chemistry, Physics, Biology, Statistics, etc.
  - ~280 application modules
- 600+ biocontainers
- Most engineering applications
  - Matlab, Ansys, Abaqus, Tecplot, Comsol, etc.
- https://www.rcac.purdue.edu/knowledge/applications

# Application Modules

- Scientific applications can be loaded using the "`module`" commands
- "`module`" is a software that updates your environment to make it easier to run applications
- Typical workflow
  - Load a module
  - Run application
  - Unload the module

# Module Commands

- Load an application module
  - `module load `<u>`appname`</u>
  - `module load matlab`
- Unload an application module
  - `module unload matlab`
- Search for an application
  - `module spider matlab`
- Find out application dependencies
  - `module spider paraview`

# Engineering Applications and Licensing

Sundeep Rao
Engineering IT
Executive Director, Information Technology

Application list: https://slic.ecn.purdue.edu/
Engineering support: https://engineering.purdue.edu/ECN/AboutUs/ContactUs

# *Resources*

- RCAC website: `www.rcac.purdue.edu`
- Cluster user guides: `https://www.rcac.purdue.edu/knowledge`
- Trainings: `https://www.rcac.purdue.edu/training`
- Coffee hour consultations: `https://www.rcac.purdue.edu/coffee`
- Purchase: `https://www.rcac.purdue.edu/purchase`
- User management: `https://www.rcac.purdue.edu/account/groups`
- PURR (data publishing): `https://purr.purdue.edu`
- RCAC facilities document: `https://docs.lib.purdue.edu/gendes/4/`

# Contacts

- **Amiya Maji**
  - Computer Science
  - amaji@purdue.edu
- **RCAC user support**
  - rcac-help@purdue.edu

# Questions

# THANK YOU

PURDUE UNIVERSITY | Rosen Center for Advanced Computing