# RSE Support for Georgia Tech's AI Makerspace

Dr. Fang (Cherry) Liu, Ronald Rahaman, Dr. Jeffrey Young

Senior Research Scientist/Manager of Research Software Engineer
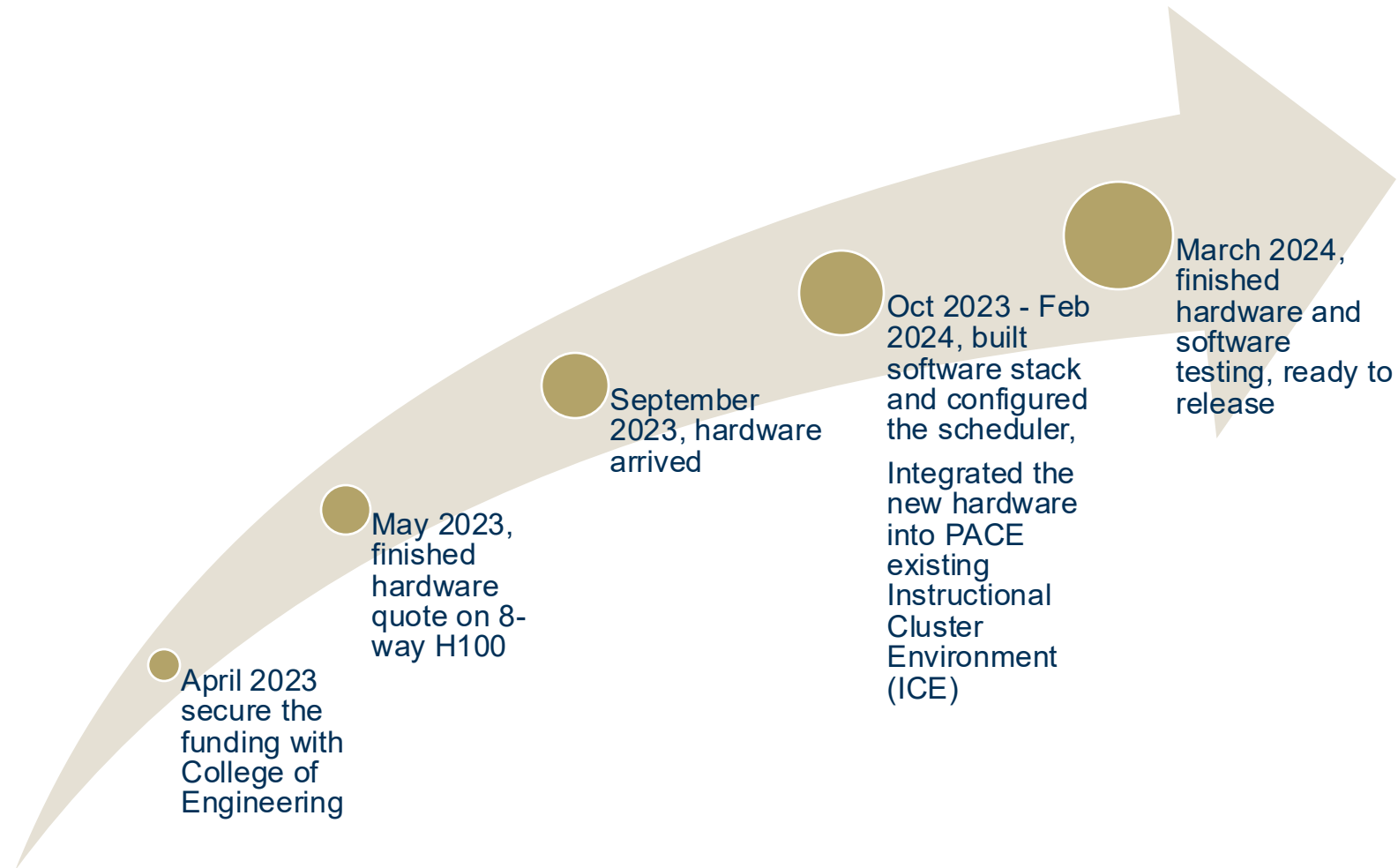PACE Research Computing Center at Georgia Tech

Georgia Tech

Digital Sandbox for students to understand and use artificial intelligence in the classroom

GEORGIA TECH
AI Makerspace

Georgia Tech
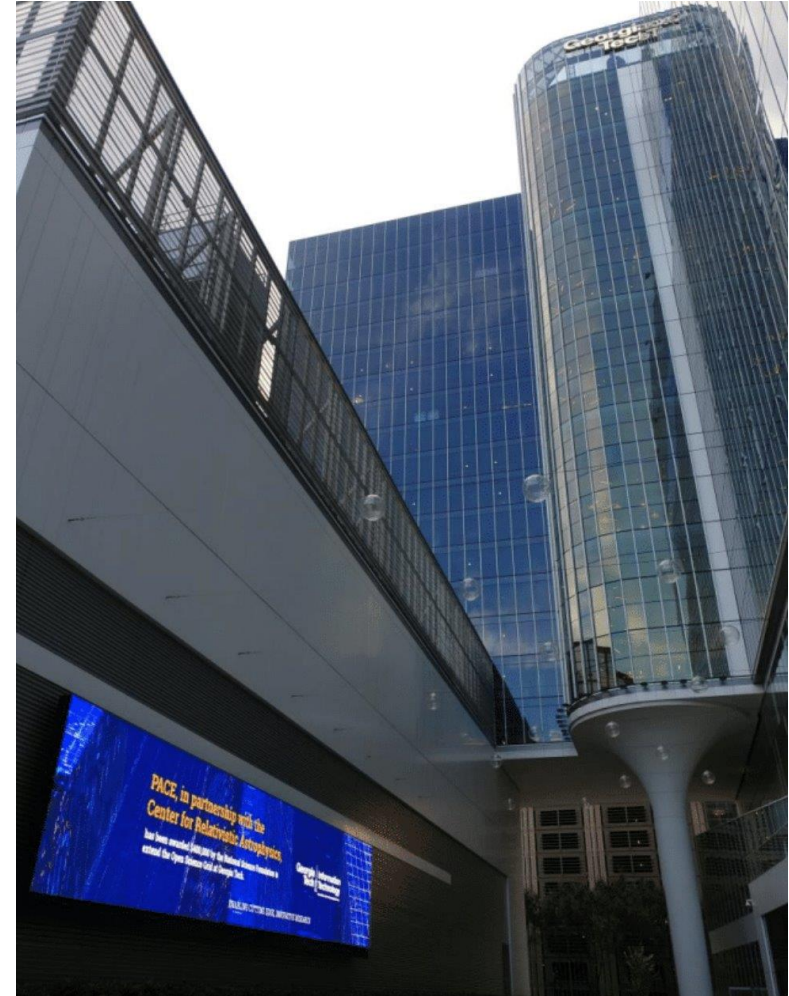
# Georgia Tech AI Makerspace

- Officially released on April 10$^{th}$ , 2024
- The first AI computing resources designed for student use within the nation
- Collaborated with Nvidia, GT College of Engineering and PACE computing center
- Phase I contains 20 8-way HGX H100 boxes, makes total 120GPUs
- Four courses are using the system at Fall 2024
- Pending Phase II contains 18 8-way HGX H200 boxes, adds 144 more powerful GPUs, system is currently under testing

Georgia Tech

# Timeline for AI Makerspace Delivery



April 2023 secure the funding with College of Engineering

May 2023, finished hardware quote on 8-way H100

September 2023, hardware arrived

Oct 2023 - Feb 2024, built software stack and configured the scheduler,

Integrated the new hardware into PACE existing Instructional Cluster Environment (ICE)

March 2024, finished hardware and software testing, ready to release

Georgia Tech

# PACE - Partnership for an Advanced Computing Environments

- Georgia Tech's PACE center (http://pace.gatech.edu) provides scalable HPC and instructional resources for Georgia Tech researchers and students

- Supports multiple clusters:
  - Phoenix – research cluster
  - Hive – NSF MRI resource
  - Firebird – CUI/ITAR-complaint research cluster
  - ICE – Instructional cluster for courses and educational workshops

- Multi-Team structure:
  - Cyber Infrastructure
  - Architecture & Platforms
  - Research Computing Facilitation & Customer Engagement
  - Research Software Engineering



Georgia Tech

# PACE Mission

## Serve & Empower

### Research

**General Research Compute** (Phoenix) : Any GT Faculty

**Controlled, Unclassified Information (CUI) Research Compute** (Firebird): As needed

**NSF-MRI cluster** (HIVE): Limited faculty, 20% allocated to **ACCESS**

### Teaching & Learning

**Instructional Cluster Environments** (CoC-ICE & PACE-ICE): Dedicated to Scientific Computing Instruction

**Technical Seminars & Tutorials**

### Democratization

**EVPR-PACESHIP**: $200k student scholarships

**Open Science Grid** (Buzzard): Funded by an NSF CC* award

**ACCESS**

**Free Tier (Compute & Storage)**

## Outreach & Engagement

Georgia Tech.

# Challenges

Students/Instructors are lacking experience on using HPC systems

Large training dataset I/O performance

Portability of Nvidia NGC containers to HPC systems

Lack of an interactive environment to develop new course material

CPU and GPU workload efficiency

Georgia Tech®

# Lowering the Access Barrier through Open OnDemand

- Funded by NSF and developed by Ohio Supercomputing Center

- Provides a web-based job submission and Jupyter Notebook interface

- Gives the freedom to add new and customized application and hides the details of resource requests

- ECE 2806 – Foundations of AI course is released on PACE Open OnDemand

# Container Support through Apptainer

- Docker is not supported in the PACE cluster and in many research computing centers
  - Podman does provide Docker-like support but requires using newer OS support features and has some NFS limitations
- Apptainer is the primary alternative to Docker on HPC systems
  - allows unprivileged users to use containers
  - prohibits escalation within the container
- To support AI/ML workflows:
  - Nvidia's NGC docker containers are converted into Apptainer containers for common scenarios, e.g. PyTorch, Tensorflow
  - Customized containers are built from base images from Nvidia
  - Integrate the container to Open OnDemand Jupyter Notebook interface
- We are in the process to enable container self-service in which everyone can build and run containers on the AI Makerspace

# Large AI Training Dataset Support

- AI training datasets usually contain many small files and a large total size

- PACE provides a central location to store AI datasets to avoid duplication issues for these common datasets

- In order to determine the best storage location, we compared the I/O performance across all available filesystems we host
  - e.g. Lustre, pNFS over RDMA, NFD over TCP, Local disk

- The comparison was done across different data sizes, data formats and filesystems

[1] Exploring Research Dataset-Sharing Strategies for Concurrent AI Workflows
https://dl.acm.org/doi/10.1145/3626203.3670597  (Best student paper at PEARC'24)

Georgia Tech

# Conclusions

## PACE Accomplishments with the AI Makerspace Include:

- Enablement for instructors to create trainings with more complex real-world problems
- Support for student teams' senior design projects in multiple semesters
- Hosting for training events (e.g. Nvidia Hackathon) for internal and external users including high school students
- Enhancement of vendor relationships – Nvidia and Penguin

## Lessons Learned

- PACE experienced a shortage of RSE-related support resources, leading to a challenging deployment
- We should also engage stakeholders ahead of time to ensure the usage of resources is fully understood

*The AI Makerspace offers students cutting-edge GPU capabilities in a classroom setting!*

Georgia Tech.

# AI Makerspace Phase II

- Hardware (with 18 8-way HGX H200 nodes) is under testing which will at least double the computing capacities

- Work to establish the student governance to make Makerspace a fully student-ran resource

- Enhance user services to accommodate the broad access to the Makerspace

- Integrate more vendor provided software solutions

# Questions

Dr. Fang (Cherry) Liu [fang.liu@gatech.edu](mailto:fang.liu@gatech.edu)