



# Programming Novel AI Accelerators for Scientific Computing

Murali Emani

Argonne Leadership Computing Facility

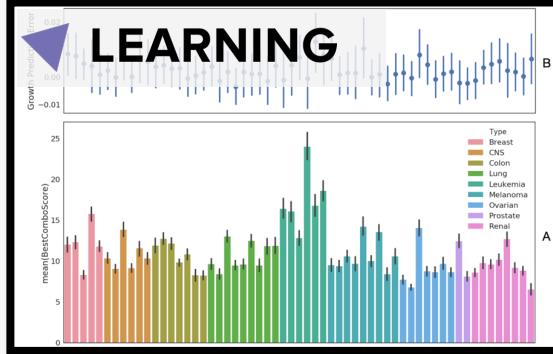
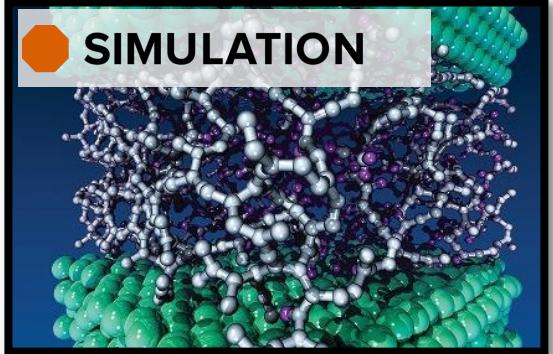
[memani@anl.gov](mailto:memani@anl.gov)

# Argonne Leadership Computing Facility



The Argonne Leadership Computing Facility provides world-class computing resources to the scientific community.

- Users pursue scientific challenges
- In-house experts to help maximize results
- Resources fully dedicated to open science



ALCF offers different pipelines based on your computational readiness. Apply to the allocation program that fits your needs.



Architecture supports three types of computing

- § Large-scale Simulation (PDEs, traditional HPC)
- § Data Intensive Applications (scalable science pipelines)
- § Deep Learning and Emerging Science AI (training and inferencing)

# ALCF AI Testbed

<https://www.alcf.anl.gov/alcf-ai-testbed>



Cerebras CS-2



SambaNova DataScale SN30



Graphcore  
Bow Pod64



Habana  
Gaudi1



GroqRack

- Infrastructure of next-generation machines with AI hardware accelerators
- Provide a platform to evaluate usability and performance of AI4S applications
- Understand how to integrate AI systems with supercomputers to accelerate science

# Recent ALCF AI Testbed Updates

ALCF AI Testbed Systems are in production and available for allocations to the research community

<https://www.alcf.anl.gov/science/directors-discretionary-allocation-program>



SambaNova SN30

SambaNova upgraded to latest 2<sup>nd</sup> generation SN30 accelerators and scaled to 8 nodes with 64 AI accelerators



Graphcore BowPod64

Graphcore upgraded to latest Bow generation accelerators and scaled to a Pod-64 configuration with 64 accelerators



Cerebras CS-2

Cerebras CS-2 upgraded to an appliance mode to include Memory-X and Swarm-X technologies to enable larger models and scaled to two CS-2 engines



GroqRack

Groq system has been recently upgraded to a GroqRack with nine nodes, each consisting of eight GroqChip v1.5 Tensor streaming processors accelerators

# Argonne Leadership Computing Facility

[ALCF Resources](#)[Science](#)[Community and Partnerships](#)[About](#)[Support Center](#)

## ALCF User Guides

[Home](#)[Account and Project Management](#)[Data Management](#)[Services](#)[Running Jobs with PBS at the ALCF](#)[Polaris](#)[Theta](#)[ThetaGPU](#)[AI Testbed](#)[Getting Started](#)[Cerebras](#)[Graphcore](#)[Groq](#)[SambaNova](#)[Data Management](#)[Cooley](#)[Aurora/Sunspot](#)[Facility Policies](#)

## ALCF AI Testbed



The [ALCF AI Testbed](#) houses some of the most advanced AI accelerators for scientific research.

The goal of the testbed is to enable explorations into next-generation machine learning applications and workloads, enabling the ALCF and its user community to help define the role of AI accelerators in scientific computing and how to best integrate such technologies with supercomputing resources.

## Table of contents

[How to Get Access](#)[Getting Started](#)[How to Contribute to Documentation](#)

# Tutorial Agenda

<https://github.com/argonne-lcf/Alaccelerators-SC23-tutorial>

Time (MST)	Topic
08.30 - 8.35	<a href="#"><u>Introduction to AI Testbed at ALCF (ANL)</u></a>
08.35 - 8.50	<a href="#"><u>Claire Zhang (Cerebras Systems)</u></a>
08.50 - 9.05	<a href="#"><u>Petro Junior Milan (SambaNova Systems)</u></a>
09.05 - 9.20	<a href="#"><u>Alex Tsyplikhin (Graphcore)</u></a>
09.20 - 9.35	<a href="#"><u>Sanjiv Shanmugavelu (Groq)</u></a>
09.35 - 9.50	<a href="#"><u>Leon Tran (Intel Habana)</u></a>
<b>10.00 - 10.30</b>	<b>Break</b>
10.30 - 12.00	Hands session on the AI Testbed (ANL)

# How to use ALCF AI Testbed

<https://github.com/argonne-lcf/Alaccelerators-SC23-tutorial>

## Request Account on AI Testbeds At ALCF

---

- Request an [ALCF Computer User Account](#) if you do not currently have one
- If you have an ALCF Account that is currently inactive, submit an [account reactivation](#) request\*.
- If you have an active ALCF account, click [Join Project](#) to submit a membership request. Specify the following in your request: Project Name: `aitestbed_tutorial`

Contact [accounts@alcf.anl.gov](mailto:accounts@alcf.anl.gov) M-F 9am to 5pm CT. Reach out to us on slack channel `#help-accounts` on [ALCF-AIAccelerator-tutorials](#) Slack.

SC23 Tutorial allocation will stay active till end of November 2023.



Director's Discretionary (DD) awards support various project objectives from scaling code to preparing for future computing competition to production scientific computing in support of strategic partnerships.

## Getting Started on ALCF AI Testbed:

**Apply for a Director's Discretionary (DD) Allocation Award**

Cerebras CS-2, SambaNova SN30, Graphcore Bow Pod64, and GroqRack are available for allocations

### Allocation Request Form

<https://www.alcf.anl.gov/science/directors-discretionary-allocation-program>

### AI Testbed User Guide

<https://www.alcf.anl.gov/alcf-ai-testbed>

# AI Testbed Talks at SC23

## Sambanova

- [RDARuntime: An OS for AI Accelerators](#)

13th International Workshop on Runtime and Operating Systems for Supercomputers (ROSS)  
Sunday, 12 November 2023 10:30am - 10:54am MST  
Location [704-706](#)

## Graphcore

- [Reducing Memory Requirements for the IPU Using Butterfly Factorizations](#)

PMBS23: The 14th International Workshop on Performance Modeling, Benchmarking, and Simulation of High-Performance Computer Systems  
Monday, 13 November 2023 11am - 11:30am MST  
Location [503-504](#)

- [Characterizing the Performance of Triangle Counting on Graphcore's IPU Architecture](#)

Tenth Workshop on Accelerator Programming and Directives (WACCPD 2023)  
Monday, 13 November 2023 11:20am - 11:40am MST  
Location [507](#)

## Cerebras

- [Scaling the “Memory Wall” for Multi-Dimensional Seismic Processing with Algebraic Compression on Cerebras CS-2 Systems](#)

ACM Gordon Bell Finalists Presentations 2  
Wednesday, 15 November 2023 4:30pm - 5pm MST  
Location [501-502](#)

## Groq

- [Exploring Converged HPC and AI on the Groq AI Inference Accelerator](#)

Exhibitor Forum: Super Intelligence II  
Thursday, 16 November 2023 10:30am - 11am MST  
Location [503-504](#)

- [From Stencils To Tensors: Running 3D Finite Difference Seismic Imaging on the Groq AI Inference Accelerator](#)

Exhibitor Forum: Super Intelligence II  
Thursday, 16 November 2023 11:30am - 12pm MST  
Location [503-504](#)

# Thank You

- This research was funded in part and used resources of the Argonne Leadership Computing Facility (ALCF), a DOE Office of Science User Facility supported under Contract DE-AC02-06CH11357.
- Venkatram Vishwanath, Michael Papka, William Arnold, Varuni Sastry, Sid Raskar, Zhen Xie, Rajeev Thakur, Bruce Wilson, Anthony Avarca, Arvind Ramanathan, Alex Brace, Zhengchun Liu, Hyunseung (Harry) Yoo, Corey Adams, Ryan Aydelott, Kyle Felker, Craig Stacey, Tom Brettin, Rick Stevens, and many others have contributed to this material.

Please reach out for further details

Venkat Vishwanath, [Venkat@anl.gov](mailto:Venkat@anl.gov)

Murali Emani, [memani@anl.gov](mailto:memani@anl.gov)