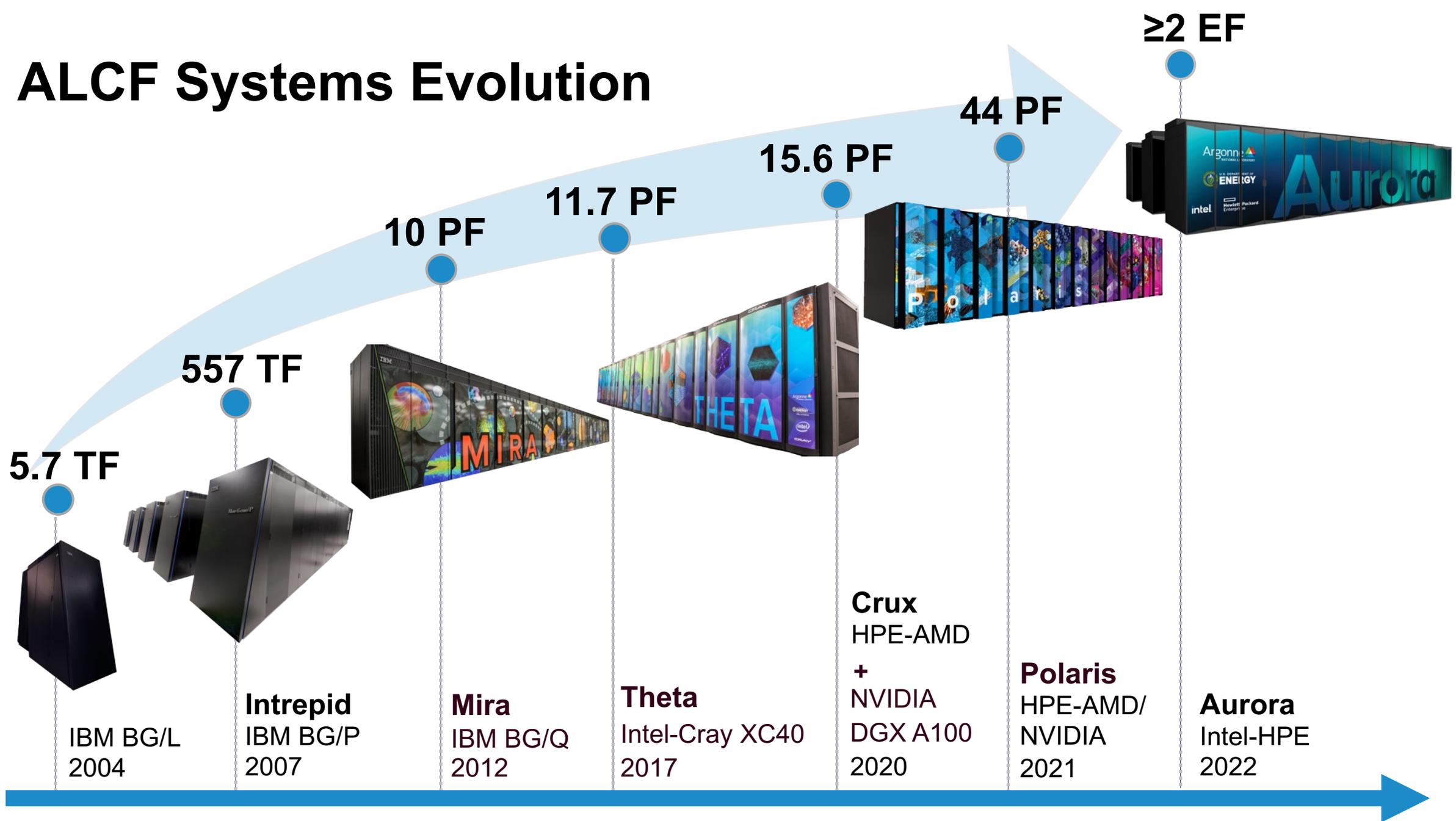# Polaris

## A Scalable Testbed Towards Aurora

**Ti Leggett**
**ALCF-Polaris Project Director, Argonne Leadership Computing Facility**

**March 30, 2022**

# ALCF Systems Evolution



**≥2 EF**

**44 PF**

**15.6 PF**

**11.7 PF**

**10 PF**

**557 TF**

**5.7 TF**

**Intrepid**
IBM BG/P
2007

IBM BG/L
2004

**Mira**
IBM BG/Q
2012

**Theta**
Intel-Cray XC40
2017

**Crux**
HPE-AMD

**+**

NVIDIA
DGX A100
2020

**Polaris**
HPE-AMD/
NVIDIA
2021

**Aurora**
Intel-HPE
2022

Argonne ▲
NATIONAL LABORATORY

# Aurora

Leadership Computing Facility
Exascale Supercomputer

PEAK PERFORMANCE

## ≧ 2 Exaflops DP

Intel GPU

## Ponte Vecchio

Intel Xeon PROCESSOR
## Sapphire Rapids wt HBM

PLATFORM
## HPE Cray-Ex

**Compute Node**
2 SPR+HBM processor;
6 PVC; Unified
Memory Architecture; 8 fabric
endpoints;

**GPU Architecture**
Xe arch-based "Ponte Vecchio"
GPU
Tile-based chiplets
HBM stack
Foveros 3D integration

**System Interconnect**
HPE Slingshot 11; Dragonfly
topology with adaptive routing

**Network Switch**
25.6 Tb/s per switch, from 64–200
Gb/s ports (25 GB/s per direction)

**Node Performance**
>130 TF

**System Size**
>9,000 nodes

**Aggregate System Memory**
>10 PB aggregate System Memory

High-Performance Storage
**220 PB** @ EC16+2**, ≧25 TB/s DAOS**

**Programming Models**
oneAPI, MPI, OpenMP, C/C++,
Fortran, SYCL/DPC++

# Polaris

Polaris will provide a platform utilizing several of the Aurora technologies and similar architectures to provide ALCF staff and users a platform for early scaling and testing purposes.

**PEAK PERFORMANCE**

## 44 Petaflop DP

**NVIDIA GPU**

## A100

**AMD EPYC PROCESSOR**

## Rome*

**PLATFORM**

## HPE Apollo Gen10+

**Compute Node**
1 AMD EPYC 7532* processor; 4 NVIDIA A100 GPUs; Unified Memory Architecture; 2 fabric endpoints; 2 NVMe SSDs

**GPU Architecture**
NVIDIA A100 GPU; HBM stack

**Processor Interconnects**
CPU-GPU: PCIe
GPU-GPU: NVLink

**System Interconnect**
HPE Slingshot 10*; Dragonfly topology with adaptive routing

**Network Switch**
25.6 Tb/s per switch, from 64–200 Gb/s ports (25 GB/s per direction)

**Programming Models**
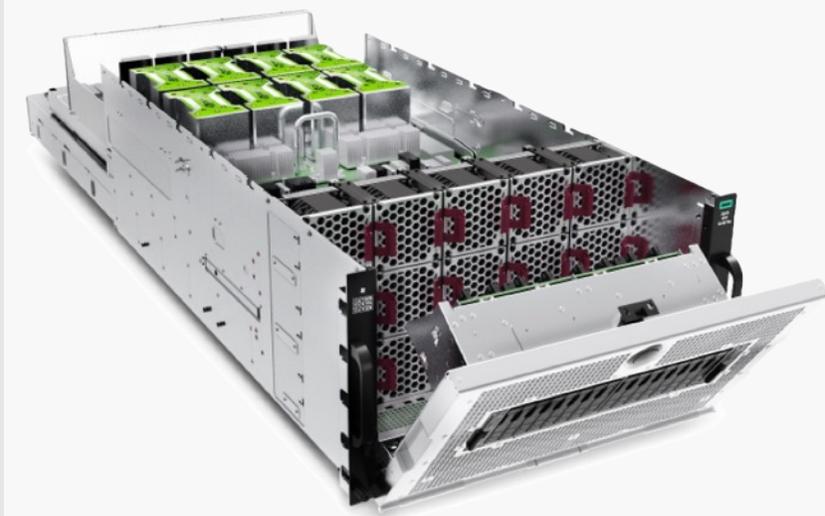CUDA, MPI, OpenMP, C/C++, Fortran, DPC++

**Node Performance**
78 TF

**Aggregate Memory**
368 TB

**System Size**
560 nodes, 1.78 MW

*Initial technology to be upgraded later

# Storage

Polaris will be connected to existing ALCF storage resources

- Grand – Global/Center-wide file system providing main project storage
  — 100 PB @ 650 GB/s
  — Accessed via Lustre LNET routers using Polaris gateway nodes
- Eagle – Community file system providing project storage that can be shared externally via Globus sharing
  — 100 PB @ 650 GB/s
  — Accessed via Lustre LNET routers using Polaris gateway nodes
- Home – shared home file system for convenience not for performance or bulk storage

Argonne ▲
NATIONAL LABORATORY

# Preparing Users for Exascale

Early Science Program (ESP)
- ALCF conducts ESP to ensure the facility's next-generation systems are ready for science on day one

- Provides research teams with critical pre-production computing time and resources
  - prepares applications for the architecture and scale of a new supercomputer
  - solidifies libraries and infrastructure for other
  - production applications to run on the system

# Bridging ESP Projects to Aurora

- To be ready for Early Science runs, projects must
  - Demonstrate INCITE level computational readiness (scaling, use GPUs, ready proposed problem in short order)
  - Complete model validations, preliminary studies, parameter-setting exercises
  - Finish integrating complex workflows for Data and Learning projects with realistic data

- Portability of applications, components, and workflows to Polaris

| Simulation components | Data components | Learning components | Workflows |
|---|---|---|---|
| ▪ OpenMP 4.5+ | ▪ Spark | ▪ TensorFlow | ▪ Containers |
| ▪ Kokkos | ▪ HDF5 | ▪ PyTorch | ▪ Balsam |
| ▪ SYCL | ▪ ADIOS | ▪ Distributed DL (eg., Horovod) | ▪ funcX/Parsl |
| ▪ PETSc, math libraries | ▪ MPI-IO | ▪ Scitkit Learn | ▪ Python-based workflows |
| ▪ *Many apps have explicit NVIDIA implementations* | ▪ Databases | ▪ JAX | |
| | ▪ Numba | ▪ Julia | |
| | ▪ Python | | |

Argonne NATIONAL LABORATORY

# Programming Models

Vendor Supported
Programming Models

ECP Provided
Programming Models

**MPI +**

OpenMP w/o target

OpenMP with target

OpenACC

CUDA

DPC++/SYCL

HIP

Kokkos

Raja

FORTRAN

C

C++

HPE

NVIDIA

LLVM

CodePlay

AMD

Argonne
NATIONAL LABORATORY

# Bridge to Aurora

| Component | Polaris | Aurora |
|---|---|---|
| System Software | **HPCM** | **HPCM** |
| Programming Models | **MPI, OpenMP, DPC++, Kokkos, RAJA, HIP,** CUDA, OpenACC | **MPI, OpenMP, DPC++, Kokkos, RAJA, HIP** |
| Tools | **PAT, gdb, ATP,** NVIDIA Nsight, cuda-gdb | **PAT, gdb, ATP,** Intel VTune |
| MPI | **HPE Cray MPI, MPICH** | **HPE Cray MPI, MPICH,** Intel MPI |
| Multi-GPU | *1 CPU : 4 GPU* | *2 CPU : 6 GPU* |
| High-Speed Network (HSN) | **HPE Slingshot** | **HPE Slingshot** |
| Data and Learning | **DL frameworks, Cray AI stack, Python, Numba, Spark, Containers,** RAPIDS | **DL frameworks, Cray AI stack, Python, Numba, Spark, Containers,** oneDAL |
| Math Libraries | cu* from CUDA | oneAPI |

Argonne ▲
NATIONAL LABORATORY

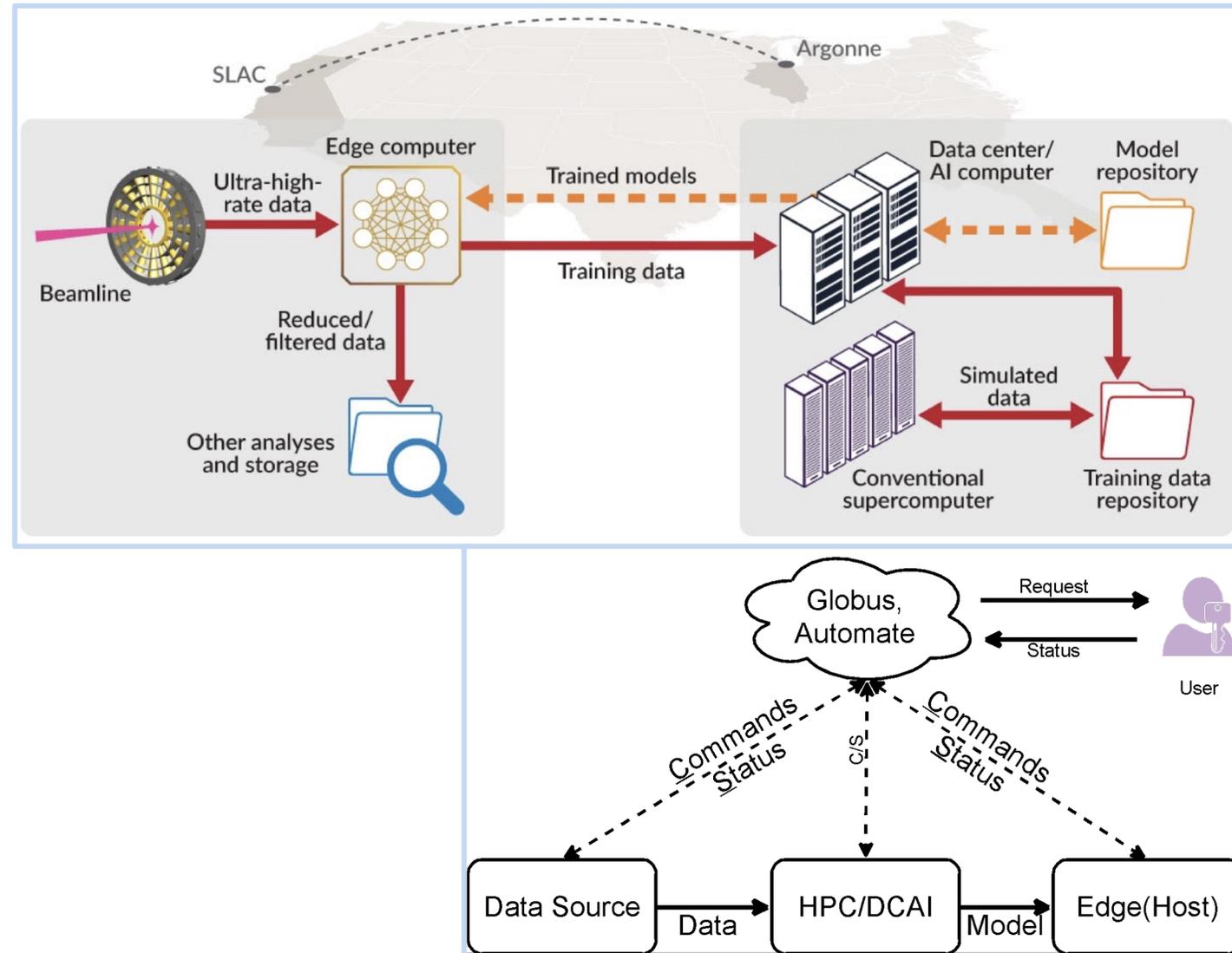# Experimental Instrument Workflows

# Experimental Instrument Workflows

# Example: Rapid Training of Deep Neural Networks using Remote Resources

- DNN at the edge for fast processing, filtering, QC

- Requires tight coupling with simulation and training with real-time data

- Near real-time steering of the experiment towards points of interest

# Upcoming

- Upgrade CPUs and HSN
  - AMD Rome → AMD Milan
  - SS-10 NICs → SS-11 NICs
  - Later this year
- Production Full User Access
  - INCITE, ALCC, ADSP, DD Allocations
  - Mid-Summer 2022

# Thanks!

- Entire Project Team
- Frank Gines & Alex Walton
- ASCR