# ALCF Operations Best Practices and Highlights

Mark Fahey

Argonne **Leadership Computing** Facility

Argonne
NATIONAL LABORATORY

# Overview
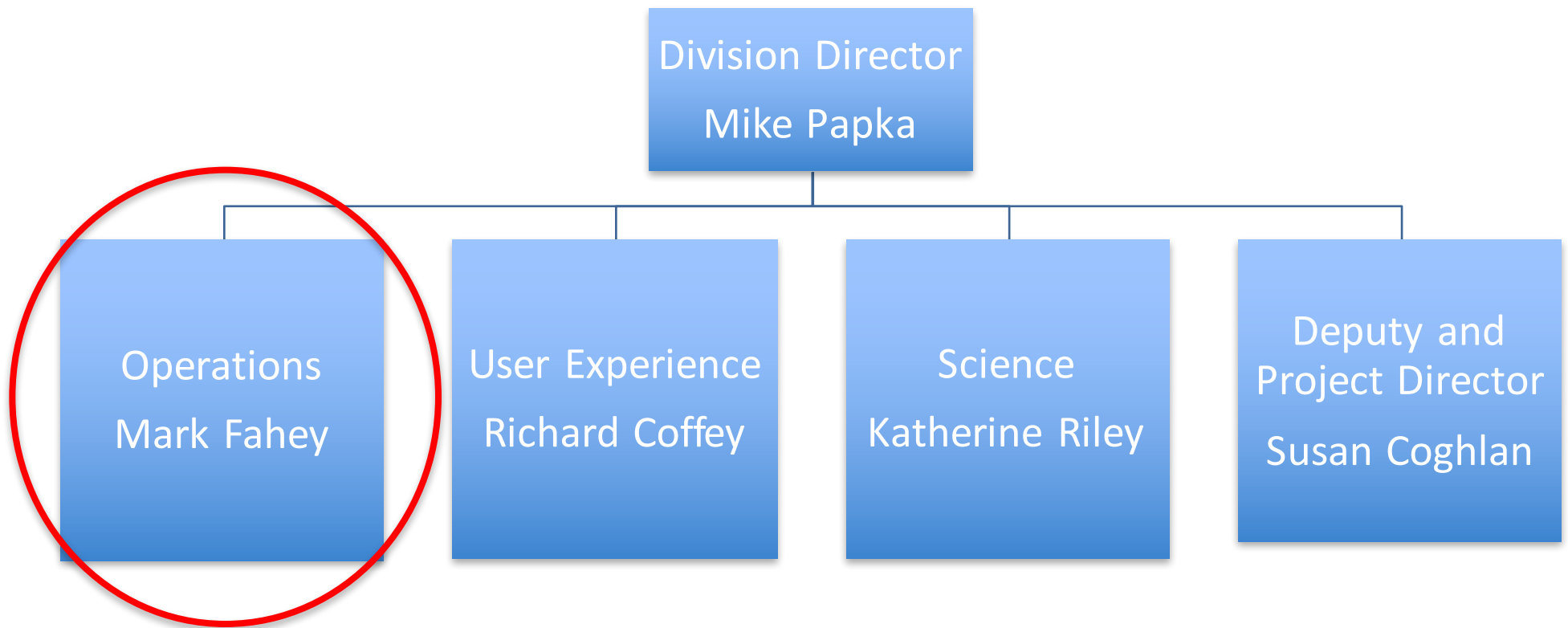
- ALCF Organization and Structure
- IBM BG/Q Mira
- Storage/Archive
- Operational highlights
  - Scheduling
  - Monitoring
  - Operational Assessment and Job Failure Analysis
- CORAL
  - Theta and Aurora
  - What's new/what stays the same
- Collaboration opportunities

Argonne **Leadership** **Computing** Facility
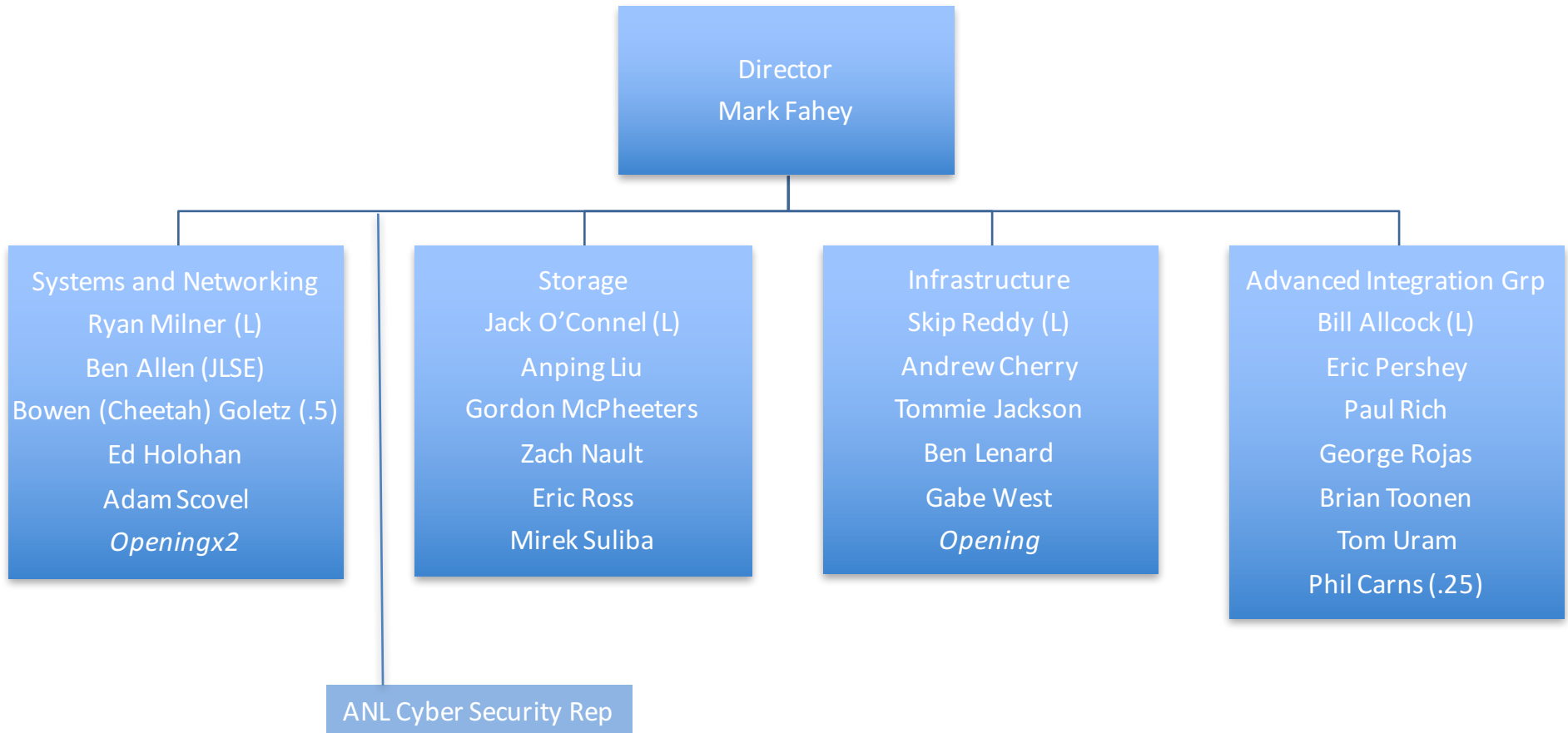
# Argonne Leadership Computing Facility

- Supported by the DOE's Advanced Scientific Computing Research program, the Argonne Leadership Computing Facility is one of two DOE Leadership Computing Facility (LCF) centers in the nation dedicated to open science.



- The LCFs deploy two diverse high-performance computer architectures that are 10 to 100 times more powerful than systems typically available for open scientific research.

- The LCF provides world-class computational capabilities to the scientific and engineering community to advance fundamental discovery and understanding in a broad range of disciplines.

Argonne **Leadership** **Computing** Facility

# ALCF Structure



Division Director
Mike Papka

Operations
Mark Fahey

User Experience
Richard Coffey

Science
Katherine Riley

Deputy and Project Director
Susan Coghlan

Argonne **Leadership Computing** Facility

# ALCF Operations Org Chart

**Director**
Mark Fahey

**Systems and Networking**
Ryan Milner (L)
Ben Allen (JLSE)
Bowen (Cheetah) Goletz (.5)
Ed Holohan
Adam Scovel
*Openingx2*

**Storage**
Jack O'Connel (L)
Anping Liu
Gordon McPheeters
Zach Nault
Eric Ross
Mirek Suliba

**Infrastructure**
Skip Reddy (L)
Andrew Cherry
Tommie Jackson
Ben Lenard
Gabe West
*Opening*

**Advanced Integration Grp**
Bill Allcock (L)
Eric Pershey
Paul Rich
George Rojas
Brian Toonen
Tom Uram
Phil Carns (.25)

ANL Cyber Security Rep

Argonne **Leadership Computing** Facility

# Mira – IBM Blue Gene/Q

- 49,152 nodes / 786,432 cores
- 768 TB of memory
- Peak flop rate: 10 PF
- Linpack flop rate: 8.1 PF

- 48 racks
- 16 1,600 MHz PowerPC A2 cores per node
- 5D torus interconnect
- 384 I/O nodes

**Mira and her cables**

# Other ALCF resources

- **Cetus** (T&D and prod.) – IBM Blue Gene/Q
  - 4,096 nodes / 65,536 cores
  - 64 TB of memory
  - 838 TF peak flop rate

- **Vesta** (T&D) – IBM Blue Gene/Q
  - 2,048 nodes / 32,768 cores
  - 32 TB of memory
  - 419 TF peak flop rate

- **Cooley** (Visualization) – Cray + NVIDIA
  - 126 nodes / 1512 x86 cores  (Haswell)
  - 126 NVIDIA Tesla K80 GPUs
  - 47 TB x86 memory / 3 TB GPU memory
  - 293 TF peak flop rate



**IBM Blue Gene/Q**

Current ALCF
GPFS File System
Infrastructure

**Mira**

(GPFS Mounts)

**Cetus**

25 DDN  SFA12KE
Couplets
w/embedded  VM
file servers

40Gb/s Mellanox IB

**Mira-home cluster**

**1PB capacity**

**22 GB/s sustained
bandwidth**

**24 embedded VM**

DDN SFA12KE

**Mira-fs0 cluster**

**19PB capacity**

**240 GB/s sustained
bandwidth**

**128 embedded VM**

DDN SFA12KE

**Mira-fs1 cluster**

**7PB capacity**

**90 GB/s sustained
bandwidth**

**48 embedded VM**

DDN SFA12KE

Argonne **Leadership**
**Computing** Facility

Coming soon
GPFS File System
Infrastructure

Integration of IBM
Elastic Storage System
as a burst buffer for
Mira and Cetus

Mira

Cetus

40Gb/s Mellanox IB

(GPFS Mounts)

**Mira-home cluster**

**Mira-fs0 cluster**

mira-fs0 remote mount

**Mira-fs1 cluster**

mira-fs1 remote mount

GPFS
Active File
Manager

**Mira-fs2 cluster**

**13PB capacity**

**400 GB/s sustained bandwidth**

**30 IBM Power 8 nodes**

**IBM ESS**

In progress

Argonne **Leadership**
**Computing** Facility

# Scheduling - Cobalt

- Orginally COBALT (Component-Based Lightweight Toolkit) was a set of component-based system software for system and resource management developed within Argonne's Mathematics and Computer Science Division

- Cobalt is a set of system software for high performance machines
  - The main piece is a set of resource management components for IBM BG systems and clusters.

- ALCF adopted the resource scheduling component and continued to enhance it for use within the facility
  - ALCF sees resource scheduling a major component of future facilities and its research/development efforts are focused on future needs

# Mira multiple rack partitions ("blocks")

**512 nodes**  **1024**  **2048**  **4096**

| Row | R00 | R01 | R02 | R03 | R04 | R05 | R06 | R07 | R08 | R09 | R0A | R0B | R0C | R0D | R0E | R0F |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| Row 0 | | | | | | | | | | | | | | | | |

| Row | R10 | R11 | R12 | R13 | R14 | R15 | R16 | R17 | R18 | R19 | R1A | R1B | R1C | R1D | R1E | R1F |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| Row 1 | | | | | | | | | | | | | | | | |

| Row | R20 | R21 | R22 | R23 | R24 | R25 | R26 | R27 | R28 | R29 | R2A | R2B | R2C | R2D | R2E | R2F |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| Row 2 | | | | | | | | | | | | | | | | |

**8192**

**16384**

http://status.alcf.anl.gov/mira/activity

**partlist** will show you if a large free block is busy due to a wiring dependency

The number of large block sizes possible is:

| # of nodes | # of blocks |
|------------|-------------|
| 49152 | 1 |
| 32768 | 3 |
| 24576 | 2 |
| 16384 | 9 |
| 12288 | 12 |
| 8192 | 6 |
| 4096 | 12 |
| 2048 | 24 |
| 1024 | 64 |
| 512 | 96 |

| Nodes | A | B | C | D | E |
|-------|-----|-----|-----|----|---|
| 512 | 4 | 4 | 4 | 4 | 2 |
| 1024 | 4 | 4 | 4 | 8 | 2 |
| 2048 | 4 | 4 | 4 | 16 | 2 |
| 4096 | 4/8 | 4 | 8/4 | 16 | 2 |
| 8192 | 4 | 4 | 16 | 16 | 2 |
| 12288 | 8 | 4 | 12 | 16 | 2 |
| 16384 | 4/8 | 8/4 | 16 | 16 | 2 |
| 24576 | 4 | 12 | 16 | 16 | 2 |
| 32768 | 8 | 8 | 16 | 16 | 2 |
| 49152 | 8 | 12 | 16 | 16 | 2 |

# Mira job scheduling

- ◉ Restrictions in queues
  - ◎ **prod-long:** restricted to the row 0.
  - ◎ **prod-short**, **prod-capability:** can run in the full machine

http://www.alcf.anl.gov/user-guides/job-scheduling-policy-bgq-systems

Mira queues by node count and walltime requested

| User Queued | Underlying Queue | Nodes | Wall-clock Time (hours) | Max. Running per User | Max. Queued per User |
|---|---|---|---|---|---|
| prod | prod-short | 512 - 4096 | 0 - ≤6 | 5 | 20 |
| | prod-long | 512 - 4096 | >6 - 12 | 5 | 20 |
| | prod-capability | 4097 - 49152 | 0 - 24 | 5 | 20 |
| | backfill (*) | 512 - 49152 | 0 - 6 | 5 | 20 |
| prod-1024-torus | prod-1024-torus | 1024 | 0 - 12 | 5 | 16 |
| prod-32768-torus | prod-32768-torus | 32768 | 0 - 24 | 1 | 20 |

(*) This queue is automatically selected based on the scheduling policy.

- ▪ **I/O to compute node ratio 1:128**

# Machine status web page



http://status.alcf.anl.gov/mira/activity

# Monitoring

- Check_MK is a comprehensive Open-Source-Solution for monitoring developed around the Nagios-core
  - Allows creating rule-based configuration using Python and offloading work from the Nagios core to make it scale better, allowing more systems to be monitored from a single Nagios server
  - Checks that consist of agent-side and server-side parts
- Check_MK is monitored using monit and MRTG
- Team members are asked to subscribe to categories of alerts. Individual subscriptions are meant to ensure that notifications remain relevant for each team member.
- Stakeholders are required to tune the monitoring of hosts and services

Argonne **Leadership Computing** Facility

# Monitoring – Slack integration

The ALCF Check_MK instance is further customized to publish alerts to a dedicated channel using the SLACK API

# Check_MK GUI

Standard UI landing

~720 - hosts monitored

~50000 - services monitored

~10 - host checks per second

~1200 - service checks per second

ALCF monitoring statistics

# Host and service views

# Operational Assessment Process

- We account for every core-second on the primary production machines
- We track the fate of every job launched on the system and classify the cause of interrupt if it does not end successfully
- Once a week, at least one member representing the major components of the machine (BG, storage, networking, infrastructure, etc.) meets to validate the previous weeks data. Most can be done automatically, but some require scrubbing logs
- All results are stored in a database and we use this information to drive where we focus our improvement efforts

# Job Failure Analysis Process

- How we track availability, outages, and failure categories
- We do weekly Job Failure Analysis (JFA)
  - We do root cause analysis on every job that ran the previous week; Considered system error unless we can find explicit proof / "probable cause" it is user error.
  - On Wednesday afternoon, the Ops team gets in a room and walks through anything that wasn't pre-classified
  - Produces Job MTTI, which internally is what we track. Also categorizes failures, which drives improvement projects.
- We account for every core-second on the machine
  - Up, scheduled down, unscheduled down; utilized or idle
  - Integers make reconciliation easy
- This software is very specific to us, but maybe someday (more on that later)
- We try to have a command for everything that also logs relevant data for later reporting
  - [begin|end]-service-action; maintman for maintenance; Scheduler reservations.
  - The Blue Gene comes with this built-in; porting to the Cray is going to be a challenge

Argonne **Leadership** **Computing** Facility

# Maintenance Manager – maintman or mm2

⦿ Script that automates our maintenance processes

⦿ In our opinion, a very nice tool

⦿ For this discussion, what is apropos is that it writes records into our database and modifies scheduler reservations that are part of the availability calculation

```
Usage: mm2 <command> [<args>]

Some useful mm2 commands are:
    binotify   Send out email to users notifying about BI maintenance related tasks
    call       Send out the call for scheduled maintenance
    defer      defers a pm reservation for a resource
    extend     extends a pm reservation already in place for a resource
    initsched  initializes a calendar based schedule with pm ticket items.
    nagios     Enable/disable nagios alerting
    notify     Send out email to users notifying about maintenance related tasks
    reserve    checks and/or sets a pm reservation for a resource
    sendsched  Emails the schedule for next maintenance to the team
    version    prints out the version string of mm2

See 'mm2 help <command>' for information on a specific command.
```

# The pre-classification script

- The script is run daily, and loads interrupts to be analyzed.
- Staff can choose to do analysis / data entry ahead of the Wed meeting
- Below is an example of the output of what the script produces.
- This email is post-JFA, so it includes the resulting analysis (the comments)

```
New interrupt: id=0, num_events=1
    2015-07-29 03:24:14: jobid=514738, mode=script user=[        part='MIR-00000-73FF1-16384', type=Unknown,
msg='UNCLASSIFIED: abnormal termination by signal 15 from rank 1080. Delivered by kill_job user         on host
miralac1'

# User error: user killed task with kill_job

New interrupt: id=1, num_events=2
    2015-07-30 03:54:01: jobid=523195, mode=script user=.        part='R22-M1-N08-J13', type=System, msg='fatal RAS
event'
        The install of a kernel image failed,  domain[0] rc=1 for image /bgsys/drivers/ppcfloor/boot/cnk
    2015-07-29 05:36:44: jobid=518826, mode=script user=.    _, part='MIR-08000-3BFF1-8192', type=System, msg='RAS
42745348 with task kill signal SIG35'

# System error: CFAM, L1P and DDR machinechecks - node was replaced
```

# JFA Web App

- This shows all the records the pre-classification couldn't automatically identify

- Each person in the room runs this (as well as it being projected) and they can select an event that they will analyze

# "Component" Analysis

- We also classify by "component"
- This allows us to see what is giving us problems and drives improvement projects.
- One of the first real wins: We discovered that GPFS was 3x the next source of failures. We investigated and discovered we were getting timeouts and moving the management functions to dedicated nodes dropped GPFS down into the noise.

```
Component fault analysis (Recent range: prior 90 days,

Component              Count (Recent) Count (All)
--------------------   -------------- -----------
machine                          25          305
gpfs                             21           59
sw/driver stack                   3           33
ddn                               2            4
facility                          1            4
```

# Incident Timeline



- Trying to figure out when an incident began and ended is non-trivial.
- This shows all the sources of data about a given incident.
- We take the "union" of all the events to determine the duration of the incident.

# Incident Timeline – drill down



- From this screen you can drill down into the details for any entry.

# Our standard availability report



Mira Availability

| ■ Scheduled Down | ■ Unscheduled Down | □ Available |

Generated on 2015-08-06 13:01:08

Scheduled        99.435652%
Overall          99.435652%
Start 2015-07-29
End   2015-08-04

Incidents:

 ID  Name
 396. Bad Node - R22-M1
 397. Bad Node - R22-M1
 398. Link Errors - R13-M1 to R1B-M1
 399. Bad Node - R06-M1-N15-J27

# The (complicated) Big Picture...

# Machine Time Overlay...



Y axis are the allocable chunks of the machine (mid-planes here, nodes on the vis cluster)

X axis is time

- Easy to see scheduling; Also helps find bugs, like two jobs running on the same resource at one time.
- There is a lot of information encoded here
- This is general; Any information you can provide that is (data, location, time) can be displayed this way; We also use this for coolant temperature, power consumption, etc..

Argonne **Leadership Computing** Facility

# CORAL– Collaboration of ORNL, Argonne, LLNL

- Provide the Leadership computing capabilities needed for the DOE Office of Science mission from 2018 through 2022
  - Capabilities for INCITE and ALCC science projects

- CORAL was formed by grouping the three Labs who would be acquiring Leadership computers in the same timeframe (2017-2018), benefits include:
  - Shared technical expertise
  - Decreases risks due to the broader experiences, and broader range of expertise of the collaboration
  - Lower collective cost for developing and responding to RFP

Argonne **Leadership**
**Computing** Facility

# CORAL Overview

**Objective -** Procure 3 leadership computers to be sited at Argonne, ORNL, and LLNL in CY17-18.

*Current DOE Leadership Computers*

**Mira (ANL)**
**2012 - 2017**

**Sequoia (LLNL)**
**2012 - 2017**

**Titan (ORNL)**
**2012 - 2017**



**Leadership Computers** RFP requests >100 PF, 2 GB/core main memory, local NVRAM, and science performance 4x-8x Titan or Sequoia

**Approach**

Competitive process – 1 RFP (issued by LLNL) leading to 2 R&D contracts and 3 computer procurement contracts

For risk reduction and to meet a broad set of requirements, 2 architectural paths will be selected – and Argonne and ORNL must choose different architectures

Once selected, multi-year lab-awardee relationship to co-design computers

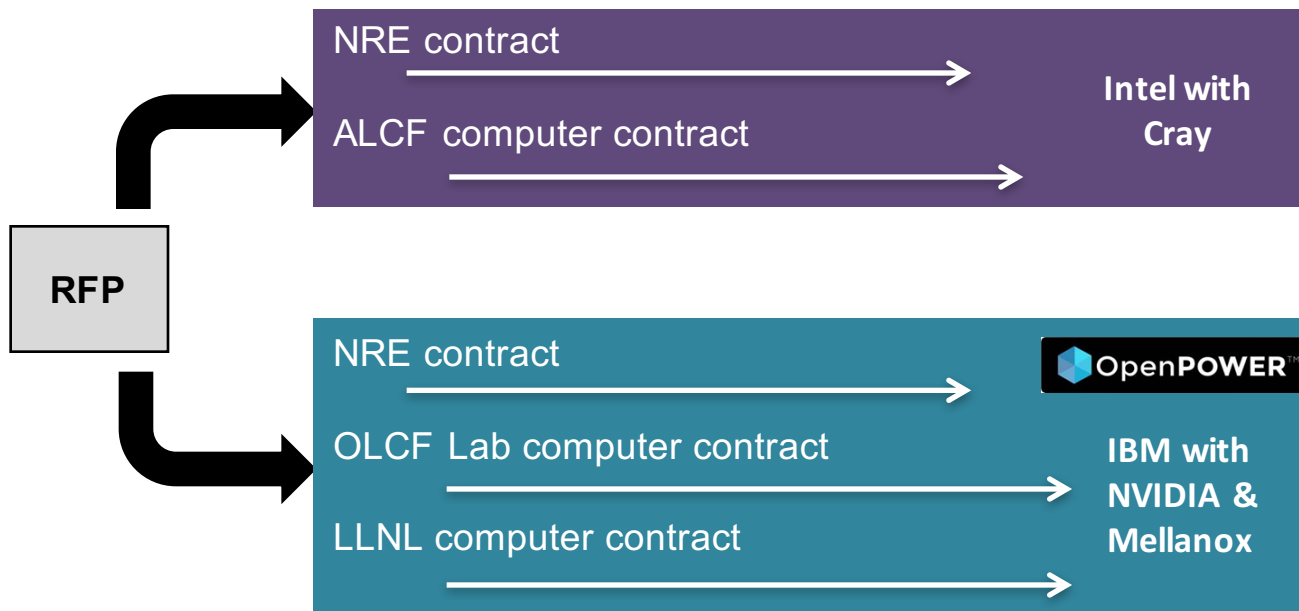Both R&D contracts jointly managed by the 3 Labs

Each lab manages and negotiates its own computer procurement contract, and may exercise options to meet their specific needs

Understanding that long procurement lead time may impact architectural characteristics and designs of procured computers

Argonne **Leadership**
**Computing** Facility

# Results of CORAL Procurement

**Two Diverse Architecture Paths**

# 2018 ALCF Leadership System

**Many Core architecture**



**System Name: Aurora**

**Vendor: Intel (Prime) / Cray (Integrator)**

**Delivery date: 2018**

- Over 13X Mira's application performance
- Over 180 PF peak performance
- More than 50,000 nodes with 3rd Generation Intel® Xeon Phi™ processor
  - Code name Knights Hill, > 60 cores
- Over 7 PB total system memory
  - High Bandwidth On-Package Memory, Local Memory, and Persistent Memory
- 2nd Generation Intel® Omni-Path Architecture with silicon photonics in a dragonfly topology
- More than 150 PB Lustre file system capacity with > 1 TB/s I/O performance

# 2016 ALCF Theta System
## Many Core architecture

**Vendor: Intel (Prime) / Cray (Integrator)**

- Transition and data analytics system
- Over 8.5 PF peak performance
- More than 2,500 nodes with 2nd Generation Intel® Xeon Phi™ processor
  - Code name Knights Landing, > 60 cores
- 192GB DDR4 memory and up to 16GB HBM on each node
- 128GB SSD on each node
- Cray Aries high speed interconnect in dragonfly topology
- Initial file system: 10PB Lustre file system, 200 GB/s throughput
- Cray XC system
- Cray software stack
- ~1.7 MW peak power

# Systems feature summary

| System Feature | Mira (2012) | Theta (2016) | Aurora (2018) |
|---|---|---|---|
| Peak Performance | 10 PF | > 8.5 PF | 180 PF |
| Number of Nodes | 49,152 | > 2,500 | > 50,000 |
| Aggregate HBM, local memory, and persistent mem | 786 TB | > 480 TB | > 7 PB |
| File system capacity | 26 PB | 10 PB (initial) | > 150 PB |
| File system throughput | 300 GB/s | 210 GB/s (initial) | > 1 TB/s |
| Peak Power Consumption | 4.8 MW | 1.7 MW | 13 MW |
| GFLOPS/watt | 2.1 | > 5 | > 13 |
| Facility Area | 1,536 sq. ft. | ~1,000 sq. ft. | ~3,000 sq. ft. |

# What changes, what doesn't

⊙Same
- ◎ many core
- ◎ GPFS
- ◎ MPI+OpenMP
- ◎ Cobalt scheduler

⊙Different
- ◎ Network
- ◎ RAS
- ◎ Lustre
- ◎ System software
- ◎ Cray Programming Environment
- ◎ On package, HBM memory
- ◎ Intel x86, not powerpc
- ◎ SSDs

Argonne **Leadership Computing** Facility

# Future/Opportunities

- We WILL port
  - Cobalt
  - Monitoring
  - JFA process
- Much of our reporting depends on it
- We are looking at abstractions and architectural improvements that will make this easier to use on general machines
- We have also been leaning on Cray and Intel to work with us to develop standardized interfaces and mechanisms for obtaining this kind of data
- Opportunities to share/codevelop tools

Argonne **Leadership Computing** Facility

# Other opportunities

- Evaluate each other's Petascale computers and software stacks
  - Scaling studies, tools, libraries, compilers
- Modeling and simulation of applications
- Community codes
- Visualization support: software, techniques
- Industry engagements

Argonne **Leadership Computing** Facility

# Acknowledgment

- This research used resources of the Argonne Leadership Computing Facility, which is a DOE Office of Science User Facility supported under Contract DE-AC02-06CH11357.

# Acknowledgements

⦿ The entire ALCF team

# Questions?

Mark Fahey <mfahey@anl.gov>

Argonne **Leadership**
**Computing** Facility