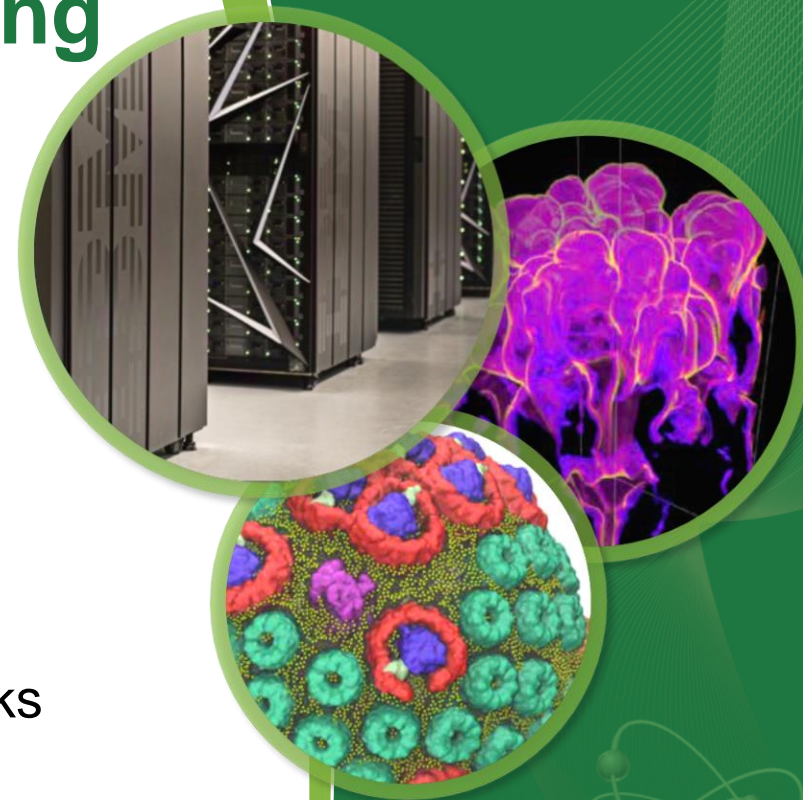# Data-Intensive Science Executed within Leadership-Scale Computing Facilities

Jack C. Wells
Director of Science
Oak Ridge Leadership Computing Facility
Oak Ridge National Laboratory

Mini-Symposium on Data over Distance: Convergence of Networking, Storage, Transport, and Software Frameworks

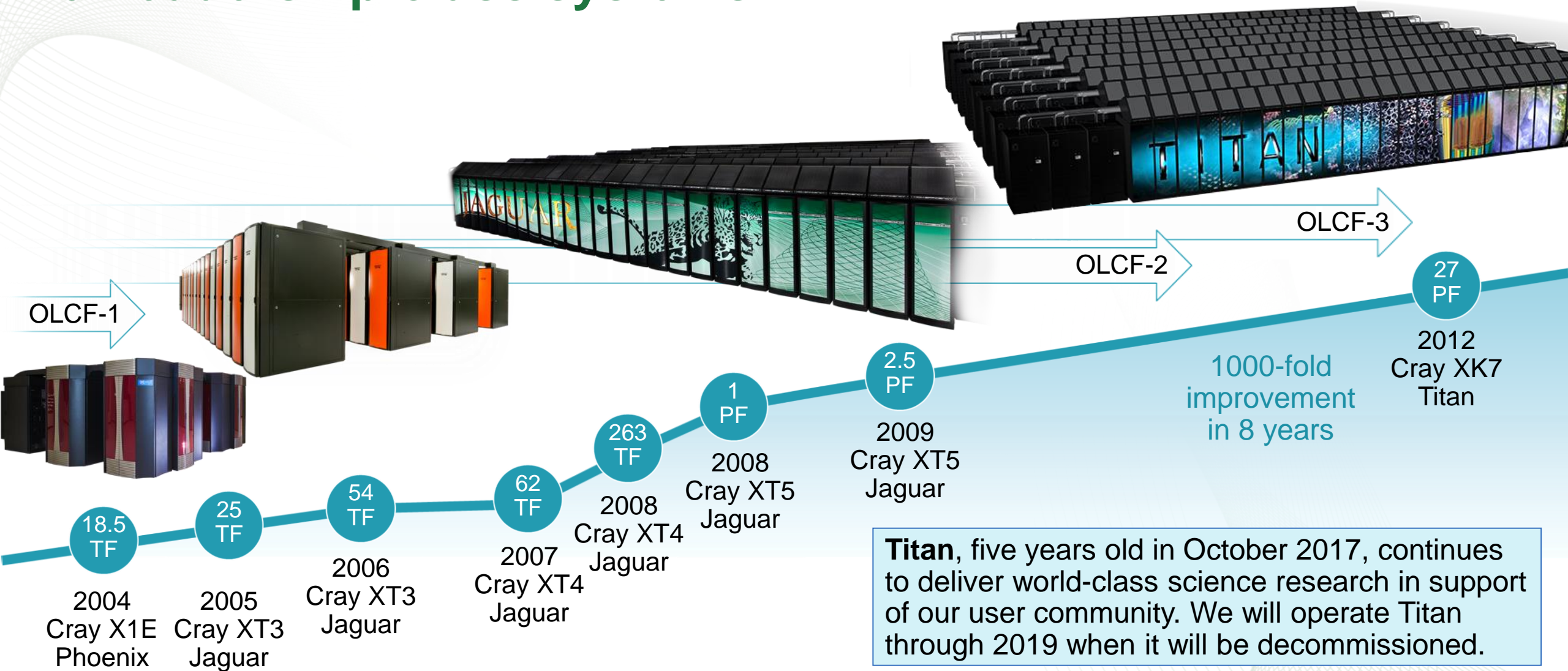19 July 2018

Hanover, MD, USA

# Outline

- Brief introduction to US DOE Leadership Computing Facility program

- Science requirements for Experimental and Observational Data (EOD) Science: highlight DOE/SC/ASCR workshop reports.

- LHC/ATLAS – OLCF integration: Big PanDA Demonstrator project at OLCF

- PanDA WMS beyond HEP: PanDA server instance at OLCF

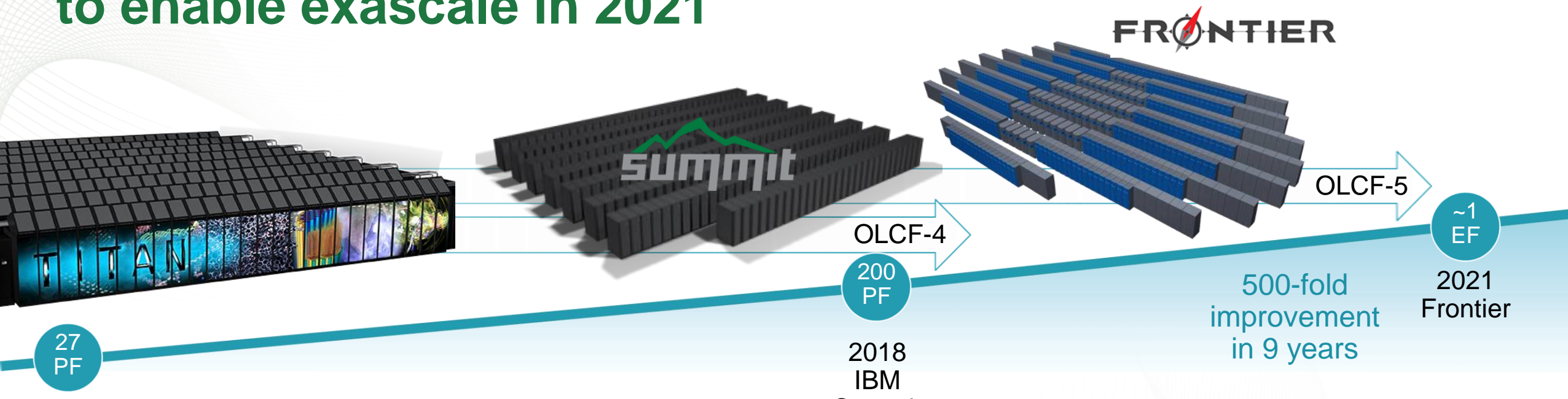- Conclusions

# What is a Leadership Computing Facility (LCF)?

- Collaborative DOE Office of Science user-facility program at ORNL and ANL

- Mission: Provide the computational and data resources required to solve the most challenging problems.

- 2-centers/2-architectures to address diverse and growing computational needs of the scientific community

- Highly competitive user allocation programs (INCITE, ALCC).

- Projects receive 10x to 100x more resource than at other generally available centers.

- LCF centers partner with users to enable science & engineering breakthroughs (Liaisons, Catalysts).

# ORNL has systematically delivered a series of leadership-class systems



OLCF-3

OLCF-2

OLCF-1

**27 PF**
2012
Cray XK7
Titan

1000-fold improvement in 8 years

**2.5 PF**
2009
Cray XT5
Jaguar

**1 PF**
2008
Cray XT5
Jaguar

**263 TF**
2008
Cray XT4
Jaguar

**62 TF**
2007
Cray XT4
Jaguar

**54 TF**
2006
Cray XT3
Jaguar

**25 TF**
2005
Cray XT3
Jaguar

**18.5 TF**
2004
Cray X1E
Phoenix

**Titan**, five years old in October 2017, continues to deliver world-class science research in support of our user community. We will operate Titan through 2019 when it will be decommissioned.

OAK RIDGE National Laboratory | LEADERSHIP COMPUTING FACILITY

# We are building on this record of success to enable exascale in 2021

**FRONTIER**

OLCF-5

OLCF-4

~1 EF

2021 Frontier

500-fold improvement in 9 years

200 PF

2018 IBM Summit

27 PF

2012 Cray XK7 Titan

June 25, 2018

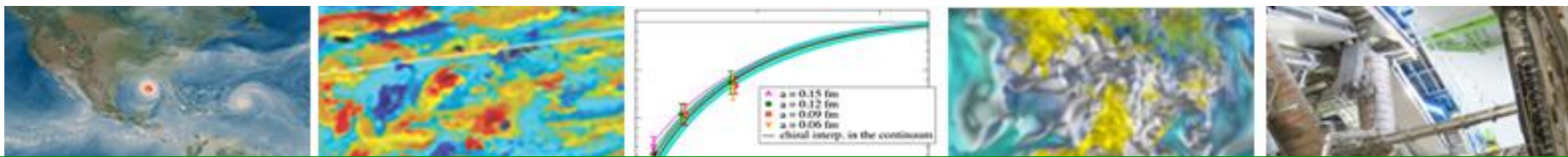**US Regains TOP500 Crown with Summit Supercomputer, Sierra Grabs Number Three Spot**

TOP500 News Team | June 25, 2018 02:37 CEST

FRANKFURT, Germany; BERKELEY, Calif.; and KNOXVILLE, Tenn.—The TOP500 celebrates its 25th anniversary with a major shakeup at the top of the list. For the first time since November 2012, the US claims the most powerful supercomputer in the

TOP 500 The List.

OAK RIDGE National Laboratory | LEADERSHIP COMPUTING FACILITY

# DOE-ASCR Exascale Requirements Reviews

ASCR facilities conducted six exascale requirements reviews in partnership with DOE Science Programs

- Goals included:
  – Identify mission science objectives that require advanced scientific computing, storage and networking in exascale timeframe
  – Determine future requirements for a computing ecosystem including data, software, libraries/tools, etc.
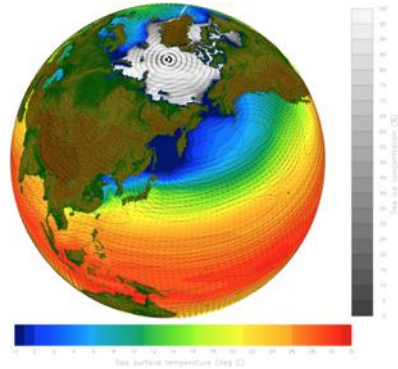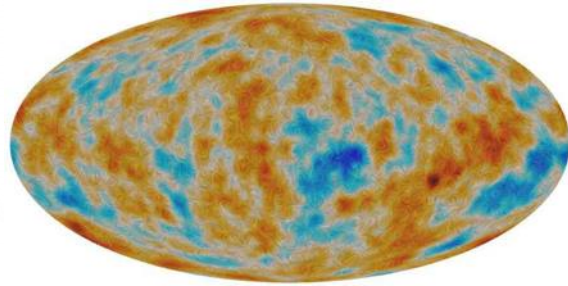
**Schedule**

| | |
|---|---|
| June 10–12, 2015 | HEP |
| November 3–5, 2015 | BES |
| January 27–29, 2016 | FES |
| March 29–31, 2016 | BER |
| June 15–17, 2016 | NP |
| Sept 27–29, 2016 | ASCR |
| March 9–10, 2017 | XCut |

All 7 workshop reports are available online: http://exascaleage.org/

# Common Themes Across DOE Science Offices
## *Data: Large-scale data storage and analysis*



Experimental and simulated data set volumes are growing exponentially. Examples: High luminosity LHC, light sources, climate, cosmology data sets ~ 100s of PBs.
Current capability is lacking.

Methods and workflows of data analytics are different than those in traditional HPC. Machine learning is revolutionizing field. Established analysis programs must be accommodated.

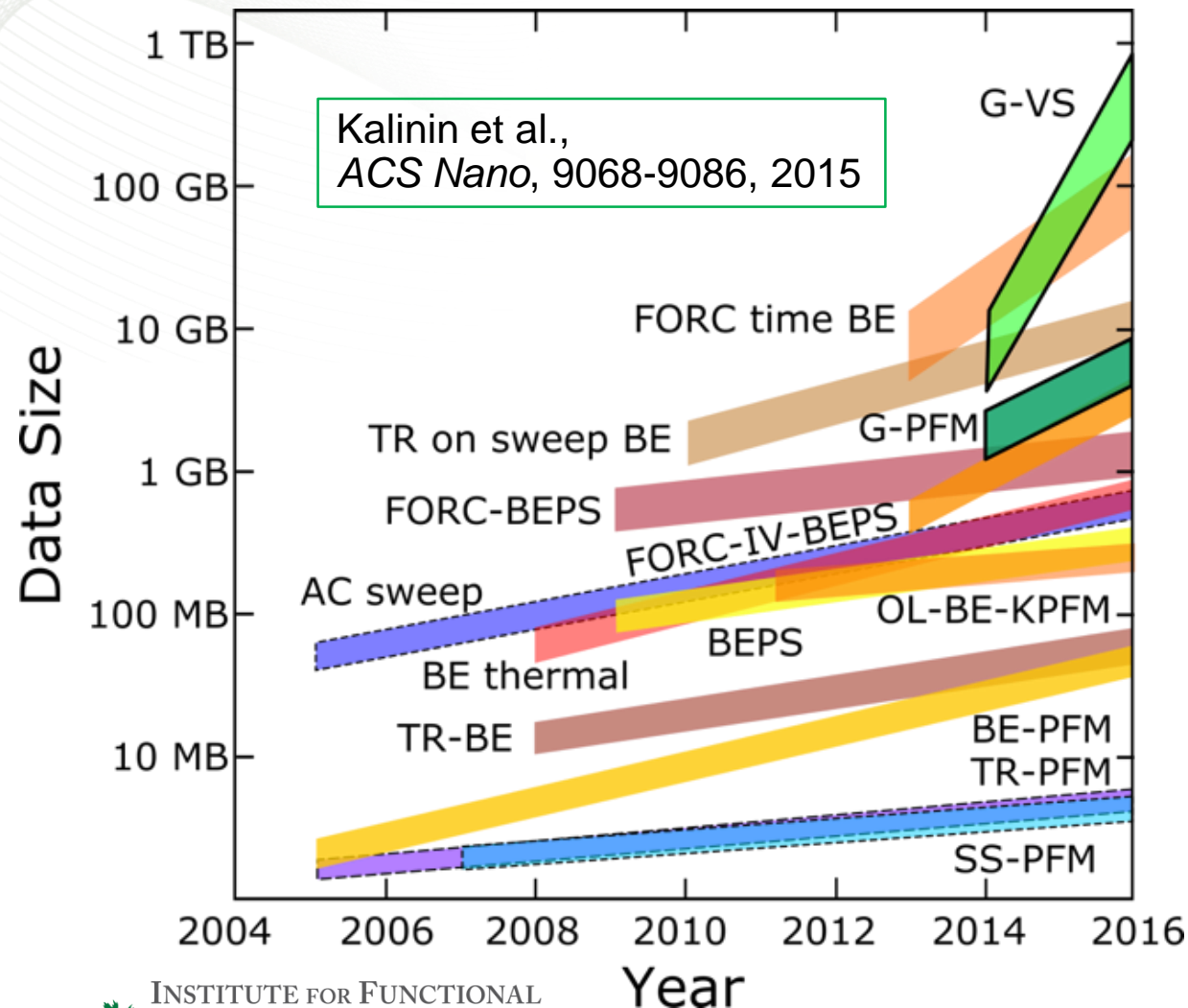# Experimental and Observational Science Data is Exploding
## *LHC Upgrade Timeline*



| | | |
|---|---|---|
| 2009 | LHC startup, √s 900 GeV | |
| 2010 | | Trigger r ~400 Hz |
| 2011 | √s=7+8 TeV, L~6x10³³cm⁻²s⁻¹, bunch spacing 50ns | Run 1 |
| 2012 | | ~25 fb⁻¹ |
| 2013 | Go to design energy, nominal luminosity - Phase 0 | |
| 2014 | LS1 | |
| 2015 | | |
| 2016 | √s=13~14 TeV, L~1x10³⁴cm⁻²s⁻¹ | k Hz |
| 2017 | | ~75-100 fb⁻¹ |
| 2018 | Injector + LHC... design luminosity | Pile-up: ~4 |
| 2019 | LS2 | |
| 2020 | | |
| 2021 | ...10³⁴cm⁻²s⁻¹, bunch spacing 25ns | Run 3 |
| 2022 | | ~350 fb⁻¹ |
| 2023 | | |
| 2024 | HL-LHC Phase II upgrade: Interaction Region, crab cavities? | |
| 2025 | √s=14 TeV, L~5x10³⁴cm⁻²s⁻¹, luminosity levelling | ~3000 fb⁻¹ |

Run1 2009 - 2013
Run2 2015 - 2018
Run3 2020-2022
Run4

A new detector

e.g. tracking, calorimeters

**In 10 years, increase by factor 10 the LHC luminosity**
➜ **More complex events**
➜ **More Computing Capacity**

# Experimental and Observational Science Data is Exploding
## *Multi-mode Scanning Probe Microscopy*

**Evolution of information volume
in multidimensional scanning probe microscopies**



Kalinin et al.,
*ACS Nano*, 9068-9086, 2015

- **Growing data sizes & complexity**
  - Cannot use desktop computers for analysis
  - ➤ **Need HPC!**
- **Multiple file formats**
  - Multiple data structures
  - Incompatible for correlation
  - ➤ **Need universal, scalable, format**

- **Disjoint and unorganized communities**
  - Similar analysis but reinventing the wheel
  - Norm: emailing each other code, data
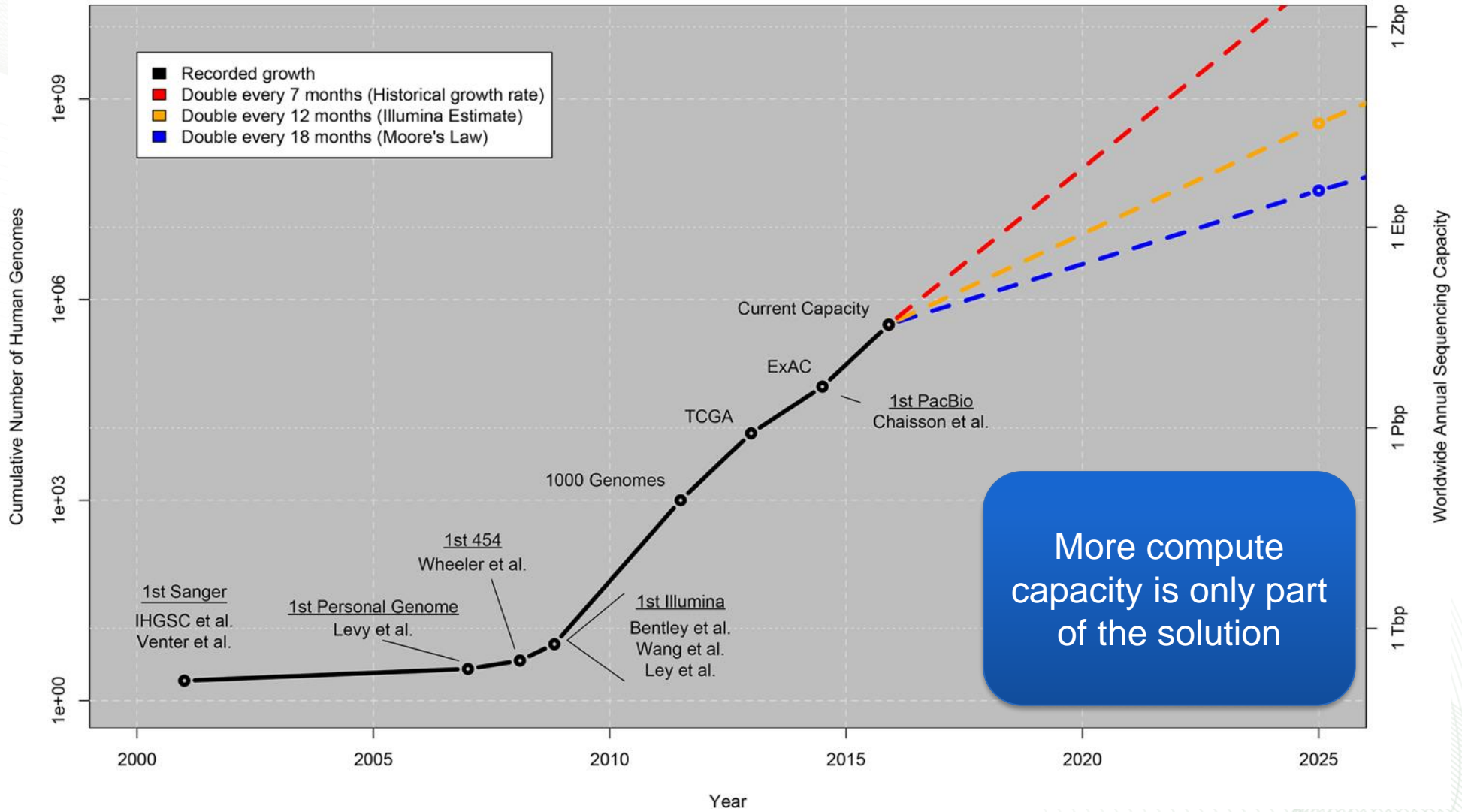  - ➤ **Need centralized repositories**
- **Proprietary, expensive software**
  - ➤ **Need robust, open, free software**

INSTITUTE FOR FUNCTIONAL
IMAGING OF MATERIALS
OAK RIDGE NATIONAL LABORATORY

OAK RIDGE
National Laboratory | LEADERSHIP COMPUTING FACILITY

# Experimental and Observational Science Data is Exploding
## *Sequence generation is outpacing Moore's law*



Growth of DNA Sequencing

**Legend:**
- Recorded growth
- Double every 7 months (Historical growth rate)
- Double every 12 months (Illumina Estimate)
- Double every 18 months (Moore's Law)

Current Capacity

ExAC

TCGA

1st PacBio
Chaisson et al.

1000 Genomes

1st 454
Wheeler et al.

1st Sanger
IHGSC et al.
Venter et al.

1st Personal Genome
Levy et al.

1st Illumina
Bentley et al.
Wang et al.
Ley et al.

More compute capacity is only part of the solution

JGI JOINT GENOME INSTITUTE
UNITED STATES DEPARTMENT OF ENERGY

OAK RIDGE National Laboratory | LEADERSHIP COMPUTING FACILITY

# Operational Demo: LHC/ATLAS-OLCF Integration



Program: DOE/SC/ASCR Next-Generation Networking for Science, Manager: Rich Carlson
Project: "BigPanDA Workflow Management on Titan for High Energy and Nuclear Physics and for Future Extreme Scale Scientific Applications,"
PI: Alexei Klimentov (BNL); Co-Pis; K. De (U. Texas-Arlington), S. Jha (Rutgers U) J.C. Wells (ORNL)

# The Opportunity for Supercomputer-Grid (HTC) Integration I

*How do we efficiently integrate supercomputing resources and <u>distributed</u> High Throughput Computing (HTC, or Grid) resources?*

- From the perspective of large supercomputer centers, how best to integrate large capability workloads, e.g., the traditional workloads of leadership computing facilities, with the large capacity workloads emerging from, e.g., experimental and observational data?

- Workflow Management Systems (WFMS) are needed to effectively integrate experimental and observation data into our data centers.

# The Opportunity for Supercomputer-Grid (HTC) Integration II

*The ATLAS experiment provides an attractive science driver, and the PanDA Workflow Management System has attractive features for capacity-capability integration*

- *The Worldwide LHC Computing Grid and a leadership computing facility (LCF) are of comparable compute capacity.*
  - *WLCG: Several 100,000's x86 compute cores*
  - *Titan: 300,000 x86 compute cores and 18,000 GPUs*

- *There is a well-defined opportunity to increase LCF utilization through backfill.*
  - *Batch scheduling prioritizing leadership-scale jobs results in ~90% utilization of available resources.*
  - *Up to 10% of Titan's cycles (~400M core hours) are available if a very large volume of capacity jobs can be run in backfill mode.*
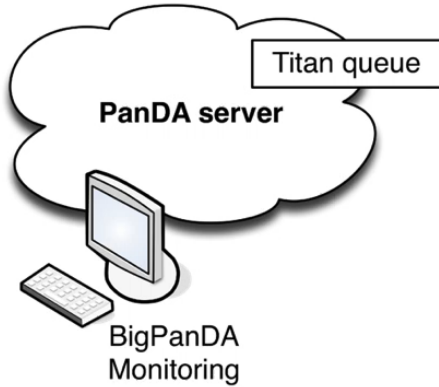
# Leadership-Job Mandate: Job Priority by Processor Count

| Bin | Min Nodes | Max Nodes | Max Walltime (Hours) | Aging Boost (Days) |
|-----|-----------|-----------|----------------------|---------------------|
| 1 | 11,250 | 18,688 | 24 | 15 |
| 2 | 3,750 | 11,249 | 24 | 5 |
| 3 | 313 | 3,749 | 12 | 0 |
| 4 | 126 | 312 | 6 | 0 |
| 5 | 1 | 125 | 2 | 0 |

Titan Hourly % Nodes Allocated by Batch System:
March 2017

# OLCF Titan Integration with ATLAS Computing

Oak Ridge

CERN

Brookhaven



Implements Dynamic Payload Shaping

# Understanding Backfill Slot Availability



Data points = 62555
x mean (red line) = 126
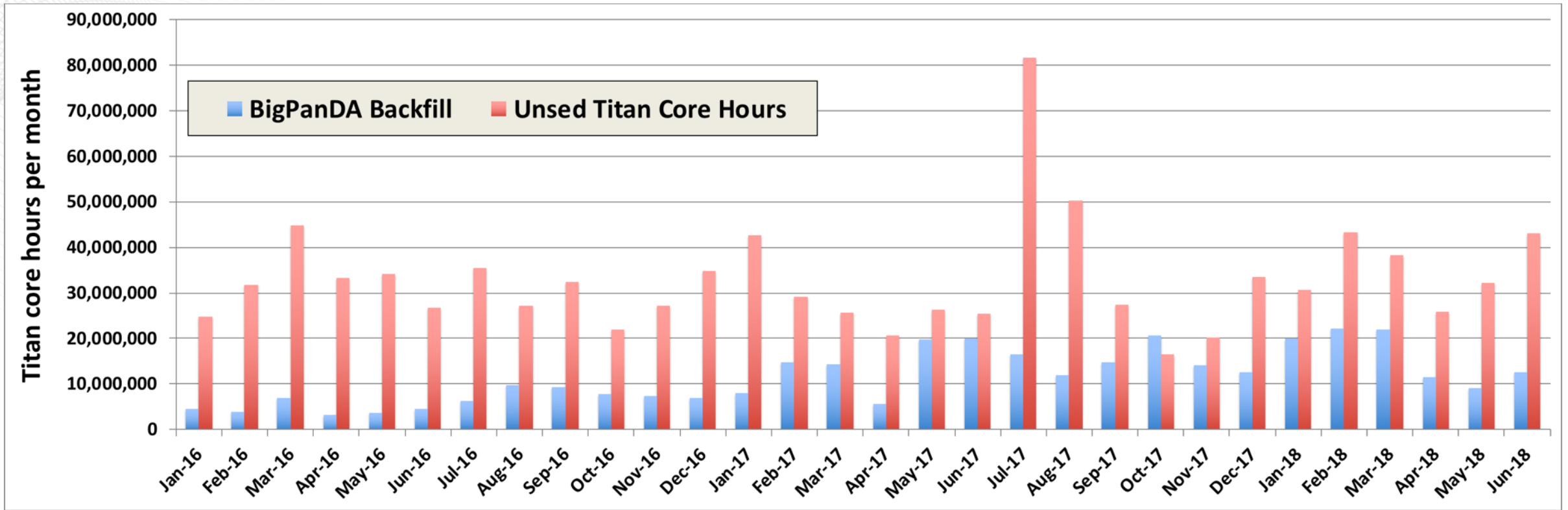y mean (orange line) = 691

2555 measures of Backfill availability on Titan during a time window.
- Mean Backfill availability: 691 worker nodes for 126 minutes.
- Up to 15K nodes for 30-100 minutes
- Large margin for optimization

OAK RIDGE
National Laboratory | LEADERSHIP COMPUTING FACILITY

- Currently shaping number of nodes per payload
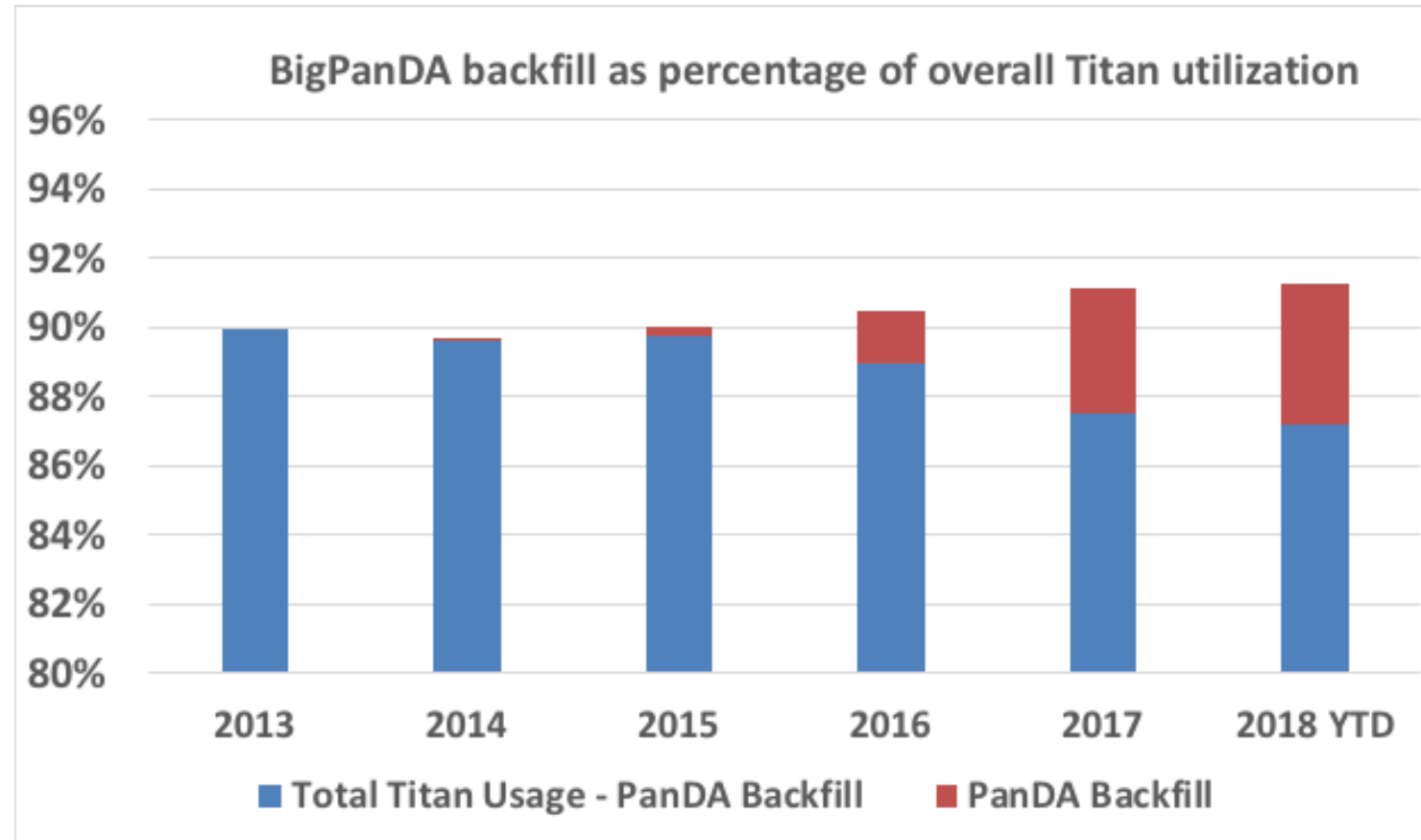- No shaping of payload duration

# Operational Demo: Scaling Up Active Backfill w/ Job Shaping



- Consumed 340 Million Titan core hours from January 2016 to present
  - This is 2.9 percent of total available time on Titan over this period
- Remaining used backfill slots are often too short or too small for assigned ATLAS payloads

# Operational Demo: Scaling Up Active Backfill

- Increased Titan's utilization by ~2 percent over historical trends
    - May have displaced ~ 2% of Titan's small jobs in the batch queue. This is currently under evaluation.
- Preemption of PanDA payloads to be evaluated
    - Checkpointing needed
    - Checkpointing will be enabled by Event Service: ability to save incremental, event-by-event results.
- Currently, only shaping payloads through number of nodes employed
    - Additional opportunity for shaping through payload duration.

BigPanDA backfill as percentage of overall Titan utilization



■ Total Titan Usage - PanDA Backfill    ■ PanDA Backfill

OAK RIDGE | LEADERSHIP COMPUTING FACILITY
National Laboratory

# US ATLAS won a large ALCC allocation at ALCF/OLCF/NERSC

**2017 ASCR Leadership Computing Challenge (ALCC) Application**

**Assigned Proposal ID:**

286

**Title:**

2017 ASCR Leadership Computing Challenge (ALCC) Application

Consortium/End-Station Proposal

**Principal Investigator:**

John T Childers, Argonne National Laboratory, Tel: 3313024647. Email: jchilders@anl.gov

**Project collaborators:**

Thomas LeCompte (Argonne National Laboratory)
Doug Benjamin (Duke University)
Radja Boughezal (Argonne National Laboratory)
Paolo Calafiura (Lawrence Berkeley National Laboratory)
Stefan Hoeche (SLAC National Laboratory)
Burt Holzman (Fermi National Accelerator Laboratory)
Alexei Klimentov (Brookhaven National Laboratory)
Jim Kowalkowski (Fermi National Accelerator Laboratory)
Frank Petriello (Northwestern University)
Vakho Tsulaia (Lawrence Berkeley National Laboratory)
Craig Tull (Lawrence Berkeley National Laboratory)
Thomas Uram (Argonne National Laboratory)
Torre Wenaus (Brookhaven National Laboratory)

**HPC resources requested:**

**Titan:** 80 x10^6 Titan-core hours, ~20 TB online storage, ~0 TB offline storage

**Mira:** 18 x10^6 Mira-core hours, ~20 TB online storage, ~0 TB offline storage
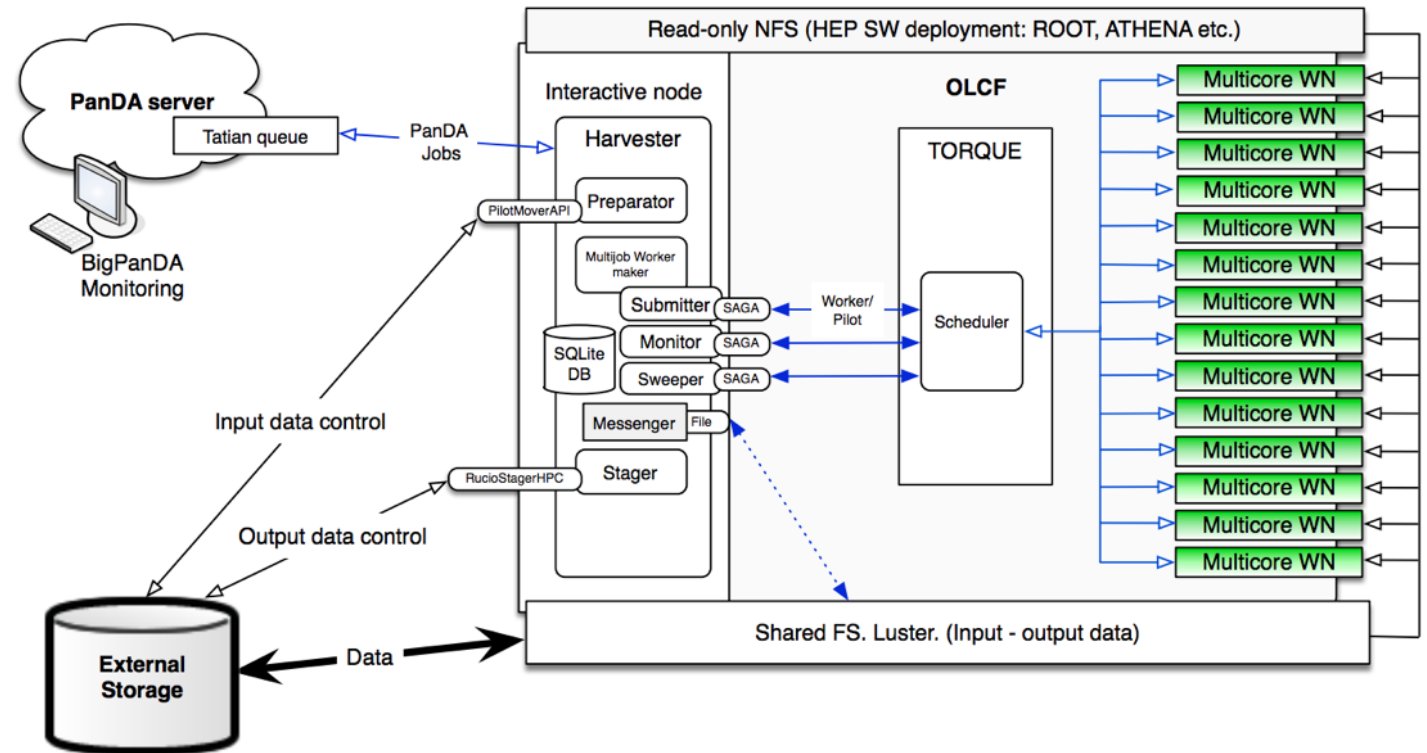
**Theta:** 45.5 x10^6 Theta-core hours, ~20 TB online storage, ~0 TB offline storage

**Cori (Cray XC40 Intel Xeon Phi KNL nodes):** 36 x10^6 NERSC-core hours, ~100 TB scratch storage, ~10 TB proje

**Cori/Edison (Cray XC40/30 Intel Xeon nodes):** 33 x10^6 NERSC-core hours, ~100 TB scratch storage, ~10 TB proj
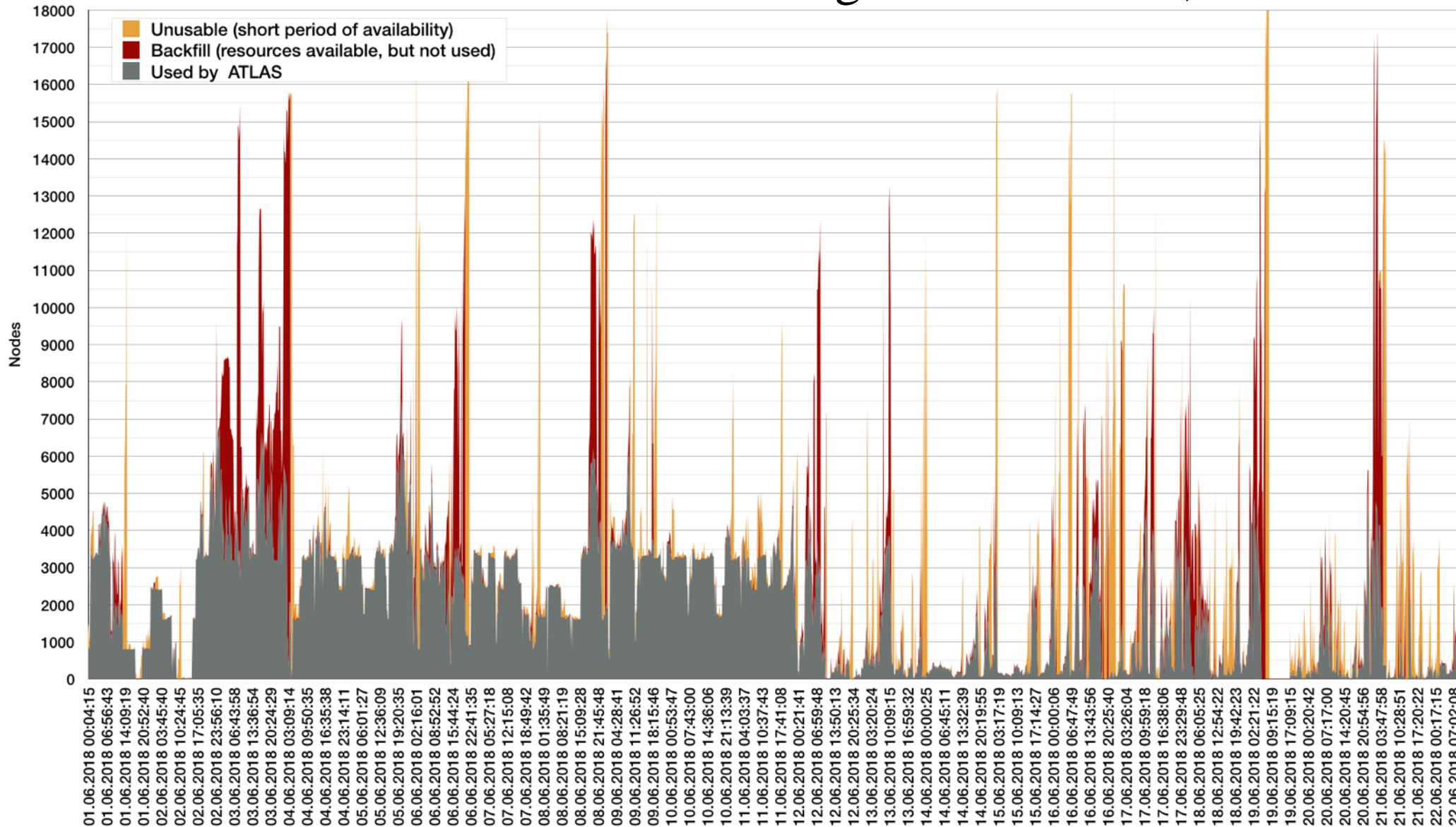
**Additional notes:** At NERSC, the request for Storage is in total. That means 100TB (Scratch) + 10TB (Project) across both Cori and Cori/Edison. Ideally these two spaces would be visible to both machines.

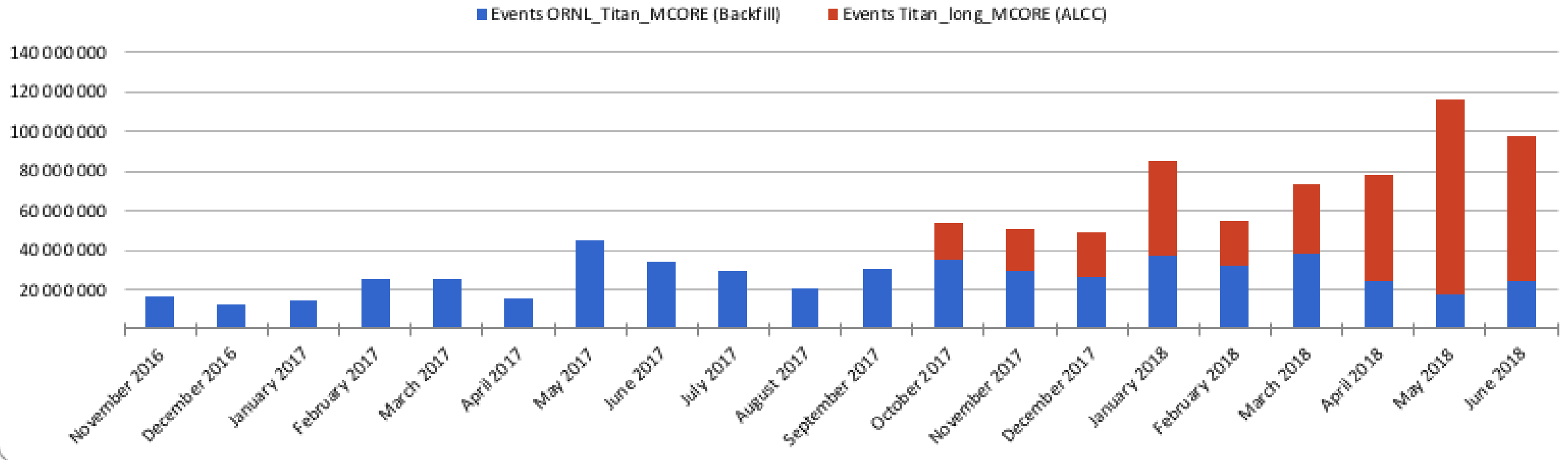Implement ALCC Project at OLCF using Harvester

# ATLAS@OLCF: Batch Queue Submission & Active Backfill

- Backfill utilization in 1 June through 22 June 2018, 10-min data frequency



Plot:
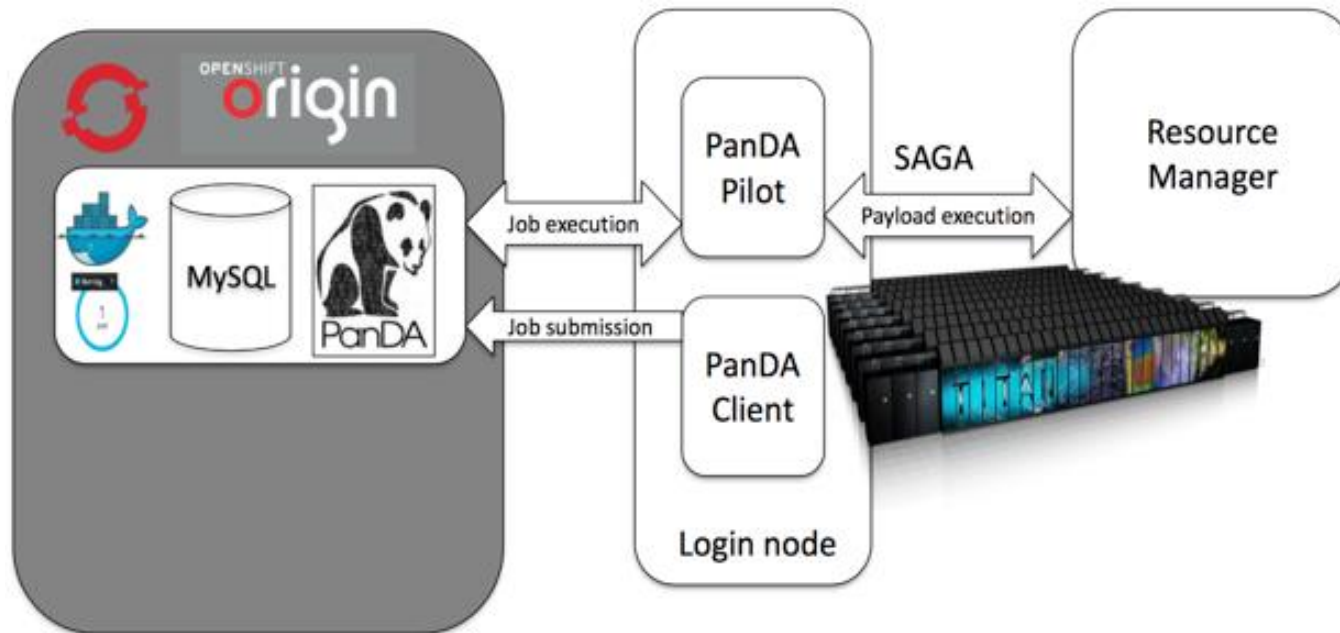Danila Oleynik

Events (Backfill), Events (ALCC) and Events

- Since Nov. 2016, 508 million ATLAS events computed via backfill
- Since Oct. 2017, 395 million TLAS events computed via "normal" batch queue
  - Increases in batch queue event generation beginning in Feb. 2018 show the impact of Harvester

# PanDA Server at OLCF: Broad application across domains

- In March 2017 a new PanDA server instance has been established at ORNL to serve various experiments. This installation the first at OLCF to demonstrate application of a container cluster management and orchestration system, Red Hat OpenShift Origin.

- OpenShift, when fully in production, will give OLCF users the ability to deploy and manage their own middleware and infrastructure services

  - https://www.olcf.ornl.gov/2017/06/05/olcf-testing-new-platform-for-scientific-workflows/
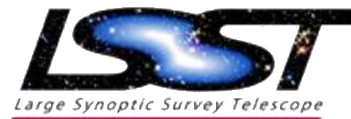
Key Contributors:
Jason Kincl (ORNL),
Ruslan.Mashinistov (BNL)
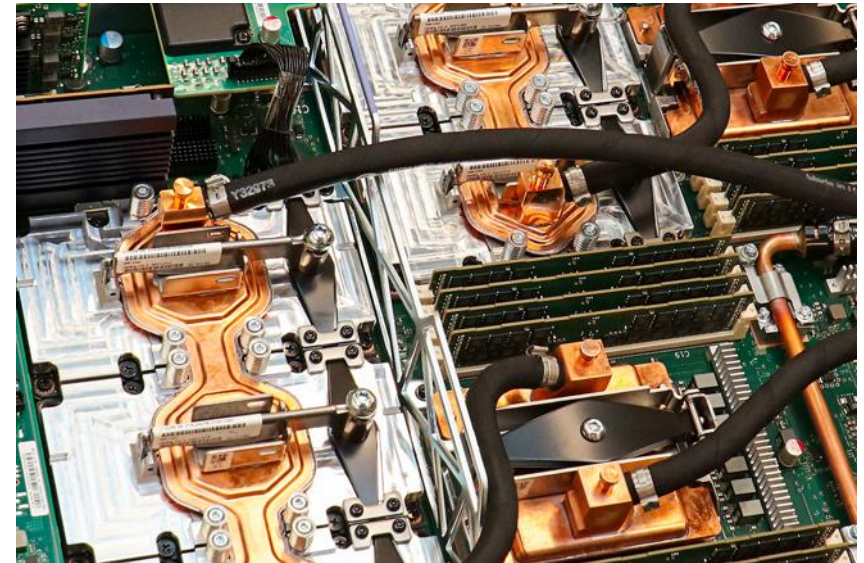
# PanDA Server at OLCF: PanDA WMS beyond HEP

- Biology / Genomics: Center for Bioenergy Innovation at ORNL
- Molecular Dynamics: Prof. K. Nam (U. Texas-Arlington)
- nEDM, (neutron Electric Dipole Moment Experiment, ORNL
- IceCube Experiment
- Blue Brain Project (BBP), EPFL
- SST (Large Synoptic Survey Telescope) project
- LQCD, US QCD SciDAC Project

# Coming in 2018: Summit will replace Titan as the OLCF's leadership supercomputer

Summit is the Department of Energy's Oak Ridge National Laboratory's newest supercomputer for open science.
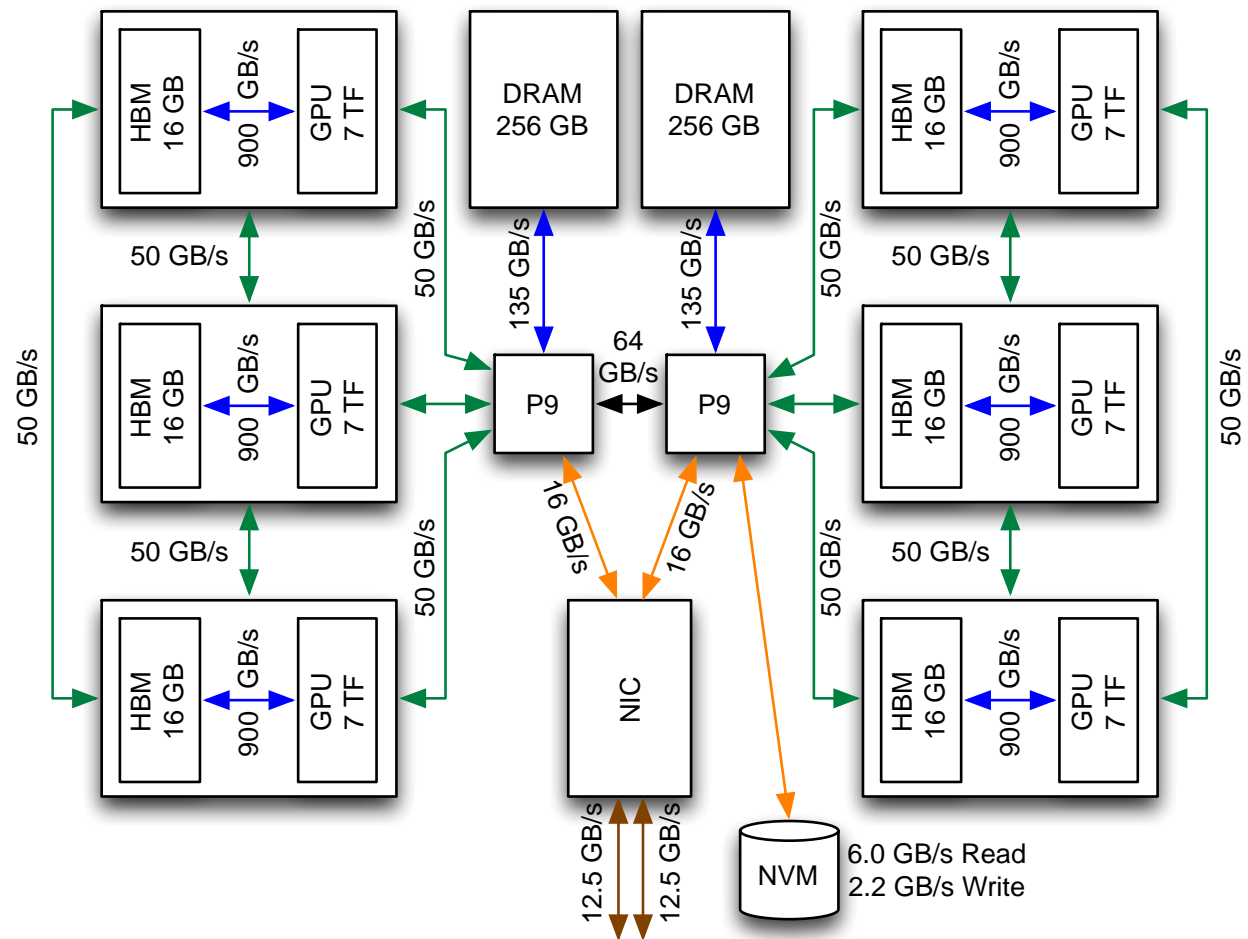
# Summit compared to Titan

- Many fewer nodes

- Much more powerful nodes

- Much more memory per node and total system memory

- Faster interconnect

- Much higher bandwidth between CPUs and GPUs

- Much larger and faster file system

| Feature | Titan | Summit |
|---|---|---|
| Peak Flops | 27 PF | 200 PF |
| Application Performance | Baseline | 5-10x Titan |
| Number of Nodes | 18,688 | ~4,600 |
| Node performance | 1.4 TF | > 40 TF |
| Memory per Node | 32 GB DDR3 + 6 GB GDDR5 | 512 GB DDR4 + 96 GB HBM |
| NV memory per Node | 0 | 1600 GB |
| Total System Memory | 710 TB (600 TB DDR3 + 110 TB GDDR5) ) | 10 PB (2.3 PB DDR4 + 0.4 PB HBM + 7.4 PB NVRAM) |
| System Interconnect (node injection bandwidth) | Gemini (6.4 GB/s) | Dual Rail EDR-IB (23 GB/s) |
| Interconnect Topology | 3D Torus | Non-blocking Fat Tree |
| Processors per node | 1 AMD Opteron™ 1 NVIDIA Kepler™ | 2 IBM POWER9™ 6 NVIDIA Volta™ |
| File System | 32 PB, 1 TB/s, Lustre® | 250 PB, 2.5 TB/s, GPFS™ |
| Peak power consumption | 9 MW | 13 MW |

# Summit Node Overview



| | | |
|---|---|---|
| TF | 42 TF (6x7 TF) | |
| HBM | 96 GB (6x16 GB) | |
| DRAM | 512 GB (2x16x16 GB) | |
| NET | 25 GB/s (2x12.5 GB/s) | |
| MMsg/s | 83 | |

- → HBM/DRAM Bus (aggregate B/W)
- → NVLINK
- → X-Bus (SMP)
- → PCIe Gen4
- → EDR IB

HBM & DRAM speeds are aggregate (Read+Write).
All other speeds (X-Bus, NVLink, PCIe, IB) are bi-directional.

# Summit will be the world's smartest supercomputer for open science
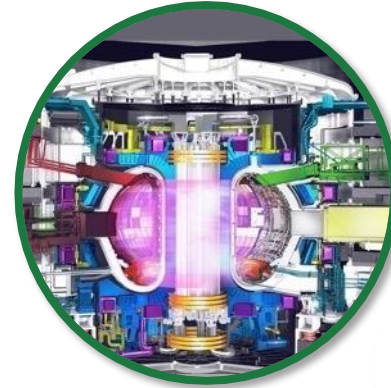*But what can a smart supercomputer do?*



## Science challenges for a smart supercomputer:

**Identifying Next-generation Materials**
By training AI algorithms to predict material properties from experimental data, longstanding questions about material behavior at atomic scales could be answered for better batteries, more resilient building materials, and more efficient semiconductors.
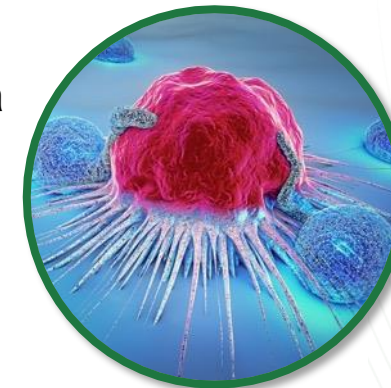
**Predicting Fusion Energy**
Predictive AI software is already helping scientists anticipate disruptions to the volatile plasmas inside experimental reactors. Summit's arrival allows researchers to take this work to the next level and further integrate AI with fusion technology.

**Deciphering High-energy Physics Data**
With AI supercomputing, physicists can lean on machines to identify important pieces of information—data that's too massive for any single human to handle and that could change our understanding of the universe.

**Combating Cancer**
Through the development of scalable deep neural networks, scientists at the US Department of Energy and the National Cancer Institute are making strides in improving cancer diagnosis and treatment.

# Emerging Science Activities:
## Selected Machine Learning Projects on Titan: 2016-2017

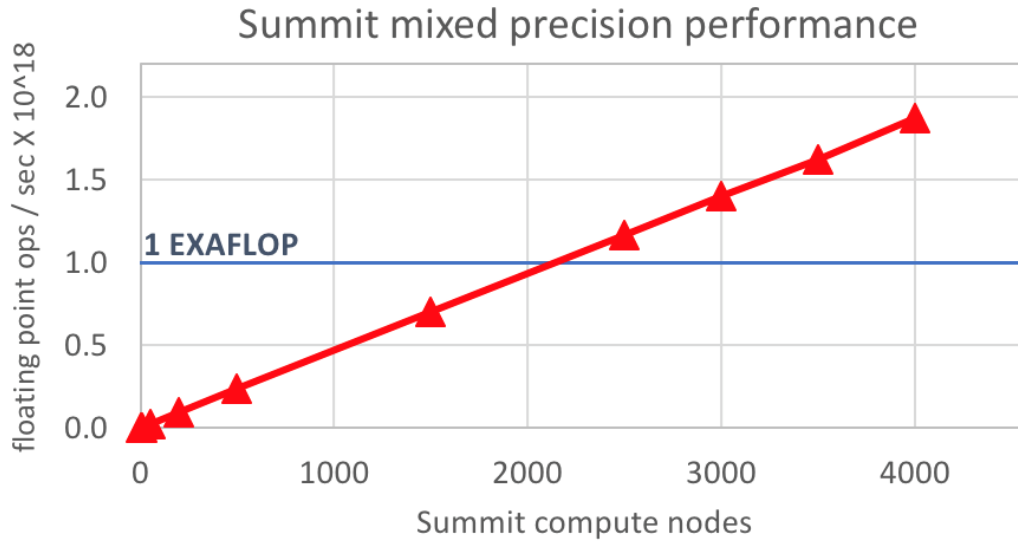| Program | PI | PI Employer | Project Name | Allocation (Titan core-hrs) |
|---|---|---|---|---|
| ALCC | Robert Patton | ORNL | Discovering Optimal Deep Learning and Neuromorphic Network Structures using Evolutionary Approaches on High Performance Computers | 75,000,000 |
| ALCC | Gabriel Perdue | FNAL | Large scale deep neural network optimization for neutrino physics | 58,000,000 |
| ALCC | Gregory Laskowski | GE | High-Fidelity Simulations of Gas Turbine Stages for Model Development using Machine Learning | 30,000,000 |
| ALCC | Efthimions Kaxiras | Harvard U. | High-Throughput Screening and Machine Learning for Predicting Catalyst Structure and Designing Effective Catalysts | 17,500,000 |
| ALCC | Georgia Tourassi | ORNL | CANDLE Treatment Strategy Challenge for Deep Learning Enabled Cancer Surveillance | 10,000,000 |
| DD | Abhinav Vishnu | PNNL | Machine Learning on Extreme Scale GPU systems | 3,500,000 |
| DD | J. Travis Johnston | ORNL | Surrogate Based Modeling for Deep Learning Hyper-parameter Optimization | 3,500,000 |
| DD | Robert Patton | ORNL | Scalable Deep Learning Systems for Exascale Data Analysis | 6,500,000 |
| DD | William M. Tang | PPPL | Big Data Machine Learning for Fusion Energy Applications | 3,000,000 |
| DD | Catherine Schuman | ORNL | Scalable Neuromorphic Simulators: High and Low Level | 5,000,000 |
| DD | Boram Yoon | LANL | Artificial Intelligence for Collider Physics | 2,000,000 |
| DD | Jean-Roch Vlimant | Caltech | HEP DeepLearning | 2,000,000 |
| DD | Arvind Ramanathan | ORNL | ECP Cancer Distributed Learning Environment | 1,500,000 |
| DD | John Cavazos | U. Delaware | Large-Scale Distributed and Deep Learning of Structured Graph Data for Real-Time Program Analysis | 1,000,000 |
| DD | Abhinav Vishnu | PNNL | Machine Learning on Extreme Scale GPU systems | 1,000,000 |
| DD | Gabriel Perdue | FNAL | MACHINE Learning for MINERvA | 1,000,000 |
| | | **TOTAL** | | **220,500,000** |

# Summit is still under construction

- We expect to accept the machine in Summer of 2018, allow early users on this year, and allocate our first users through the INCITE program in January 2019.
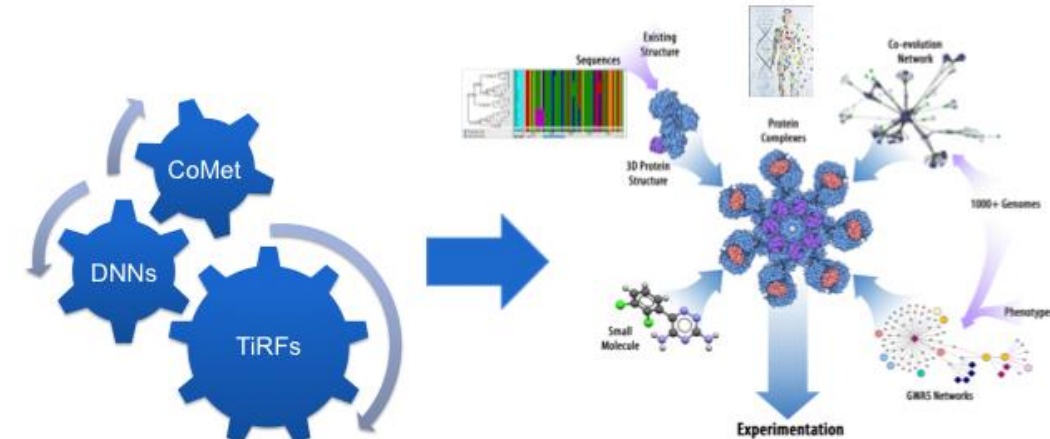
# CoMet: ExaOp Comparative Genomics on Summit

**Dan Jacobson, Wayne Joubert (ORNL)**

- Modified 2-way CCC algorithm uses NVIDIA Volta Tensor Cores   and cuBLAS library to compute counts of bit values

- Near-ideal weak scaling to 4000 nodes (87% of Summit) – **1.8 EF (FP16)** mixed precision performance reached; 234 quadrillion element comparisons / sec attained

- **4.5X faster** than previous optimized bitwise CCC/sp code on Summit

- **80 TF (FP16)** achieved per GPU for full algorithm – cuBLAS performance per GPU nearly **100 TF (FP16)**

- Expect **2+ EF (FP16) mixed precision achievable** on full Summit system

- **Summit allows us to:**
  – Discover co-evolutionary relationships across a population of genomes at an unprecedented scale
  – Discover epistatic interactions for Opioid Addiction
    - 2018 Gordon Bell Prize finalist



Summit mixed precision performance

W. Joubert, J. Nance, S. Climer, D. Weighill, D. Jacobson, "Parallel Accelerated Custom Correlation Coefficient Calculations for Genomics Applications," arxiv 1705.08213 [cs], *Parallel Computing,* accepted.

# Summary & Conclusions

- Over the past 5 years, Titan has delivered the DOE Leadership Computing Mission at ORNL: delivering science and constraining energy consumption.

- ORNL is advancing hybrid-accelerated supercomputing based on success of Titan; Summit is our next step taking place this year.

- DOE Office of Science is advancing an integrated vision for exascale computing ecosystem, including data-intensive and distributed applications dependent on networks and storage, e.g., experimental and observational data.

- ATLAS/PanDA deployment of Titan shows the potential of distributed, high-throughput computing to be integrated with high-performance computing infrastructure.
  – Offers significant value for other projects beyond HEP

- Summit's high-bandwidth, GPU-accelerated architecture should be very effective for data analytics and machine learning.

# Acknowledgements

- DOE Office of Science, Advanced Scientific Computing Research
  - Next-Generation Networking for Science, Manager: Rich Carlson
  - Oak Ridge Leadership Computing Program, Manager: Christine Chalk
- Collaboration with:
  - U.S. ATLAS (DOE/SC/HEP)
  - nEDM (DOE/SC/NP)
  - Plant-Microbe Interfaces (DOE/SC/BER)
  - Center for Bioenergy Innovation, CBI (DOE/SC/BER)
  - US QCD SciDAC Project (DOE/SC)
  - LSST (DOE/SC/HEP)
  - ICECUBE (NSF)
  - Blue Brain Project (EPFL)

OAK RIDGE National Laboratory | LEADERSHIP COMPUTING FACILITY

**Questions?  Jack Wells, wellsjc@ornl.gov**

OAK RIDGE National Laboratory | LEADERSHIP COMPUTING FACILITY