

WHPC Lightning Talks



Programme Committee

- Elsa Gonsiorowski, Lawrence Livermore National Laboratory, USA
- Raquell Holmes, Improvscience
- Elizabeth Bautista, NERSC, Lawrence Berkeley National Laboratory, USA
- Jo Adegbola, Amazon Web Services, USA
- Mozghan Kabiri, University of Sheffield, UK
- Karen Devine, Sandia National Laboratory, USA
- Zhiling Lan, Illinois Institute of Technology, USA
- Hadia Ahmed, Lawrence Berkeley National Laboratory, USA
- Lavanya Ramakrishnan, Lawrence Berkeley National Laboratory, USA
- Debbie Bard, Lawrence Berkeley National Laboratory, USA
- Rosa Filgueira, EPCC, UK
- Danielle Sikich, Intel, USA
- Mahwish Arif, CAM, UK
- Baiou Shi, PSU, USA
- Neelofer Banglawala, EPCC, UK
- Catherine Schumann, Oak Ridge National Laboratory, USA
- Shubbhi Taneja, Sonoma State University, USA

Thanks to the mentors for volunteering!



FIRST EVER WOMEN-IN-HPC SUMMIT!

In partnership with Simon Fraser University

Call For Participation (Papers, Tutorials & Posters) Open!

<https://womeninhpc.org/events/summit-2020>

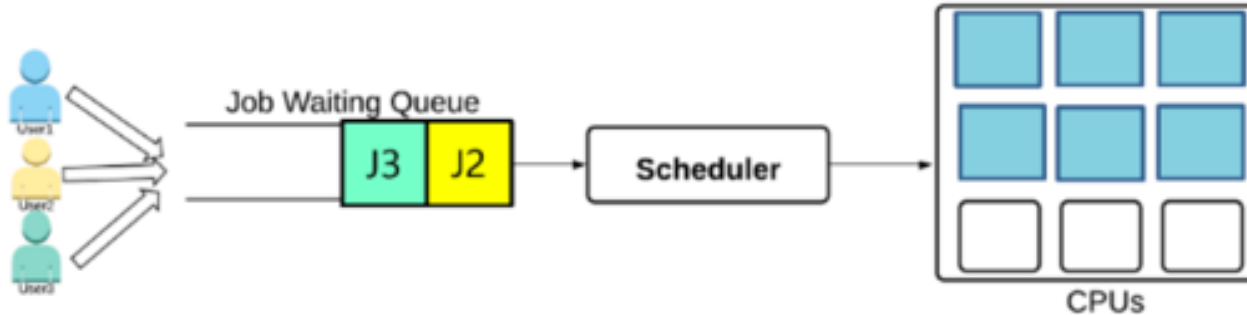
Multi-Resource Scheduling in HPC

Yuping Fan | yfan22@hawk.iit.edu
Illinois Institute of Technology



Job Scheduling in HPC

- Decides when and where to execute jobs
 - Policies: FCFS, SJF
 - Examples: Slurm, PBS, Mesos
- Traditionally, focuses on CPU utilization



- Multiple Resources in HPC
 - Local resources: CPU, GPU, SSD
 - Shared resources: burst buffer, RAN
- **How can we efficiently use of multiple HPC resources?**

Existing Methods

- Naïve method: no optimization; first job in the queue can run when there are sufficient resources
- Optimize the utilization of one resource (e.g., CPU)
- Optimize the weighted sum of multiple resource utilization
- **Problem:** They are single-objective and can only provide one solution, but multiple optimal solution exist when scheduling multiple resources

BBSched: a Novel Multi-Resource Scheduling Framework

- Formulates the multi-resource scheduling problem in to multi-objective optimization problem
 - Maximize multiple independent objectives
 - Node utilization: $f_1(x) = \sum_{i=1}^w n_i \times x_i$
 - Burst Buffer utilization: $f_2(x) = \sum_{i=1}^w b_i \times x_i$
- Rapidly solves the problem by genetic algorithm
- Provides multiple scheduling options for system administrators
- Flexible to embrace new resources added to HPC systems

Moment Representation in the Lattice Boltzmann Method on Massively Parallel Hardware

Madhurima Vardhan | madhurima.vardhan@duke.edu
PhD Candidate | Biomedical Engineering
Duke University



CHANGING
THE FACE
OF HPC

Why is the lati... LBM algorithm an active area of research?

LBM is a widely used in CFD algorithm

- Explicit and second order accurate
- Highly scalable
- Handles complex geometries

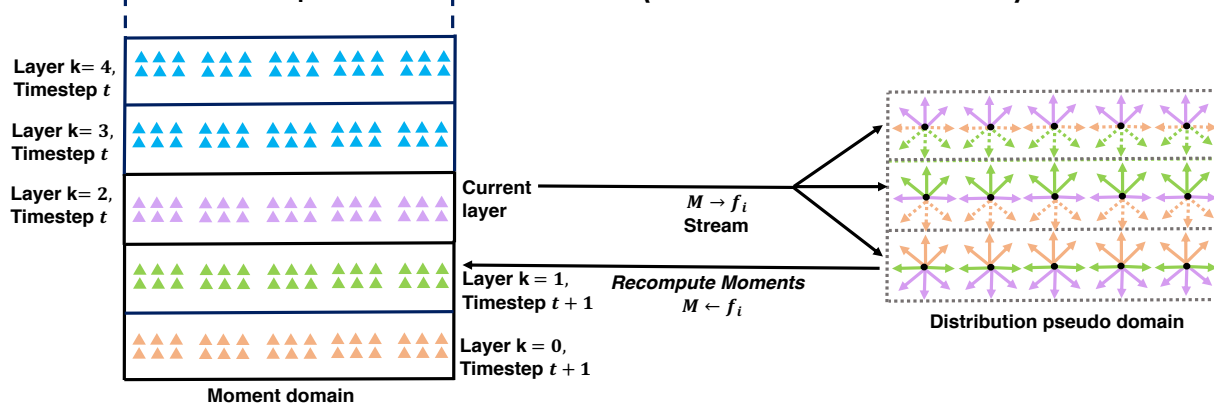
However, LBM is

- Memory-expensive, 38 doubles per lattice site
- Imposes limit on the resolution

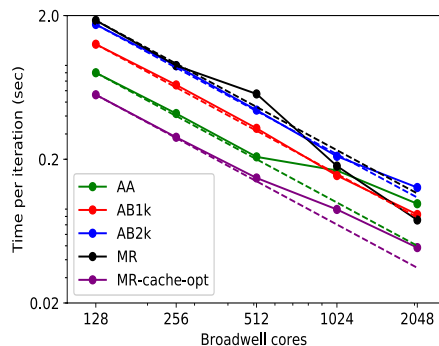
What do we propose to reduce memory requirement?

Adapt Regularized LBM

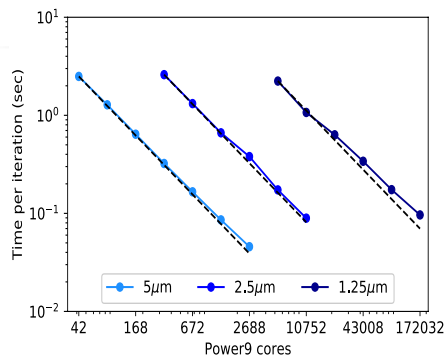
- First three order moments can accurately solve fluid dynamic equations¹
- 10 doubles per lattice site (~74% reduction)



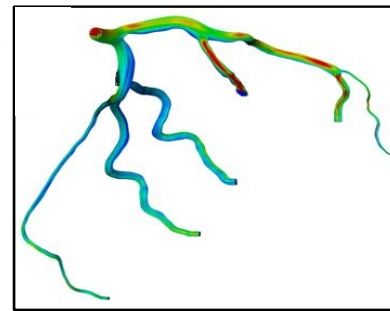
What have our efforts yielded so far?



Strong scaling comparison on Intel Broadwell Cores



Strong scaling on Summit Power9 cores



Patient-specific hemodynamics in left coronary arterial circulations

- Relative to other methods, best time to solution
- Near ideal scaling on Summit Power9 cores
- Applicable to real-world complex problems
- Challenges – Accelerators and GPUs

libCEED:

Lightweight High-Order Finite Elements Library

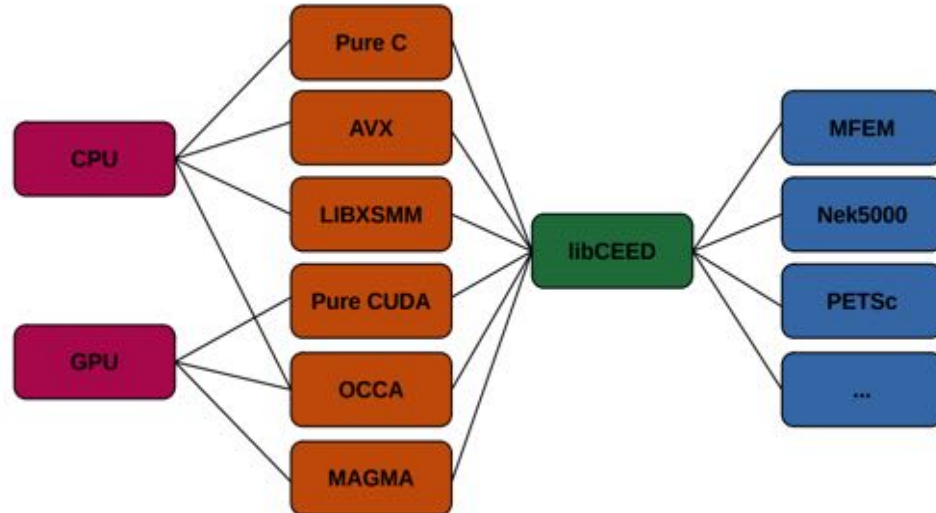


Valeria Barra | valeria.barra@colorado.edu
University of Colorado Boulder

What is the point of having a **Ferrari**,
if you drive it stuck in second gear?

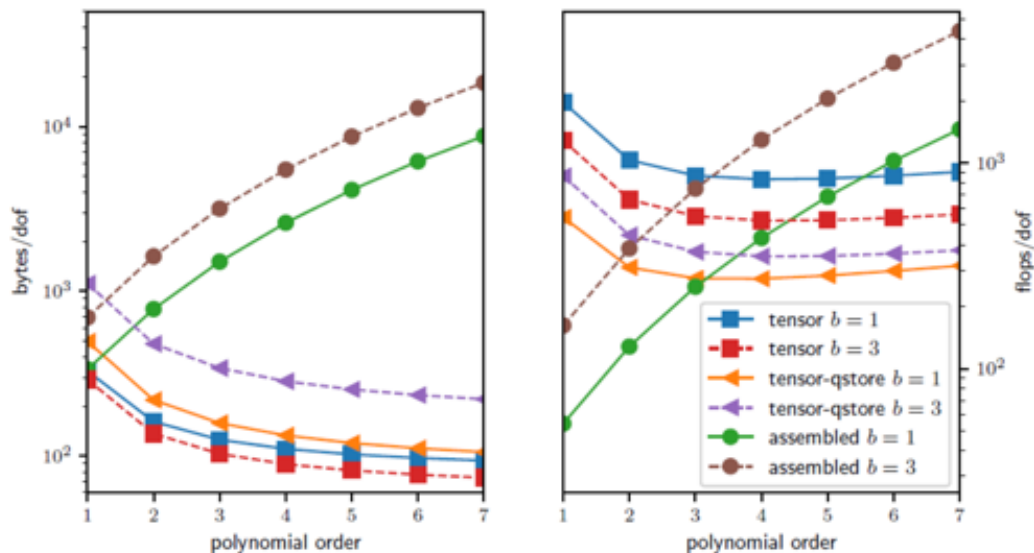
libCEED is a low-level API for achieving high-performance scientific computing
on different architectures.

libCEED supports run-
time selection of
implementations tuned
for a variety of
computational device
types (e.g., CPUs,
GPUs, etc).



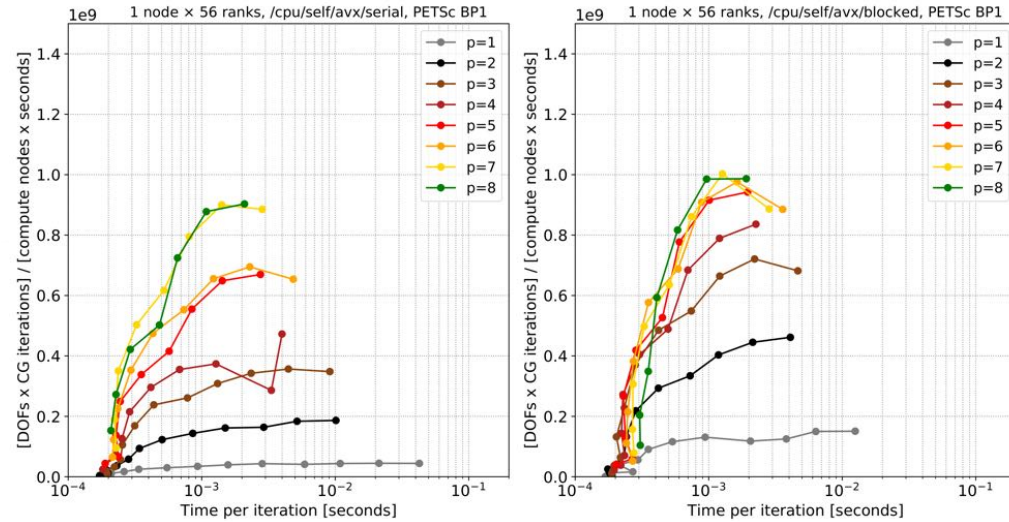
Why matrix-free?

A sparse matrix is no longer a good representation for high-order operators.
libCEED has a purely algebraic interface for matrix-free operator representation

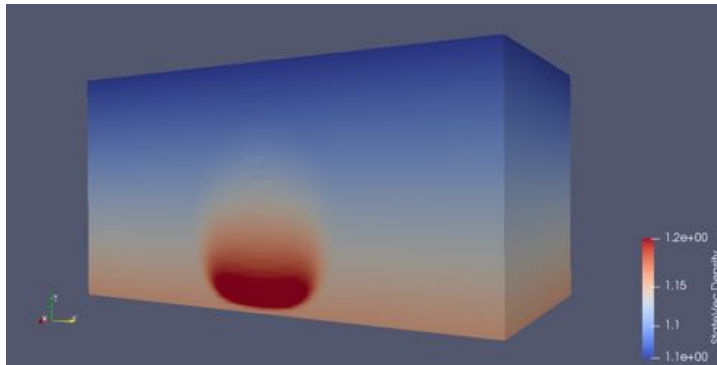


Memory bandwidth (left) and flops per dof (right) to apply a Jacobian matrix, obtained from discretization of a b -variable PDE system.

Performance



Application Example



Not just toy problems: A fast and efficient **Navier-Stokes** solver.

Polynomial order of spectral elements: $p=10$

Computational domain:
[0,6000] x [0,6000] x [0,3000] m

Elem. Resolution: 500 m

893101 Nodes

Analysing Digital Historical Textual Data in HPC clusters and Cloud using Apache Spark and Jupyter Notebooks

Rosa Filgueira | rosa.filgueira@ed.ac.uk
EPCC, University of Edinburgh



Context and Motivation

Working with

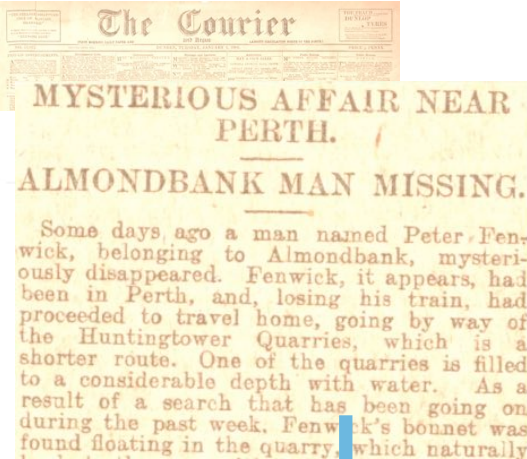
- Historians, Humanities and computational linguistics researchers
- Large digital collections been available for research

Motivation – toolbox for Historians & Humanities communities

- Hunger for large scale text mining facilities
- Limited capacity and/or skills to use:
 - HPC/Cloud environments
 - analytic frameworks to create applications

Challenges

- Several large digital collections (semi-structured data)
- Different levels of quality of data – OCR
- Data with different physical representations and schemas



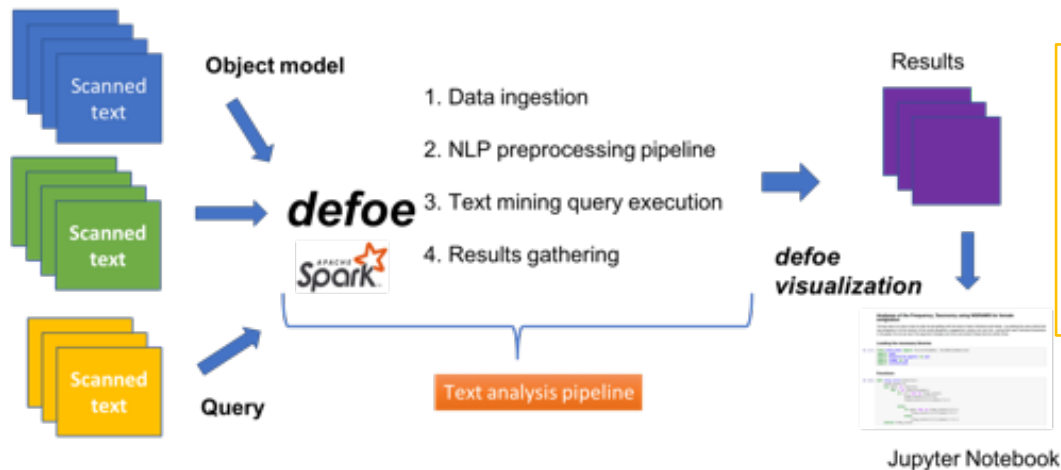
Digitalized Newspaper issue

```
...  
<text.title>  
  <pg pgref="5" clipref="1"  
    pos="4069,3036,4949,3154"/>  
  <p>  
    <wd pos="4069,3036,4949,3154">MYSTERIOUS AFFAIR NEAR  
PERTH.</wd>  
  </p>  
</text.title>  
<text.cr>  
  <pg pgref="5" clipref="1"  
    pos="4039,3191,4987,4235"/>  
  <p>  
    <wd pos="4041,3192,4496,3241">ALMONDBANK</wd>  
    <wd pos="4523,3200,4663,3246">MAN</wd>  
    <wd pos="4696,3198,4976,3250">MISSING.</wd>  
    <wd pos="4085,3290,4189,3323">Some</wd>  
    <wd pos="4214,3290,4312,3329">days,</wd>  
  </p>  
...
```

XML schema

defoe: new toolbox for historical research

Digital Collections



- Extracting knowledge from historical data.
- Running parallel text analyses across large collections.
- Rich set of text mining queries.
- Scalable and distributed analyses.
- NLP preprocessing techniques to mitigate OCR errors.
- Portability - computing environments and collections.

<https://github.com/alan-turing-institute/defoe>

https://github.com/alan-turing-institute/defoe_visualization

Developing the QuEST Library for Quantum Circuit Simulation

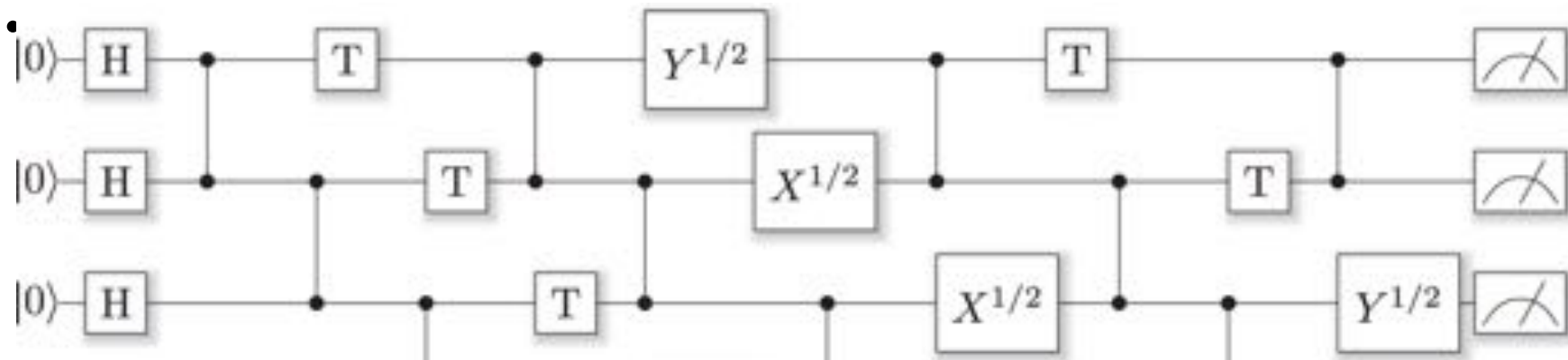
Lessons learnt concerning good software practice in HPC

Anna Brown | anna.brown@oerc.ox.ac.uk
University of Oxford; University of Southampton



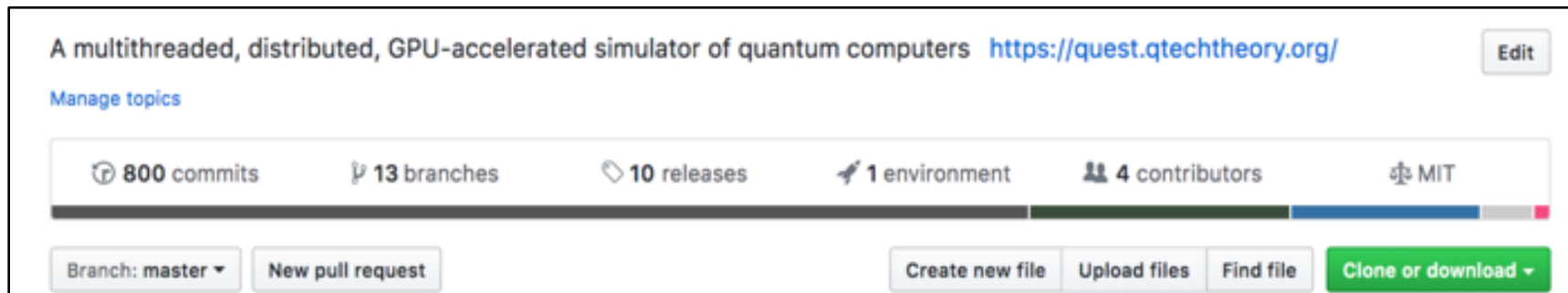
QuEST

- Simulate quantum circuits on classical HPC
- Develop quantum algorithms that are robust to noise
- Memory requirements and run time scale exponentially in number of qubits



Good software practices

- Single CPU, MPI and GPU versions available through same API
- Complexity handled by build system
- Minimum dependencies - cmake, C99, Python3
- Github, doxygen, unit testing framework, CI



Visualization for Dynamic Fans Speed Control

Meagan Mak

meagan.mak@alumni.ubc.ca | meagan@lbl.gov

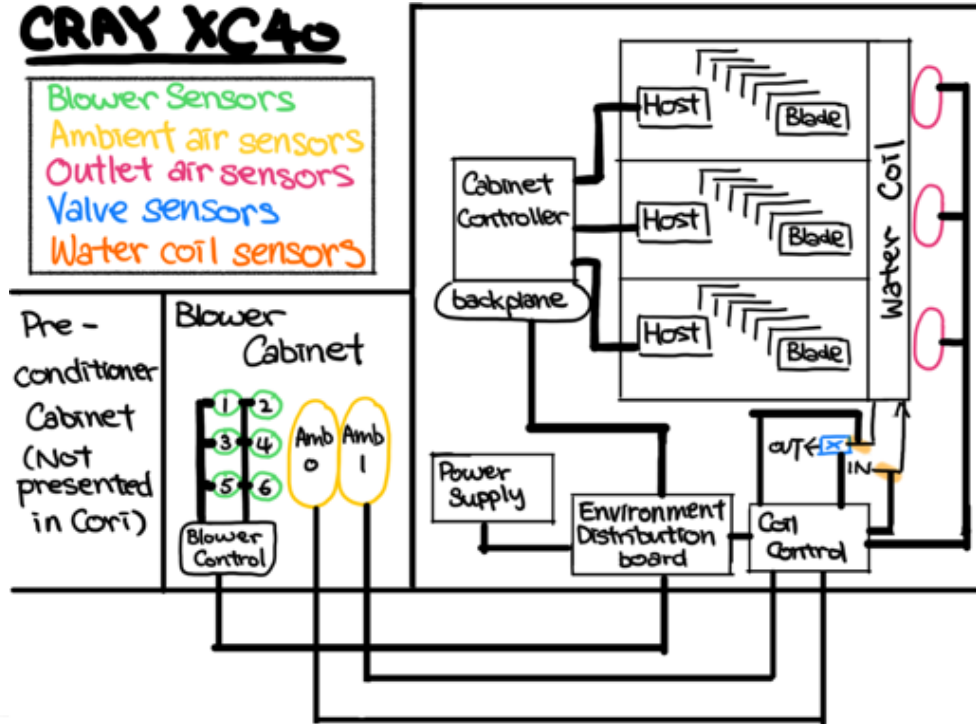
University of British Columbia | Lawrence Berkeley National Laboratory



CHANGING
THE FACE
OF HPC

Dynamic Fans Speed Control

- Put in place in April 2018
- Blower speeds are adjusted according to processor temperatures
- Reply on a variety of sensors (both internal and external)



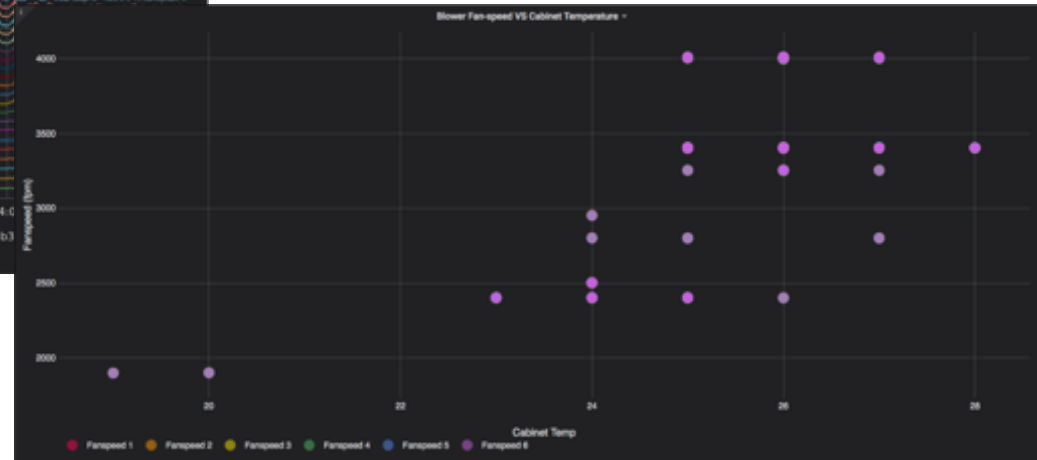
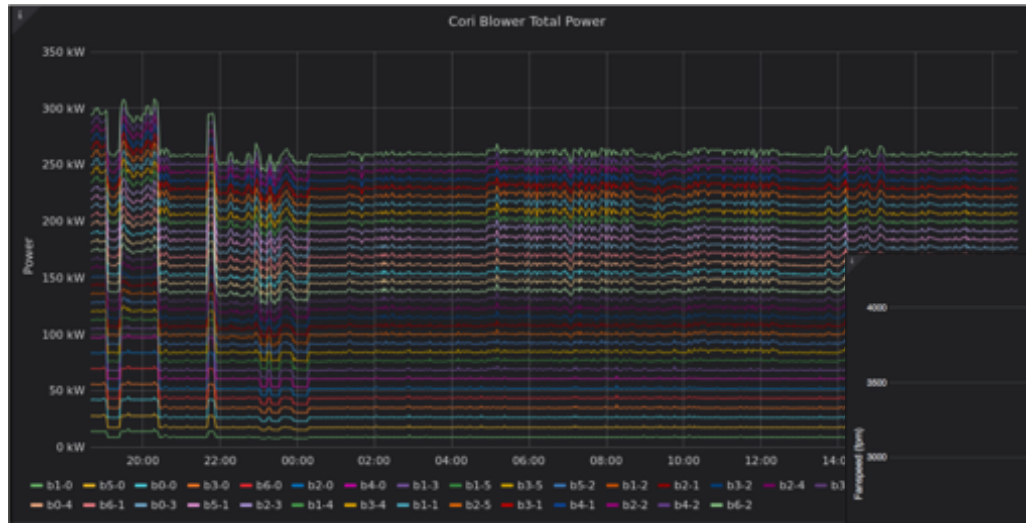
Real-time Monitoring

- Data collected from different sensors are processed through Elasticsearch and visualized with Grafana (Open-sourced)
- Provide a general picture of the conditions of the DFSC system to SREs
- Def
- pur



Analysis and Optimization of DFSC

- Used at evaluating the efficiency of the system
- Provide data on the implementation of the dynamic temperature setpoint feature, which allows the system to be further optimized





Backbone Network Design for a Sustainable Data Warehouse at NERSC

Meriam Gay Bautista | mgbautista@lbl.gov
Graduate Intern

Thomas Davis | tadavis@lbl.gov
System Architect Engineer

Operations Technology Group
NERSC
Lawrence Berkeley National Laboratory



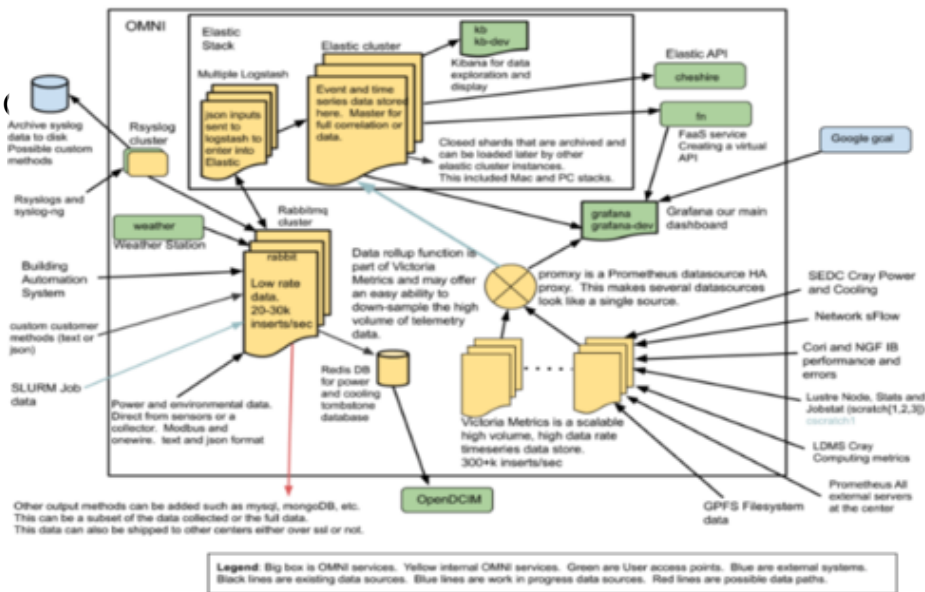
CHANGING
THE FACE
OF HPC

National Energy Research Scientific Computing Center (NERSC)

- non-classified computational facility for Department of Energy

Instrumented a facility that gathers (from heterogeneous source;

- Syslogs
- I/O Disks
- Sflow network
- Cyber logs
- Building management data (power, temperature, humidity, etc.)

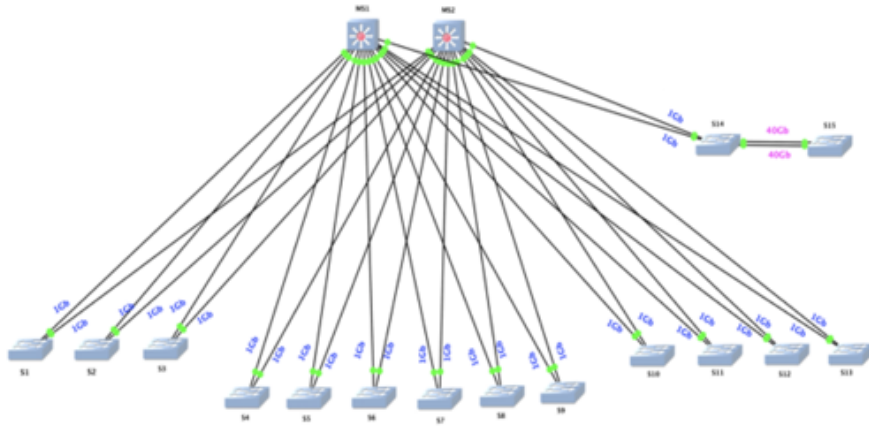


OMNI Architecture (Operations Monitoring and Notification Infrastructure)



Current OMNI network

2 level Fat-Tree based network topology

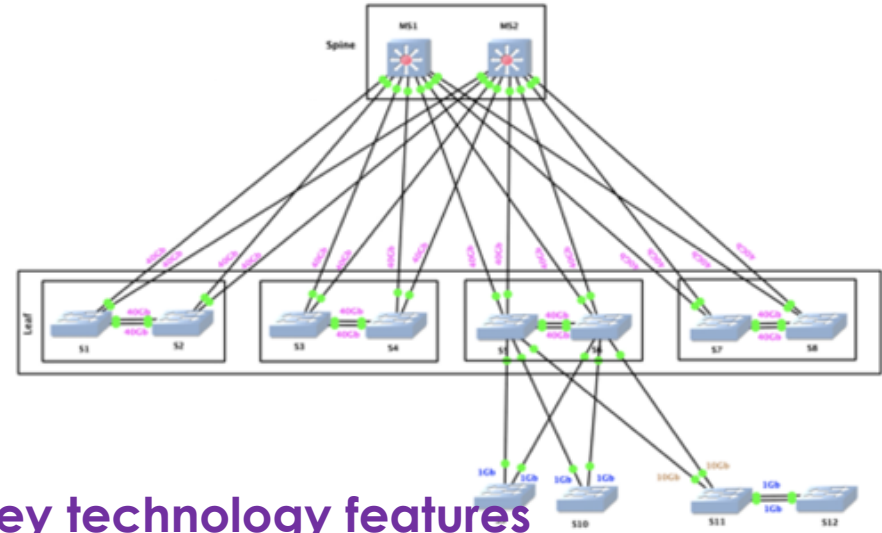


Why we need to upgrade ?

- The top of rack (TOR) switches are in End-of-life
- Cori - Compute node switch is also in End-of life
- No control plane for single management
- Leverage Ultra POE to implement new IoT technologies
- SPOF (single point of failure) configuration in network
- Reduce 3 switch OS to 1 single OS

Upgrading OMNI Network

Spine-Leaf based Network Topology



Key technology features

- Use of VXLAN for SDN based networking
- Cumulux Linux: leverage automation, EVPN, MLAG, Ansible
- Use of Border Gateway Protocol (BGP): flexible, scalable routing, improves latency, minimize downtime
- Much more compatible to Perlmutter, collect more data



Modernizing Supercomputer Monitoring via Artificial Intelligence

Elisabeth (Lissa) Moore | lissa@lanl.gov
Ultrascale Systems Research Center
Los Alamos National Laboratory





















CHANGING
THE FACE
OF HPC

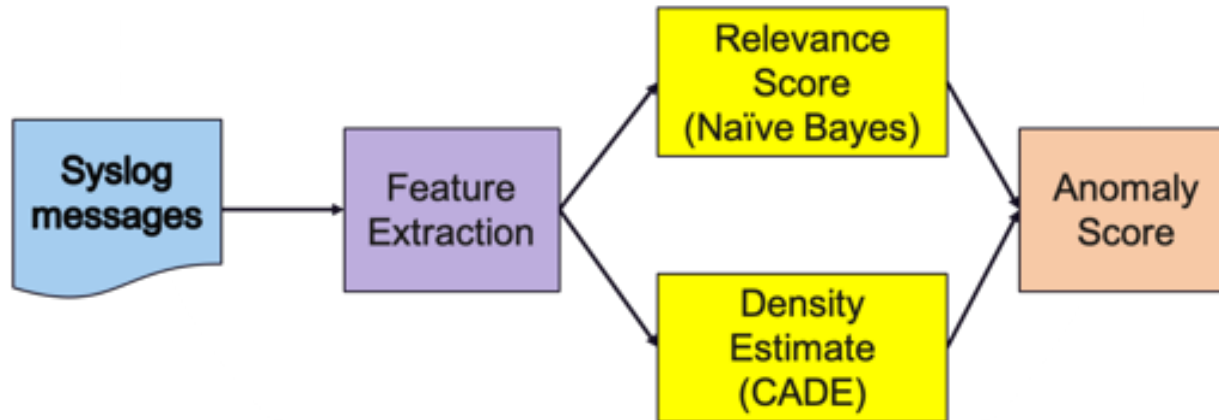
Anomaly Detection in Text Logs

Mix density estimation, user feedback, and explainable ML to finding interesting syslog messages.

```
Jun 21 14:39:54 centostest1 kernel: task_pid 1760-key switches prio exec-rolling sub-44c
c oom-kill
Jun 21 14:39:54 centostest1 kernel: bash 8686 28738.348788 215 128 28738.348788 73.6365
87 320852.286947 /
Jun 21 14:39:54 centostest1 kernel:
Jun 21 14:40:01 centostest1 CROND[8791]: (root) CMD (/usr/lib/ana/sa/1.1)
Jun 21 14:40:08 centostest1 kernel: Sysrq : Emergency Sync
Jun 21 14:40:08 centostest1 kernel: Emergency Sync complete
Jun 21 14:40:16 centostest1 kernel: Sysrq : Manual OOM execution
Jun 21 14:40:16 centostest1 kernel: events/0 invoked oom-killer: gfp_mask=0xd0, order=0, oom_adj=0, oom_score_adj=0
Jun 21 14:40:16 centostest1 kernel: events/0 cgroup=/msg_allowed#
Jun 21 14:40:16 centostest1 kernel: Pid: 7, comm: events/0 Not tainted 2.6.32-584.1.3.el6.i686 #1
Jun 21 14:40:16 centostest1 kernel: call Trace:
Jun 21 14:40:16 centostest1 kernel: [c084f8d9a] ? dump_header+0x84/0x190
Jun 21 14:40:16 centostest1 kernel: [c084f113b] ? oom_kill_process+0x68/0x283Jun 21 14:40:59 centostest1 kernel: linklog 5.0
18, log source = /proc/kmsg started
Jun 21 14:41:07 centostest1 kernel: eth2: no IPd routers present
Jun 21 14:41:13 centostest1 kernel: mdadmrun: failed to make dump initrd
Jun 21 14:41:14 centostest1 acpid: starting up
Jun 21 14:41:14 centostest1 acpid: 1 rule loaded
Jun 21 14:41:14 centostest1 acpid: waiting for events: event logging is off
Jun 21 14:41:15 centostest1 acpid: client connected from 82B9(68:68)
Jun 21 14:41:15 centostest1 acpid: 5 client rule loaded
Jun 21 14:41:16 centostest1 automount[8389]: lookup_read_master: lookup(nisplus): couldn't locate nis+ table auto.master
Jun 21 14:41:16 centostest1 sshd[8388]: Server listening on 0.0.0.0 port 22.
Jun 21 14:41:16 centostest1 sshd[8388]: Server listening on :: port 22.
Jun 21 14:41:16 centostest1 xinetd[8348]: Reading included configuration file: /etc/xinetd.d/chargen-dgram [file/etc/xinet
d.conf] [line49]
Jun 21 14:41:16 centostest1 xinetd[8348]: Reading included configuration file: /etc/xinetd.d/chargen-stream [file/etc/xine
t.d/chargen-stream] [line67]
Jun 21 14:41:16 centostest1 xinetd[8348]: Reading included configuration file: /etc/xinetd.d/daytime-dgram [file/etc/xinet
d.d/daytime-dgram] [line67]
Jun 21 14:41:16 centostest1 xinetd[8348]: Reading included configuration file: /etc/xinetd.d/daytime-stream [file/etc/xine
```

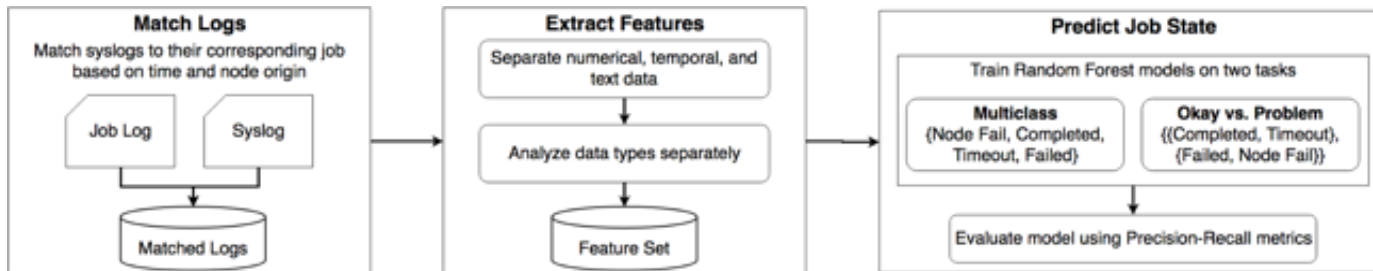


Login						
List of Entries		User Recorded Rules		User: user1	Filters: Host Filter + Alert Filter +	
	Score	Host	Time	Ident	Message	
 	1.0000	gr-fc3	7/15/2019, 7:41:01 AM	CROND	Lorem ipsum dolor sit amet, ne malis possit splendide eos, debet dolores cu nec. Et pro legimus copiosae scripserit, duo an	
 	1.0000	gr-fc3	7/15/2019, 7:41:01 AM	CROND	Lorem ipsum dolor sit amet, ne malis possit splendide eos, debet dolores cu nec. Et pro legimus copiosae scripserit, duo an	
 	1.0000	gr-fc3	7/15/2019, 7:41:01 AM	CROND	Lorem ipsum dolor sit amet, ne malis possit splendide eos, debet dolores cu nec. Et pro legimus copiosae scripserit, duo an	
 	1.0000	gr-fc3	7/15/2019, 7:41:01 AM	CROND	Lorem ipsum dolor sit amet, ne malis possit splendide eos, debet dolores cu nec. Et pro legimus copiosae scripserit, duo an	
 	1.0000	gr-fc3	7/15/2019, 7:41:01 AM	CROND	Lorem ipsum dolor sit amet, ne malis possit splendide eos, debet dolores cu nec. Et pro legimus copiosae scripserit, duo an	
 	1.0000	gr-fc3	7/15/2019, 7:41:01 AM	CROND	Lorem ipsum dolor sit amet, ne malis possit splendide eos, debet dolores cu nec. Et pro legimus copiosae scripserit, duo an	
 	1.0000	gr-fc3	7/15/2019, 7:41:01 AM	CROND	Lorem ipsum dolor sit amet, ne malis possit splendide eos, debet dolores cu nec. Et pro legimus copiosae scripserit, duo an	
 	1.0000	gr-fc3	7/15/2019, 7:41:01 AM	CROND	Lorem ipsum dolor sit amet, ne malis possit splendide eos, debet dolores cu nec. Et pro legimus copiosae scripserit, duo an	
 	1.0000	gr-fc3	7/15/2019, 7:41:01 AM	CROND	Lorem ipsum dolor sit amet, ne malis possit splendide eos, debet dolores cu nec. Et pro legimus copiosae scripserit, duo an	

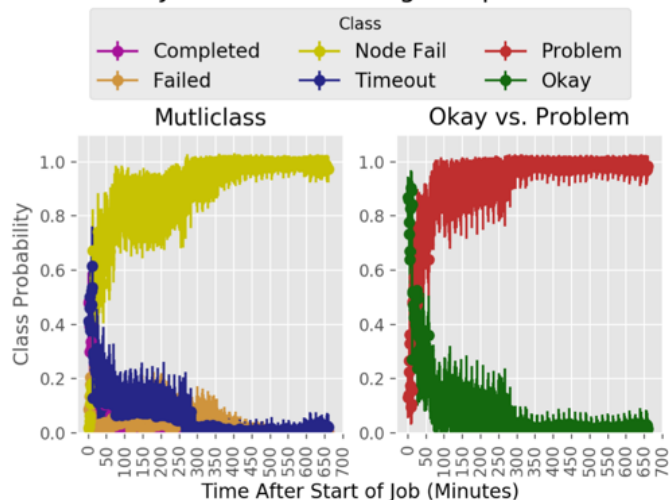


Job Outcome Prediction

Extract features from syslog messages for early detection of failing, timeout, and successful jobs.



Class Probabilities over Time for
Node Fail Job wf-404963: Tag Temporal Numerical

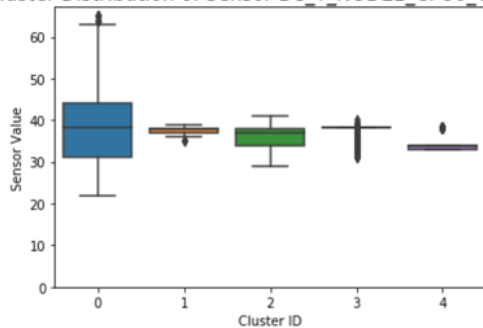


For more detail, also come to
DAAC Workshop on Friday!

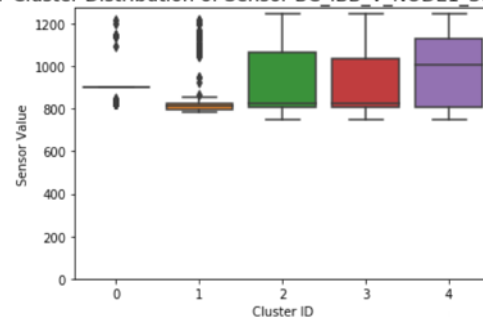
Telemetry Analysis

Characterize telemetry data from HPC systems to detect signals correlated with node failures.

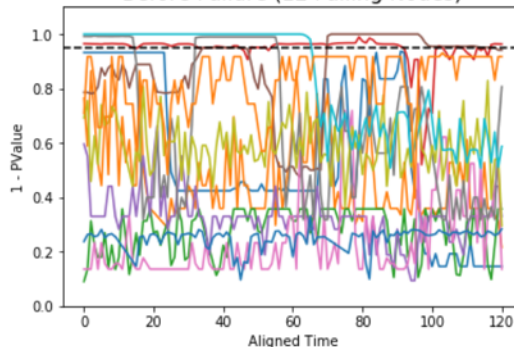
Per-Cluster Distribution of Sensor BC_T_NODE1_CPU0_CH0_DIMM0



Per-Cluster Distribution of Sensor BC_IBB_V_NODE1_S0_VCC_OUT



1 - PValue for Sensor BC_T_NODE1_CPU0_TEMP
Before Failure (12 Failing Nodes)



Workflow Visualization for Maintaining NERSC Data Center

Jameelah N. Mercer

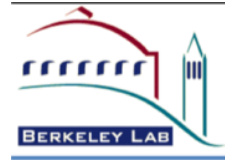
JMercer@lbl.gov

Lawrence Berkeley National Laboratory



CHANGING
THE FACE
OF HPC

Significance of Efficient Data Centers



Introduction:

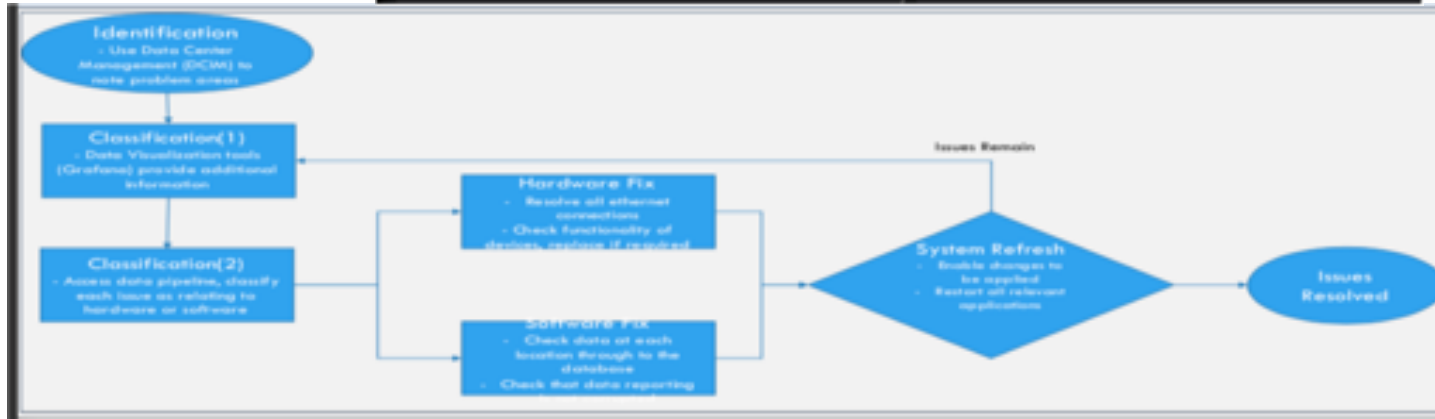
This study has two main goals:

1. First, to troubleshoot NERSC supercomputing sensor functionality
 - Analyzing and debugging the temperature, humidity, and power sensor data to ensure complete workability of the NERSC supercomputing environment.
2. To create an instance go OMNI infrastructure using Kubernetes and Docker so that we can run a data set through the new version of Elastic Stack.
 - Specifically by updating the Elastic Stack by implementing K3's/Rancher and installing an ES-Operator.



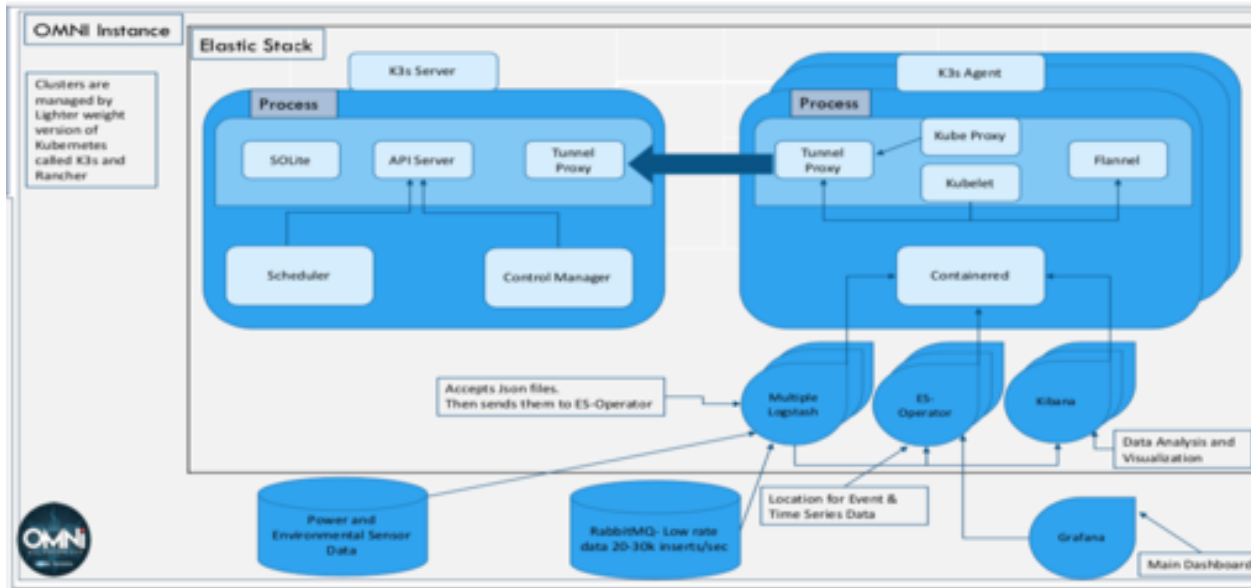
Daemon Environment

1. Consistent troubleshooting of environmental sensors.
2. Periodic analysis of PDU/Power consumption.



Operation Monitoring and Notification Infrastructure

1. OMNI is built using open source technologies.
2. OMNI contains over two years of operational data, accumulating over 125 of data.
3. An instance of ONMI is created to limit unforeseen problems and to increase reliability.



WHPC Fellows



CHANGING
THE FACE
OF HPC

Scalable Assembly of Large Genomes

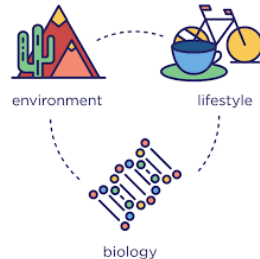
Priyanka Ghosh | Priyanka.ghosh@pnnl.gov
Pacific Northwest National Laboratory



CHANGING
THE FACE
OF HPC

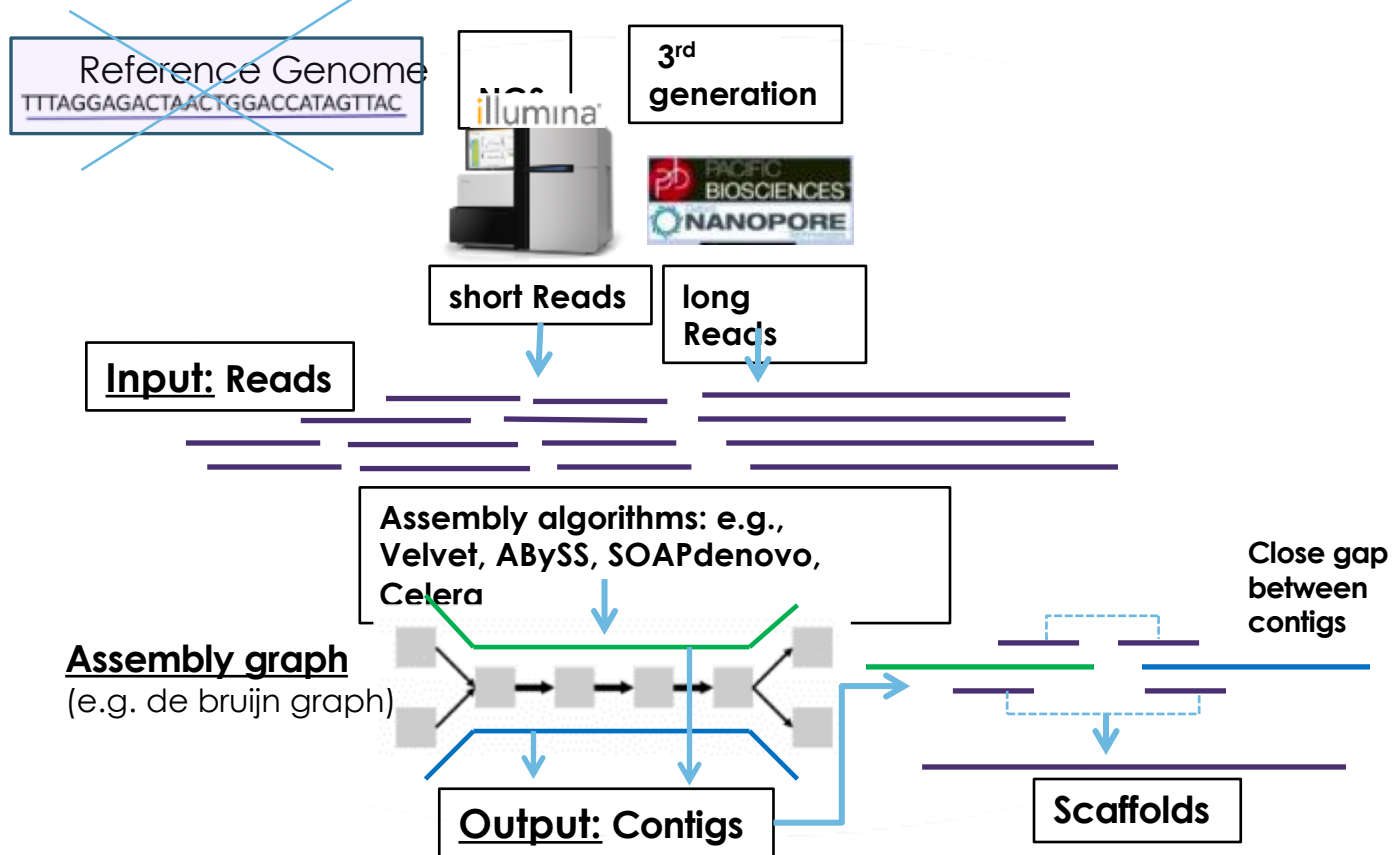
Motivation

- Precision Medicine Initiative – research effort for disease prevention and treatment
 - ❑ Take into account individual differences in people's genes, environments, and lifestyles
 - ❑ Study and analyze genetic variants/mutations in diseased tissues (such as tumors) - facilitate development of targeted therapeutics
- Genome and Metagenome assembly recognized as one of the key applications in the DOE Exascale Project
 - ❑ Develop highly scalable algorithms/software to overcome high computational demands of assembling millions of (meta)genomes
 - ❑ Reduce assembly time by orders of magnitude and make feasible the assembly of larger complex genomes



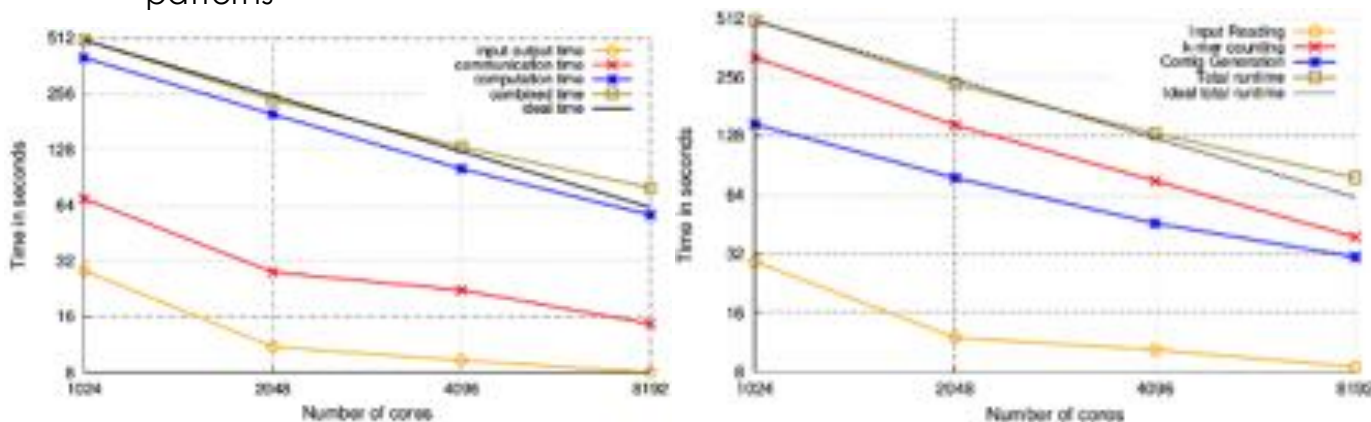
De Novo Genome Assembly: Problem Statement

Goal: Assemble the DNA sequence of an unknown target genome from numerous fragments (or 'reads') obtained from it



Distributed-memory approach (*PaKman*)

- Scalability Challenge: Typical read dataset comprises of billions of reads
 - Several hundred billions of vertices in the graph
 - Computationally demanding with respect to memory and time
- PaKman*: scalable algorithm tackling assemblies of large genomes at extreme scale
 - novel distributed-memory graph data-structure (PaK-Graph) that enables minimal communication during contig enumeration
 - novel contig generation algorithm with simplified I/O and communication patterns



PaKman assembles a complete set of contigs for full human genome in 78.4 secs on 8k cores

Technical Support in Configuring HPC Systems

Raksha Roy | raksha.roy@icimod.org
ICIMOD, Nepal



Objectives

- Get proper training on configuration and use of HPC Systems
- Install and configure scientific applications required for the Supercomputing facility
- High Performance Computing Benchmarks
- Train Students and Scientists on use of HPC Systems

Impact

Installation and configuration of OpenHPC, Lustre

Accomplishments

- A full fledged HPC Production System
- Introduction on HPC, Parallel Programming Techniques, Public Awareness on SuperComputing

Enabling HPC for neuroimaging science

- *CT template creation using nonlinear
image registration for TBI analysis*

Zhe Bai | zhebai@lbl.gov

Computational Research Division

Lawrence Berkeley National Laboratory



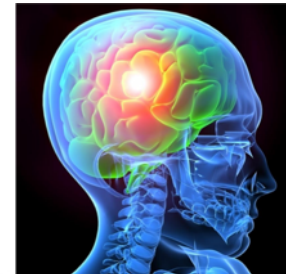
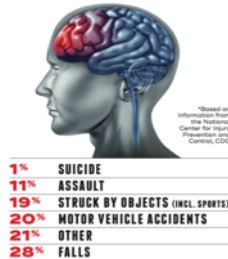
CHANGING
THE FACE
OF HPC

Background: Traumatic Brain Injury (TBI)

- Number of TBI cases occurring in the U.S. every year: ~1.7 million.
- Influences in life: physical, psychological, social, and spiritual.
- Complexities: *multi-modal data, large volume image, statistical varieties.*

**data science + computer vision +
high-performance computing**

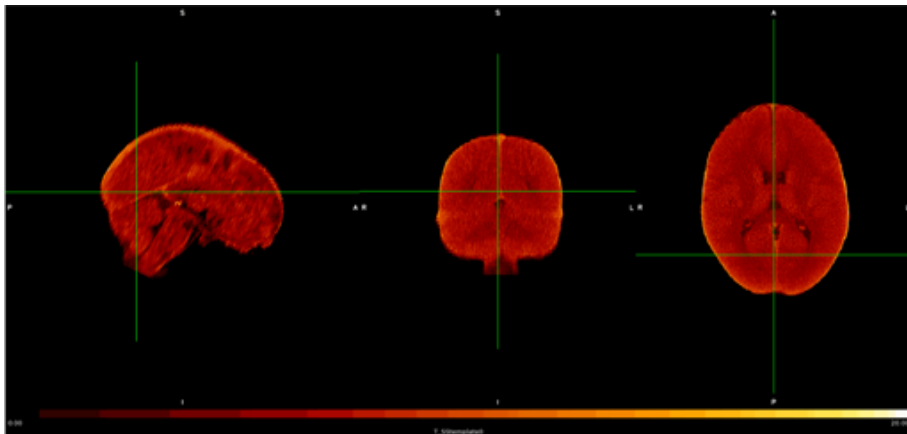
MAJOR CAUSES OF
TRAUMATIC BRAIN INJURIES*



CT template creation

- Subgroup study based on physiological features.
- Iterative algorithm: rigid + affine + nonlinear transformation.
- High performance: image similarities are optimized in parallel.

Created template based on 12 subjects (shown in MNI space)



Computational time vs. # cores

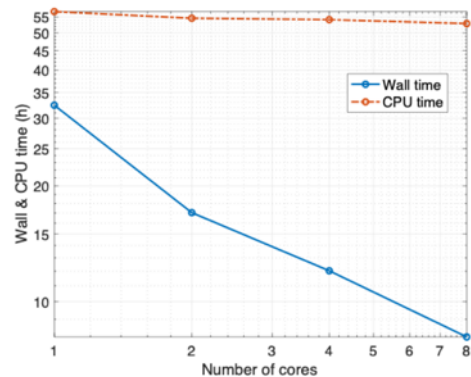
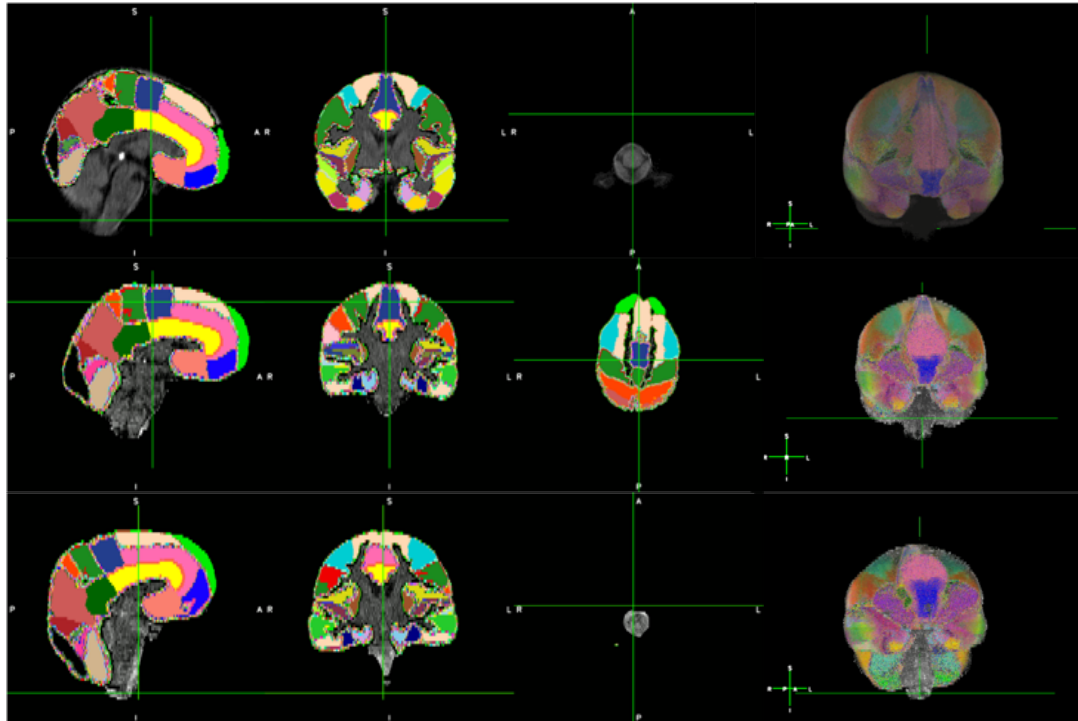


Image segmentation for TBI patients

- Segmented template & patients' CT scans.
- Automatic parcellation: 48 structural areas < 1 min.



Atlas: HarvardOxford-Cortical
Bottom: Group skull-stripped CT
template

Patient GCS: 15
Bottom: Patient's skull-stripped CT

Patient GCS: 4
Bottom: Patient's skull-stripped CT