# Acceptance Testing the Chicoma HPE Cray EX Supercomputer

Presented by:  Jennifer K. Green, Scientist

Prepared by:   K. Everson, P. Ferrell, J. Green, F. Lapid, D. Magee, J. Ogas, C. Seamons, N. Sly
**Programming & Runtime Environments Team
High Performance Computing Environments
Los Alamos National Laboratory**

May 3, 2021
Presented to the CUG 2021 Virtual Conference
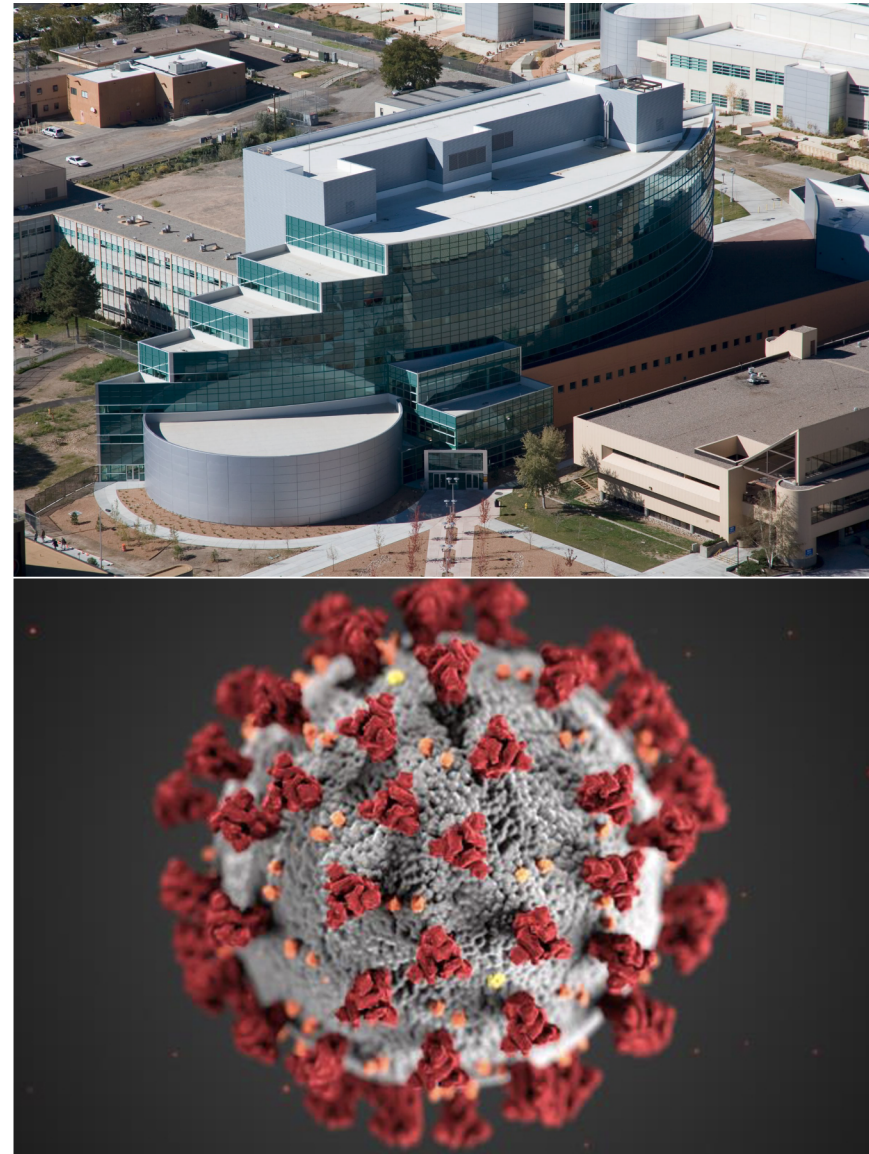
LA-UR-21-24209

Managed by Triad National Security, LLC., for the U.S. Department of Energy's NNSA.

1

# Background

- Los Alamos National Laboratory (LANL) remains at the forefront of addressing global crises using state-of-the-art computational resources to accelerate scientific innovation and discovery
- LANL is supplying high-performance computing (HPC) resources to contribute to the recovery from the impacts of SARS-CoV-2 (Coronavirus Pandemic)
- Chicoma is an HPE Cray EX Supercomputer recently installed at LANL to specifically serve as a platform to supply molecular dynamics simulation computing cycles for epidemiological modeling, bioinformatics, and chromosome/RNA simulations as part of the 2020 Coronavirus Aid, Relief, and Economic Security (CARES) Act



https://www.lanl.gov/discover/news-release-archive/2020/October/1020-hpc-to-fight-against-covid19.php

## Overview of Presentation

- Chicoma HPE Cray EX System Description
- Acceptance Testing Approach
- Testing Tools Description
- Test Suite Contents
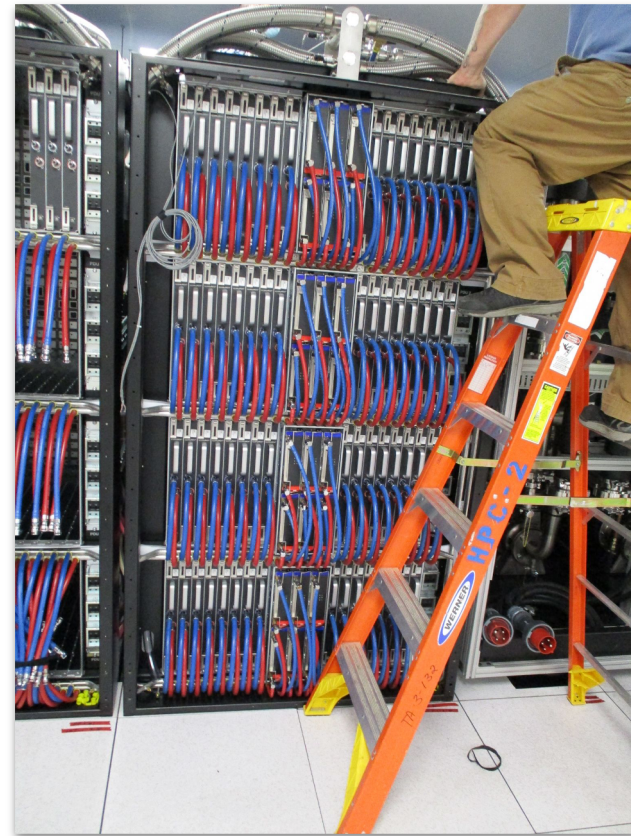- Results
- Conclusions
- Future Work

# The Chicoma Supercomputer

# System Details

- Chicoma is an early deployment of HPE Cray EX
- Has a large-scale system architecture
- Brand new **Shasta** system software stack
- Features direct-to-chip liquid cooling
- HPE Slingshot interconnect
- AMD EPYC 7H12 processor
- In total
  - more than 73,000 cores
  - 300 TB of system memory

# Testing Motivation and Context

- Preparing an Acceptance Testing Plan involves:
  - Develop an understanding of the intended workloads for the system
  - Identify the specifications and expectations of performance and reliability for supporting the science
  - Develop a testing plan to ensure that the final installation has met those requirements
- Chicoma is among the earliest installations of the HPE Cray EX system, running the Shasta Architecture
  - Integration and testing activities continue to date
  - Continued testing ensures that the system can support a synthetic workload representative of the science for which it is intended

# Acceptance Testing Approach

- Integration Testing
  - New architecture and progression of system software while developing the test suite required integration testing

- Functionality testing was accomplished during this phase
  - Evaluating the readiness of the Cray Programming Environment (CPE) to support workloads
  - Testing viability of containerized FEs to host the harness
  - Unprivileged container testing using Charliecloud
  - Usability of the supplied GROMACS application with COVID-19 study .tpr file
  - Scaling tests for MPI applications
  - Setting up Pavilion configurations and developing Acceptance Test Suite

- Seven Weeks from Plan Draft to Running Acceptance Tests
  - Drafted Testing Plan - *July 14, 2020*
  - Implemented Plan - *September 3, 2020*

**Los Alamos**
NATIONAL LABORATORY

# Test Suite

DGEMM - single node performance

ExaMiniMD - proxy MD application

GROMACS Covid-19 problem - real world application

HPCG - full system benchmark

HPL (8 nodes, full system, single node) - various sized benchmark

LULESH - proxy application

MILC7 - Mini app QCD problem

QuickSilver - CTS Mini App

Stream - Memory benchmark

SystemConfidence - network latency benchmark

VPIC - Kinetic plasma modeling simulation

Intel MPI Benchmarks - MPI-1 benchmark suite

## Pavilion2 HPC Test Harness

Pavilion is a Python 3 (3.5+) based framework for running and analyzing tests targeting HPC systems

- Maintained by LANL's High Performance Computing Environments Group and is open-sourced for community contributions & usage
- Supplies a framework for creating sophisticated YAML configurations to automate the workflow of running jobs on HPC systems
- Plugin components include those for gathering system data, adding additional schedulers, parsing test results, and more
- Pavilion outputs results of every test in a json log file, which then is able to be processed by a number of analysis utilities
- https://github.com/hpc/pavilion2/

# Splunk - Data Visualization Tool for Test Results

- Splunk is a flexible data analytics platform that enables searching, monitoring and analyzing machine-generated data
- Captures and correlates data in searchable database for supporting dashboard and graphical displays of data
- The monitoring infrastructure is supported by a distributed Splunk instance on every network, gathering large temporal data for analysis
- LANL's Splunk distributed monitoring infrastructure indexed Chicoma test result logs to enable automated results analysis and correlation of system events to test underperformance/failures

# Acceptance Testing Results Summary

**PASSING TESTS COUNT guaje 1599228000 - 1599400800**

## 13,998

TEST SUCCESSES

**FAILING TESTS COUNT guaje 1599228000 - 1599400800**

## 61

TEST FAILURES

**SUCCESS PERCENT guaje**

## 99.45%

Success Percentage

*{ 1 of 61 failures was legitimate }*

*{ HPCG test mis-configuration led to 60 failures }*

*{ 225+ TFLOP/s - 256 nodes }*

Pie chart labels: vpic, sysconfidence, stream, quicksilver, milc7, lulesh_single_node, imb, hpl-full, hpl, hpcg-guaje, hello_mpi, gromacs, examini-guaje, dgemm

HPL GFLOP/s Full System

**HPL Full System Performance guaje**

Bar chart values: 129,000 (MAX), 113,700 (MIN), 121,350 (AVG), 460800; 224,200 (MAX), 224,200 (MIN), 224,200 (AVG), 614400

GFLOP/s axis: 300,000 / 200,000 / 100,000

valn

Legend: MAX GFLOP/s, MIN GFLOP/s, AVG GFLOP/s

**Los Alamos** NATIONAL LABORATORY

11

# Load Testing



## Table Count Tests 1599400800 - 1599699600

| | testname ⇕ | count ⇕ |
|---|---|---|
| 1 | dgemm | 427 |
| 2 | examini-guaje | 8 |
| 3 | gromacs | 8 |
| 4 | hello_mpi | 8 |
| 5 | hpcg-guaje | 23 |
| 6 | hpl | 19 |
| 7 | hpl-8 | 248 |
| 8 | hpl-full | 11 |
| 9 | imb | 18 |
| 10 | lulesh_single_node | 7 |
| 11 | milc7 | 6 |
| 12 | quicksilver | 18 |
| 13 | stream | 18 |
| 14 | sysconfidence | 13 |
| 15 | vpic | 41 |

| name ⇕ |
|---|
| hpcg-guaje.base.128-4-threads-true-sockets-128 |
| hpcg-guaje.base.64-4-threads-true-sockets-128 |
| hpl-8.rome.HPL-ROME |
| hpl-full.rome.HPL-ROME |

*Only failures were due to oversubscribing or running too large a problem*

- *59 Hours* with no hardware failures or system related test failures

**25,132**
TEST SUCCESSES

**99.10%**
Success Percentage

**48**
TEST FAILURES

# LULESH

# HPL-8 Node Benchmark

HPL-8 Node

N Val (Size)

| 80000 | ▼ |

PPN (Processes Per Node)

| 64 | ▼ |

threads per node

| 8 | ▼ |

*62 GFLOP/s diff between fastest/slowest*



8-Node Linpack guaje



8-Node Linpack guaje

*Problem Specs Selected to Match Integration Team's 8-Node HPL Size and Runtime Configs*
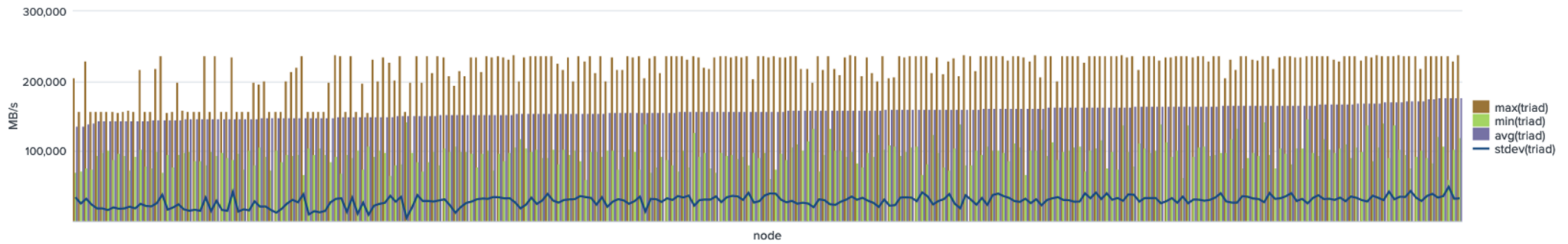
# HPL Single Node Benchmarks

*{ Some interesting HPL spikes in std dev }*
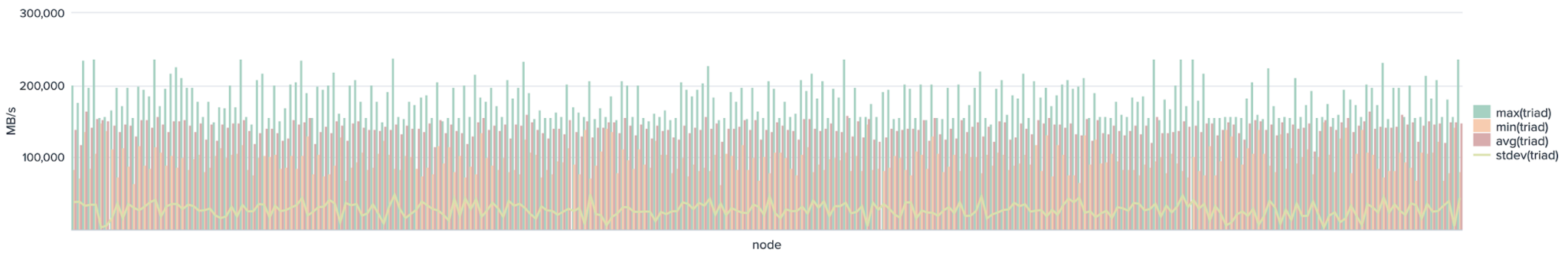


HPL Single Node Performance guaje

# Stream Memory Benchmark

**Stream Single Node guaje triad Rate**



**Stream Single Node guaje triad Rate**

# Full System
## High Performance Conjugate Gradient (HPCG)



17

# QuickSilver Proxy Application



Quicksilver Figure of Merit - guaje | sort stdev(fom)

## Quicksilver

- MPI/MPI-OMP proxy application
- Included in the CTS Benchmarks Suite
- Solves a simplified dynamic Monte Carlo particle transport problem.
- Its performance is bound by poor vectorization potential, latency bound table look-ups and a heavily branching or divergent code path.

# VPIC LPI 3D Deck (Lyin-Sequoia)

## Vector Particle-In-Cell (VPIC)

- Simulation code for modeling kinetic plasmas on one, two, or three dimensions.
- It employs a second-order, explicit, leapfrog algorithm to update charged particle positions and velocities in order to solve relativistic kinetic equations.
- The input deck, a modified version of lyin_sequoia problem conducted, exercises the problem that Lawrence Livermore National Laboratory used to evaluate their Sequoia system's potential to model the interaction of realistic fast-ignition-scale lasers with dense plasmas in three dimensions.

### VPIC Lyin-Sequoia
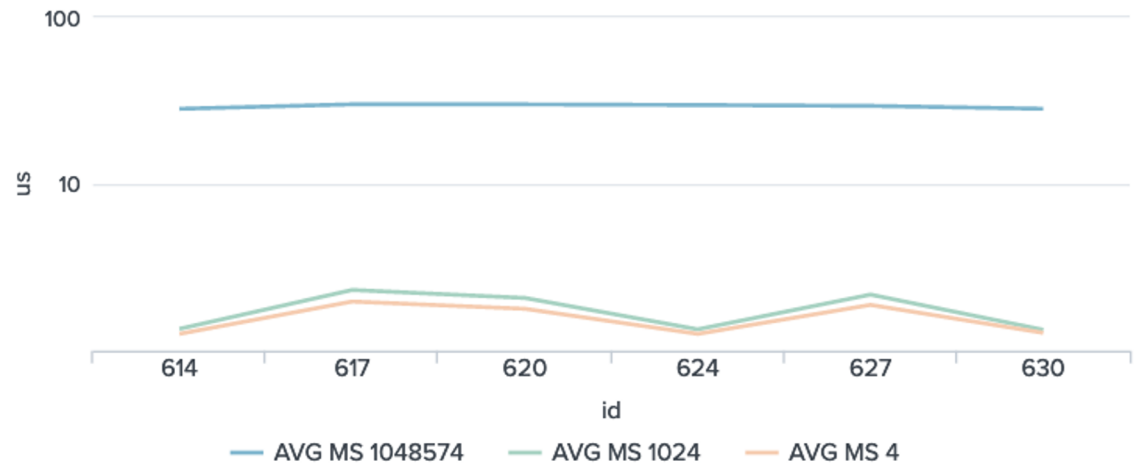
**VPIC LPI 3D Deck: time to completion with 256 nodes.**
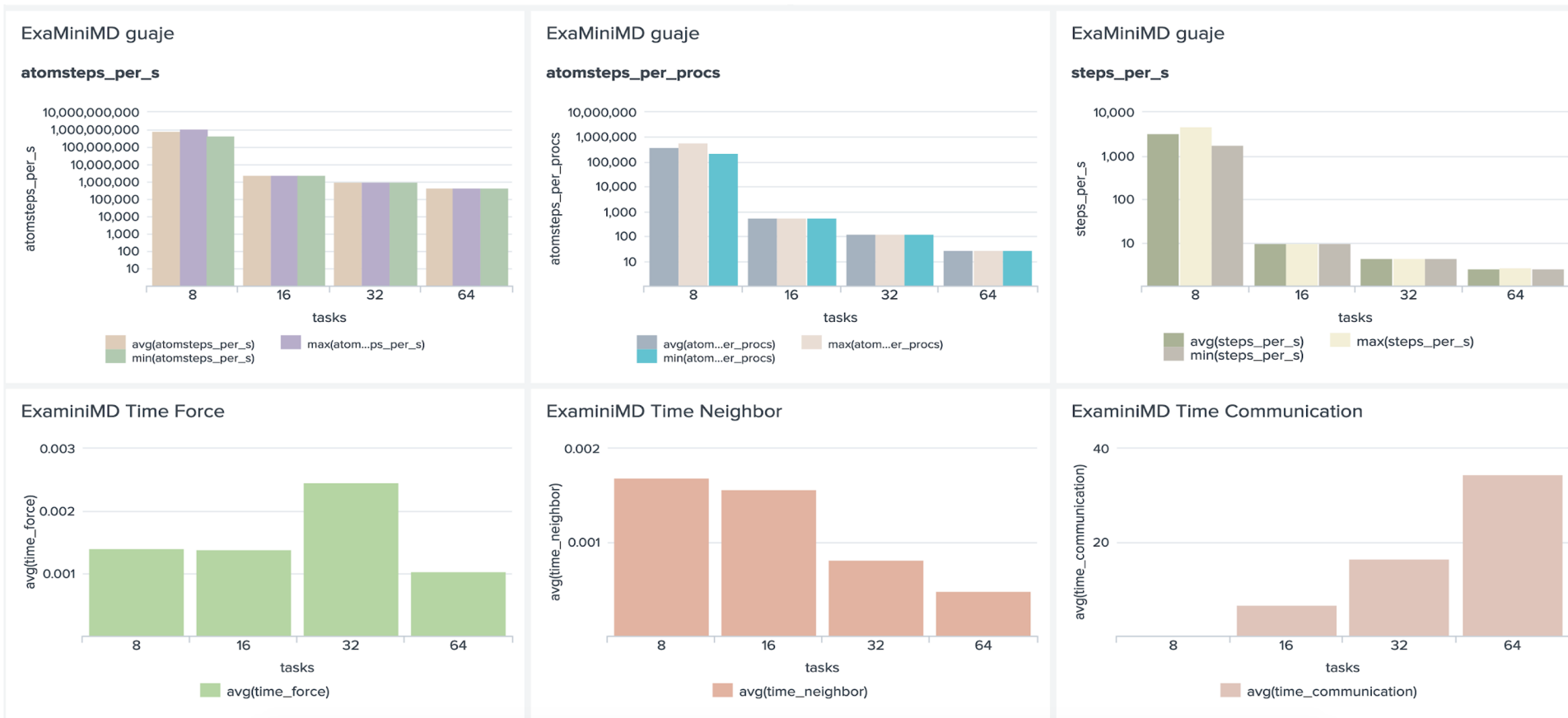
# Intel MPI Benchmarks

## Intel MPI Benchmarks



IMB Alltoall guaje gcc_cray



IMB PingPong guaje gcc_cray

# ExaMiniMD

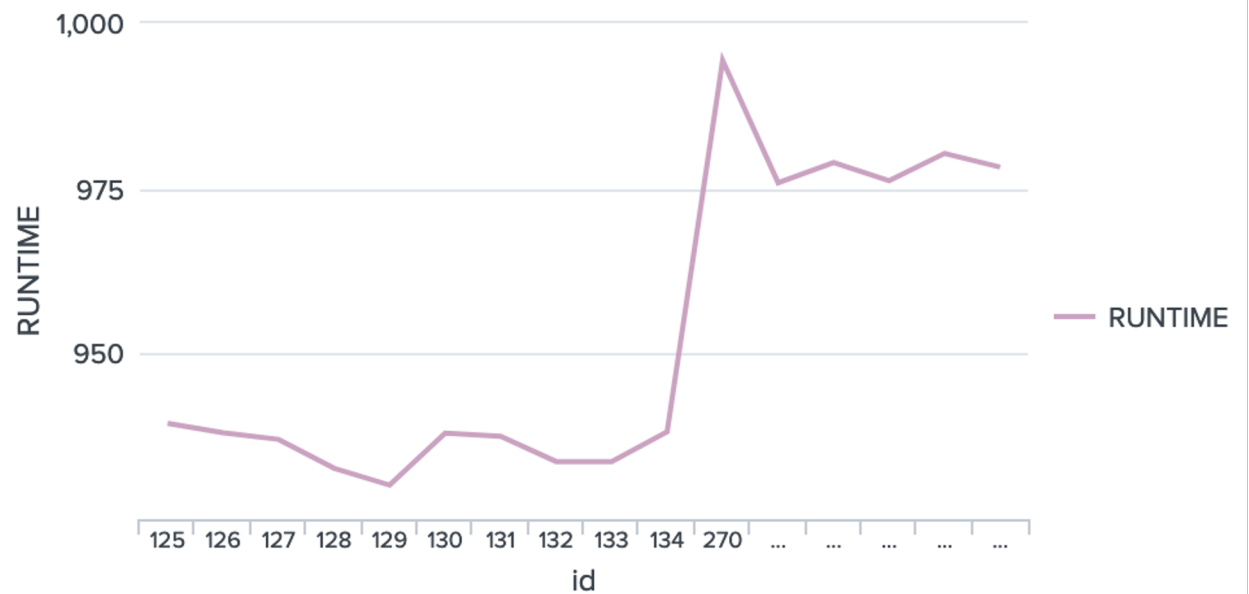## ECP Proxy Application (Kokkos) Simplified MD Simulation

# MILC

The MILC Code is a body of high performance research software written in C for doing SU(3) lattice gauge theory on high performance computers
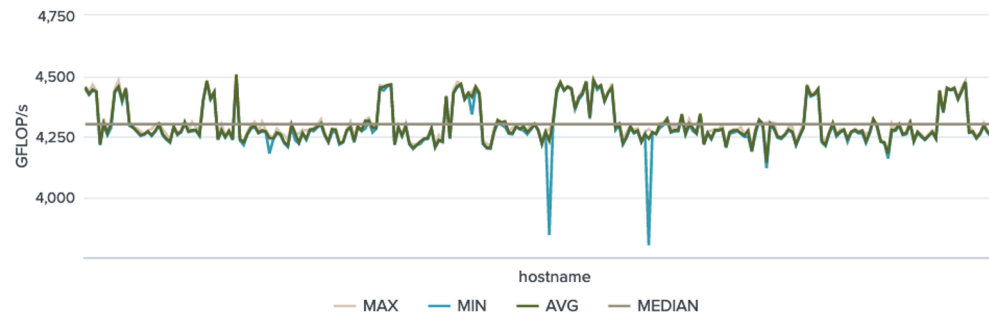
MILC7

**MILC7 Runtime (LESS IS BETTER)**
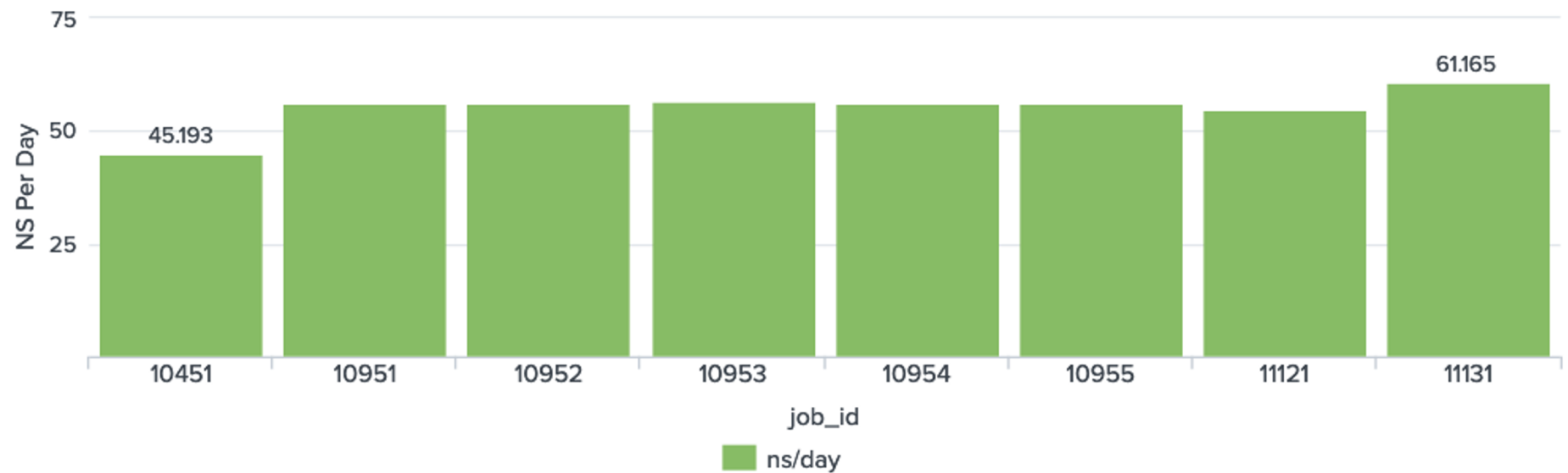
# DGEMM



## DGEMM

- DGEMM was built with the Cray Libsci package and we weren't able to get the job to spread to both sockets of the nodes.
- We're working on a DGEMM built with OpenBLAS to see if we can overcome this issue, for improved performance.

# GROMACS - COVID-19 Simulation

One major motivation for this effort was to ensure MD problems would run on the system. The repeatability of successful runs with a real GROMACS simulation proves it.

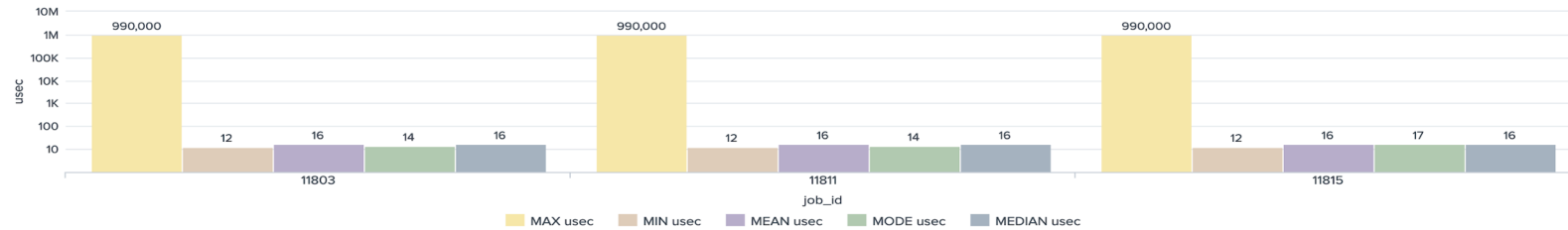GROMACS Simulation

**COVID-19 MD Sim**

# System Confidence Network Latency Test

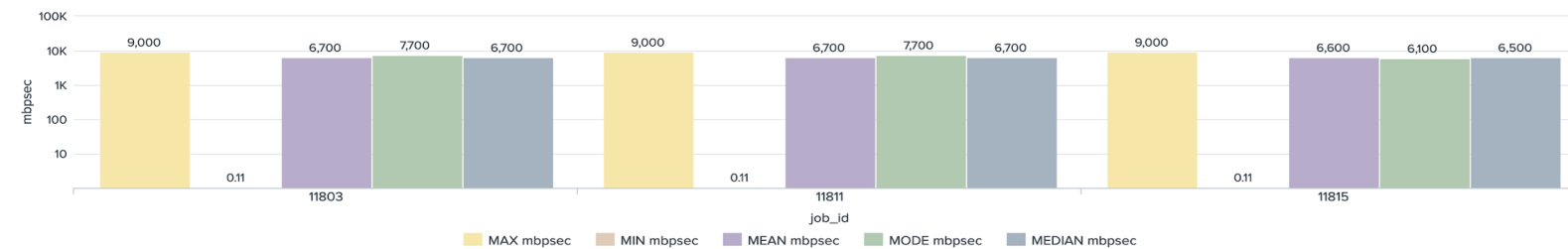Sysconfidence guaje Pair-wise Latency Test

Unit
[ U secs ▾ ] [ × ]

Buffer Length
[ 52376 ▾ ] [ × ]

**Sysconfidence guaje usec 52376**



Unit
[ MB/sec ▾ ]

Buffer Length
[ 52376 ▾ ] [ × ]

**Sysconfidence guaje mbpsec 52376**



## System Confidence

- Captures latency (usecs) and rate (MB/s)
- Consistent ~1 sec latency max measured for all buffer lengths
- Could be a test anomaly

**Los Alamos**
NATIONAL LABORATORY

# Conclusions

- Chicoma was "accepted" by LANL after demonstrating that it was capable of sustaining a workload and measuring acceptable performance
- Chicoma was constructed to serve as the IC Program's Platform for supporting COVID-19 studies
- Chicoma is currently undergoing an upgrade to Shasta v1.4
- Tests will be repeated after that upgrade to ensure continued stability and performance of the machine
- Chicoma is currently running in a pre-production mode at LANL while efforts to fully integrate into production environments are underway
- Users are using the pre-production system to conduct their research for the IC Program

# Future Work

- Pavilion tests are being developed to target unprivileged containerized runtimes on HPC resources at LANL
- This effort proved that Pavilion was able to satisfy the requirements to conduct Acceptance testing of future procurements
- Test implementations under Pavilion for Chicoma acceptance will be re-run during the course of transitioning the Chicoma system to full-production
- Results comparison of the initial baselined results will be rerun with upgrades, including the upgrade to Shasta v1.4, and conducted over the life-cycle of the machine to
    - identify any performance degradation
    - support optimization of configurations
    - feed into future procurements

PAVILION
HPC Test Harness

SOURCE CODE

READ THE DOCS

Los Alamos
NATIONAL LABORATORY