# HIGH PERFORMANCE COMPUTING

## 2023 ANNUAL REPORT
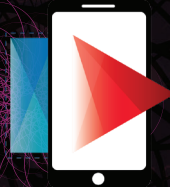
# TACKLING THE NATION'S TOUGHEST CHALLENGES

## HIGH PERFORMANCE COMPUTING IN 2023

**Sandia National Laboratories**

**BUT WAIT, THERE'S MORE**
Watch innovation come to life using **SNLSimMagic**, an augmented reality application developed at Sandia National Laboratories.

**Download SNLSimMagic on your iOS device.**

Use your phone to scan images with an icon to watch content come to life.

# TABLE OF CONTENTS

# Director's message



LABS DIRECTOR,
*Dr. James S. Peery*

High performance computing (HPC) and the brilliant researchers who use it have a long history at Sandia National Laboratories. The HPC Annual Report highlights the exceptional research being done at the Labs in support of our mission to innovate and discover new technologies to strengthen the nation's technological superiority.

In this year's report we showcase how Sandia uses HPC to address the nation's most interesting and pressing challenges, including improving networks supporting exascale computing and working to stretch the thermodynamic limits of HPC efficiency.

You'll read how our teams are studying machine learning used in many ways, from investigating its credibility, to predicting how trillions of different ionic liquids will behave without having to test them all, to investigating if we can design structural joints from patterned features without relying on human intuition. You will learn how HPC-driven digital engineering supports unique mission needs and how Sandia engineers are using computerized tomography, or CT scans and computing power to create and test digital twins of explosive tools.

Our researchers tap HPC for studies into Arctic coastal erosion and permafrost demise and the impact on national security posed by climate change and threats to critical infrastructure in the Arctic. These topics and the other innovative research explored in these pages demonstrate Sandia's excellence in high performance computing.

With this 2023 HPC Report, we pay tribute to a lost colleague and leader, Scott Collis, who was integral to Sandia's HPC legacy. Scott passed away in September 2022 at the age of 55 from a rare and aggressive cancer. Though his life was cut short, his impact on Sandia and the nation will endure.

I have no doubt that Scott would be proud of the work portrayed in this report. I hope you are fascinated and enriched by its contents and share my excitement about HPC at Sandia as we aim for excellence while addressing the greatest challenges facing our nation and world today.

# Tribute to Scott Collis

## HONORING A LOST COLLEAGUE AND LEADER, SCOTT COLLIS

CONTRIBUTING WRITER | *Sheina MacCormic*

Scott came to Sandia National Laboratories after an Assistant Professorship at Rice University. Though he started his studies in aerospace engineering, he finally settled on and obtained his PhD in mechanical engineering and followed a path into computational research. He spent nearly 20 years of distinguished service at Sandia in Albuquerque, New Mexico. A conference room in Sandia's Computer Science Research Institute, located in the Sandia Science and Technology Park was dedicated to Scott in March 2023 in honor of his work and service at Sandia.

Scott's ability to lead was quickly recognized by leadership and he rose through the management ranks to Director of Sandia's Center for Computing Research. He also served as Director of Sandia's DOE/NNSA Advanced Simulation and Computing program and as Program Executive for Sandia's DOE/Office of Science Advanced Scientific Computing Research program.

"Staying on the forefront of research is critical for a national laboratory," Scott asserted in a rare video meeting chat with friends and colleagues in the months before his death. "Historically, it is really important to have an institution that helps build bridges for the laboratory, to academia and to industry. For us, this has been the Computer Science Research Institute that we use to host visitors and foster collaboration."

Scott fostered a highly collaborative environment and helped create a wide range of programs where computing was key, including computing for energy and climate applications; quantum, neural-inspired, and other non-conventional computing concepts; and exascale supercomputing architectures and implementations. Under his directorship he saw the development of two new applications critical to the weapons program that will take advantage of El Capitan, the exascale computing platform for NNSA that will be sited at Lawrence Livermore National Laboratory.

Scott maintained an attitude of people first. He was a visionary and a leader in forming teams and partnerships that made it possible for those he worked with to accomplish extraordinary work, such as that which is represented in this report. He is greatly missed, but his proud legacy of excellence continues.

*Scott passed away in September 2022 at the age of 55 from a rare and aggressive cancer. Though his life was cut short, the impact of his life to Sandia and to the nation is far-reaching.*

# Frontiers in artificial intelligence and machine learning

CONTRIBUTING WRITER | *Susan Jean-Pierre*

FIGURE 1. *Artificial intelligence, machine learning and deep learning.*

## WHAT IS MACHINE LEARNING?

Humans and animals are capable of perceiving, synthesizing what they perceive and applying that information to their lives. Artificial intellig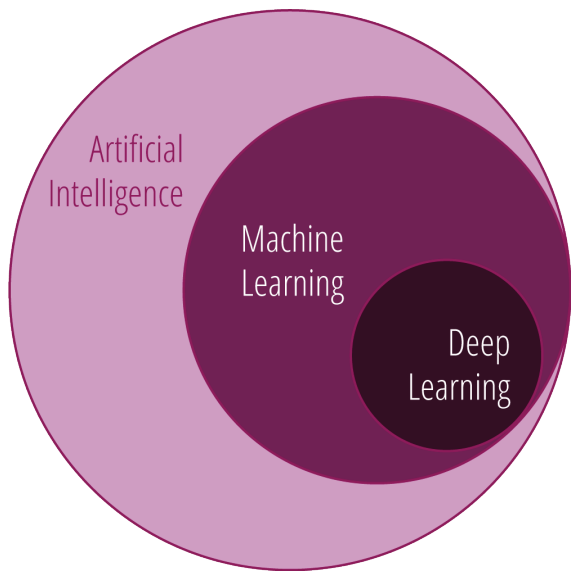ence, or AI, is a branch of science with the goal of developing machines that can function in the same way. Humans and animals use their five senses to gather information, their brains synthesize those details and from there, the information is shared with others. However, machines are more limited and depend on information that is shared with them. If you use the Google search function or Amazon's Alexa, you've used AI applications designed to help you find what you need.

Machine learning is a subset of AI focused on using algorithms to build models from training data. These models then can use this data to make decisions or predictions even though they are not programmed to do so. Basically, the model enables the machine to program itself to solve a specific set of problems. If you use an email spam filter, then you've used machine learning.

Artificial neural networks are a subset of machine learning that mimic the way a biological brain behaves and interconnects by connecting and weighting different computer nodes so they act in a coordinated and complementary manner. Deep neural networks, or DNNs, have multiple layers of nodes and inputs/outputs. Self-driving cars use DNNs to recognize objects such as traffic signs and pedestrians.

Today, the broader scientific computing research community has gone beyond classical computational simulation modeling and has started leveraging machine learning models in lieu of, to complement or as surrogates for classic scientific computing models. This emergent modeling paradigm is known as scientific machine learning, or SciML, and has the potential to revolutionize computational modeling.

Here are four of the many examples of how machine learning is used at Sandia to investigate and solve complex scientific problems.

**FIGURE 2.** *Computer simulation credibility model.*

## ENSURE MACHINE LEARNING RESULTS ARE ACCURATE

TEAM | *Erin Acquesta, Bill Rider*

Computational simulation, or CompSim, models are digital prototypes of something a researcher is interested in building or testing. These prototypes make it easier and faster for the researcher to design a product that accomplishes a specific goal without first having to build and test a long list of physical prototypes.

Machine learned models are used in lieu of, complementary to or as surrogates for CompSim models. Machine learned models are useful for situations like pandemic tracking, climate modeling and wherever research and traditional CompSim models are unrealistic or inadequate.

However, a key aspect of using either type of model is knowing whether their simulations are credible. The existing CompSim credibility process provides a proven method for gathering evidence to support the credibility of CompSim predictions (see Figure 2).

Machine learning models are fundamentally different from CompSim models in their relationship to the data. In CompSim, the models are derived from physics or other domain knowledge, and then data is used to calibrate some model parameters or to perform model validation. In machine learning models, there is a primary dependence on the data itself, so developing a machine learning credibility framework requires an initial proof of credibility for the data used to train, test and validate the machine.

The team has extended the standard CompSim credibility process (Figure 2) and adapted it specifically for SciML (see Figure 3).

- Data representation: Does the data provide a representative population for training/testing/validation?

- Domain-aware: What physical phenomena need to be preserved in the model? Is it captured?

- Code correctness: Can the software quality be assured?

- Solution explainability: Are the model solutions explainable to model customers?

- Validation: Do the model predictions agree with the ground truth data that was not used during training?

- Uncertainty quantification: What sources of uncertainty cannot be mitigated or eliminated? What sources can be mitigated or eliminated?

This SciML framework has been systematically tested and refined through a broad set of prototype SciML examples derived from analysis of test and experimental data. HPC computers were used to discover and calibrate SciML models from large-scale datasets and examine the data provenance, machine learning algorithm influence and impact of expert domain knowledge. All this work is helping refine the SciML modeling and credibility framework to provide increased knowledge and enable trustworthy decision making.



**FIGURE 3.** *Scientific machine learning credibility model.*

Comprehensive Evidence Basis
*(Plan, Execute, Organize & Analyze)*

UQ

Data Representation

Validation

**SciML Model**

Domain-Aware

Explainability

Code Correctness /SQA

Application Context
*(Application Requirements, Test-SciML Integration, Derived SciML Requirements)*

SciML Deliverables
*(Plausible Prediction Bounds*

**SciML Model**

PREDICTION

Assess & Communicate
*(Customer, Peer Reviews)*

**LITHIUM-ION BATTERY**

DISCHARGE

CATHODE (+)
ALUMINIUM
CURRENT
COLLECTOR

ELECTROLYTE

ANODE (-)
COPPER
CURRENT
COLLECTOR

LI-METAL
CARBON

LITHIUM ION

SEPARATOR

ELECTRON

LI-METAL
OXIDES

CHARGE

CATHODE (+)
ALUMINIUM
CURRENT
COLLECTOR

ELECTROLYTE

ANODE (-)
COPPER
CURRENT
COLLECTOR

LI-METAL
CARBON

LITHIUM ION
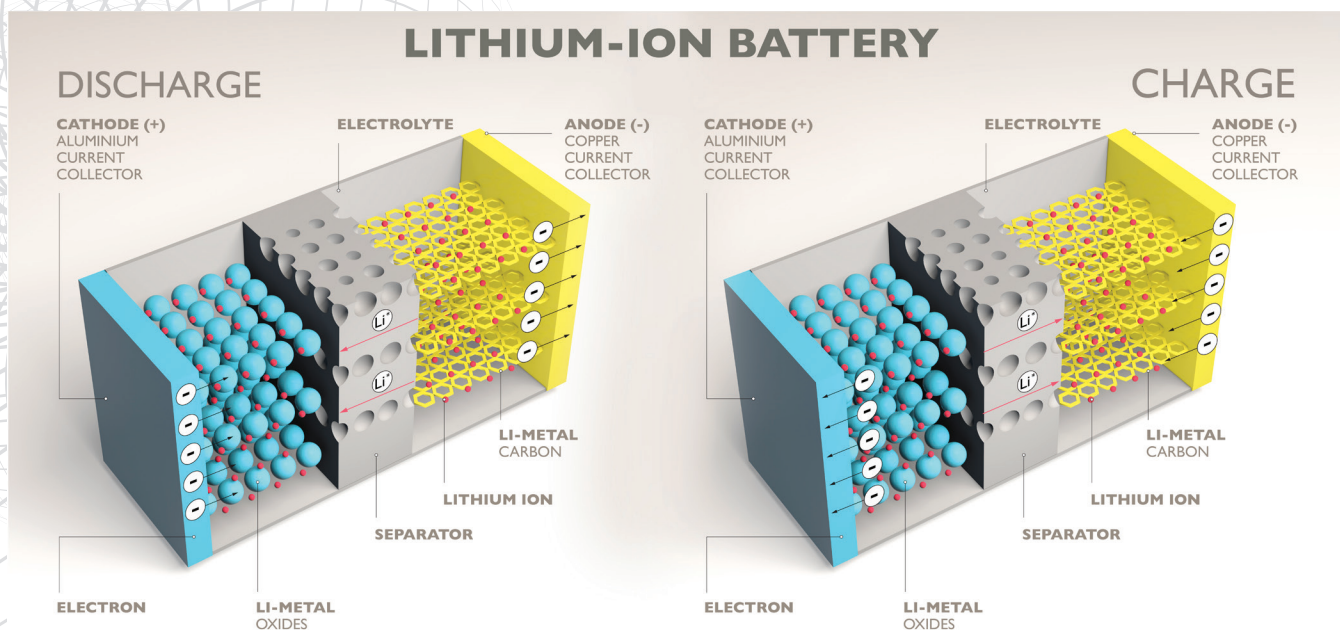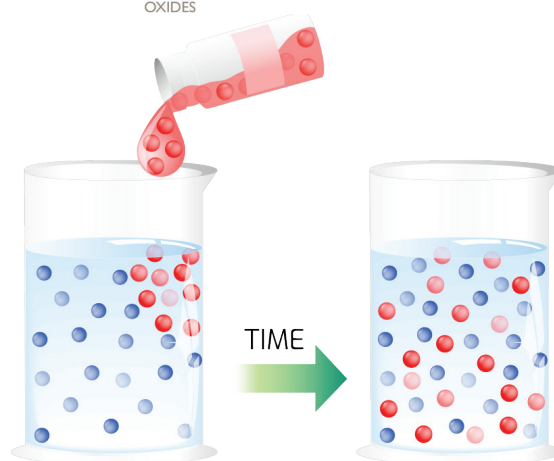
SEPARATOR

ELECTRON

LI-METAL
OXIDES

**FIGURE 4.** *Lithium-ion battery; Diffusion of ink droplets in water.*

TIME

# PREDICT IONIC LIQUID DIFFUSION

TEAM | *N. Scott Bobbitt, Joshua P. Allers, Jacob A. Harvey, Derrick Poe, Jordyn D. Wemhoner, Jane Keth, Jeffery A. Greathouse*

Ionic liquid, or liquid salt, has a wide variety of uses, such as for solar power, batteries, carbon dioxide capture and pharmaceuticals. In addition to the sodium chloride we use in our favorite dishes, chemical salts include a vast array of other possible chemical combinations, such as cobalt nitrate and nickel chloride.

Diffusion is the intermingling of substances by the natural movement of their particles, typically shifting from a higher concentration in a small area to a lower concentration in a large area, for example, a drop of ink in a cup of water (see Figure 4). In liquid salts, diffusion describes the movement between positively charged molecules (cations) and negatively charged molecules (anions). This movement creates energy, for example, a lithium-ion battery powering a smartphone (see Figure 4).

Still, up to a trillion forms of ionic liquids have yet to be investigated to learn more about how they could be used to solve problems not only now but in the future. An enormous hurdle in these investigations is understanding how new types of ionic liquids will behave in various internal and external conditions. Since the behavior of ionic liquids is complex and non-linear, it's extremely challenging to predict.

Sandia researchers began to tackle the problem by having the HPC computers use a combination of molecular dynamics simulations and machine learning methods to investigate 29 different ionic liquids in a controlled environment. They found that machine learning could accurately predict ion diffusion based on six to eight descriptive inputs, such as ion size, mass and geometry. This has proved to be a fast and inexpensive way to accurately predict diffusion behavior for ionic liquids, an important step forward in harnessing their full potential.

# RECOGNIZE RADAR TARGETS QUICKLY AND ACCURATELY

**TEAM** | *Craig Vineyard, William Severa, Ryan Melzer*

Radar bounces radio waves off distant objects and uses the reflected radio waves to display those objects on a screen, enabling us to perceive things normally beyond our sight. The larger the antenna dish receiving the reflected radio waves, the more details we can discern from the target. Synthetic Aperture Radar, or SAR, creates the equivalent of a huge receiving antenna dish by tasking orbiting satellites or a radar attached to a moving airplane to take multiple radar images of a target from different angles (see Figure 5). The controlled movement of the satellite or airplane enables a series of precise images to be captured and then integrated into a single high-resolution image of the target.

However, gathering the data is only the beginning of the work: it also must be analyzed and interpreted. Since high-resolution radar images have so much data to interpret, automating SAR data analysis and interpretation has been a challenge ever since SAR was initially developed in the 1950s. The early attempts to automatically analyze and interpret SAR data with DNNs in the 1980s, 1990s and later did not work well. Although development of DNNs began in the 1970s, they became more powerful and high-performing in 2012, ushering in a multitude of neural network advances. A shortcoming of more recent DNN SAR research has been the focus on only a few questions or variables.



FIGURE 5. *Synthetic Aperture Radar from satellite and airplane*.



Sandia researchers wanted to dig deep and broad in assessing DNNs for SAR automatic target recognition. They used HPC to evaluate thousands of deep neural networks to assess different approaches to automatic target recognition of SAR images. The team focused on assessing multiple DNN architectures for accuracy and speed (see Figure 6), then used data augmentation to close the gap between training datasets and current real-world scenarios, making the DNN results far more accurate.



FIGURE 6. *Training and testing multiple neural network architectures for automatic target recognition.*

Input

GradCAM

Vanilla Gradient Saliency Maps

The Sandia team rigorously applied an extensive analysis exploring the various DNN approaches that have advanced over the last decade. In doing so, they investigated whether neural networks could perform well and what computational structures enabled high accuracy. They investigated deeper/bigger networks, wider netwo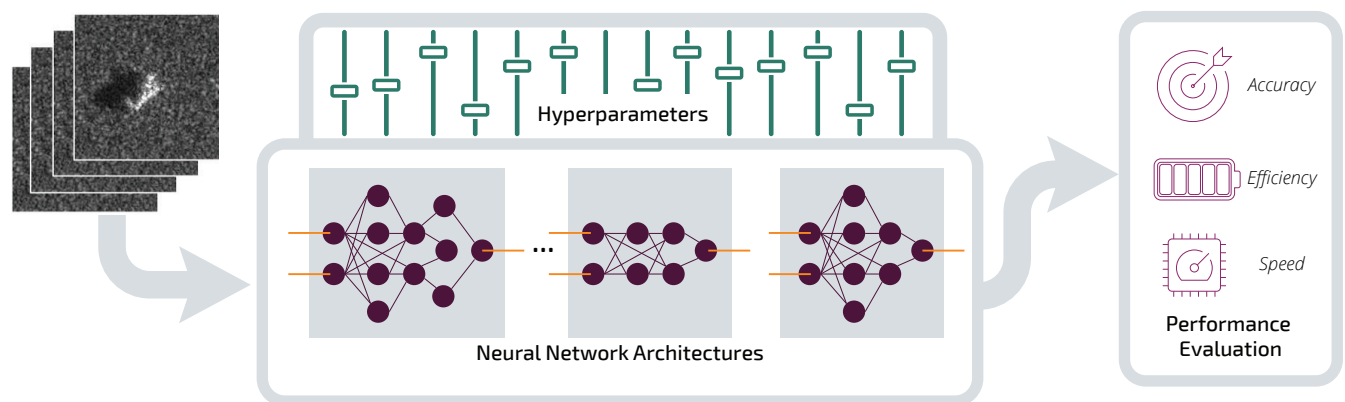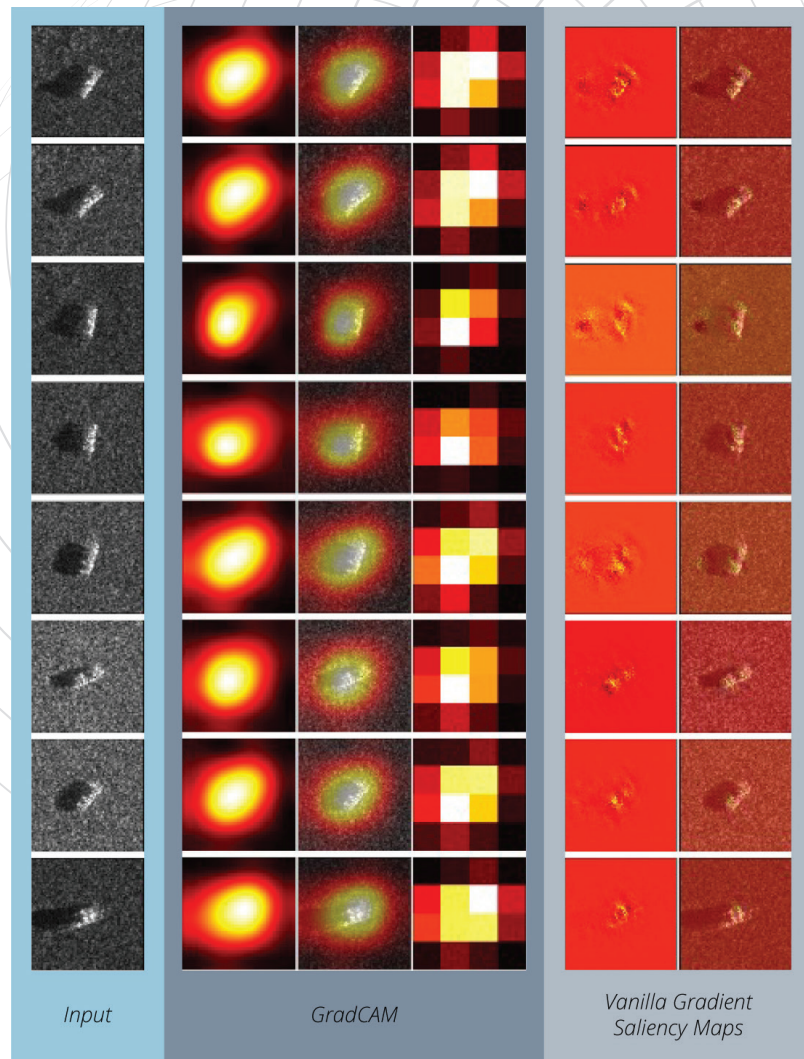rks, sophisticated connections between neural network layers, various optimization methods and more. This involved training many thousands of networks and attaining a state-of-the-art benchmark for accuracy.

Beyond looking at accuracy, the team wanted to learn why the neural networks are making their decisions and so they integrated model explainability methods. To investigate this further, they developed heatmap images showing the parts of the input image that potentially influence the DNN's decision making (see Figure 7). For example, is the DNN making decisions based upon the shadow of a vehicle, the target itself, the background or something else?

Another important area of investigation was reproducibility. Can a given DNN's performance be reproduced? Is one particular computational structure of the DNN actually better than another, or did that training instance just happen to work out better one time? The team investigated this question thoroughly, which contributed to their understanding of how these facets of DNNs apply to SAR automatic target recognition.

The team's extensive model analysis, application of augmentations, reproducibility studies and work on explainability together are building a deep understanding of using neural network approaches for SAR automatic target recognition.

FIGURE 8. *LEGO bricks® and a stackable building constructed from hexagonal interlocking metasurfaces illustrate the specific interlock details..*

# DESIGN STRONG AND FLEXIBLE INTERLOCKING METASURFACES

**TEAM:** *Nathan Brown, Benjamin Young, Ophelia Bolmin, Brad L. Boyce, Philip J. Noell*

Interlocking metasurfaces are a recent Sandia-invented joining technology that may replace other structural joining methods, such as bolts and adhesives. Much like LEGO® bricks (see Figure 8 and Figure 9), they are patterned, modular surfaces that constrain motion between two bodies when assembled into a larger structure. Unlike LEGO bricks, they are designed to sustain high mechanical loads. Each metasurface component precisely attaches and aligns with other metasurface components. The specific component shapes can increase the strength and stability of the individual unit cell as well as the assembled structure. Interlocking metasurfaces can serve as an impermanent joining technology for complex shapes, dissimilar materials, extreme environments in aerospace, civil engineering, micro-robotics and prosthetics.



FIGURE 9. *Intuition-based T-slot unit cell feature design.*

FIGURE 10.
*Performance of computer designs against T-slot design.*

Sandia researchers decided to use machine learning to improve on the human-expert designed unit cells, such as the T-slot, an intuitive design that served as the team's baseline to compare against their computer designs (see Figure 9).

The HPC computers created training data, which were used to teach machine learning models. The models were then tasked with redesigning this interlocking structure with a focus on vertical strength (tensile loading), horizontal strength (shear loading) or a focus on both. The most promising computer designs were 3D printed and then tested for vertical and horizontal strength.

The vast majority of computer-generated designs performed significantly better than the T-slot baseline. Tall, branching unit cells generally had increased vertical strength and shorter unit cells had increased horizontal strength. Designs that combined both features (such as MOGAND in Figure 10) had both increased vertical and horizontal strength.

The team's data-driven design approaches led to radical, organic-looking designs for interlocking metasurface unit cells. The vast majority of these designs performed far better than prior designs based on expert human intuition. These insights will guide future high-performance interlocking metasurface design.

FIGURE 11. *Performance of computer designs against T-slot design.*

# Mission- *focused modeling & simulation*

Time: 0.000000 days since 01JUL2018

# *An innovative HPC approach for modeling Arctic permafrost degradation and coastal erosion*

TEAM: *Irina Tezaur, Diana L Bull, R. Charles Choens, Jennifer M Frederick, Alejandro Mota*

CONTRIBUTING WRITER: *Laura Alice Sowko*

The Arctic is warming at a rate of four times the rest of the globe. The increasing rates of Arctic coastal erosion, driven by climate change, have endangered essential infrastructure and indigenous communities. Despite the Arctic accounting for a third of the world's coastline, existing tools cannot accurately explain unpredictable, storm-induced incidents. Funded by Laboratory Directed Research and Development (LDRD) and Biological and Environmental Research (BER) programs, and in collaboration with University of Texas at Austin, University of Alaska Fairbanks, US Geological Survey (USGS) and Integral Consulting, scientists from Sandia have established an innovative, multi-physics finite element simulation framework, known as the Arctic Coastal Erosion (ACE) model, for numerically modeling Arctic coastal permafrost. The ACE model enhances comprehension of coastal erosion processes, contributes to predictions of erosion-driven biogeochemical and sediment burdens, and can be used to help assess infrastructure vulnerability.

**Legend:**
- Ocean
- Ice wedge
- Permafrost
- Boundary Conditions

**FIGURE 1.** *Figure 1. Aerial photograph of a typical Arctic coastal permafrost geomorphology (top) compared with an analogous computational geometry modeled in ACE (bottom). Ice wedges are enhanced with blue, ocean with purple, and permafrost with green.*

# EXISTING STATE-OF-THE-ART PERMAFROST MODELS

When the Energy and Homeland Security LDRD project that led to the creation of ACE began in 2017, no existing tools could accurately predict Arctic coastal erosion. Many state-of-the-art permafrost models were based on trend projections and/or empirical relationships, lacking consideration for the governing partial differential equations, or PDEs. While a few limited PDE-based models were available, these were predominantly thermal models that ignored the crucial role of mechanics and deformation in the permafrost degradation process. Most of these models assumed a specific erosion type (such as block failure) and failed to account for permafrost geophysics or geomorphologies.

Further, these models tended to not realistically account for the changing conditions driving the erosion or the atmospheric and oceanographic changes; they simply relied on trend projections or aggregate environmental changes incapable of resolving the storm-driven episodic events by which significant amounts of erosion are forced. In contrast to earlier models, ACE is a systems model that ingests accurately modeled boundary conditions at relevant spatio-temporal scales into a 3D PDE-based model that incorporates interlinked thermo-mechanical effects on a discretized permafrost geomorphology, consisting of an ice-wedge polygon morphology (see Figure 1).

**FIGURE 2.** *The ACE model has been calibrated and used to simulate an observed block collapse event at Drew Point, AK on Sept. 1, 2018 (right panel), when forced with realistic oceanic and atmospheric boundary conditions (bottom panel). Permafrost, ice and peat are shown in orange, yellow, and green, respectively.*

# A COUPLED THERMO-MECHANICAL FINITE ELEMENT FORMULATION WITH ELEMENT EROSION/FAILURE

The ACE model's core terrestrial modeling component comprises two main elements: (1) a solid mechanics model, which calculates the 3D stress, strain and displacement fields of the evolving permafrost based on a frozen water content-dependent plasticity model, and (2) an innovative thermal model that regulates the heat conduction and phase transformation within the permafrost. The coupling of these two physics sets occurs through a sequential thermo-mechanical coupling scheme, developed within the Albany-Laboratory for Computational Mechanics, or LCM, open-source finite element code (https://github.com/sandialabs/LCM).

The mechanical component of ACE includes plasticity, modeling the mechanical strength of a permafrost bluff as a function of its thermal state through calculations by the thermal component of ACE. More specifically, the strain and displacement fields of the underlying permafrost are coupled to a novel thermal model governing the 3D heat conduction and solid-liquid phase change occurring within the permafrost.

The terrestrial state is forced using realistic atmospheric and oceanic boundary conditions. An oceanographic modeling suite (consisting of WAVEWATCH III, Delft3D-FLOW, and Delft3D-WAVE) produced time-dependent surge and run-up boundary conditions for the terrestrial model. This suite utilizes historical and projected atmospheric conditions from Earth System Models, or ESMs, to boundary conditions; these ESMs also provide the atmospheric conditions to the terrestrial model. The material and thermal parameter for frozen soil within ACE are informed by experiments on natural cores collected from Drew Point, Alaska, conducted by Sandia's Geomechanics Laboratory. Compressive and tensile experiments were performed at near thaw conditions to define failure and elastic moduli-based ice saturation levels.

# UNIQUE CAPABILITIES OF THE ACE MODEL

The ACE model offers several unique capabilities. First, it can simulate permafrost erosion events by dynamically removing elements from the underlying finite element mesh based on physically motivated failure criteria. Moreover, because the ACE thermo-mechanical model captures permafrost deformation according to ice saturation-dependent stresses, it naturally captures permafrost subsidence as heat flow to and from permafrost. The ACE model hence differs significantly from previous approaches, as it enables failure from any permissible deformation (block failure, thermo-denudation, thermo-abrasion), and defines failure modes based on constitutive relationships innate in the underlying physics, instead of predetermined failure planes, boundary condition approximations, critical niche depth or empirical relationships.

# UTILIZATION OF SANDIA HPC

The ACE thermo-mechanical model is implemented within the Albany-LCM HPC finite element code, which is developed and maintained on a variety of Sandia HPC machines, including Common Engineering Environment (CEE) and Skybridge. Nightly testing is performed on all targeted architectures to ensure code quality and performance are maintained. The ACE team is currently in the process of calibrating and validating their model using observational data collected in northern Alaska during a 2018-2019 field campaign and Sandia HPC (Figure 2). The projected oceanographic boundary conditions through 2050 were also generated on Sandia HPC.

# FUTURE EXTENSIONS AND APPLICATIONS

Once the ACE thermo-mechanical model is calibrated and validated, the ACE team intends to expand the model toward addressing national security issues linked with climate change and permafrost decline, such as threats to crucial Arctic infrastructure. As part of the BER InteRFACE project (https://arcticinterface.org/), techniques to upscale model outputs are being pursued toward studying permafrost demise along the entire Alaskan coastline, among other locations.

Because the ACE thermo-mechanical model is able to naturally capture permafrost, together with its implementation within the Albany-LCM finite element HPC code, the model is particularly well-suited to serve as the foundation for a groundbreaking framework that enables the simulation of critical infrastructure placed on top of permafrost. The team is currently looking for funding opportunities to create a unique simulation tool that would couple the ACE permafrost model to a mechanics-based model of critical infrastructure placed on top of permafrost. The new model would enable assessments of climate change-induced risk to a variety of infrastructure-permafrost combinations, toward predicting likely failure mechanisms and proposing possible solutions to decrease risk. The team is additionally interested in developing, through a series of laboratory experiments involving synthetic soil samples, a universal constitutive model for frozen soils that would be broadly applicable across the Artic.

## REFERENCES

[1] J. Frederick, A. Mota, I. Tezaur, D. Bull. "A thermo-mechanical terrestrial model of Arctic coastal erosion", J. Comput. Appl. Math. 397 113533, 2021.

[2] M. Thomas, A. Mota, B. Jones, C. Choens, J. Frederick, D. Bull. Geometric and material variability influences stress states relevant to coastal permafrost bluff failure. Frontiers in Earth Science, 8:143, 2020.

[3] D. Bull, C. Flanary, C. Jones, J. Frederick, A. Mota, I. Tezaur, J. Kasper, E. Brown, B. Jones, M. Jones, E. Bristol, C. Choens, C. Connoly, J. McClelland. "Arctic Coastal Erosion: Modeling and Experimentation". Sand No. 2020-10223. Sandia National Laboratories, Albuquerque, NM, 2020.

# Exploring plasticity and damage through high-fidelity modeling and novel X-ray techniques

TEAM | *Kyle Johnson, Philip Noell, Hojun Lim, John M Emery, Andrew Polonsky, Matthew Vaughan, Robert A Buarque de Macedo, Demitri Maestas, Carianne Martinez, Kevin Matthew Potter, Aniket Pant*

CONTRIBUTING WRITER | *Laura Alice Sowko*

Predicting plasticity and damage evolution in metal components during mechanical loading environments is critical to ensuring the nation's nuclear stockpile remains safe, secure and reliable. However, doing so remains extremely challenging due to the myriad processes occurring at the microstructural level. As part of a Laboratory Directed Research and Development project focused on predicting material failure using deep learning, researchers evaluated damage evolution in the 2219 aluminum alloy using diffraction contrast tomography, or DCT, and in-situ X-ray computed tomography, or CT, in conjunction with high-fidelity modeling.

In the experimental effort, using a Zeiss Xradia Versa system, pre-test, nondestructive DCT scans delivered 3D reconstructions of crystallographic size, shape and orientation. During subsequent tensile testing, researchers used in situ CT scans to quantify initial $Al_2Cu$ second-phase particle and void distributions, as well as emerging void nucleation, growth and coalescence leading to fracture. After testing, the researchers performed serial sectioning within an electron microscope using the state-of-the-art TriBeam system, resulting in 3D reconstructions of deformed grain structure and orientation and providing insight into microstructural evolution and damage accumulation.

## Model Fracture Predictions



von Mises (Pa)

5.0e+08

4e+8

3e+8

2e+8

1e+8

0.0e+00

## CT Fractured Sample



**FIGURE 1.** *Predicted crack path from continuum finite element simulation (left) and fractured sample from in situ CT test (right).*

**VIDEO.** *View video on SimMagic App available on iOS*

a)

b)

**von Mises (Pa)**



c)



1.2e+03
1150
1100
1050
1000
950
900
850
800
750
700
650
600
550
500
450
400
350
300
250
2.0e+02

**VIDEO.** *View simulation on SimMagic App available on iOS*

In the computational portion of the project at the continuum scale, the initial pre-test CT scan was segmented, reconstructed and converted to hexahedral finite elements using Cubit/Sculpt, resulting in a mesh containing over 8.1M elements. Researchers calibrated continuum plasticity and damage models to macroscale stress-strain behavior for the aluminum matrix, while second-phase particles were treated as elastic based on nano-indentation results. Sierra/SolidMechanics was then employed to predict the plastic deformation, damage evolution and eventual failure. Crack initiation was modeled through element death based on a critical damage value. This simulation utilized 1,656 central processing units, or CPUs, on the Attaway cluster for 48 hours. When compared to post-test CT scans, the predictions for crack path and fracture surface appear qualitatively similar. At the mesoscale, DCT data was converted to a finite element mesh and directly used in a crystal plasticity simulation, allowing for the study of stresses and strains in individual metal grains. Combined with DCT and in situ CT, these simulation results allow modelers to study the effects of crystallographic orientation on void nucleation. Notably, in addition to the effects of initial grain orientation, the crystal plasticity results uncovered relationships between void nucleation and grain rotation during loading, which allowed for the discovery of correlations that were not present in initial orientations.

200 μm

The results of this combined experimental and computational approach have provided great insight into microstructural deformation mechanisms as well as improved workflows. The researchers developed a streamlined approach to go directly from CT scan to finite element model while still considering different material phases present in the sample. Analyzed in situ, CT data has provided unparalleled detail into void formation and evolution, which will allow future researchers to calibrate damage models by using measured information on void evolution rather than stress-strain data alone. This improvement will enable better predictions of mechanical performance of structural components. Combined DCT, in situ CT and post-test TriBeam characterization has enabled for the first time at Sandia a direct comparison between mesoscale model predictions and experiments as well as insight into the local, grain-level stress states and their effect on damage evolution. Finally, these simulations are currently serving as training data for the development of Deep Learning algorithms in an effort to rapidly predict failure of components based on loading conditions and microstructural features.



CT Lower Frature Surface          Model Lower Fracture Surface

$y_{CT}$
$x_{CT}$
$z_{CT}$

$y_{CT}$
$x_{CT}$
$z_{CT}$

Height (micron)
-500.0  -450.0  -400.0  -350.0  -300.0  -250.0  -200.0  -150.0  -100.0  -50.0   0.0

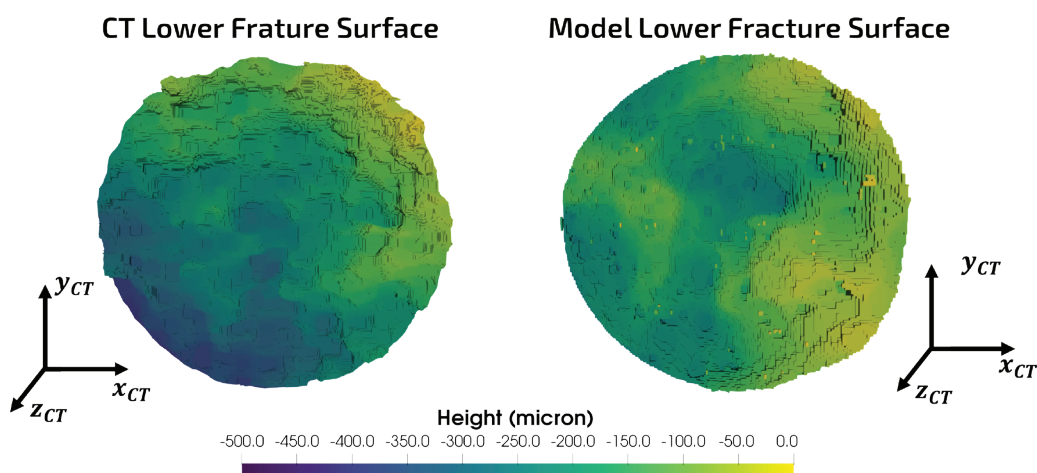**FIGURE 3:** *Comparison of topography in experimental fracture surface (left) and predicted fracture surface from simulation (right).*

# Keeping nuclear weapons safe and reliable by predicting wear in critical safety mechanisms
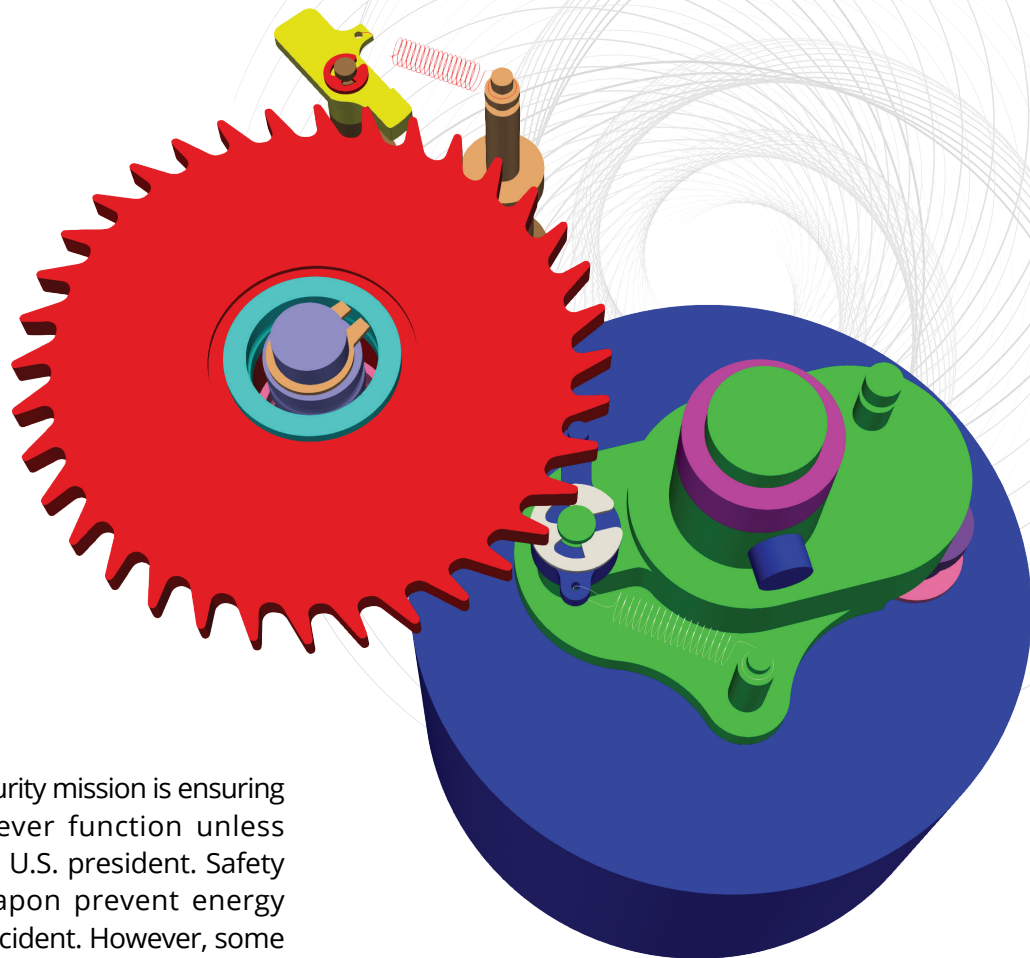
**TEAM |** *Jennifer Fort*

**CONTRIBUTING WRITER |**
*Kristen Meub*

FIGURE 1. *The gears and other components inside a watch are similar in complexity to safety mechanisms.*

**FIGURE 2.** *Sandia can now predict wear on safety mechanisms using computational mechanics models on HPC platforms.*

Part of Sandia's national security mission is ensuring that nuclear weapons never function unless authorized to do so by the U.S. president. Safety mechanisms in each weapon prevent energy from passing through by accident. However, some conditions a safety mechanism may experience, such as vibration, could cause parts to rub together, leading to wear and reduced reliability.

A team of researchers at Sandia built HPC models that can predict how much wear nuclear weapon components, including the subcomponents that make up the safety mechanisms, will experience over time in a specified environment.

"Mechanisms are an important part of a weapon's operation," said Jennifer Fort, project lead and mechanical engineer. "If certain parts within a component wear too much, they may not function as intended, which means the reliability is not what it should be. This is similar to how components in a watch can rub together and cause damage and wear."

The team used Sierra Solid Mechanics to simulate the wear rate on various parts in a safety mechanism, modeling contact between various sub-components and springs. This was done for multiple orientations in a specified environment. The simulation modeled a compressed one-second acceleration time history, including acceleration in the x, y and z directions.

Each simulation used 720 processors and ran for more than 75 days.

"We simulated the wear rates so that we could avoid future testing of actual hardware," Fort said. "We gathered testing data on wear rates for one environment and then compared it with the simulations using www.paraview.org. We observed similar wear patterns and more wear on the same parts in both the test and simulation data."

The team then simulated wear rates for a new environment that had not been tested, again running simulations for 75 days on 720 processors. They compared the wear rates and patterns to the previous simulations and found similarity between the two.

"This will help us evaluate future environments," Fort said. "The results showed the ability to predict the wear rates and compare future environments that a mechanism could experience. This work would not have been possible without HPC due to the size and complexity of the models."

# Edge computing for explosives

**TEAM** | *Jered Wade Mitchell, Theo Stangebye, Sam Bowie, John Korbin, Paul Sanchez, Anita Watson and Dustin Romero*

**CONTRIBUTING WRITER** | *Mollie Rappe*

## CREATING DIGITAL TWINS OF EXPLOSIVE TOOLS PROVIDES CONFIDENCE

In a metal high bay approximately 10 miles south of Sandia's main New Mexico campus, Sandia engineers are using computerized tomography, or CT, scans and computing power to create digital twins of explosive tools.

Using HPC algorithms, Sandia engineers can detonate precise digital replicas of explosive tools over and over again, unlike real-life explosives which can only be used once, according to Theo Stangebye, Sandia software systems engineer involved in the project. This capability provides Sandia's national security customers the confidence that these devices will always perform as expected.

Over the past 18 months, Sandia has built the HPC infrastructure and refurbished several HPC clusters to enable this important work at the far edge of Sandia's information technology infrastructure, where the internet is slow and the engineers can only devote 12 kW of power for computing. These small clusters aren't intended to replace Sandia's corporate HPC resources, but rather provide an agile complement to them. They mimic, at a much smaller scale, the technologies that make HPC possible. That's why it's called mini-HPC at the edge.

While the engineers only have a few small clusters, they use the same scalable provisioning technology such as diskless booting and Simple Linux Universal Resource Manager queuing.

## FORECASTING EXPLOSIVE BEHAVIOR FOR CONFIDENCE

John Korbin and Sam Bowie, Sandia mechanical engineers, use CT scans, or computed tomography, and commercial image processing frameworks to turn the X-ray images of an explosive tool into a digital twin of the device. The process is somewhat similar to how CT scans in medicine can help locate a tumor. However, the team's CT scanner is larger and uses more intense X-rays.

Then Korbin and Bowie use Sandia's shock physics modeling code to create a high-reliability forecast of the behavior of the digital explosive. Having the CT machine allows them to create accurate digital representations of each individual explosive. This allows them to do inspections and take measurements to forecast how the explosive will perform, affording Sandia's customers the confidence that the tool they hold in their hands will perform as expected.

By being able to automate the end-to-end construction of digital twins, the team can ensure the quality of each tool much quicker and examine more devices in a day. Some days the CT team produces up to a terabyte of data.
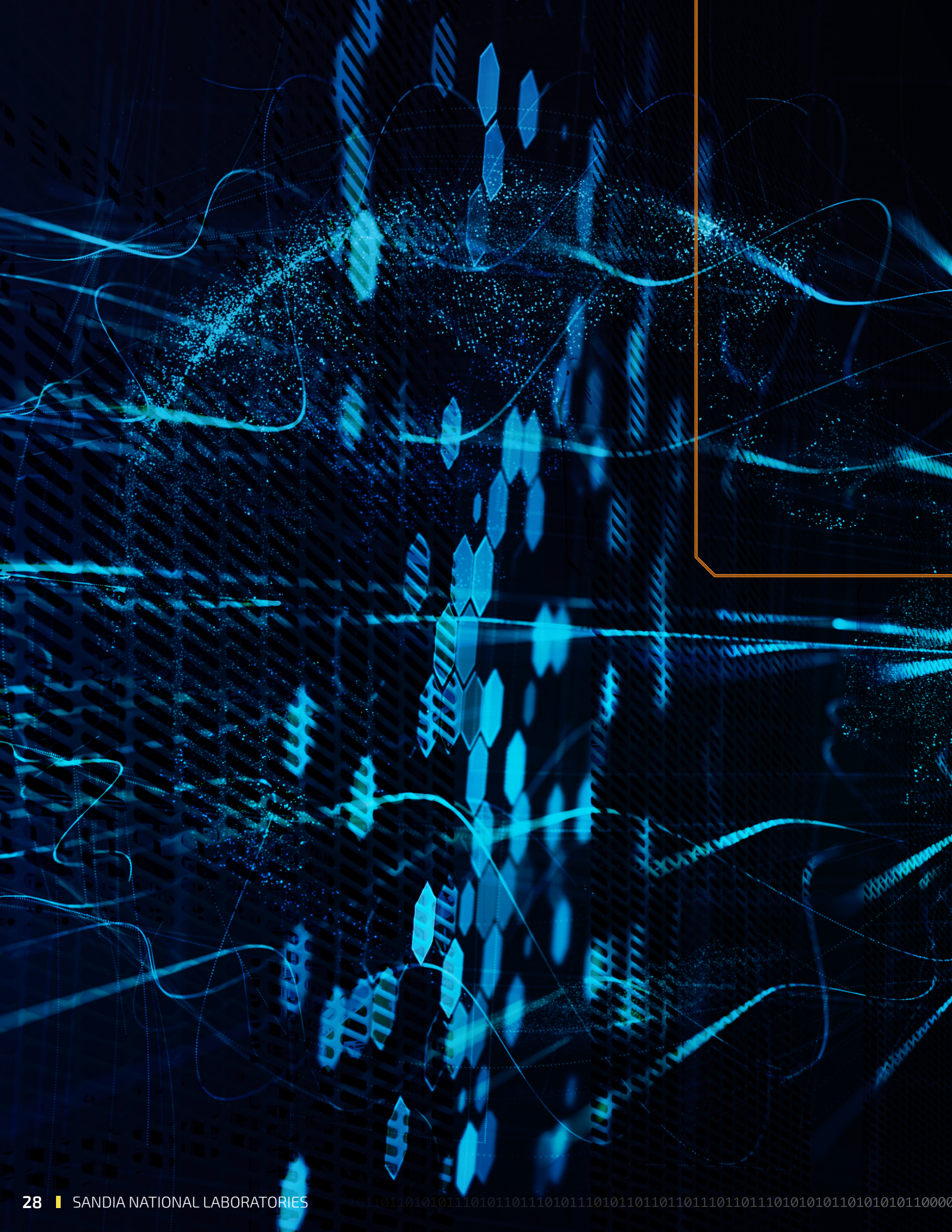
## LEAN, AGILE COMPUTING TO AID THE MISSION

Using this agile approach has allowed the team to find manufacturing defects in prototype tools, as well as smaller variances in devices due to the computer-aided design process. Using these methods and having a dedicated computing resource without an extensive queue has greatly increased the agility of the explosives group. Integrating the computing side with the hardware side has changed the culture of how the team thinks about engineering explosive tools.

Explosives experts create digital twins of every explosive tool before testing them. Sandia doesn't use HPC forecasting for every device, but the digital twin is available for analysis if the device performs oddly or particularly well. And since the cutting-edge resources are conveniently located, the explosives experts can ask questions and run simulations they would never have thought of a decade ago.

Next steps for the project include incorporating more machine learning into the processes. The challenge with machine learning is voluminous training data sets, in this case X-ray images, need to be created and labeled. However, the team has developed new algorithms that can perform crude materials segmentations, which might be used as training data for machine learning.

The team hopes to share their paradigm of lean, agile, scalable HPC at the edge with other groups at Sandia which would benefit the Labs and the nation as a whole.

*High Performance Computing looking forward*

# Stretching the thermodynamic limits of HPC efficiency

TEAM | *Michael P Frank*

CONTRIBUTING WRITER | *Neal Singer*

*As supercomputing systems evolve toward ever higher levels of performance density, their continued advance is threatened by looming fundamental thermal limitations. Historical performance improvements were largely enabled by increases in energy efficiency, which will soon slow down as the low-hanging fruit will have been plucked. Continued progress beyond this point would require the use of an unconventional approach called reversible computing, to conserve energy in digital systems by conserving information. Research at Sandia is investigating the viability of this approach.*
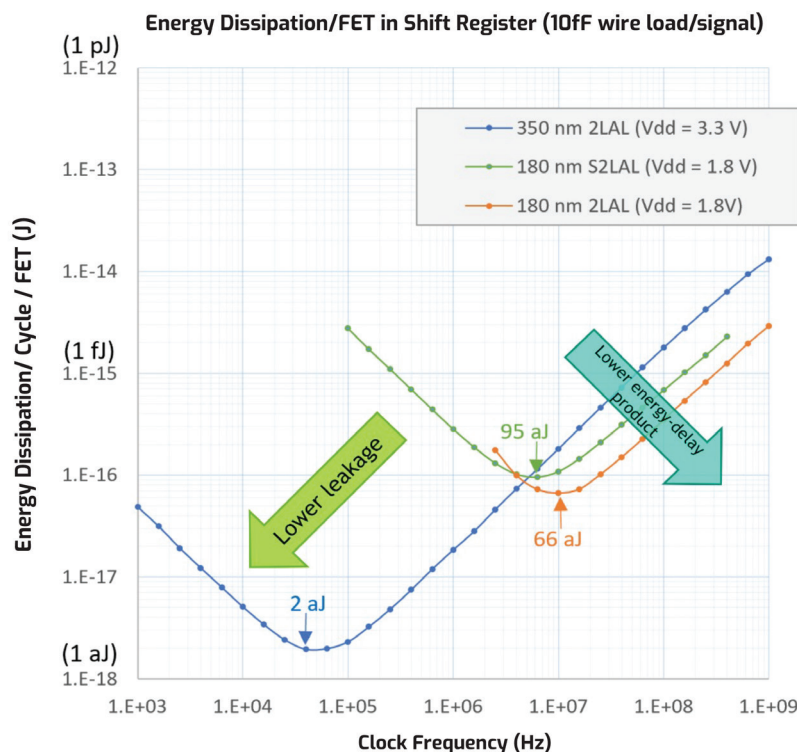
**Energy Dissipation/FET in Shift Register (10fF wire load/signal)**

Legend:
- 350 nm 2LAL (Vdd = 3.3 V)
- 180 nm S2LAL (Vdd = 1.8 V)
- 180 nm 2LAL (Vdd = 1.8V)

Lower leakage

Lower energy-delay product

95 aJ

66 aJ

2 aJ

Y-axis: Energy Dissipation/ Cycle / FET (J); (1 pJ) 1.E-12, 1.E-13, 1.E-14, (1 fJ) 1.E-15, 1.E-16, 1.E-17, (1 aJ) 1.E-18

X-axis: Clock Frequency (Hz); 1.E+03, 1.E+04, 1.E+05, 1.E+06, 1.E+07, 1.E+08, 1.E+09

**FIGURE 1.** *Results from Cadence simulations of a simple sequential circuit in the fully adiabatic 2LAL (two-level adiabatic logic) style as well as its fully static S2LAL variant, in two different MESA processes. Interconnect loads were modeled at 10 fF. In the 350 nm process, energy dissipation per device-cycle in a conventional CMOS circuit would have been more than 10 fJ; whereas in 2LAL it could be made as low as 2 aJ, corresponding to a raw efficiency boost of up to ~5000x. Even accounting for circuit complexity overheads, we can save >99% of the energy that would be dissipated by a conventional circuit for the same function.*

Over the last several decades the overall performance rating of leading HPC systems in floating-point operations per second (FLOPS) has increased steadily, from teraflops in the 90s to over an exaflop today. This increase would not have been feasible if not for significant improvements in the energy efficiency of digital technologies. The astounding millionfold performance increase from the ASCI Red supercomputer in 1997 to Frontier today was achieved while only increasing power consumption by 25×; this corresponds to about a 40,000× improvement in performance per watt over the last quarter-century. About 200× of this relates to a reduction in the scale of digital signal energies, with minimum transistor gate energies having gone from about a femtojoule in the 1997 250 nanometer (nm) technology to around five attojoules in today's leading-edge "three nm" fin field-effect transistors, or FinFETs, while the other 200× came from optimizations at higher levels in system design (architecture, packaging, memory, interconnects and software).

But now, further improvements across all these levels threaten to soon enter an era of diminishing returns. In particular, at the device level, transistor gate energies are only projected to decrease by around another 2× through further process improvements. These energies are limited by fundamental thermal noise considerations to no more than about a 50× reduction from today's level, even if aggressive subthreshold voltage scaling techniques can be applied. The room left for optimization at higher levels in the technology stack is similarly limited. Hence, it seems unlikely that the rates of performance improvement that we have been accustomed to in digital computing can be sustained much longer if we persist in using conventional methods.

Among unconventional approaches to computing, only one offers the potential to increase the efficiency of low-level digital compute far beyond the fundamental thermal limits of conventional technology mentioned above: Namely, reversible computing. This approach, first conceived in the 1960s and '70s by pioneering researchers Rolf Landauer and Charles Bennett of IBM Research, stems from the observation that there is only a fundamental thermodynamic minimum on the energy that must be consumed to perform a digital operation if logical information is lost in the course of carrying out that operation. When logical information is conserved, the digital signal energies used to represent that information can, in principle, be recovered and reused as well; that is, we can theoretically approach the ideal of computing in a manner that is both logically and thermodynamically reversible.
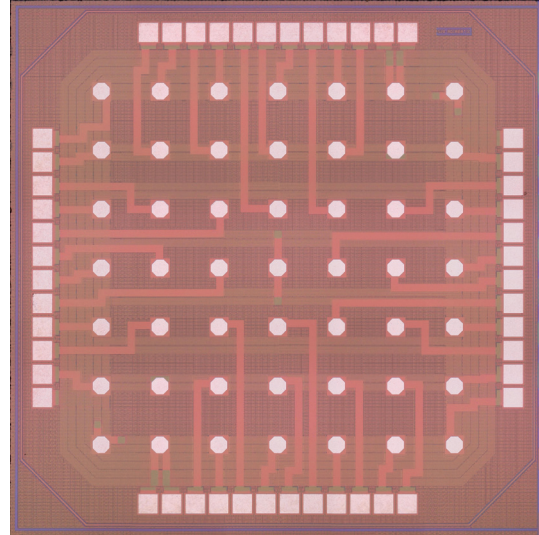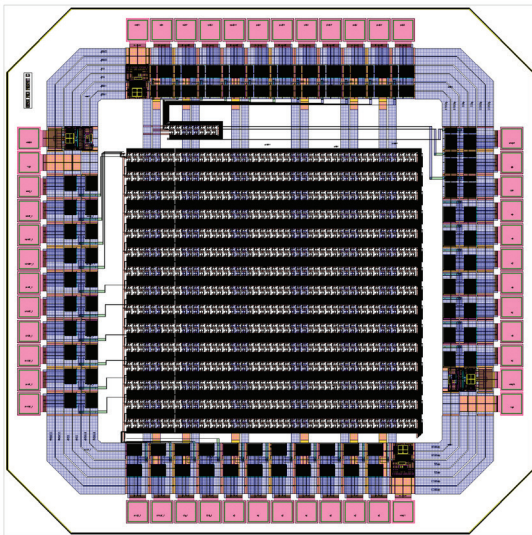
**FIGURE 2**. *Layout and photo of 2LAL test chip fabricated in MESA's 180nm process. This small (2 mm) die contains a 720-stage 2LAL shift register suitable for testing functionality and energy dissipation. It is currently undergoing probe station testing.*

The engineering challenges that would need to be addressed to approach this ideal in practice are formidable. Given the National Strategic Computing Initiative's mandate to find a path forward beyond the limits of current technology, NNSA's Advanced Simulation and Computing, or ASC, program has supported an effort at Sandia since 2017 to carefully assess the potential for further development of the reversible approach to feasibly carry digital computing efficiency beyond the end of the semiconductor roadmap. In parallel, a Laboratory Directed Research and Development, or LDRD, project started that year led to subsequent Strategic Intelligence Partnership Program, or SIPP, sponsorship of a complementary line of work to expand the scope of the team's effort to include early-stage research into novel reversible superconducting technologies for digital compute.

A foundational aspect of this work has been to solidify the theoretical basis for reversible computing. In a series of papers, Sandia showed how the concepts of reversible computing can be generalized from traditional models to a more adaptable form suitable for practical computing hardware. Sandia clarified why the physical motivation for reversible computing

follows rigorously from fundamental concepts of statistical mechanics and information theory. Briefly, because fundamental physics is reversible, discarded information cannot be destroyed, and instead it inevitably manifests as new physical entropy. Currently, in collaboration with university partners, the team investigating the ultimate physical limits of reversible computing efficiency using theoretical tools from non-equilibrium quantum thermodynamics.

Building on the theoretical foundations, the central thrust of this work is to develop and analyze engineering implementations of these concepts. Sandia can implement reversible computing in existing CMOS technology via adiabatic switching based on alternative logic gate designs and clocking disciplines. While the basic principles behind this approach have long been recognized, past implementations were far from perfect. However, the team's focus is to truly push the limits of this approach. Sandia's major contributions in this area to date include: (1) the invention of a new complementary metal-oxide semiconductor, or CMOS, logic family that is perfectly adiabatic; (2) the design and fabrication of a test chip using Microsystems Engineering, Science and Applications, or MESA,
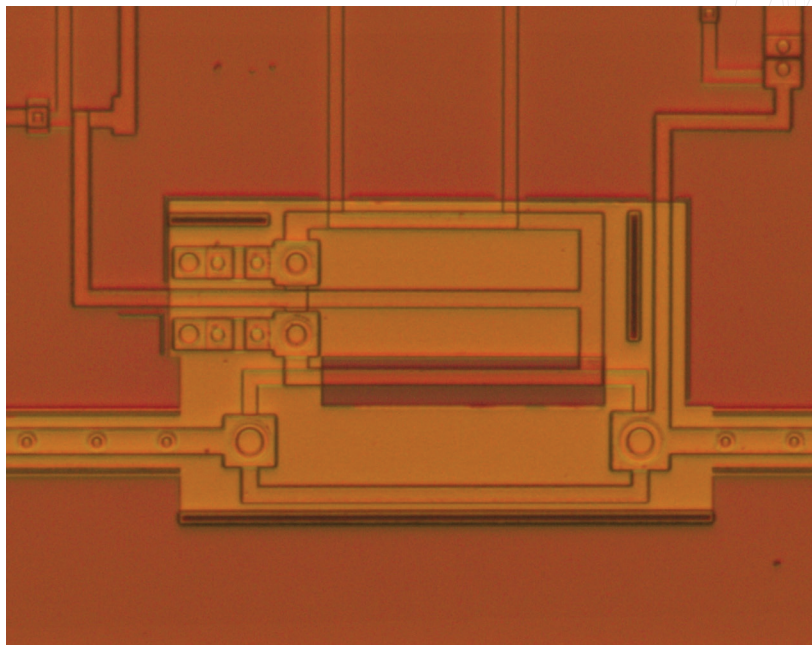
**FIGURE 3.** *Closeup photo of a test circuit for our ballistic reversible memory (RM) cell fabricated in SeeQC's niobium superconducting process. Shown are the storage loop, discretized long Josephson junction (dLJJ) based interconnects, larger junctions for exchange and reset of flux quanta, and an RF SQUID circuit for flux state measurement. This circuit is also still under test.*

180nm process; (3) circuit simulations showing that raw efficiency gains of up to three orders of magnitude are achievable in MESA's processes (beating even end-of-roadmap CMOS!); (4) the invention of a novel resonant oscillator circuit designed to recover >99.9% of the signal energy from an adiabatic CMOS chip; (5) a detailed analysis of the theoretical gains in raw efficiency and throughput density that could be achieved across future CMOS technology nodes, showing that >100× boosts in raw throughput density versus conventional CMOS should be attainable. Sandia's work in this area continues, with results to date suggesting that, with further development, reversible computing technology could significantly benefit future HPC systems.

In a complementary line of work funded by external agencies, Sandia is looking farther ahead to investigate a novel approach to reversible computing. This approach is based on the ballistic transport and asynchronous elastic interaction of flux solitons, or fluxons, conveying single magnetic flux quanta in superconducting circuits based on Josephson junctions, or JJ$_s$. Superconducting approaches are interesting for their potential utility in applications at cryogenic temperatures, such as for the control and readout of superconducting qubits in quantum computing. This line of work is still at a very early stage of development, but important advances to date include a patented design concept for a single-JJ reversible memory, or RM, cell, a fabricated niobium test chip for this circuit, and an enumeration of possible

primitive operations in our new Ballistic Asynchronous Reversible Computing in Superconductors, or BARCS, computing paradigm. The team's future plans involve utilizing Sandia HPC resources for AI-enhanced automated discovery of additional BARCS circuits suitable for general-purpose computation.

In closing, the team's work exemplifies Sandia's dedication to innovation and addresses the national strategic need for groundbreaking efficiency in digital computation. This research, while challenging, helps to bolster the long-term contribution of Sandia National Laboratories to the national interest, and has enjoyed support from NNSA's ASC program, Sandia's LDRD program, and SIPP partners. Looking forward, the team anticipates continued progress and new discoveries in the quest to push the limits of energy efficiency in digital computation. The team will continue to leverage Sandia's HPC resources as to flesh out this novel approach to reversible computing, aiming to contribute not only to the field of HPC, but also to advanced applications and a more sustainable future. Consistent results and ongoing advancements in this field serve to reinforce confidence in this research and to affirm Sandia's role in leading technological innovation.

# 'Sandia Inside' HPCs and their impact on exascale systems

**TEAM** |
*Ronald Brightwell*

**CONTRIBUTING WRITER** |
*Johann Snyder*

## portals

In 2008 the Defense Advanced Research Projects Agency, or DARPA, released a report entitled "ExaScale Computing Study: Technology Challenges in Achieving ExaScale Systems." This report outlined several obstacles in achieving the next three orders of magnitude performance increase beyond petascale systems. At that time Sandia developed a strategy and infrastructure for doing hardware/software co-design, which is a process for exploring hardware design tradeoffs and the necessary changes to algorithms, applications and system software to motivate and exploit new hardware capabilities. As the focus of the HPC community shifted toward addressing the many challenges in providing a usable exascale system for the Department of Energy, or DOE, one of the primary goals for Sandia's co-design work was to put a "Sandia Inside" logo on a DOE exascale system. Playing off the popular "Intel Inside" marketing campaign, the ultimate objective was to have Sandia's influence on a DOE exascale system be clearly evident.

The Frontier system at Oak Ridge National Laboratory, from supercomputer vendor Hewlett Packard Enterprise, or HPE, became the first machine to break the exaflops barrier on the Top500 benchmark in June 2023, and two of Sandia's co-design technologies were instrumental in designing the Slingshot-11 high-performance interconnect, which began initially at Cray, Inc. and continued at HPE when they acquired Cray in 2019. Sandia's Portals 4 network programming interface and Structural Simulation Toolkit, or SST, both played a key role in the development of the HPE's Slingshot-11 fabric, which is composed of two hardware components, the network interface controller, or NIC, called Cassini and the network switch called Rosetta. Frontier is the first DOE exascale system to be deployed, but the next two DOE exascale systems, El Capitan at Lawrence Livermore National Laboratory and Aurora at Argonne National Laboratory, will also use the Slingshot-11 network, so all three DOE exascale systems will have "Sandia Inside."

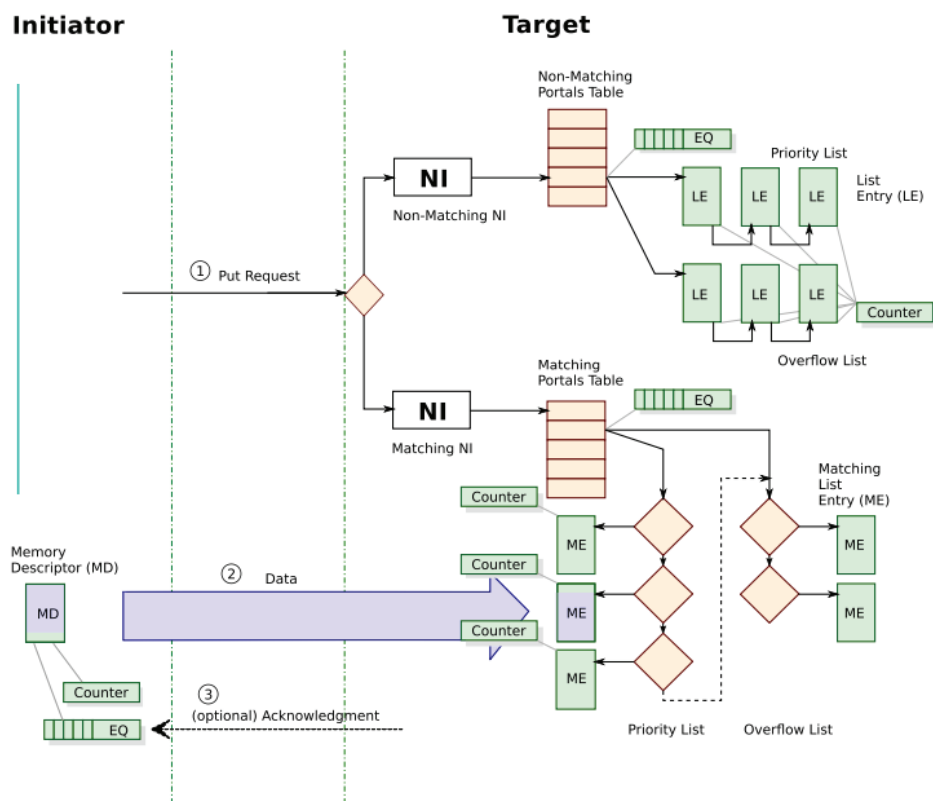The Portals 4 network programming interface was specifically designed to enable co-design and drive the development of network hardware to meet the stringent performance and scalability requirements of applications and services running on today's HPC platforms. The network of an HPC machine has stringent performance and scalability requirements that are unique to scientific computing applications and system services. For HPC systems, the interconnection network determines how large the system can grow and how fast applications that span the machine can execute.

Prior versions of the Portals network programming interface developed primarily by Sandia were simply portability layers that translated low-level network functionality to meet the needs of HPC applications and system services. However, Portals 4 encapsulates the capabilities needed by all HPC applications and system services in a way that allows NIC designers to implement these capabilities efficiently in hardware, reducing the semantic gap between what the hardware provides and what the software needs. Unlike other low-level network programming interfaces, the focus of Portals 4 is not on a software implementation, but rather on a detailed specification of objects and functions. This approach allows hardware architects the freedom to continue to innovate but does so in a way that reduces the need for software portability layers.

"Portals 4 provided inspiration for the Cassini line of NICs in the HPE Cray Slingshot network, the interconnect enabling the first exascale system on the Top500 list. As a network Application Programming Interface, Portals 4 does not specify a network interface architecture or a network transport protocol; rather, it provides insight into how to build a network transport for high performance computing," HPE Vice President and General Manager Mike Vildibill said. "The visionary work done by the Portals team significantly accelerated the architecture definition phase of Cassini."

In addition to HPE, several vendors have designed and deployed network hardware based on Portals 4. The Bull eXascale Interconnect (BXI, BXI-2 and BXI-3) networks produced by European

**FIGURE 1.** *This figure shows a Portals Put operation that sends data from the initiator to the target. This figure illustrates several concepts in the Portals programming interface.*
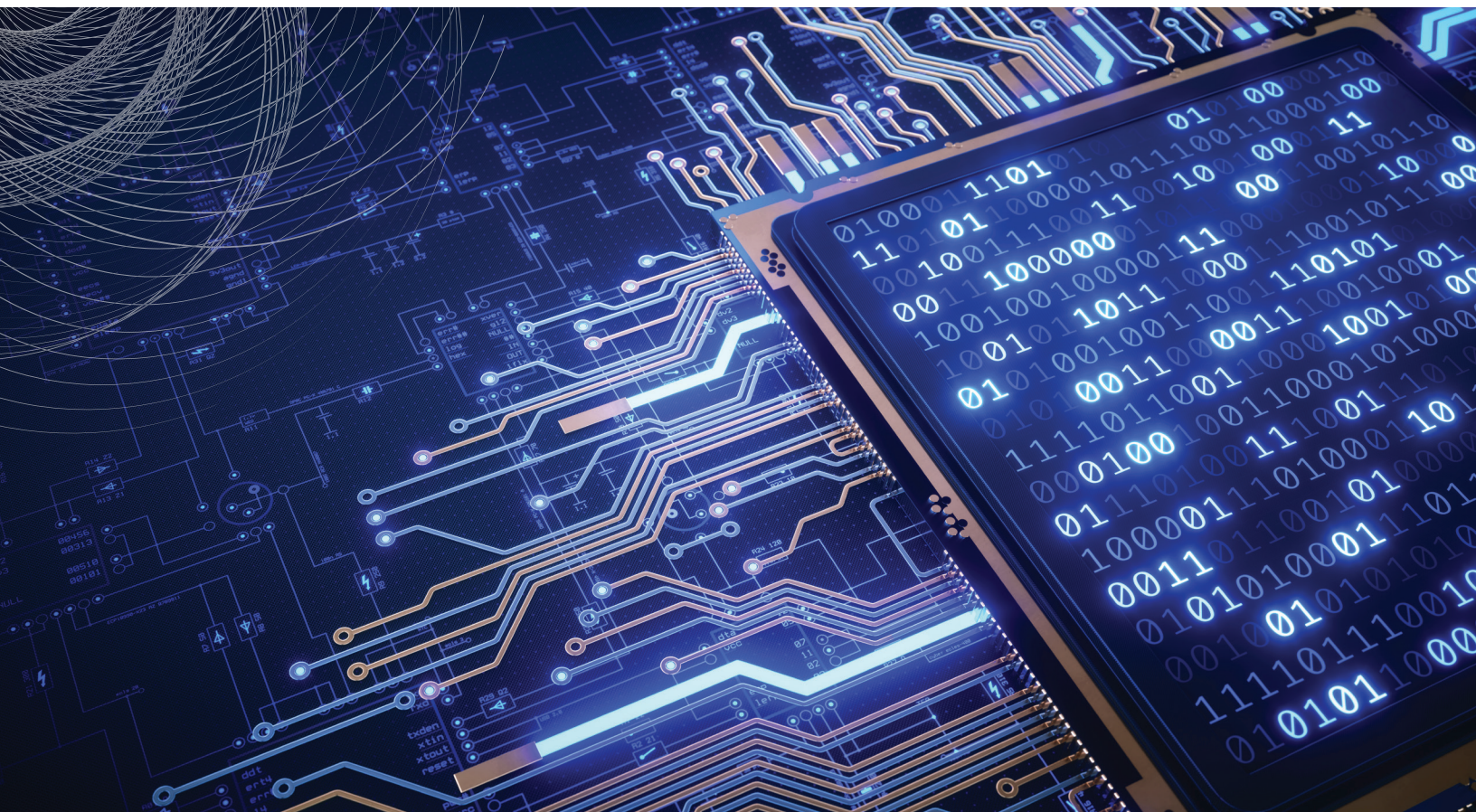
supercomputing vendor Atos are also based on Portals 4. Portals 4 is also a key driver in the development of Intel's OmniPath Architecture, which was recently taken over by Cornelis Networks as OmniPath Express.

Additionally, Portals 4 has influenced software design. The Open Fabrics Interface is an extension of Portals 4 that Intel spearheaded to help transition the deployment of their OmniPath Architecture series of HPC networks from their initial product based on InfiniBand exposing the Verbs interface to later generations based on Portals 4.

"In my previous role as Fellow & Chief Technologist for at Intel, the Portals 4 specification was a key technology that informed and guided the evolution of Intel's OmniPath Architecture," Chief Technologist for HPC at Google Cloud, Bill Magro, said. "Further, Portals 4 was very influential in the design of the network interface hardware and the semantics of the Open Fabrics Interface API, or OFI. OFI was created with two goals in

mind: first, create a much-improved semantic match between the needs of HPC and Machine Learning application software and, second, create a flexible and extensible interface to unlock faster innovation and greater competition in the interconnect fabric space. OFI has been very successful in both of these regards, as it has seen widespread adoption by multiple network vendors. Portals 4 was a critical proof point that the industry could develop tailored networking technology while maintaining a vendor-neutral software interface, and it is arguably the most successful hardware/software co-design tool for high performance networking I encountered during my years of technical leadership."

OFI has seen widespread commercial adoption by several vendors, including Amazon's Elastic Fabric Adapter. Portals 4 has also been extremely influential in the supercomputing research community, spurring further explorations in network software and hardware design for cloud and data center computing.

Sandia's SST also played an important role in the development of the Slingshot-11 network. SST is a parallel discrete event simulation framework that simulates the low-level interactions of computing devices at the architectural level. SST has a modular design that enables extensive exploration of individual components, such as memory configurations and processor instructions, without intrusive changes to the simulator. In particular, it allows vendors to easily plug their own device simulators into the framework. Since SST is parallel, it is able to run complex simulations more quickly. SST was used extensively by many component and system vendors, including Cray, for design-space exploration as part of the PathForward program under the DOE Exascale Computing Project.

HPE Senior Distinguished Technologist Duncan Roweth explains, "Cray had always performed extensive system scale simulation using home grown tools. As the Slingshot project started, the Cray team decide to adopt an open-source simulator. We selected SST for two main reasons. Firstly, we knew that we would need to scale to large system sizes, thousands of switches with tens of thousands of endpoints or NICs. Secondly, the DOE Labs community was developing motifs and miniApps that characterized their applications and integrating them into SST. The Cray team worked with developers at Sandia introducing a device interface that allowed for interchange of open source and proprietary device models. Cray implemented models of the Rosetta switch and Cassini NIC that used this interface. The resulting devices are being used in all three of the US exascale systems. We regard the work on SST that Cray undertook with Sandia as an excellent example of vendor-labs codesign, one that we will build upon in future projects."

Sandia's efforts in developing co-design tools, such as Portals 4 and SST, using those tools to explore the impact and interaction of hardware and software and partnering with HPC system and component vendors like HPE to transition design concepts into reality have seen many successes. While the HPE logo will eventually appear on all three DOE exascale machines, all three could just as easily display the "Sandia Inside" logo too.

# 'Always on' performance monitoring for HPC applications systems

TEAM | *James M. Brandt, Ben Allan, Jeanine Cook, Ann Gentile, Vivek Kale, Stephen Olivier, Mark Schmitz, Benjamin Schwaller, Vanessa Surjadidjaja, M. Scot Swan, Sara Walton*

CONTRIBUTING WRITER | *Steve Scott*

The power, scale and efficacy of today's HPC systems makes them turbocharged vehicles for innovation and discovery. A new capability developed by Sandia researchers seeks to maximize efficiency and performance for every one of those vehicles and the application developers who drive them.

The capability, called AppSysFusion, gives HPC users and administrators an "always on" data collection and analysis framework to investigate, diagnose and efficiently address performance issues at runtime or post-run. The framework functions something like the on-board diagnostics in a modern car, going beyond simple alerts, check engine for example, to analyze detailed performance data and indicate potential solutions.

In HPC terms, this rapid diagnostic capability provides a means for reducing downtime and the number of debugging, development and other non-science runs, thereby increasing science throughput for HPC systems as they continue addressing the nation's most significant scientific and security challenges.

Toward this end, the project team that developed AppSysFusion seeks to enhance the computing performance available to every team that uses an HPC system. The project improves and broadens access to HPC monitoring capabilities, providing users and administrators with real-time insights into both application and systems performance. The capability provides resource utilization views, application progress metrics and root cause diagnosis capabilities. This enables application users and developers, system administrators and HPC architects to make data-driven decisions about how shared HPC resources are provisioned and utilized. The project team also sees it as a key step toward a more advanced self-driving system that would further automate performance enhancement.

AppSysFusion derives its name from a signature innovation – the ability to rapidly merge performance data gathered from running applications with data pulled from the system itself. Correlating these data sets is often crucial to fully understanding performance issues and identifying their source
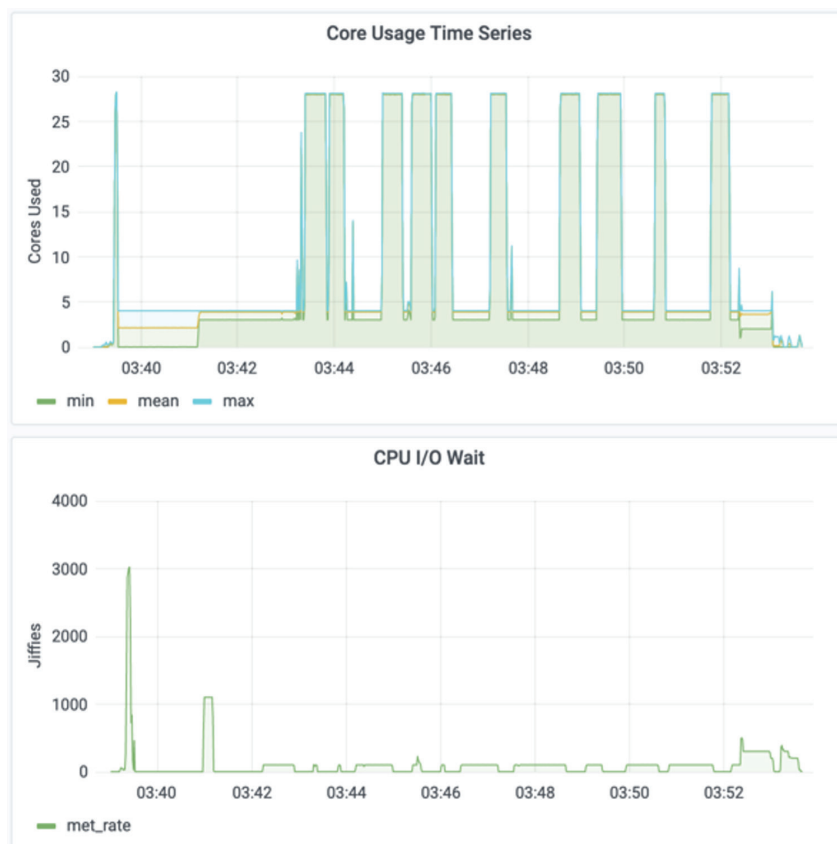
FIGURE 1. *This eight-node simulation usually runs in about five minutes, but with considerable variability (up to 3x slowdown). The AppSysFusion dashboard shows that the I/O is the cause for this slowdown.*

but, without automation, the task can take weeks of human effort. By making the combined data easily available for use, AppSysFusion avoids that detour and helps untangle previously intractable performance issues.

The AppSysFusion capability is built around the Lightweight Distributed Metric Service, or (LDMS), a Sandia-developed monitoring application recognized in 2015 with an R&D 100 award for its ability to regularly collect fine-grained detail about HPC system resources. Lightweight, in this instance is crucial; it means that the data gathering adds no significant processing burden to the HPC resources it is monitoring.

An LDMS capability called Streams also enables scalable transport of information from application data sources. AppSysFusion currently integrates LDMS with three application monitoring frameworks:

- **Kokkos** – an application portability layer developed at Sandia and used heavily by code teams. AppSysFusion samples Kokkos function timers to obtain lightweight application performance information.

- **Darshan** – an I/O profiling tool developed at Argonne National Laboratory.

- **Caliper** – an application profiler from Lawrence Livermore National Laboratory, used at multiple HPC sites. Caliper provides additional application performance information.

The application and system data gathering capabilities of AppSysFusion collect thousands of metrics per compute-node. In its current implementations on Sandia HPC systems, AppSysFusion generates roughly 15 terabytes of system and application performance data per day. That data is stored in a custom database created to handle the size and complexity of HPC monitoring data.

**VIDEO**: *The AppSysFusion dashboard can be used to visualize performance data on a job that is currently running. View this video on SimMagic App available on iOS*

To make it all work, the project team developed two new components. To assist in collecting the application data, the team developed connectors to link the application monitoring frameworks with the lightweight data-collection mechanisms. The team also set up the engine for analyzing and visualizing the time-aligned application and system data.

This engine enables users to create custom modules for data transformation and interpretation which can range from simple statistical measurements, such as rates and averages to things as complex as machine learning modeling for anomaly detection. The results are viewable, during runtime, in meaningful and actionable dashboards created by the project team in collaboration with application development teams and system administrators. Taken together, the capabilities of AppSysFusion enable other potential uses:

- Improving system design and acquisition by identifying common performance bottlenecks.

- Aiding application development with application-related timing information gathered on production runs at scale. This can alert developers to optimizations that can improve performance.

- Using artificial intelligence and machine learning to detect and diagnose anomalous and poor performance and to identify opportunities for improving performance.

- Users can potentially reconfigure or rebalance applications in response to detected performance problems.
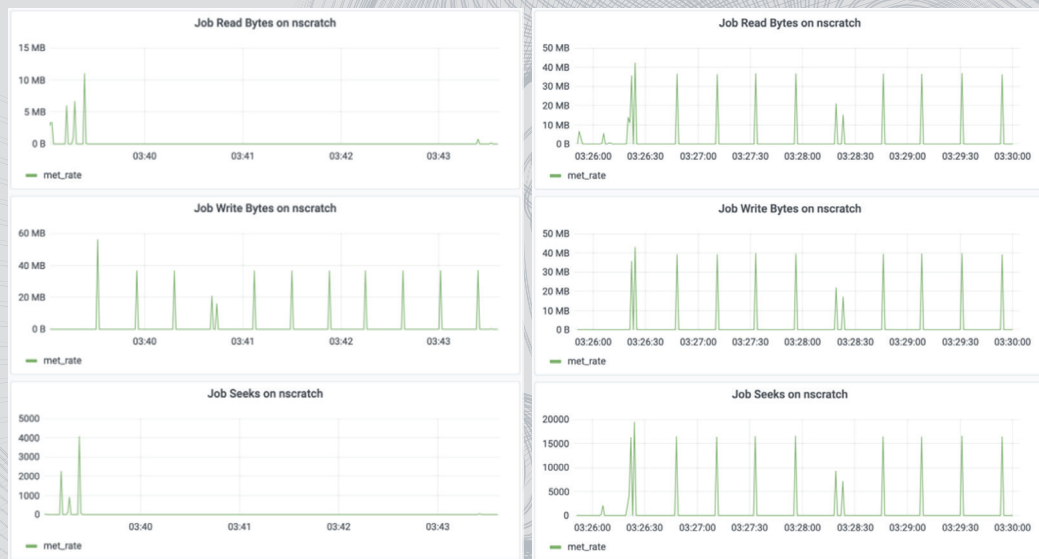
FIGURE 2. *By seeing the allocation's read, write and seek rate for the scratch filesystem, it was found that the application's implementation of NetCDF3[KM20] had unnecessary reads and seeks while writing binary output (left). By changing the format to NetCDF4[KM21], those unnecessary reads and seeks were removed.*

- Using run time feedback for scheduling and resource management, to avoid killing jobs before completion when run time has been underestimated and when delaying other jobs is acceptable.

- Using historical job resource utilization and performance profiles to potentially suggest better node counts, run times and application placements.

Currently, AppSysFusion is in use at Sandia by multiple application teams' developers and analysts for diagnostics and for resource allocation decisions. System administrators are using the capability for getting resource performance and utilization information.

Project team members plan to take the system's capabilities even further. The team's ultimate goal is to develop a self-driving, capability in which monitoring data would be gathered, analyzed and acted upon automatically by the system to optimize performance instantaneously.

# Revolutionary speedups in SIERRA structural dynamics enhance mission impact

TEAM | *Johnathan Vo*

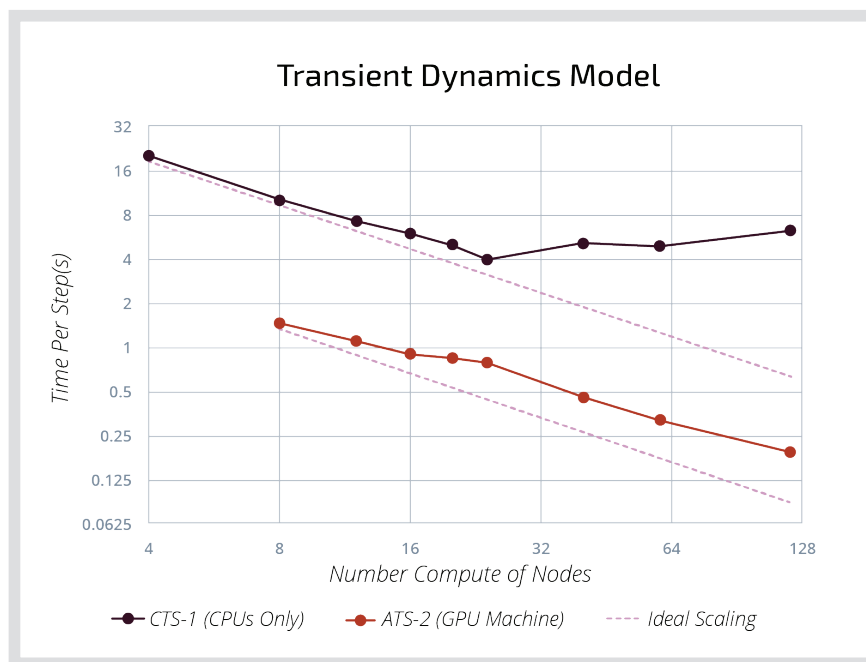CONTRIBUTING WRITE | *Kristen Meub*

**Transient Dynamics Model**

*Figure 1. Strong scaling of a transient dynamics acceptance test for CPU-only machine CTS-1 and GPU-based machine ATS-2. The GPU machine continues to scale up to 128 compute nodes while the CPU-only machine begins to plateau around 24 compute nodes. GPU machines shown to be faster, more scalable, and high throughput.*

A Sandia team is performing computer simulations for nuclear deterrence missions 10 to 20 times faster after making a series of improvements, ultimately enabling them to produce better information about how a nuclear weapon will function in various environments and wear throughout its lifetime.

Sandia researchers have been actively developing a massively parallel structural dynamics finite element code, called Sierra Structural Dynamics, or SD, for more than 25 years. The code is used for system-level analyses and design of nuclear weapons and is part of the Sierra Engineering Mechanics code suite for mechanical, fluid and thermal modeling. During the last several years, the team made software and hardware updates to speed up performance.

Code developers integrated the latest graphical processing units, or GPUs, and updated default code parameters so that users would not need to make modifications to the simulation files to run simulations on machines with or without GPUs.

As a primarily linear code, a majority of the runtime in structural dynamics analyses are in the linear system solves and dense matrix computations. The GPU is well suited for these types of computations.

Developers enabled computation on the GPU by using other code packages developed at Sandia, including Trilinos, Kokkos and Tacho. The team used packages in Trilinos that implement GPU-ready operations via Kokkos, and the sparse-direct linear solvers provided by Tacho helped improve GPU performance.

Most of SD's performance-critical operations are now built upon Trilinos objects. Through the sustainable component architecture, most of the GPU-related complexity and maintenance are hidden from the SD application. This approach to central processing unit, or CPU, GPU portability allows developers to focus on the big picture as they work to expand GPU support to SD algorithms and will be the approach used to migrate to next generation machines moving forward.

Along with algorithmic optimizations, better use of memory on the GPU platforms also provided better scalability. In a strong scaling study performed for a transient dynamics acceptance test model, the time per time step for the CPU-only machine CTS-1 begins to plateau after approximately 24 compute nodes while that for the GPU-based machine ATS-2 continues to scale up to 128 compute nodes.

Additionally, analysts used Lawrence Livermore National Laboratory's new advanced high performance computing system Sierra, known as ATS-2, to run simulations. ATS-2 offers immense throughput performance and scalability. The faster runtimes and shorter queues have enabled more analyses within a typical project timeline, higher fidelity simulations and routine heroic simulations.



FIGURE 2. *LLNL's Sierra supercomputer, known as ATS-2*



CTS-1 (Serrano)

45min

=

1x
Modify Model Parameters
Run Eigen Simulation
Check against Test Modes

ATS-2 (Sierra)

5min

=

1x
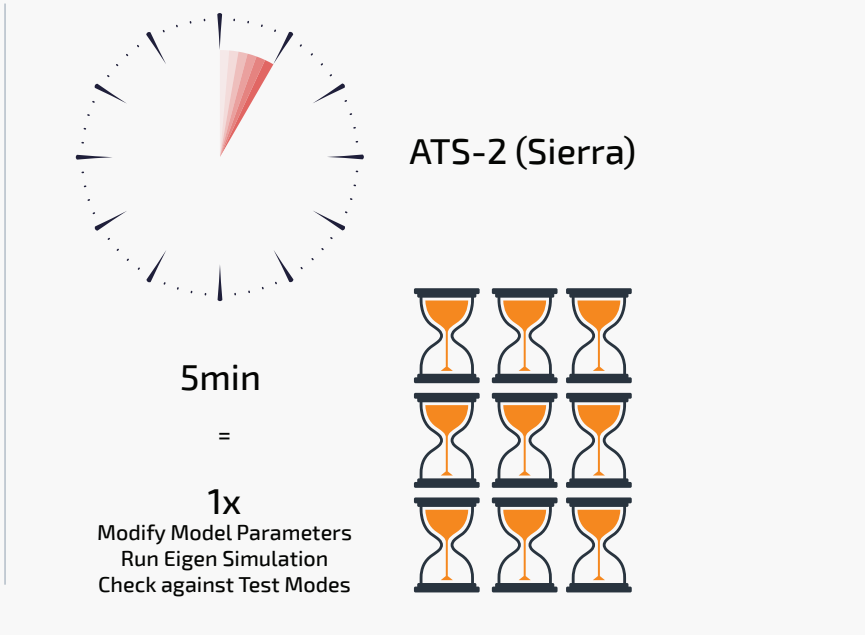Modify Model Parameters
Run Eigen Simulation
Check against Test Modes

FIGURE 3. *Enhance Existing Simulations. Faster iterations for model calibration provide better insight on the effects of parameter changes.*

# ENHANCE EXISTING ANALYSES: COMPONENT MODEL UPDATING

The performance improvements in SD have enhanced existing analyses by improving turnaround times for model calibration. In one example, an analyst iteratively performed eigen simulations to calibrate model parameters until frequencies and mode shapes match between test measurements and simulation. Each iteration

on the CTS-1 required 45 minutes to perform the modification of parameters, run the simulation, and compare results to measurements. Using the GPU-based machine ATS-2, the analyst was able to complete each iteration in five minutes. The speedups provided closer to real-time feedback about parameter changes and more freedom to see what happens with each iteration.
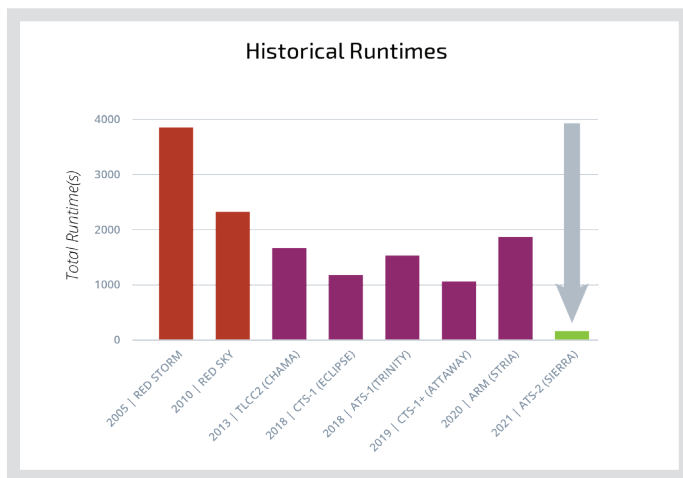
Historical Runtimes

**FIGURE 5.** *Historic runtimes of acceptance test model show that the multiyear GPU development yielded dramatic runtime reduction.*
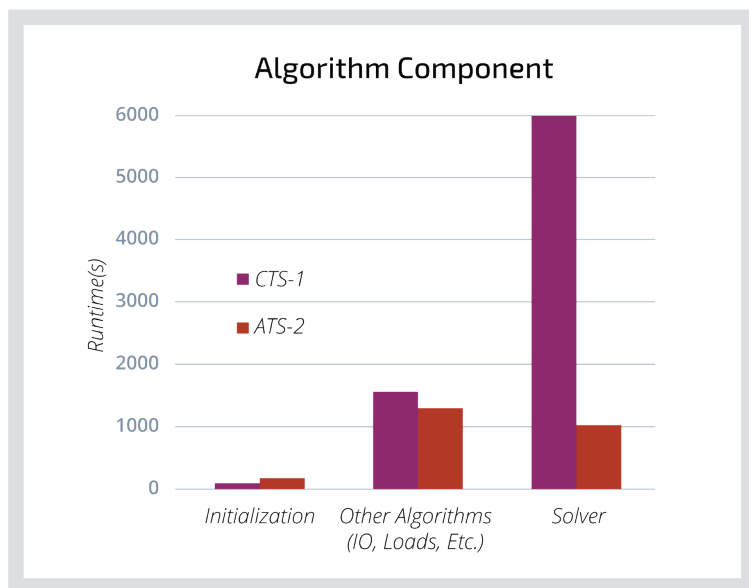
# HIGH FIDELITY EXPERIMENTAL TEST SUPPORT

Analysts often work to support experimental design through simulations. In this second example, an analyst performed system model simulations to inform the design of an Impedance-match multi-axis test. IMMAT testing is used to represent the response of a structure to complex vibration loads by attaching multiple shakers to excite the structure at discrete measurement points. Simulations can help determine the location of the shakers to produce representative response in the structure. Simulation speedups from 20 hours on legacy machines to three hours using ATS-2 resulted in an approximately one-week turnaround from request to results. Such completion times were previously unobtainable in support of experimental test design.

# ROUTINE "HEROIC" SIMULATIONS

Heroic simulations are long, complex, high-fidelity simulations that are rarely run. For example, a high-fidelity, full-system model may need to be simulated when sub-component analyses raise questions about the full-system response. When such simulations are attempted on legacy machines, the run could take weeks to complete due to queue times and the need to perform restarts. Access to ATS-2 has enabled high-fidelity, full-system dynamics to be simulated in as little as 10 hours. Routine "heroic" simulations can now be performed overnight. The ability to perform heroic simulations and produce conclusive results overnight allows analysts to provide better information to customers, resulting in more informed decision making.

**FIGURE 4.** *Example SIERRA/SD runtime breakdown for CPU–only machine (CTS–1) vs. GPU–based machine (ATS–2) showing the drastic reduction in solver time.*



Algorithm Component

# Acknowledgements

Sandia National Laboratories

U.S. DEPARTMENT OF ENERGY

NNSA
National Nuclear Security Administration